

RESEARCH ARTICLE

Distilling identifiable and interpretable dynamic models from biological data

Gemma Massonis¹, Alejandro F. Villaverde^{2,3*}, Julio R. Banga^{1*}

1 Computational Biology Lab, MBG-CSIC (Spanish National Research Council), Pontevedra, Galicia, Spain, **2** CITMAga, Santiago de Compostela, Galicia, Spain, **3** Universidade de Vigo, Department of Systems and Control Engineering, Vigo, Galicia, Spain

* afvillaverde@uvigo.gal (AFV); j.r.banga@csic.es (JRB)

Abstract

Mechanistic dynamical models allow us to study the behavior of complex biological systems. They can provide an objective and quantitative understanding that would be difficult to achieve through other means. However, the systematic development of these models is a non-trivial exercise and an open problem in computational biology. Currently, many research efforts are focused on model discovery, i.e. automating the development of interpretable models from data. One of the main frameworks is sparse regression, where the sparse identification of nonlinear dynamics (SINDy) algorithm and its variants have enjoyed great success. SINDy-PI is an extension which allows the discovery of rational nonlinear terms, thus enabling the identification of kinetic functions common in biochemical networks, such as Michaelis-Menten. SINDy-PI also pays special attention to the recovery of parsimonious models (Occam's razor). Here we focus on biological models composed of sets of deterministic nonlinear ordinary differential equations. We present a methodology that, combined with SINDy-PI, allows the automatic discovery of structurally identifiable and observable models which are also mechanistically interpretable. The lack of structural identifiability and observability makes it impossible to uniquely infer parameter and state variables, which can compromise the usefulness of a model by distorting its mechanistic significance and hampering its ability to produce biological insights. We illustrate the performance of our method with six case studies. We find that, despite enforcing sparsity, SINDy-PI sometimes yields models that are unidentifiable. In these cases we show how our method transforms their equations in order to obtain a structurally identifiable and observable model which is also interpretable.

OPEN ACCESS

Citation: Massonis G, Villaverde AF, Banga JR (2023) Distilling identifiable and interpretable dynamic models from biological data. PLoS Comput Biol 19(10): e1011014. <https://doi.org/10.1371/journal.pcbi.1011014>

Editor: Marc R Birtwistle, Clemson University, UNITED STATES

Received: March 13, 2023

Accepted: October 3, 2023

Published: October 18, 2023

Copyright: © 2023 Massonis et al. This is an open access article distributed under the terms of the [Creative Commons Attribution License](https://creativecommons.org/licenses/by/4.0/), which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

Data Availability Statement: All data and code used for running experiments is available on <https://doi.org/10.5281/zenodo.7713047>.

Funding: This research has received support from grant PID2020-117271RB-C22 (BIODYNAMICS) funded by MCIN/AEI/10.13039/501100011033; from the CSIC intramural project grant PIE 202070E062 (MOEBIUS); from grant PID2020-113992RA-I00 funded by MCIN/AEI/10.13039/501100011033 (PREDYCTBIO); from grant ED431F 2021/003 funded by Consellería de Cultura, Educación e Ordenación Universitaria,

Author summary

Dynamical models provide a quantitative understanding of complex biological systems. Since their development is far from trivial, in recent years many research efforts focus on obtaining these models automatically from data. One of the most effective approaches is based on implicit sparse regression. This technique is able to infer biochemical networks with kinetic functions containing rational nonlinear terms. However, as we show here,

Xunta de Galicia; and from grant RYC-2019-027537-I funded by MCIN/AEI/10.13039/501100011033 and by “ESF Investing in your future”. The funding bodies played no role in the design of the study, the collection and analysis of data, or in the writing of the manuscript.

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

one limitation is that it may yield models that are unidentifiable. These features may lead to inaccurate mechanistic interpretations and wrong biological insights. To overcome this limitation, we propose an integrated methodology that applies additional procedures in order to ensure that the discovered models are structurally identifiable, observable, and interpretable. We demonstrate our method with six challenging case studies of increasing model complexity.

Introduction

Mathematical models are increasingly used to describe, monitor, analyze and predict the behavior of complex biological systems. One of the major benefits of using mathematical models to study biology is that they can provide an objective and quantitative understanding that would be difficult to achieve through any other means. In systems biology, dynamical models (typically sets of ordinary differential equations, ODEs) are widely used to provide mechanistic insights into the functioning of biological systems [1, 2].

The use of dynamical systems theory originated in Newtonian mechanics is now pervasive in all the natural and engineering sciences [3]. Dynamic models are highly versatile, enabling researchers to study complex biosystems from a range of different perspectives, such as (i) analyzing the effect of changes in conditions and scenarios different from those studied experimentally, (ii) guiding research by identifying key aspects that need to be further investigated, (iii) helping to generate new testable hypotheses, or (iv) guiding the design of interventions. However, the systematic development of mechanistic dynamic models is a non-trivial exercise. In the case of biological systems, the situation is particularly difficult due to the fact that we cannot rely on first principles in the same way as in e.g. physics. As a consequence, model development is one of the key open problems in mathematical biology [4].

Can we automate the development of mechanistic models? This question of model discovery (in the sense of symbolic reconstruction of equations) from data was already addressed by pioneering attempts in the field of artificial intelligence several decades ago [5–7]. However, the data-driven automatic identification of nonlinear dynamic models has only been addressed more recently. In this area, several different statistical and machine learning frameworks have been considered, including symbolic regression [8, 9], grammar-based methods [10, 11], sparse regression [12], neural networks [13–15], Gaussian process regression [16, 17] and Bayesian approaches [18–20]. More detailed reviews can be found in [21–25]. The sparse identification of nonlinear dynamics (SINDy) algorithm [12] has been particularly successful, and a number of extensions have been developed (see review in [26]).

In the case of biological systems, a large amount of research has been devoted to different classes of subproblems with different simplifying assumptions (such as e.g. static networks, non-mechanistic dynamic networks, linear dynamics, etc.), as reviewed by [27–29]. In this work, we consider the more general problem of fully reconstructing interpretable (mechanistic and parameterized) nonlinear dynamic models from time-series data. Recently, several approaches using methods based on sparse regression, Bayesian identification or symbolic regression have appeared [18, 30–36]. In this context, SINDy-PI [37] is an especially interesting parallel implicit version of SINDy because it allows the incorporation of implicit dynamics and rational nonlinear terms, thus enabling the discovery of kinetic functions (such as Michaelis-Menten) common in biochemical networks.

Many of these SINDy-based methods pay special attention to the recovery of parsimonious models, usually penalizing model complexity [38] or evaluating performance on a validation

data-set [37]. The objective is to find the simplest model which can explain the data, in agreement with the well known principle of Occam's razor. These strategies help to discard more complex models which would be indistinguishable (i.e. would explain the data equally well but adding spurious terms). Besides enforcing simplicity, a related key aspect in model discovery is ensuring structural identifiability and observability (SIO). The property of structural identifiability refers to the theoretical possibility of inferring the unknown parameters of a given model (assuming that its equations are known, except for the numerical values of the parameters) from observations of the model output, which typically consist of time-resolved measurements of its state variables, or of a subset of them [39]. Likewise, observability is the possibility of inferring all the state variables of a model at a given time from future observations of a subset of them. Since lack of SIO makes it impossible to uniquely infer parameters and state variables, it can compromise the usefulness of the model [40–46]. The analysis of these properties can be performed with symbolic computation tools [47], and numerical approaches have also been proposed for their study [48, 49]. However, to the best of our knowledge, ensuring SIO has not been considered in dynamic model discovery yet.

Here we present a methodology that ensures SIO in automatic model discovery in two possible scenarios: with and without prior knowledge. In both cases the end product is a dynamic model of a biological system consisting of (typically nonlinear) ODEs. The equations may contain rational terms, such as Michaelis-Menten kinetics, thus being suitable for the description of many biochemical processes. If there is no prior knowledge about the model structure, the methodology performs equation discovery with the SINDy-PI approach, and incorporates a SIO analysis as a post-processing stage. If there is prior knowledge (i.e. we have a candidate model), another SIO analysis is added as a pre-processing step. If the analyses reveal structural unidentifiabilities, a reparameterization step is carried out to ensure that the resulting model is fully identifiable and observable. Furthermore, equivalent model reformulations are generated to facilitate its interpretation in a mechanistic sense.

Using representative case studies, we illustrate how ignoring these structural properties can lead to wrong conclusions or poorly identified models. Although we demonstrate the use of the methodology with SINDy-PI, it is straightforward to apply it in combination with other automatic discovery methods. In particular, it could easily be adapted to future methods capable of considering partially-observed systems.

Overall, our study presents a novel and non-trivial integrated methodology to ensure that the discovered models are structurally identifiable and interpretable. To the best of our knowledge, this is the first study to address these questions in model discovery. Further, our method involves an original and non-obvious combination of algorithmic steps regarding structural identifiability analysis (SIO), reparameterization, reformulation and interpretability analysis. While the concepts of SIO and reparameterization draw on recent ideas developed in our group, the remaining steps and their integration represent fresh and innovative contributions to the field.

Methods

In this section we describe the methodology, which can be used in two different scenarios. Both of them entail performing model discovery (using SINDy-PI or a similar approach) and performing SIO analysis. If a model is structurally identifiable and observable, we say that it is FISPO (full input, state, and parameter observability). If the SIO analysis reveals that the model is not FISPO, our method suggests a reparameterization step. The two scenarios and their procedures are as follows:

- Scenario (I): full model discovery from time-series data with no prior knowledge. Since we assume zero prior knowledge, we use SINDy-PI to discover a candidate model (CM). We

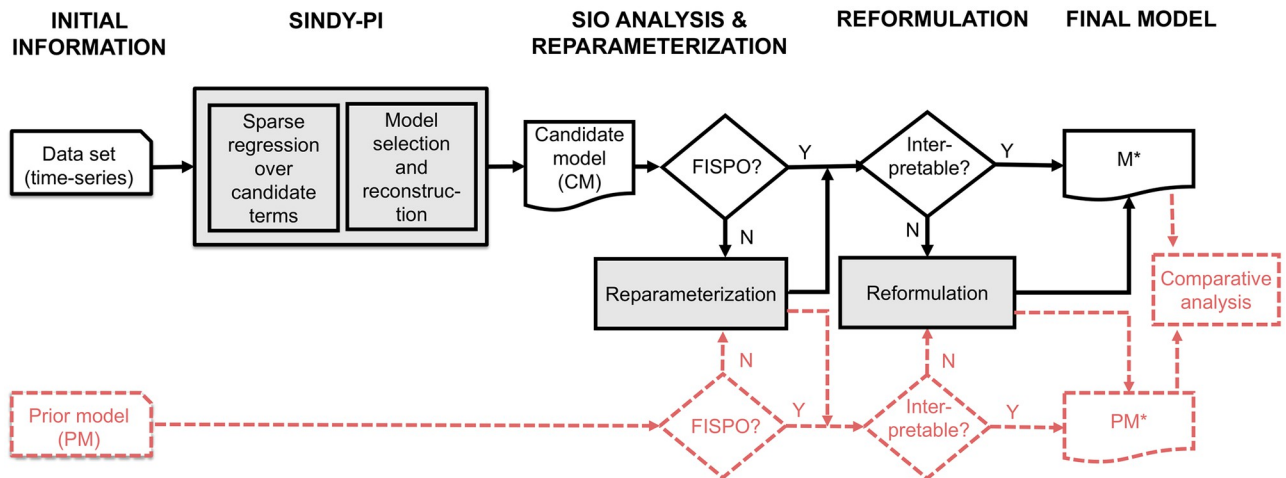


Fig 1. Workflow of the methodology. Scenario (I) (solid black lines only): data-driven full model discovery from (time-series) data with no prior knowledge. We apply SINDy-PI and test the SIO of the discovered candidate model (CM). If the CM is not FISPO, we reparameterize it. Next, we check if the model is interpretable; if not, we reformulate it via symbolic manipulation. The result is a FISPO interpretable model, M^* . Scenario (II) (solid black lines + dashed, dark orange lines in the lower part): model discovery from (time-series) data with good prior knowledge. In this scenario we seek model (in)validation and/or refinement. We have a prior model (PM) which we want to compare with an alternative candidate discovered from data (CM). To this end, we check the SIO of the PM and reparameterize if needed. In parallel, we apply SINDy-PI to the data to obtain a CM, and we make sure it is FISPO (using reparameterisation if not). Then, we use model reformulation techniques to ensure interpretable versions (M^* and PM^*) if needed. Lastly, we perform a comparative analysis.

<https://doi.org/10.1371/journal.pcbi.1011014.g001>

then analyse its SIO. If it is not FISPO, we reparameterize it in order to obtain an equivalent model which is FISPO. Finally, we check if the model is *interpretable*, in the sense that it contains monomials and simple rational terms which belong to a dictionary of mechanistic kinetic terms. If not, we apply a symbolic reformulation step in order to render it interpretable.

- **Scenario (II):** model discovery from time-series data with prior knowledge. This scenario corresponds to situations where we seek model (in)validation and/or refinement. We assume good prior knowledge and time-series data, that is we are reasonably confident that our prior model (PM), which already represents the data quite well, is close to the ‘true’ one. Here the motivation to use SINDy-PI is to compare this PM with an alternative candidate obtained via model discovery (CM). To this end we check the SIO of the PM, obtaining a reparameterized version if needed. In parallel, we apply SINDy-PI to the data, obtaining a CM, and we make sure that it is FISPO (using reparameterisation if not). If needed, we use model reformulation techniques to obtain interpretable versions of the CM and the PM. Finally, we perform a comparative analysis of these latter models.

An schematic diagram of our method considering these two scenarios is depicted in Fig 1. In the remainder of this section we describe in detail each of the steps.

Automatic model discovery using sparse regression

We assume that the dynamical system is governed by classical reaction-rate nonlinear ordinary differential equations with the following form:

$$M = \begin{cases} \dot{x}(t) & = f(x(t), p) \\ y(t) & = g(x(t), p) \\ x_0 & = x(t_0, p) \end{cases}$$

where $\dot{x}(t) \in \mathbb{R}^n$ is the state vector, $p \in \mathbb{R}^{n_p}$ is the parameter vector, the function $f(x(t), p)$ represents the dynamics, $y(t)$ is the measurable output, and x_0 is the vector of initial conditions. SINDy [12] assumes a fully observed system, $y(t) = x(t)$. In the remainder of this section, we will consider $\dot{x}(t) = f(x(t))$ to simplify the notation. SINDy also assumes that $f(x(t))$ can be expressed as the product of a suitable library function, $\Theta(x(t))$, and a sparse vector ξ (indicating the active library terms), where each entry in the library function is a candidate term:

$$\Theta(x) = [\theta_1(x) \ \theta_2(x) \ \theta_3(x) \ \dots \ \theta_p(x)] , \tag{1}$$

By arranging the time-series data as a matrix, $X = [x(t_1), \dots, x(t_m)]$, and its associated derivative matrix $\dot{X} = [\dot{x}(t_1), \dots, \dot{x}(t_m)]$, $\dot{x}(t)$ can be expressed as:

$$\dot{x}(t) \approx \Theta(x(t))\Xi , \tag{2}$$

where Ξ corresponds to the sparse matrix of active terms. When the system includes rational terms, $f(x)$ can be rewritten as:

$$\dot{x}(t) = f(x) = N(x)/D(x) \tag{3}$$

leading to the implicit problem formulation [31]:

$$\dot{x}(t)D(x) = N(x) . \tag{4}$$

Eq 4 has a different kind of term in each side of the equality: the *Left Hand Side (LHS)*, in which there are combinations of term involving the derivative data and the candidate library, and the *Right Hand Side (RHS)*, in which we only have library terms. When $f(x)$ includes rational terms, model complexity can be viewed as the number of terms in the LHS, as they will involve the denominator degree too.

The generalized function library $\Theta(X)$ allows the inclusion of X and \dot{X} . Under this consideration, the implicit problem formulation can be rewritten as:

$$\dot{x}(t)D(x) - N(x) = 0 \rightarrow \Theta(X, \dot{X})\Xi = 0 . \tag{5}$$

For example, if a model has two states and the chosen degree is 2, the function library for the first state, i.e. x_1 , will be:

$$\Theta(X, \dot{X}) = [1, x_1, x_2, x_1^2, x_1x_2, x_2^2, \dot{x}_1, \dot{x}_1x_1, \dot{x}_1x_2, \dot{x}_1x_1^2, \dot{x}_1x_1x_2, \dot{x}_1x_2^2] \tag{6}$$

It should be noted that the function library for state x_2 will differ from Eq 6 as it will include \dot{x}_2 instead of \dot{x}_1 . The design of the library of candidate functions is a critical aspect of SINDy-PI. However, in the context of dynamic modelling of biochemical and biological phenomena, by including monomials up to an order of 5 or 6, we can accommodate the vast majority of nonlinear terms, such as e.g. mass-action kinetics, or feedback regulatory loops. Furthermore, common nonlinear rational terms such as Michaelis-Menten in enzyme kinetics, or Monod in microbial growth, can be inferred from such library due to the implicit nature of SINDy-PI. The case studies considered here cover a wide range of nonlinear terms, illustrating the capabilities of SINDy-PI in biological modelling.

The implicit form of Eq 5 admits the sparsest trivial solution $\Xi = 0$. The implicit-SINDy algorithm [31] surmounted this issue by using nonconvex optimization to find the sparsest vector ξ in the null space. However this particular formulation is very sensitive towards noise levels, thereby affecting the robustness of the method.

In an effort to address this issue, Kaheman et al. [37] introduced SINDy-PI, a novel method that solves the problem using a sequence of convex relaxations of the non-convex optimization

problem. By doing so, the algorithm can utilize the same noise levels as those employed in the original SINDy algorithm. The authors achieved this by assuming the knowledge of at least one term on the *LHS*, specifically of the form $\dot{x}D(x)$. Consequently, Eq 5 can be rewritten as:

$$\theta_j(X, \dot{X}) = \Theta(X, \dot{X}|\theta_j(X, \dot{X}))\xi_j, \tag{7}$$

where $\Theta(X, \dot{X}|\theta_j(X, \dot{X}))$ denotes the library $\Theta(X, \dot{X})$ without the θ_j element. SINDy-PI proceeds iteratively by examining each term within the library. In order to balance accuracy and complexity, the method performs Pareto optimal model selection.

To find the sparsest vector ξ , SINDy-PI considers the problem:

$$\xi_j = \operatorname{argmin}_{\xi_j} \|\theta_j(X, \dot{X}) - \Theta(X, \dot{X}|\theta_j(X, \dot{X}))\xi_j\|_2 + \beta\|\xi_j\|_0 \tag{8}$$

SINDy-PI solves this non-convex optimization problem using a sequentially thresholded least-squares (STLSQ) approach. This method proceeds by iteratively solving the least squares term in the cost function, zeroing out elements of ξ that are below a certain threshold λ . This threshold must be fine-tuned to select the model that provides the best trade-off between accuracy and efficiency. To discover the model, SINDy-PI considers a finite set of λ values and proceeds by sweeping the library terms for each value of the threshold λ , obtaining a family of possible candidate models. Next, SINDy-PI performs model selection choosing the best trade-off, i.e. the Pareto-optimal model. The Pareto front is obtained by considering a model complexity metric (such as the Akaike information criterion, AIC), and the score for each candidate model. Details of this process, illustrated with an example, are given in the Supporting Information.

Structural identifiability and observability analysis and reparameterization

Once a candidate model structure (CM) has been discovered, the next step is to analyse its structural identifiability and observability (SIO) [50]. This test assesses the possibility of determining the values of the model parameters and state variables, respectively, from output measurements. These properties are *structural* (i.e. they depend only on the model equations) and hence they can be analysed *a priori* (i.e. before taking experimental measurements) using symbolic computation. They should not be confused with the so-called *practical* versions of these properties, which depend on the features of the experimental data and are analysed *a posteriori*, i.e. after performing measurements [39].

We can provide a mathematical definition of *structural local identifiability* (SLI) as follows. Let us denote by $y(t, p)$ the output vector obtained with a parameter vector p at time t . (For fully observed systems $y(t, p) = x(t, p)$, while for partially observed systems y typically consists of a subset of x .) We say that a parameter p_i (which is the i^{th} element of the parameter vector $p \in \mathbb{R}^{n_p}$) is structurally locally identifiable (SLI) if, for almost any parameter vector $p^* \in \mathbb{R}^{n_p}$, there is a neighbourhood $\mathcal{N}(p^*)$ such that:

$$\hat{p} \in \mathcal{N}(p^*) \quad \text{and} \quad y(t, \hat{p}) = y(t, p^*) \Rightarrow \hat{p}_i = p_i^*. \tag{9}$$

The definition of structural *global* identifiability is similar, but with the neighbourhood $\mathcal{N}(p^*)$ extending to the whole parameter space. In this paper we focus on SLI.

There are several approaches for determining structural local identifiability and observability. We apply a differential geometry approach, which we explain briefly in these paragraphs. In this framework, parameters are treated as state variables that happen to be constant, i.e. the state vector is augmented with the parameters, $\tilde{x} = (x^T, p^T)^T$, and has dimension $n_{\tilde{x}} = n + n_p$.

The augmented dynamic equations are $\dot{\tilde{x}} = \tilde{f}(\tilde{x})$, and the output function is $y = g(\tilde{x})$, omitting the dependence on time for ease of notation.

Thus, SLI is considered as a particular case of a more general property, observability, which describes the possibility of inferring the internal state of a model by observing its output vector—hence the use of the term FISPO for “full input, state, and parameter observability” (note that this concept also allows for the treatment of unknown inputs as additional state variables, a possibility that we will not consider in this paper).

We analyse SIO by building an observability-identifiability matrix and computing its rank. The matrix is built with Lie derivatives of the output function. The zero-order Lie derivative is $L_f^0 g(\tilde{x}) = g(\tilde{x})$, and for $i \geq 1$ the i -order Lie derivatives are obtained as:

$$L_f^i g(\tilde{x}) = \frac{\partial L_f^{i-1} g(\tilde{x})}{\partial \tilde{x}} \tilde{f}(\tilde{x}).$$

The observability-identifiability matrix \mathcal{O}_I is:

$$\mathcal{O}_I(\tilde{x}) = \frac{\partial}{\partial \tilde{x}} \left(L_f^0 g(\tilde{x})^T L_f g(\tilde{x})^T L_f^2 g(\tilde{x})^T \dots L_f^{n_{\tilde{x}}-1} g(\tilde{x})^T \right)^T, \tag{10}$$

A model is FISPO around a point \tilde{x}_0 if the rank of its observability-identifiability matrix equals the number of its states and parameters, $\text{rank}(\mathcal{O}_I(\tilde{x}_0)) = n_{\tilde{x}} = n_x + n_p$. If the rank is smaller, the model contains structurally unidentifiable parameters. By performing additional tests it is possible to determine which specific parameters are structurally identifiable, and which state variables are observable.

If a model is not FISPO, its calibration will almost surely produce wrong parameter estimates. Furthermore, structural unidentifiability is often linked with non-observability, in which case the simulations of some state variables will also be wrong. Thus, structural non-identifiability and non-observability are undesirable features of a model’s structure, which compromise its reliability as a source of biological insight. These features are caused by symmetries in the differential equations of the model that make its output invariant with respect to certain changes in their parameters and/or state variables [51–53]. Said symmetries can be studied in the framework of Lie group theory. We say that a mapping of the form

$$x^* = X(x, \varepsilon), \tag{11}$$

is a one-parameter Lie group of transformations (with ε being the parameter) if it has the following properties: it is smooth in x and analytic in ε , it satisfies the four group axioms (closure, associativity, and existence of an identity and an inverse), and the identity element can be chosen as $\varepsilon = 0$. The transformation above is also called a *symmetry transformation*, or a Lie symmetry. Examples of the simplest and possibly most common symmetries in biological modelling include the following:

Translation:

$$x_i^* = x_i + \varepsilon, \quad X = \frac{\partial}{\partial x_i} \tag{12}$$

Scaling:

$$x_i^* = e^\varepsilon x_i, \quad X = x_i \frac{\partial}{\partial x_i} \tag{13}$$

Moebius:

$$x_i^* = \frac{x_i}{1 - \varepsilon x_i}, \quad X = x_i^2 \frac{\partial}{\partial x_i} \quad (14)$$

It is sometimes possible to remove or ‘break’ these symmetries by transforming the model equations via a suitable reparameterization. To this end, we first search for the symmetry transformations admitted by the model. If a model has such symmetries, it is overparameterized and therefore structurally unidentifiable. Then, we express the ε of those transformations in terms of other parameters, thus setting the value of one of the transformed parameters to one and removing it from the equations. The end result of the reparameterization is a FISPO model that has exactly the same dynamic behaviour as the original one. In previous work [54] we presented a methodology to perform such reparameterizations automatically, which has been integrated in the workflow described here.

In summary, if the SIO analysis of the CM reveals structural unidentifiability and/or non-observability, our methodology applies a symmetry-breaking reparameterization that makes it FISPO.

Model reformulation for interpretability

The dynamic model obtained in the previous step supports the experimental data and is structurally identifiable and observable. However, the rational expressions in 3 may lack a clear biological interpretation. In the case of biological networks, we will need to reformulate expressions of the form $N(x)/D(x)$ into terms that belong to a dictionary of interpretable terms.

Our model reformulation procedure seeks to transform it into simple monomials and rational terms that have a mapping with the dictionary of kinetic and regulatory terms compatible with the specific type of biological reaction network under study. Typically, this dictionary will include mass-action kinetics and simple rational functions (e.g. Michaelis-Menten for enzyme kinetics, or Hill for cooperative binding). However, care should be taken in order to ensure that the reformulation does not destroy identifiability and observability. Further, as shown in the case studies below, sometimes these rational terms can have high degrees, complicating model discovery.

Our reformulation procedure makes use of symbolic manipulation and involves the following steps:

1. Obtain the list of p non-trivial divisors of the denominator:

$$dd(x) = [dd_1(x), \dots, D(x)] \quad (15)$$

2. Obtain a family A of expressions composed of monomials (interpretable as e.g. mass-action kinetics), minimizing the number of rational terms and their degree, by obtaining the quotients and the residuals:

$$\frac{N(x)}{dd_i(x)} = dd_i(x)q_i(x) + r_i(x) . \quad (16)$$

3. If any residuals $r_i(x)$ lack interpretability, factorize and simplify $N(x)$ by means of the nested Horner form:

$$N(x) = a_0 + x(a_1 + x(a_2 + \dots + x(a_{n-1} + a_n x))) , \quad (17)$$

obtaining a family of coupled and factorized equations with the same degree, but with monomials involving different state combinations:

$$N_{x_1} = a_0 + x_1(a_1 + \dots) + x_2(\dots), \quad (18)$$

$$N_{x_2} = b_0 + x_2(b_1 + \dots) + x_1(\dots). \quad (19)$$

Thus, the Horner nested form gives different possible decompositions of the numerator. Next, obtain a family B of reformulations by simplifying the rational terms using the divisors of Eq 15 and the Horner form of the numerator.

As an example, consider Eq 20 below, where we can decompose the fraction in the left into a monomial plus a simpler rational term as follows:

$$\frac{a_0x + a_1x^2}{b_2 + b_3x} = \frac{b_0b_2x + x^2(b_1 + b_3b_0)}{b_2 + b_3x} = b_0x + \frac{b_1x^2}{b_2 + b_3x}. \quad (20)$$

4. Match the monomials and simplified rational terms in families A and B with elements in the dictionary of canonical kinetic and regulatory expressions (or by inspection by a human domain expert), finding members that are fully interpretable.
5. Ensure that the resulting interpretable model is FISPO. If not, reparameterize and repeat until an interpretable and identifiable model is obtained.

Implementation

We implemented our methodology, as depicted in Fig 1, using Matlab and the Symbolic Math Toolbox, integrating the following components:

- Sparse regression using SINDy-PI [37] with some modifications, as detailed in the Supporting Information.
- Structural identifiability and observability (SIO) analysis using the algorithm FISPO [55], plus reparameterization using the algorithm AutoRepar [54], as implemented in STRIKE-GOLDD 4.0 and later releases [56].
- Reformulation for interpretability, implementing the algorithm described above using symbolic manipulation.

The resulting code is available at <https://doi.org/10.5281/zenodo.7713047>. In order to facilitate reproducibility and illustrate the results at each step of the workflow, we have included interactive notebooks (Matlab live scripts) and reports (in HTML format) for each of the case studies described below. More details are given in S1 File.

Results

Below, we apply our methodology to a set of challenging case studies. (Table 1 summarises their main features). These examples are presented in order to illustrate the performance of our method for a variety of situations of increasing complexity, from models without rational terms and fully identifiable and observable (FISPO) structure, to larger (in terms of number of parameters and states), non-FISPO models with more difficult non-linearities, as indicated by the different maximum degrees in their rational terms.

Table 1. Main features of the case studies: Relevant references and main characteristics of the models considered in the case studies. The fourth and fifth rows show the maximum degree of $N(x)$ and $D(x)$ in Eq 4. The last row indicates if the original (ground truth, GT) model is fully identifiable and observable (FISPO).

| ID (short name) | Lorenz | Immunity | Bacterial | Microbial | Crypt | Glycolysis |
|-------------------|----------|----------|-----------|-----------|-------|--------------|
| References | [12, 57] | [58] | [31, 59] | [60] | [61] | [31, 37, 62] |
| # states | 3 | 2 | 2 | 2 | 3 | 7 |
| # parameters | 3 | 8 | 5 | 4 | 11 | 13 |
| Max degree $N(x)$ | 2 | 3 | 6 | 2 | 3 | 6 |
| Max degree $D(x)$ | 1 | 2 | 6 | 2 | 4 | 4 |
| FISPO | Y | N | Y | Y | N | Y |

<https://doi.org/10.1371/journal.pcbi.1011014.t001>

In order to illustrate all the steps and capabilities of our workflow, we consider Scenario (II) in all these examples. For each problem, a ground truth (GT) model is defined and subsequently considered as prior model (PM) for the sake of simplicity but without loss of generality. This GT model is used to generate training data sets in all the case studies. After confirming the identifiability and interpretability of the final discovered model M^* , we also assess its structural, parametric and predictive accuracy. The predictive power is evaluated taking into account conditions different from those used for generating the training data. Details regarding the training data generation and the conditions to evaluate predictive accuracy are given in [S1 File](#).

Lorenz system (Lorenz)

This case study involves the well-known Lorenz system [57], which is a classical example of dynamic model with chaotic behaviour. This model was previously used in [12] to demonstrate the original SINDy algorithm. The governing equations describe the dynamics of a fluid layer warmed from below and cooled from above:

$$\dot{x}_1 = a(x_2 - x_1) , \quad (21a)$$

$$\dot{x}_2 = x_1(b - x_3) - x_2 , \quad (21b)$$

$$\dot{x}_3 = x_1x_2 - cx_3 . \quad (21c)$$

where x_1 is proportional to the rate of convection, x_2 to the horizontal temperature variation and x_3 to the vertical one. For certain values of parameters a , b , and c , the system exhibits chaotic dynamics.

We consider the ideal Scenario (II) case where the prior model (PM) is the same as the nominal (or ground truth, GT) model. We generate a synthetic training data set using the GT model and settings similar to [12] (details in Supporting Information). Following the workflow in [Fig 1](#), we perform structural and identifiability analysis and confirm that the PM is fully identifiable and observable (FISPO). We then apply SINDy-PI to the training data, obtaining the following candidate model (CM):

$$\dot{x}_1 = p_1x_2 + p_2x_1 , \quad (22a)$$

$$\dot{x}_2 = p_3x_1 + p_4x_2 + p_5x_1x_3 , \quad (22b)$$

$$\dot{x}_3 = p_6x_1x_2 + p_7x_3 . \quad (22c)$$

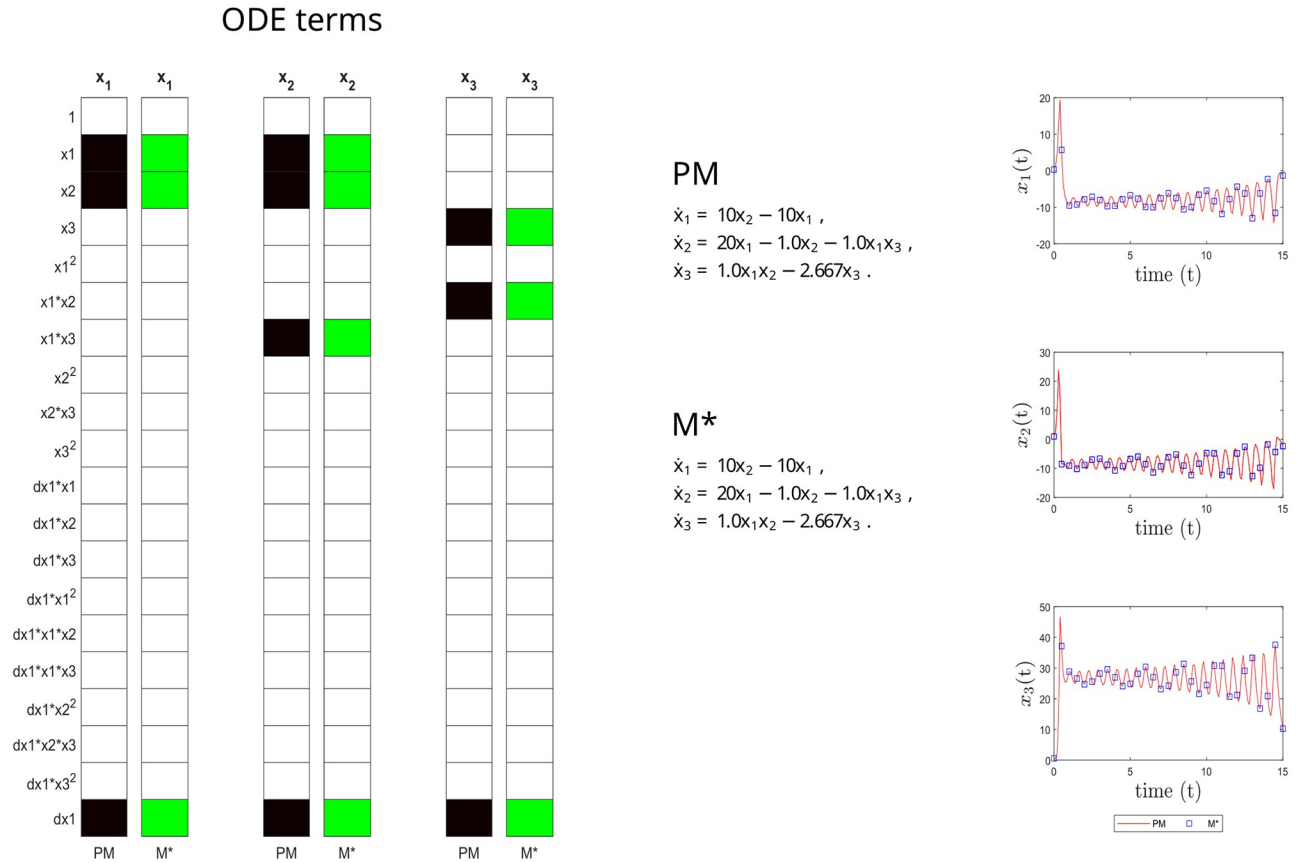


Fig 2. Lorenz case study. Structural accuracy: on the left, active terms in ξ (non-zero terms of the prior model PM in black, and of the inferred model M* in green). Parameter accuracy: center, matching parametric ODEs for PM and M*. Predictive accuracy: on the right, time evolution of the different states (x_1 , x_2 and x_3) of the PM and M* models.

<https://doi.org/10.1371/journal.pcbi.1011014.g002>

Our algorithm then confirms that this inferred CM model is FISPO and interpretable (thus, it corresponds to M*). Further, it is fully equivalent to the expanded ground truth model in terms of structural, parametric and predictive accuracy, as shown in Fig 2. In summary, in this case study we have the ideal situation where both the nominal and the inferred models are fully observable and identifiable. As we will see below, this situation might change as soon as we consider rational terms in the dynamics.

Competition between bacteria and the immune system (Immunity)

This model describes the influence of quorum sensing signaling molecules (QSSM) on the competition between bacteria and the immune system, as studied in [58]. The following differential equations depict the dynamics for the concentrations of bacteria (\dot{x}_1) and immune cells (\dot{x}_2):

$$\dot{x}_1 = a(1 - x_1/k)x_1 - ex_1x_2 - \frac{\beta\gamma x_1x_2^2}{\gamma x_2 + \alpha x_1} \tag{23a}$$

$$\dot{x}_2 = S + dx_1 - \delta x_2 \tag{23b}$$

where it is assumed that bacteria grow logistically at rate a with an effective carrying capacity

of the environment given by parameter k , and that they are cleared by the immune system following a mass action term ex_1x_2 . The rational term at the end of Eq 23a represents the modulation of QSSM in the competition between bacteria and the immune system. We consider Eqs 23a and 23b as the GT model.

We consider Scenario (II) again, assuming that the prior model (PM) is the same as the ground truth (GT) or nominal model. Considering that the unknown parameters are $a, k, e, \beta, \gamma, \alpha, S, d$ and δ , the structural identifiability analysis of the PM indicates that two of the three parameters involved in the rational term are unidentifiable. Specifically, there is a scaling symmetry between γ and α , which is probably the most common type of symmetry in biological models [63]. Our reparameterization step indicates that this issue can be solved by dividing the numerator and denominator by one of the unidentifiable parameters; for example, if we choose α , Eq 23a will be:

$$\dot{x}_1 = a(1 - x_1/k)x_1 - ex_1x_2 - \frac{\beta\gamma x_1x_2^2}{\alpha} = a(1 - x_1/k)x_1 - ex_1x_2 - \frac{\beta\gamma^*x_1x_2^2}{\gamma^*x_2 + x_1} \tag{24}$$

where $\gamma/\alpha = \gamma^*$. Thus our new reference model will be the following PM*:

$$\dot{x}_1 = a(1 - x_1/k)x_1 - ex_1x_2 - \frac{\beta\gamma^*x_1x_2^2}{\gamma^*x_2 + x_1} \tag{25a}$$

$$\dot{x}_2 = S + dx_1 - \delta x_2 \tag{25b}$$

Next, our workflow proceeds by applying SINDy-PI to a data set generated with the GT model, obtaining the following candidate model (CM):

$$\dot{x}_1 = p_1x_1 + p_2x_1^2 + p_3x_1x_2 + \frac{p_4x_1^3}{p_5x_2 + p_6x_1} \tag{26a}$$

$$\dot{x}_2 = p_7 + p_8x_1 + p_9x_2 \tag{26b}$$

Interestingly, our method then finds that this CM is not FISPO due to three structurally unidentifiable parameters: p_4, p_5, p_6 . The reformulation step is then able to find structurally identifiable reformulations of the form:

$$\dot{x}_1 = p_1x_1 + p_2x_1^2 + p_3x_1x_2 + \frac{\frac{1}{p_j}p_4x_1^3}{(p_5x_2 + p_6x_1)} , \text{ for } j \in [4, 5, 6], \tag{27}$$

We chose $j = 6$, but the same result can be obtained with $j = 4, 5$. Denoting as

$p_j^* = \frac{p_j}{p_6}$, $j = 4, 5$, the resulting dynamic system becomes identifiable. Eq 27 is re-arranged as:

$$\dot{x}_1 = p_1x_1 + p_2x_1^2 + p_3x_1x_2 + \frac{p_4^*x_1^3}{p_5^*x_2 + x_1} . \tag{28}$$

The resulting model is fully identifiable, but the rational term does not match the one in the GT describing the modulation of QSSM. However, the reformulation step in our workflow

CM

$$\begin{aligned} \dot{x}_1 &= 0.18x_1 - 0.018x_1x_2 + 0.00325x_1^2 - \frac{(0.0065x_1^3)}{1x_1 + 2x_2} \\ \dot{x}_2 &= 0.001x_1 - 0.105x_2 + 1.1 \end{aligned}$$

CM2

$$\begin{aligned} \dot{x}_1 &= 0.18x_1 - 0.018x_1x_2 + 0.00325x_1^2 - \frac{(0.00325x_1^3)}{0.5x_1 + 1x_2} \\ \dot{x}_2 &= 0.001x_1 - 0.105x_2 + 1.1 \end{aligned}$$

M*

$$\begin{aligned} \dot{x}_1 &= 0.18x_1 - 0.005x_1x_2 - 0.003x_1^2 - \frac{(0.026x_1x_2^2)}{x_1 + 2x_2} \\ \dot{x}_2 &= 0.001x_1 - 0.105x_2 + 1.1 \end{aligned}$$

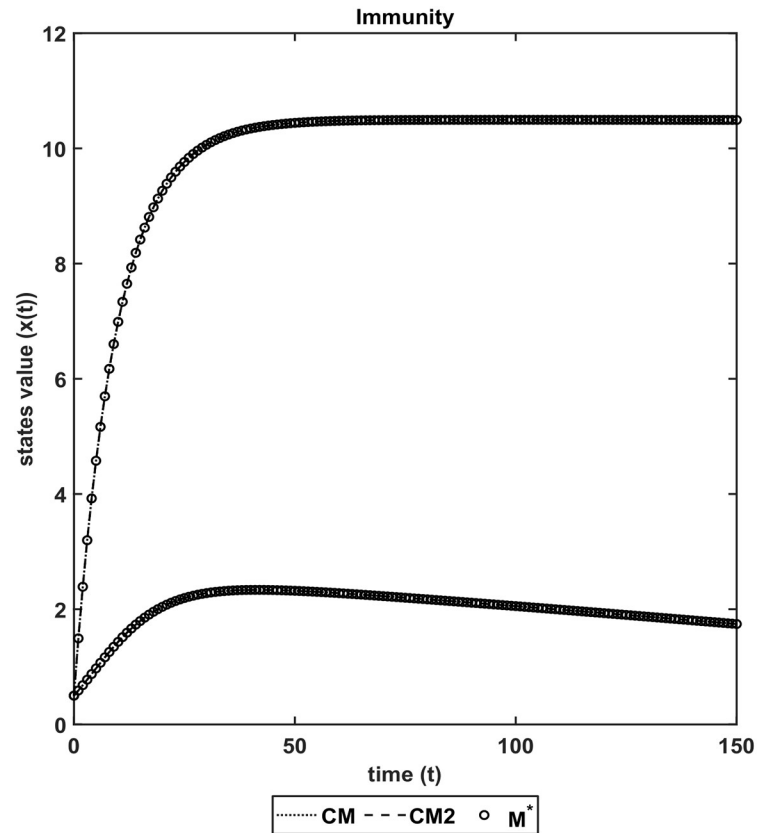


Fig 3. Immunity model. Structural unidentifiability in CM (unidentifiable parameters in red, identifiable parameters in blue) leads to the same output dynamics when different parameterizations are considered, as can be seen in CM2. In contrast, the reformulation M* is FISPO and therefore there is a unique set of parameters compatible with the output measurements.

<https://doi.org/10.1371/journal.pcbi.1011014.g003>

produces an interpretable form M* which is FISPO:

$$\dot{x}_1 = p_1x_1 + p_2x_1^2 + p_3x_1x_2 + \frac{p_4^*x_1x_2^2}{p_5^*x_2 + x_1}, \tag{29a}$$

$$\dot{x}_2 = p_7 + p_8x_1 + p_9x_2 \tag{29b}$$

Fig 3 illustrates the importance of ensuring structural identifiability and observability. The CM found by SINDy-PI is not FISPO, so there are other different parameter realizations producing exactly the same output, as shown by CM2. This means that if this CM structure is used for parameter identification, the estimated parameters will not be unique, i.e. there exist different parameterizations of CM in full agreement with the same output measurements. However, the reformulated model M* is FISPO, i.e. there is a unique set of parameter values compatible with the output. Finally, in Fig 4 we confirm the structural, parametric and predictive accuracy of M*.

Stress response in bacteria (Bacterial)

This model describes the stress response in *Bacillus subtilis* [59]. It was used by Mangan et al [31] to illustrate how an implicit SINDy approach was able to infer biological nonlinear

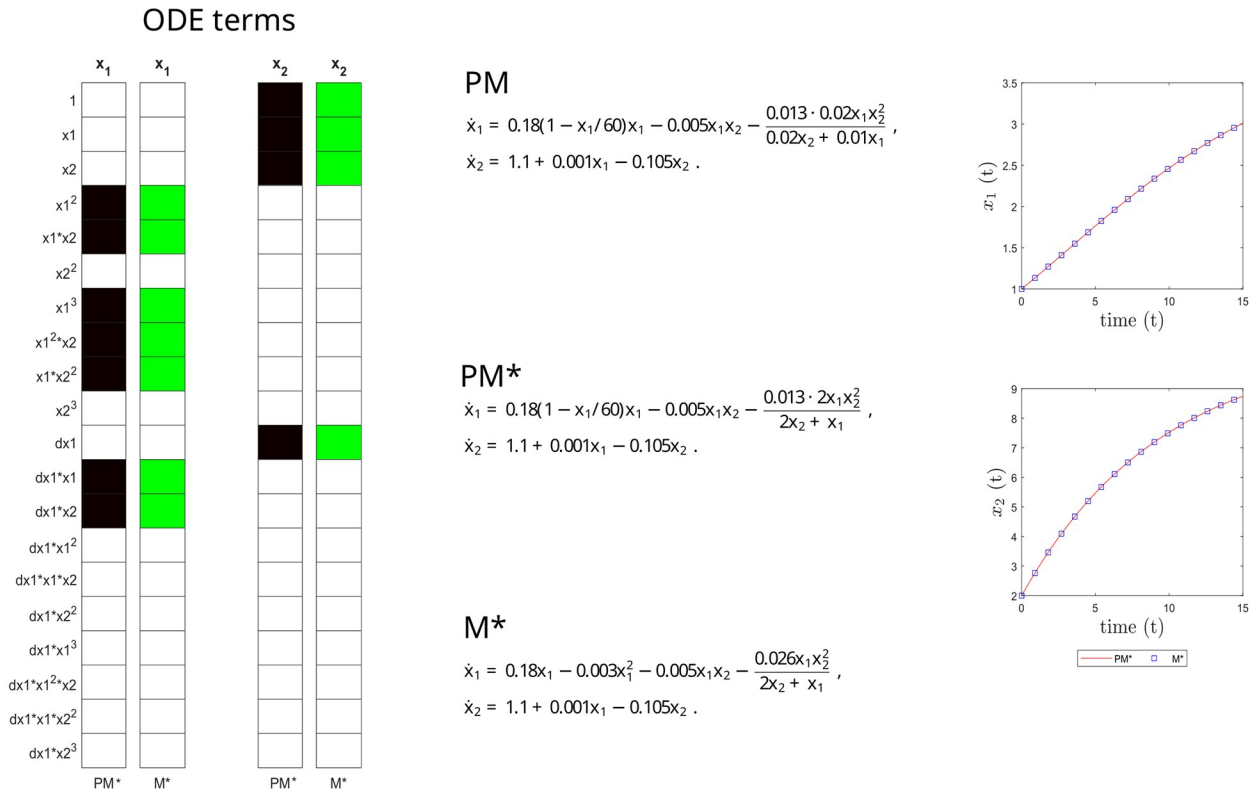


Fig 4. Immunity case study. Structural accuracy: on the left, active terms in ξ (non-zero terms of the prior model PM in black, and of the inferred model M* in green). Parameter accuracy: center, matching parametric ODEs for PM and M*. Predictive accuracy: on the right, time evolution of the different states (x_1 and x_2) of the PM and M* models.

<https://doi.org/10.1371/journal.pcbi.1011014.g004>

dynamics. Under nutrient limitation, the majority of *B. subtilis* cells switch to sporulation, but a small fraction switch to an alternative behaviour, the so called state of competence, in which they are capable of taking up extracellular DNA. This latter fraction might subsequently return to vegetative growth. Süel et al [59] described the regulatory system of this mechanism using a dynamic model with two states. In dimensionless form, the ground truth (GT) model for this example is:

$$\dot{x}_1 = a_1 + \frac{a_2x_1^2}{a_3 + x_1^2} - \frac{x_1}{1 + x_1 + x_2}, \tag{30a}$$

$$\dot{x}_2 = \frac{b_1}{1 + b_2x_1^5} - \frac{x_2}{1 + x_1 + x_2}. \tag{30b}$$

where x_1 and x_2 represent the concentration levels of the ComK and ComS proteins. The rational terms arise from time-scale separation assumptions about the regulation: an autoregulatory positive feedback loop of ComK plus and indirect negative feedback loop mediated by ComS. In Eq 30a, a_1 corresponds to the minimal rate of ComK production. The second term describes the autoregulation (via a positive feedback loop) of ComK activating its own production, where a_2 is the fully activated rate of ComK generation. The first term in Eq 30b describes the negative feedback loop regulating the repression of ComS, where b_1 is the maximum rate of ComS expression. Both the auto-activation of ComK and the repression of ComS follow Hill

kinetics where the exponent indicates the level of cooperativity (2 and 5, respectively). The last term in both equations represents the degradation of both ComK and ComS.

We again consider PM = GT and check the identifiability of PM. Considering unknown parameters a_1, a_2, a_3, b_1 and b_2 , the model is fully identifiable, thus $PM^* = GT$.

Next, SINDy-PI is applied to the training data generated using GT, obtaining the following candidate model (CM):

$$\dot{x}_1 = p_1 + \frac{p_2}{p_3 + p_4 x_1^2} + \frac{p_5 x_2 + p_6}{p_7 + p_8 x_1 + p_9 x_2}, \tag{31a}$$

$$\dot{x}_2 = \frac{p_{10}}{p_{11} + p_{12} x_1^5} + \frac{p_{13} x_2}{p_{14} + p_{15} x_1 + p_{16} x_2}, \tag{31b}$$

This CM has 16 parameters, $p_j, j = 1, \dots, 16$, and the SIO analysis reveals that all of them are non-identifiable, with the exception of p_1 . The reformulation step indicated that we can obtain an identifiable model with four scaling transformations, one per rational term, i.e. the second term in Eq 31a is scaled by $p_j, j \in [2, 3, 4]$, and the third term by $p_k, k \in [5, \dots, 9]$. In Eq 31b the same strategy is applied for $p_l, l \in [10, 11, 12]$, and $p_m, m \in [13, \dots, 16]$. That is:

$$\dot{x}_1 = p_1 + \frac{\frac{1}{p_j} p_2}{\frac{1}{p_j} (p_3 + p_4 x_1^2)} + \frac{\frac{1}{p_k} (p_5 x_2 + p_6)}{\frac{1}{p_k} (p_7 + p_8 x_1 + p_9 x_2)}, \tag{32a}$$

$$\dot{x}_2 = \frac{\frac{1}{p_l} p_{10}}{\frac{1}{p_l} (p_{11} + p_{12} x_1^5)} + \frac{\frac{1}{p_m} p_{13} x_2}{\frac{1}{p_m} (p_{14} + p_{15} x_1 + p_{16} x_2)}. \tag{32b}$$

Choosing $j = 4, k = 7, l = 11, m = 14$, the resulting structurally identifiable model is:

$$\dot{x}_1 = p_1 + \frac{p_2^*}{p_3^* + x_1^2} + \frac{p_5^* x_2 + p_6^*}{1 + p_8^* x_1 + p_9^* x_2} \tag{33a}$$

$$\dot{x}_2 = \frac{p_{10}^*}{1 + p_{12}^* x_1^5} + \frac{p_{13}^* x_2}{1 + p_{15}^* x_1 + p_{16}^* x_2} \tag{33b}$$

where * denotes a reparameterized parameter.

This reformulated model is now fully identifiable, but no longer directly interpretable: Eq 33a does not explicitly have the term involving the autoregulation of ComK. By means of the reformulation procedure, we are able to recover the autoregulation and degradation terms as in Eq 30a:

$$\dot{x}_1 = p_1 + \frac{p_2^* x_1^2}{p_3^* + x_1^2} + \frac{p_5^* x_1}{1 + p_8^* x_1 + p_9^* x_2} \tag{34a}$$

$$\dot{x}_2 = \frac{p_{10}^*}{1 + p_{12}^* x_1^5} + \frac{p_{13}^* x_2}{1 + p_{15}^* x_1 + p_{16}^* x_2} \tag{34b}$$

This reformulated model M^* is structurally identifiable and interpretable, and equivalent to the PM. Fig 5 shows the assessment of the structural, parametric and predictive accuracy of the inferred model M^* .

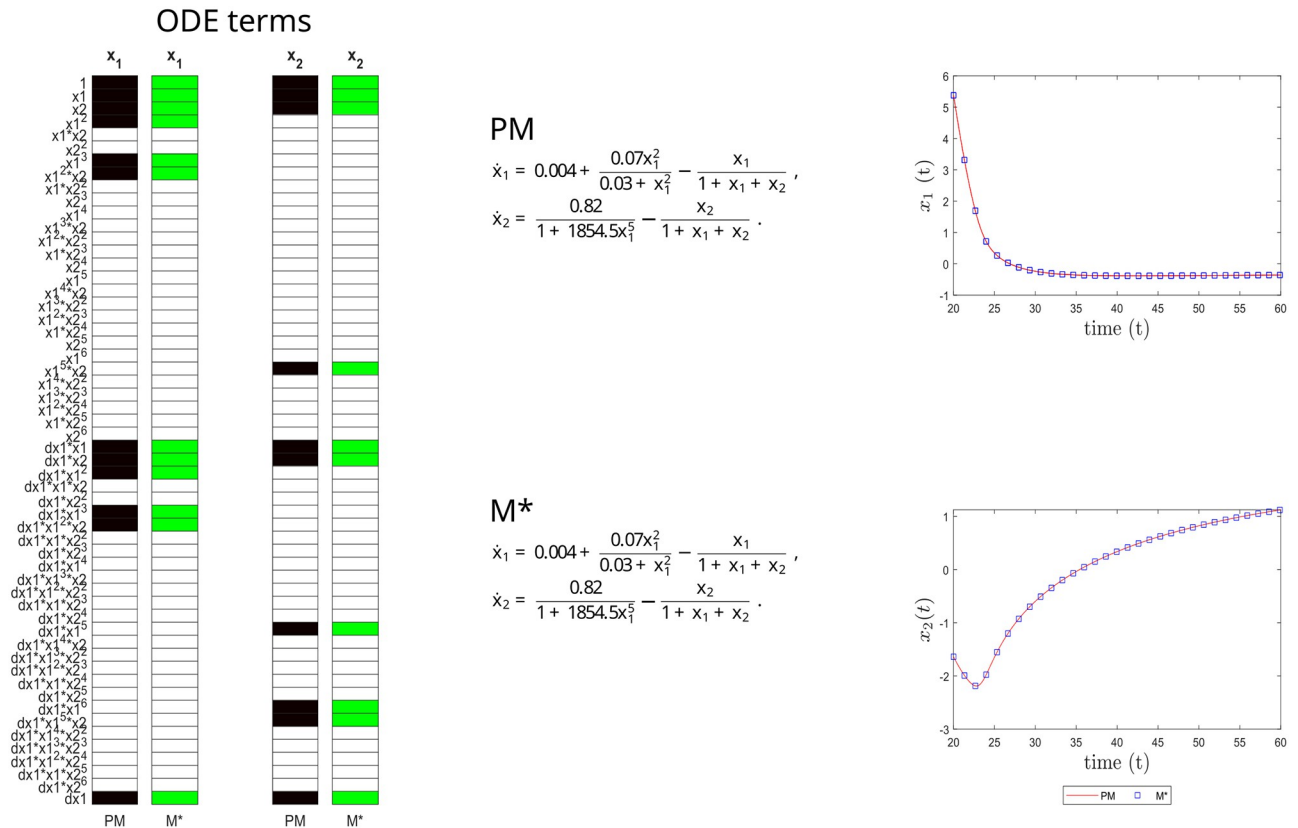


Fig 5. Bacterial case study. Structural accuracy: on the left, active terms in ξ (non-zero terms of the prior model PM in black, and of the inferred model M* in green). Parameter accuracy: center, matching parametric ODEs for PM and M*. Predictive accuracy: on the right, time evolution of the different states (x_1 and x_2) of the PM and M* models.

<https://doi.org/10.1371/journal.pcbi.1011014.g005>

Microbial growth (Microbial)

This case study considers microbial growth in a batch reactor, as presented by [64] and later used by [60] to study identifiable reparameterizations of unidentifiable systems. The following model describes the dynamics of microbial and substrate concentrations assuming Monod kinetics (similar in functional form to Michaelis-Menten enzyme kinetics):

$$\dot{x}_1 = \frac{\mu x_2 x_1}{K_s + x_2} - K_d x_1, \tag{35a}$$

$$\dot{x}_2 = -\frac{\mu x_2 x_1}{\gamma(K_s + x_2)}. \tag{35b}$$

where x_1 and x_2 represent the concentrations of microorganisms and growth-limiting substrate, respectively. The rational term in Eq 35a is the Monod kinetic term, where μ is the maximum growth velocity and K_s the substrate concentration corresponding to $\frac{1}{2}\mu$. In Eq 35a, the same rational term appears scaled by γ (the yield coefficient) to represent the depletion of substrate. The last term in Eq 35a describes the death of microorganisms assuming first order kinetics where K_d is the decay rate.

We consider a prior model (PM) that matches the ground truth (GT) model, Eqs 35a and 35b. When the initial conditions are known and different from zero, our algorithm confirms that this PM is structurally identifiable. Next, our workflow discovers the following dynamics using SINDy-PI:

$$\dot{x}_1 = p_1 x_1 + \frac{p_2 x_1}{p_3 + p_4 x_2}, \tag{36a}$$

$$\dot{x}_2 = \frac{p_5 x_1 x_2}{p_6 + p_7 x_2}. \tag{36b}$$

Next, the FISPO step finds that only p_1 is identifiable, i.e. parameters $p_i, i = 2, \dots, 7$ are unidentifiable. The reformulation step finds that it is possible to find an identifiable form by scaling each rational term by the same unidentifiable parameter. For simplicity, we have chosen that $p_2^* = \frac{p_2}{p_3}, p_4^* = \frac{p_4}{p_3}, p_5^* = \frac{p_5}{p_7}$ and $p_7^* = \frac{p_7}{p_6}$. Then, the resulting structurally identifiable model is:

$$\dot{x}_1 = p_1 x_1 + \frac{p_2^* x_1}{1 + p_3^* x_2}, \tag{37a}$$

$$\dot{x}_2 = \frac{p_5^* x_1 x_2}{1 + p_7^* x_2}. \tag{37b}$$

However, Eqs 37a and 37b are not directly interpretable because they do not contain the expected Monod kinetics terms explicitly. Next, the reformulation step finds an equivalent structure which is both interpretable and identifiable (M^*):

$$\dot{x}_1 = p_1 x_1 + \frac{p_2 x_1 x_2}{1 + p_4^* x_2}, \tag{38a}$$

$$\dot{x}_2 = \frac{p_5^* x_1 x_2}{1 + p_7^* x_2}. \tag{38b}$$

This inferred model (M^*) is compared to the ground truth in Fig 6, confirming its structural, parametric and predictive accuracy.

Cell cycle in the colonic crypt (Crypt)

This example considers a cell population model describing the cell renewal cycle in the colonic crypt [61]. This cycle is heavily regulated and the model was used to explain the rupture of homeostasis and the initiation of tumorigenesis. The equations describing the dynamics are:

$$\dot{x}_1 = (a_3 - a_1 - a_2)x_1 - \frac{k_0 x_1^2}{1 + m_0 x_1}, \tag{39a}$$

$$\dot{x}_2 = (b_3 - b_1 - b_2)x_2 + a_2 x_1 - \frac{k_1 x_2^2}{1 + m_1 x_2} + \frac{k_0 x_1^2}{1 + m_0 x_1}, \tag{39b}$$

$$\dot{x}_3 = -g x_3 + b_2 x_2 + \frac{k_1 x_2^2}{1 + m_1 x_2}. \tag{39c}$$

where the state variables represent the populations of stem cells (x_1), semi-differentiated cells (x_2), and fully-differentiated cells (x_3). Stem cells have first order kinetics for renewal (rate

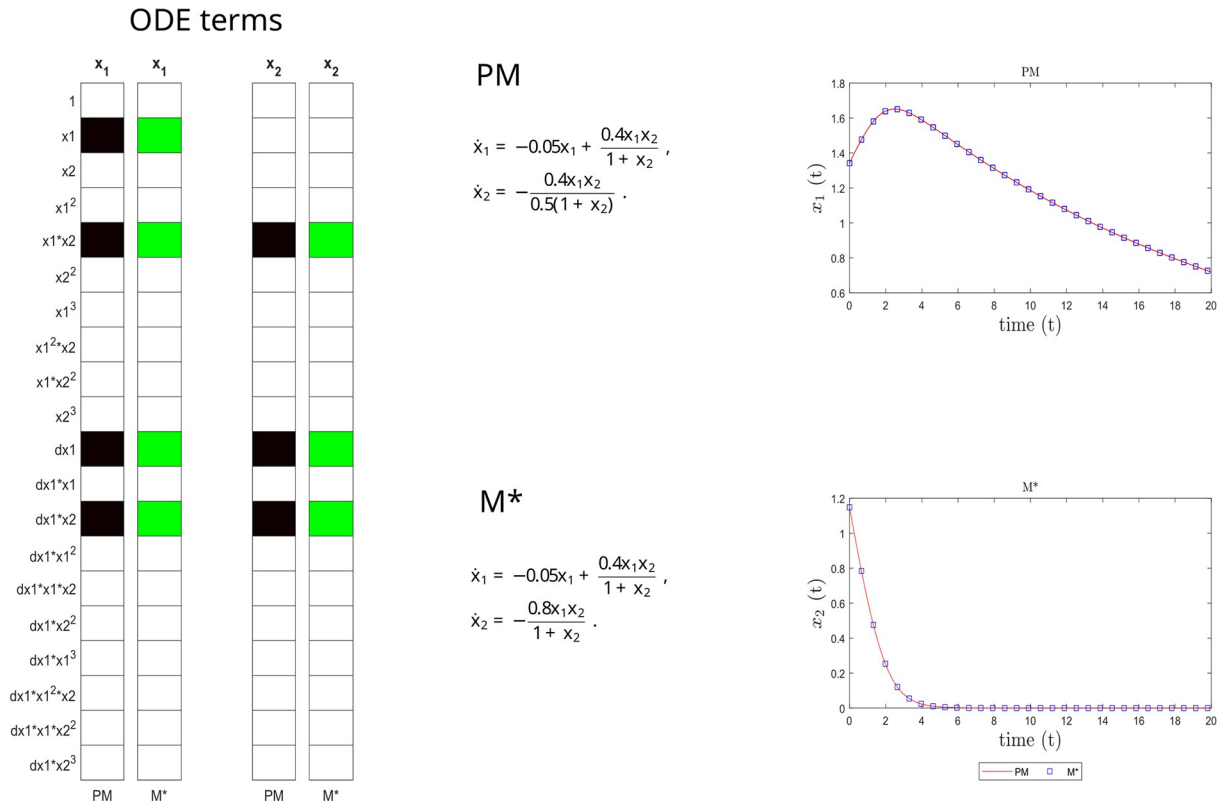


Fig 6. Microbial case study. Structural accuracy: on the left, active terms in ξ (non-zero terms of the prior model PM in black, and of the inferred model M* in green). Parameter accuracy: center, matching parametric ODEs for PM and M*. On the right, predictive accuracy: time evolution of the different states (x_1 and x_2) of the PM and M* models.

<https://doi.org/10.1371/journal.pcbi.1011014.g006>

given parameter a_3), differentiation (parameter a_2), and death (parameter a_1). Semi-differentiated cells have similar renewal, differentiation and death kinetics (with parameters b_i), plus a source term due to the differentiation of stem cells. Fully differentiated cells are generated from semi-differentiated cells with first order rate b_2 and removed with a rate modulated by parameter g . The rational terms correspond to saturating feedback mechanism in the differentiation rates.

We take the above model as GT, and PM = GT. Our algorithm finds that this PM is not structurally identifiable: it is not possible to uniquely infer a_1, a_3, b_1 and b_3 due to the presence of a translation symmetry. Next, the reformulation step finds a reparameterized prior model (PM*):

$$\dot{x}_1 = (a_3^* - a_2)x_1 - \frac{k_0x_1^2}{1 + m_0x_1}, \tag{40a}$$

$$\dot{x}_2 = (b_3^* - b_2)x_2 + a_2x_1 - \frac{k_1x_2^2}{1 + m_1x_2} + \frac{k_0x_1^2}{1 + m_0x_1}, \tag{40b}$$

$$\dot{x}_3 = -gx_3 + b_2x_2 + \frac{k_1x_2^2}{1 + m_1x_2}; \tag{40c}$$

where $a_3^* = a_3 - a_1$ and $b_3^* = b_3 - a_1$. Next, SINDy-PI is applied to the training data, obtaining

the following candidate model (CM):

$$\dot{x}_1 = p_1 + p_2 x_1 + \frac{p_3}{p_5 + x_1 p_4}, \tag{41a}$$

$$\dot{x}_2 = \frac{p_6 x_2 + p_7 x_1 + p_8 x_1 x_2 + p_9 x_2^2 + p_{10} x_1^2 + p_{11} x_1^2 x_2}{p_{12} + p_{13} x_1 + p_{14} x_2 + p_{15} x_1 x_2}, \tag{41b}$$

$$\dot{x}_3 = p_{16} x_2 + p_{17} x_3 + \frac{p_{18}}{p_{19} x_2 + p_{20}} + p_{21}. \tag{41c}$$

Considering $p_i, i = 1, \dots, 21$ as unknown parameters, the FISPO algorithm indicates that only p_1, p_2, p_{16}, p_{17} and p_{21} are structurally identifiable. The reformulation step finds the following structurally identifiable alternative:

$$\dot{x}_1 = p_1 + p_2 x_1 + \frac{p_3^*}{1 + x_1 p_4^*}, \tag{42a}$$

$$\dot{x}_2 = \frac{p_5^* x_2 + p_6^* x_1 + p_7^* x_1 x_2 + p_8^* x_2^2 + p_9^* x_1^3 + p_{10}^* x_1^2 + p_{11}^* x_1^2 x_2}{1 + p_{13}^* x_1 + p_{14}^* x_2 + p_{15}^* x_1 x_2}, \tag{42b}$$

$$\dot{x}_3 = p_{16} x_2 + p_{17} x_3 + \frac{p_{18}^*}{p_{19}^* x_2 + 1} + p_{21}. \tag{42c}$$

The above model is not directly interpretable, but the reformulation process is able to find the following interpretable and identifiable reformulation M^* :

$$\dot{x}_1 = p_1 x_1 + \frac{p_2 x_1^2}{1 + x_1 p_4}, \tag{43a}$$

$$\dot{x}_2 = p_5 x_2 + p_6 x_1 + \frac{p_7 x_2^2}{1 + p_8 x_2} + \frac{p_9 x_1^2}{1 + p_{10} x_1}, \tag{43b}$$

$$\dot{x}_3 = p_{11} x_2 + p_{12} x_3 + \frac{p_{13} x_2^2}{p_{14} x_2 + 1}. \tag{43c}$$

This discovered model M^* is fully equivalent to the identifiable version of the ground truth model in terms of structural, parametric and predictive accuracy, as shown in Fig 7. This example reinforces the importance of checking the identifiability of both the ground truth and the inferred model.

Oscillations in yeast glycolysis (Glycolysis)

Glycolysis is the transformation (in a series of reactions catalyzed by enzymes) of glucose into smaller molecules to produce energy for the cell. In many cell types, glycolysis exhibits oscillations in the concentrations of many intermediate metabolites. This phenomena has been particularly well studied in yeast cells. Wolf and Heinrich [62] studied the oscillatory dynamics of a simplified reaction scheme for yeast glycolysis under anaerobic conditions, where alcoholic

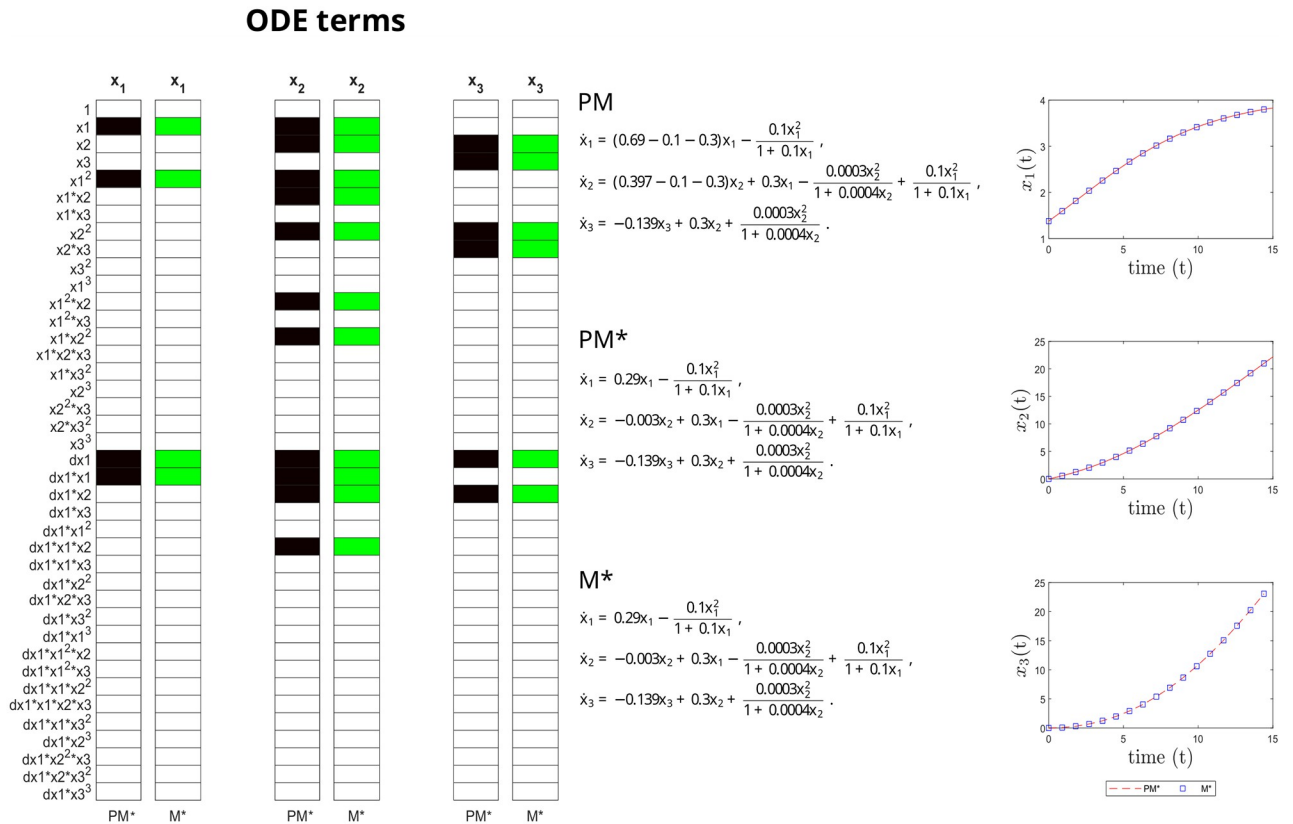


Fig 7. Crypt case study. Structural accuracy: on the left, active terms in ξ (non-zero terms of the prior model PM in black, and of the inferred model M* in green). Parameter accuracy: center, matching parametric ODEs for PM and M*. Predictive accuracy: on the right, time evolution of the different states (x_1 , x_2 and x_3) of the PM and M* models.

<https://doi.org/10.1371/journal.pcbi.1011014.g007>

fermentation takes place, proposing the following mathematical description:

$$\dot{x}_1 = c_1 + \frac{c_2 x_1 x_6}{1 + c_3 x_6^4}, \tag{44a}$$

$$\dot{x}_2 = \frac{d_1 x_1 x_6}{1 + d_2 x_6^4} + d_3 x_2 - d_4 x_2 x_7, \tag{44b}$$

$$\dot{x}_3 = e_1 x_2 + e_2 x_3 + e_3 x_2 x_7 + e_4 x_3 x_6, \tag{44c}$$

$$\dot{x}_4 = f_1 x_3 + e_2 x_4 + f_3 x_5 + f_4 x_3 x_6 + f_5 x_4 x_7, \tag{44d}$$

$$\dot{x}_5 = g_1 x_4 + g_2 x_5, \tag{44e}$$

$$\dot{x}_6 = h_3 x_3 + h_5 x_6 + h_4 x_3 x_6 + \frac{h_1 x_1 x_6}{1 + h_2 x_6^4}, \tag{44f}$$

$$\dot{x}_7 = j_1 x_2 + j_2 x_2 x_7 + j_3 x_4 x_7. \tag{44g}$$

where the state variables represent the concentrations in the cell of glucose (x_1), the pool of

triose phosphates (x_2), 1,3-bisphosphoglycerate (x_3), pool of pyruvate and acetaldehyde (x_4), NADH (x_5), ATP (x_6), and x_7 represents the pool of pyruvate and acetaldehyde in the external solution. We consider here the same formulation and parameter values as in [31, 37].

We take the above as GT, and assume PM = GT. Considering all parameters as unknown ($c_i, i = 1, 2, 3; d_i, i = 1, \dots, 4; e_i, i = 1, \dots, 4; f_i, i = 1, \dots, 5; g_i, i = 1, 2; h_i, i = 1, \dots, 5$ and $j_i, i = 1, 2, 3$), our algorithm confirms that the model is structurally identifiable and observable, i.e. $PM^* = PM$.

This problem is quite challenging for SINDy-PI due to its large number of states and parameters, and the large degree in several terms, leading to a very large library of candidate functions (over 3000 terms). However, it is able to correctly recover the following candidate model (CM):

$$\dot{x}_1 = p_1 + \frac{p_2 x_1 x_6}{p_4 + p_3 x_6^4}, \tag{45a}$$

$$\dot{x}_2 = p_5 x_2 + p_6 x_2 x_7 + \frac{p_7 x_1 x_6}{p_8 x_6^4 + p_9}, \tag{45b}$$

$$\dot{x}_3 = p_{10} x_2 + p_{11} x_3 + p_{12} x_2 x_7 + p_{13} x_3 x_6, \tag{45c}$$

$$\dot{x}_4 = p_{14} x_3 + p_{15} x_4 + p_{16} x_5 + p_{17} x_3 x_6 + p_{18} x_4 x_7, \tag{45d}$$

$$\dot{x}_5 = p_{19} x_4 + p_{20} x_5, \tag{45e}$$

$$\dot{x}_6 = p_{21} x_3 + p_{22} x_6 x_3 + p_{23} x_6 + \frac{p_{24} x_6 x_1}{p_{25} + p_{26} x_6^4}, \tag{45f}$$

$$\dot{x}_7 = p_{27} x_2 + p_{28} x_2 x_7 + p_{29} x_4 x_7. \tag{45g}$$

This model has 29 inferred coefficients, $p_i, i = 1, \dots, 29$. Our algorithm analyzes their identifiability and finds that the parameters with indices $i = 2, 3, 4, 7, 8, 9, 24, 25, 26$ are unidentifiable. Next, the reformulation step finds possible reparameterizations by scaling the rational terms. That is, considering the first rational term scaled by p_4 , then $p_2^* = \frac{p_2}{p_4}$ and $p_3^* = \frac{p_3}{p_4}$; scaling the second term with p_9 produces $p_8^* = \frac{p_8}{p_9}$ and $p_7^* = \frac{p_7}{p_9}$; and using p_{25} yields $p_{24}^* = \frac{p_{24}}{p_{25}}$ and $p_{26}^* = \frac{p_{26}}{p_{25}}$ for the last rational term. The end result is an interpretable and identifiable model M^* :

$$\dot{x}_1 = p_1 + \frac{p_2^* x_1 x_6}{1 + p_3^* x_6^4}, \tag{46a}$$

$$\dot{x}_2 = p_5 x_2 + p_6 x_2 x_7 + \frac{p_7^* x_1 x_6}{p_8^* x_6^4 + 1}, \tag{46b}$$

$$\dot{x}_3 = p_{10} x_2 + p_{11} x_3 + p_{12} x_2 x_7 + p_{13} x_3 x_6, \tag{46c}$$

$$\dot{x}_4 = p_{14} x_3 + p_{15} x_4 + p_{16} x_5 + p_{17} x_3 x_6 + p_{18} x_4 x_7, \tag{46d}$$

$$\dot{x}_5 = p_{19} x_4 + p_{20} x_5, \tag{46e}$$

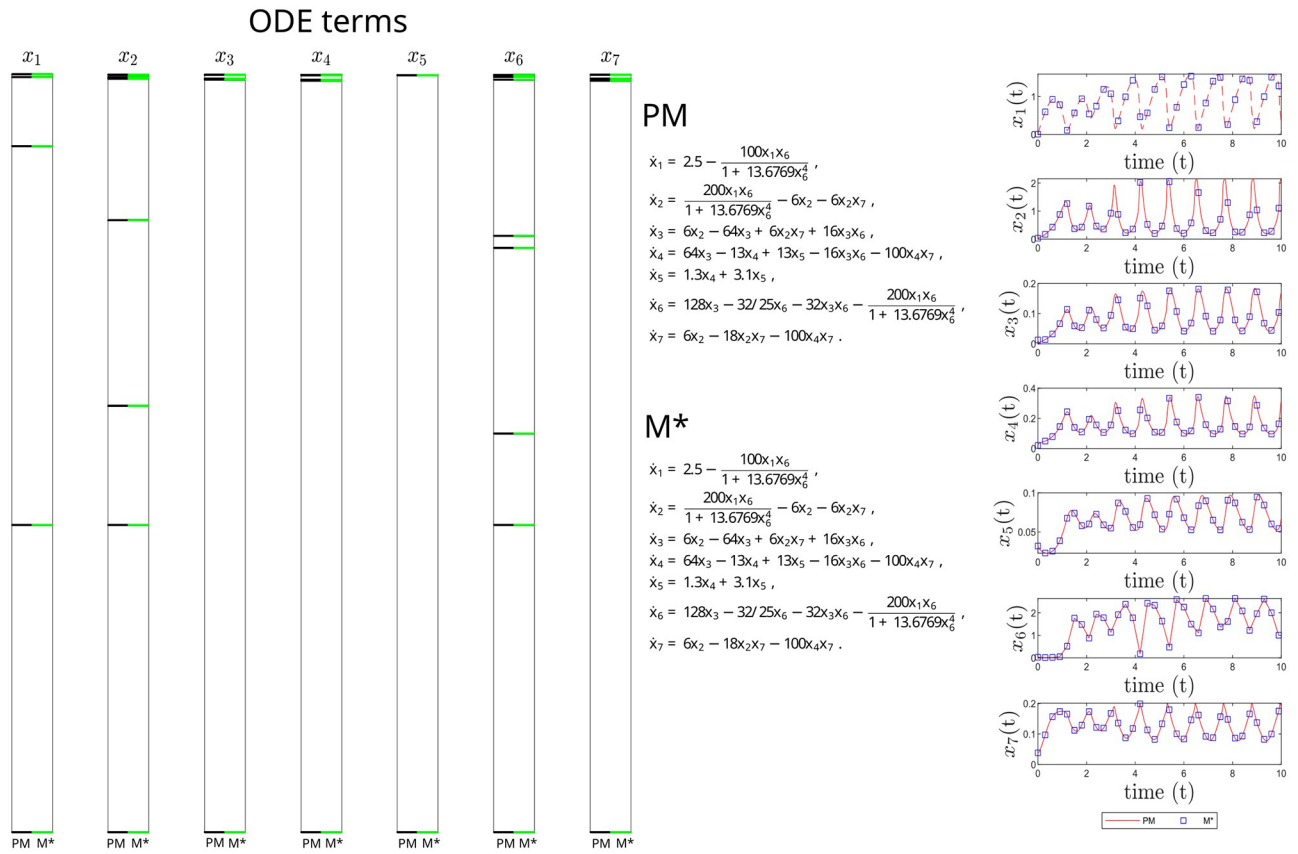


Fig 8. Yeast-Glycolysis case study. Structural accuracy: on the left, active terms in ξ (non-zero terms of the prior model PM in black, and of the inferred model M* in green). Due to the large number of terms in ξ , the candidate functions are not shown. Parameter accuracy: center, matching parametric ODEs for PM and M*. Predictive accuracy: on the right, time evolution of the different states ($x_1, x_2, x_3, x_4, x_5, x_6$ and x_7) of the PM and M* models.

<https://doi.org/10.1371/journal.pcbi.1011014.g008>

$$\dot{x}_6 = p_{21}x_3 + p_{22}x_6x_3 + p_{23}x_6 + \frac{p_{24}x_6x_1}{1 + p_{26}^*x_6^4}, \tag{46f}$$

$$\dot{x}_7 = p_{27}x_2 + p_{28}x_2x_7 + p_{29}x_4x_7 \tag{46g}$$

Fig 8 illustrates the excellent structural, parametric and predictive accuracy of the inferred model.

Discussion

In recent years, innovations in numerical methods and machine learning have been combined to improve our ability to understand complex systems. Currently, the three main classes of methods to learn equations from data are symbolic regression [65], neural-network approaches [13], and library-based sparse regression [12]. Recent reviews of these categories and their overlaps can be found in [66, 67].

In particular, data-driven model discovery methods for nonlinear dynamic systems have seen very significant advancements [22, 24]. The field has seen growth in terms of both sophistication and the range of applications. The fundamental aim remains the same: to discern the

underlying mathematical models that govern the behaviour of complex, possibly high-dimensional and nonlinear systems, using measurement data. In this study we have investigated certain aspects of automatic model discovery techniques to derive mechanistic models of biological systems from time-series data. Specifically, we have focused on possible structural deficiencies of their end result, the inferred model equations. As a reference method we have chosen SINDy-PI [37], a recent sparse regression-based methodology that is particularly suited for computational biology due to its ability to capture complex nonlinearities and rational terms.

However, it should be noted that our approach can be combined with other model discovery methods. In any case, SINDy-PI has several advantages over other model discovery methods. First, it is several orders of magnitude more robust to noise than previous approaches based on sparse regression. This means that it can learn implicit ordinary and partial differential equations and conservation laws from limited and noisy data. Second, it can discover models with very complex structure, including implicit dynamics and rational nonlinearities (such as e.g. Michaelis-Menten kinetics), which are common in biological applications. Third, it is still quite computationally efficient thanks to its parallel nature and the exploitation of a library of canonical nonlinear terms. Such a library is particularly attractive when modelling the dynamics of biological networks based on mechanistic assumptions, such as mass-action kinetics.

Since by design SINDy-PI enforces parsimonious models (with the lowest complexity to support the data), it usually produces interpretable equations with excellent predictive power. However, we have shown that sometimes these models lack structural identifiability, which means that using the discovered model structure for parameter estimation might give wrong estimates, compromising its usefulness and reliability.

To address this issue we have presented a methodology that, combined with SINDy-PI, facilitates the inference of identifiable and interpretable dynamic models. Our method integrates symbolic algorithms that analyse a model's structural identifiability and observability (SIO), reparameterize it to achieve SIO if needed, and reformulate it to make it biologically interpretable. We have illustrated its use in two scenarios, with and without prior knowledge, using six challenging case studies corresponding to different kinds of biological systems, including complex regulatory mechanisms.

Our results highlight additional challenges due to non-obvious issues in the relationship between model reformulation, identifiability and interpretability, and show how our approach is able to successfully surmount them. Importantly, our method is modular and can be easily integrated with other model discovery strategies. While we have demonstrated its application in combination with SINDy-PI, other methods could have been used as well. Furthermore, its calculations are entirely symbolic, i.e. they are not affected by numerical issues caused by insufficient or noisy data (which do however limit the application of the accompanying model discovery method).

Future work will be devoted to model discovery in partially observed systems, where the structural identifiability problem will surely be exacerbated, and observability issues—i.e. the impossibility of inferring some of the unmeasured state variables—are to be expected. It should be noted that, as a matter of fact, our methodology is applicable to partially observed systems in its present form. However, model discovery for such systems is still in its infancy (see the recent work by [68]), hence in this study we have considered fully observed systems. Another possible area of improvement is computational efficiency. While our pipeline can be applied to systems with several states and a few dozen parameters, as demonstrated with the Glycolysis example, scaling up to larger models is challenging. The main bottleneck is currently the model reparameterization step performed with AutoRepar, which involves

symbolic computations that can be very memory-consuming. We are working on improving the efficiency of the algorithms in order to alleviate its computational cost.

Other important avenues of research which are currently being explored include (i) improved approaches for the design of the library of candidate functions [69], (ii) better incorporation of partial prior knowledge [70], and (iii) taking into account noisy and missing data, uncertainty quantification and applications to real-world data-sets [71–73]. Since identifiability and observability play a major role in these scenarios, we believe that our methodology will be a useful tool in these explorations.

Supporting information

S1 File. Supporting material document.
(PDF)

Author Contributions

Conceptualization: Alejandro F. Villaverde, Julio R. Banga.

Data curation: Gemma Massonis.

Formal analysis: Alejandro F. Villaverde, Julio R. Banga.

Funding acquisition: Alejandro F. Villaverde, Julio R. Banga.

Investigation: Alejandro F. Villaverde, Julio R. Banga.

Methodology: Gemma Massonis, Alejandro F. Villaverde, Julio R. Banga.

Project administration: Alejandro F. Villaverde, Julio R. Banga.

Resources: Julio R. Banga.

Software: Gemma Massonis.

Supervision: Alejandro F. Villaverde, Julio R. Banga.

Validation: Gemma Massonis.

Visualization: Gemma Massonis.

Writing – original draft: Gemma Massonis, Alejandro F. Villaverde, Julio R. Banga.

Writing – review & editing: Alejandro F. Villaverde, Julio R. Banga.

References

1. DiStefano JJ. *Dynamic Systems Biology Modeling and Simulation*. Academic Press; 2015.
2. Ingalls BP. *Mathematical Modeling in Systems Biology: An Introduction*. MIT Press; 2022.
3. Strogatz SH. *Nonlinear Dynamics and Chaos: With Applications to Physics, Biology, Chemistry, and Engineering*. Westview Press; 2014.
4. Vittadello ST, Stumpf MPH. Open problems in mathematical biology. *Math Biosci*. 2022; 354:108926. <https://doi.org/10.1016/j.mbs.2022.108926> PMID: 36377100
5. Langley P. Data-Driven Discovery of Physical Laws. *Cognitive Science*. 1981; 5(1):31–54. <https://doi.org/10.1111/j.1551-6708.1981.tb00869.x>
6. Crutchfield JP, McNamara B. Equations of motion from a data series. *Complex systems*. 1987; 1:417–452.
7. Koza J, Keane MA, Rice JP. Performance improvement of machine learning via automatic discovery of facilitating functions as applied to a problem of symbolic system identification. In: *IEEE International Conference on Neural Networks*. IEEE; 1993. p. 191–198.

8. Bongard J, Lipson H. Automated reverse engineering of nonlinear dynamical systems. *Proc Natl Acad Sci U S A*. 2007; 104(24):9943–9948. <https://doi.org/10.1073/pnas.0609476104> PMID: 17553966
9. Udrescu SM, Tegmark M. AI Feynman: A physics-inspired method for symbolic regression. *Sci Adv*. 2020; 6(16):eaay2631. <https://doi.org/10.1126/sciadv.aay2631> PMID: 32426452
10. Džeroski S, Langley P, Todorovski L. Computational discovery of scientific knowledge. In: *Computational discovery of scientific knowledge*. Springer; 2007. p. 1–14.
11. Brencic J, Todorovski L, Džeroski S. Probabilistic grammars for equation discovery. *Knowledge-Based Systems*. 2021; 224:107077. <https://doi.org/10.1016/j.knosys.2021.107077>
12. Brunton SL, Proctor JL, Kutz JN. Discovering governing equations from data by sparse identification of nonlinear dynamical systems. *Proc Natl Acad Sci U S A*. 2016; 113(15):3932–3937. <https://doi.org/10.1073/pnas.1517384113> PMID: 27035946
13. Raissi M, Perdikaris P, Karniadakis GE. Physics informed deep learning (part i): Data-driven solutions of nonlinear partial differential equations. *arXiv preprint arXiv:1711.10561*. 2017;.
14. Raissi M, Yazdani A, Karniadakis GE. Hidden fluid mechanics: Learning velocity and pressure fields from flow visualizations. *Science*. 2020; 367(6481):1026–1030. <https://doi.org/10.1126/science.aaw4741> PMID: 32001523
15. Rackauckas C, Ma Y, Martensen J, Warner C, Zubov K, Supekar R, et al. Universal differential equations for scientific machine learning. *arXiv preprint arXiv:200104385*. 2020;.
16. Bhourri MA, Perdikaris P. Gaussian processes meet NeuralODEs: a Bayesian framework for learning the dynamics of partially observed systems from scarce and noisy data. *Philos Trans A Math Phys Eng Sci*. 2022; 380(2229):20210201. PMID: 35719075
17. VandenHeuvel DJ, Drovandi C, Simpson MJ. Computationally efficient mechanism discovery for cell invasion with uncertainty quantification. *PLoS Comput Biol*. 2022; 18(11):e1010599. <https://doi.org/10.1371/journal.pcbi.1010599> PMID: 36383612
18. Pan W, Yuan Y, Gonçalves J, Stan GB. A sparse Bayesian approach to the identification of nonlinear state-space systems. *IEEE Transactions on Automatic Control*. 2015; 61(1):182–187. <https://doi.org/10.1109/TAC.2015.2426291>
19. Zhang S, Lin G. Robust data-driven discovery of governing physical laws with error bars. *Proc Math Phys Eng Sci*. 2018; 474(2217):20180305. <https://doi.org/10.1098/rspa.2018.0305> PMID: 30333709
20. Guimerà R, Reichardt I, Aguilar-Mogas A, Massucci FA, Miranda M, Pallarès J, et al. A Bayesian machine scientist to aid in the solution of challenging scientific problems. *Sci Adv*. 2020; 6(5):eaav6971. <https://doi.org/10.1126/sciadv.aav6971> PMID: 32064326
21. Džeroski S, Todorovski L. Equation discovery for systems biology: finding the structure and dynamics of biological networks from time course data. *Current Opinion in Biotechnology*. 2008; 19(4):360–368. <https://doi.org/10.1016/j.copbio.2008.07.002> PMID: 18672061
22. North JS, Wikle CK, Schliep EM. A Review of Data-Driven Discovery for Dynamic Systems. *arXiv preprint arXiv:221010663*. 2022;.
23. Willard J, Jia X, Xu S, Steinbach M, Kumar V. Integrating physics-based modeling with machine learning: A survey. *arXiv preprint arXiv:200304919*. 2020;1(1):1–34.
24. Brunton SL, Kutz JN. *Data-driven science and engineering: Machine learning, dynamical systems, and control*; 2nd edition. Cambridge University Press; 2022.
25. Ghadami A, Epureanu BI. Data-driven prediction in dynamical systems: recent developments. *Philosophical Transactions of the Royal Society A*. 2022; 380(2229):20210213. <https://doi.org/10.1098/rsta.2021.0213> PMID: 35719077
26. Naozuka GT, Rocha HL, Silva RS, Almeida RC. SINDy-SA framework: enhancing nonlinear system identification with sensitivity analysis. *Nonlinear Dyn*. 2022; 110(3):2589–2609. <https://doi.org/10.1007/s11071-022-07755-2> PMID: 36060282
27. Villaverde AF, Banga JR. Reverse engineering and identification in systems biology: strategies, perspectives and challenges. *Journal of the Royal Society Interface*. 2014; 11(91):20130505. <https://doi.org/10.1098/rsif.2013.0505> PMID: 24307566
28. Kirk P, Silk D, Stumpf MP. Reverse engineering under uncertainty. In: *Uncertainty in biology*. Springer; 2016. p. 15–32.
29. Mercatelli D, Scalambra L, Triboli L, Ray F, Giorgi FM. Gene regulatory network inference resources: A practical overview. *Biochimica et Biophysica Acta (BBA)-Gene Regulatory Mechanisms*. 2020; 1863(6):194430. <https://doi.org/10.1016/j.bbagr.2019.194430> PMID: 31678629
30. Sunnåker M, Zamora-Sillero E, Dechant R, Ludwig C, Busetto AG, Wagner A, et al. Automatic generation of predictive dynamic models reveals nuclear phosphorylation as the key Msn2 control mechanism. *Science signaling*. 2013; 6(277):ra41–ra41. PMID: 23716718

31. Mangan NM, Brunton SL, Proctor JL, Kutz JN. Inferring biological networks by sparse identification of nonlinear dynamics. *IEEE Transactions on Molecular, Biological and Multi-Scale Communications*. 2016; 2(1):52–63. <https://doi.org/10.1109/TMBMC.2016.2633265>
32. Daniels BC, Ryu WS, Nemenman I. Automated, predictive, and interpretable inference of escape dynamics. *Proc Natl Acad Sci U S A*. 2019; 116(15):7226–7231. <https://doi.org/10.1073/pnas.1816531116>
33. Choi K. In: Rajewsky N, Jurga S, Barciszewski J, editors. *Robust Approaches to Generating Reliable Predictive Models in Systems Biology*. Cham: Springer International Publishing; 2018. p. 301–312.
34. Hoffmann M, Fröhner C, Noé F. Reactive SINDy: Discovering governing reactions from concentration data. *J Chem Phys*. 2019; 150(2):025101. <https://doi.org/10.1063/1.5066099> PMID: 30646700
35. Yeung E, Kim J, Yuan Y, Gonçalves J, Murray RM. Data-driven network models for genetic circuits from time-series data with incomplete measurements. *J R Soc Interface*. 2021; 18(182):20210413. <https://doi.org/10.1098/rsif.2021.0413> PMID: 34493091
36. Jiang R, Singh P, Wrede F, Hellander A, Petzold L. Identification of dynamic mass-action biochemical reaction networks using sparse Bayesian methods. *PLoS Comput Biol*. 2022; 18(1):e1009830. <https://doi.org/10.1371/journal.pcbi.1009830> PMID: 35100263
37. Kaheman K, Kutz JN, Brunton SL. SINDy-PI: a robust algorithm for parallel implicit sparse identification of nonlinear dynamics. *Proc Math Phys Eng Sci*. 2020; 476(2242):20200279. <https://doi.org/10.1098/rspa.2020.0279> PMID: 33214760
38. Mangan NM, Kutz JN, Brunton SL, Proctor JL. Model selection for dynamical systems via sparse regression and information criteria. *Proc Math Phys Eng Sci*. 2017; 473(2204):20170009. <https://doi.org/10.1098/rspa.2017.0009> PMID: 28878554
39. Wieland FG, Hauber AL, Rosenblatt M, Tönsing C, Timmer J. On structural and practical identifiability. *Current Opinion in Systems Biology*. 2021; 25:60–69. <https://doi.org/10.1016/j.coisb.2021.03.005>
40. Szederkényi G, Banga JR, Alonso AA. Inference of complex biological networks: distinguishability issues and optimization-based solutions. *BMC systems biology*. 2011; 5(1):1–15. <https://doi.org/10.1186/1752-0509-5-177> PMID: 22034917
41. Chin SV, Chappell MJ. Structural identifiability and indistinguishability analyses of the Minimal Model and a Euglycemic Hyperinsulinemic Clamp model for glucose–insulin dynamics. *Computer Methods and Programs in Biomedicine*. 2011; 104(2):120–134. <https://doi.org/10.1016/j.cmpb.2010.08.012> PMID: 20851494
42. Jánzén DLI, Bergenholm L, Jirstrand M, Parkinson J, Yates J, Evans ND, et al. Parameter Identifiability of Fundamental Pharmacodynamic Models. *Front Physiol*. 2016; 7:590. <https://doi.org/10.3389/fphys.2016.00590> PMID: 27994553
43. Villaverde AF, Banga JR. Dynamical compensation and structural identifiability of biological models: Analysis, implications, and reconciliation. *PLoS Comput Biol*. 2017; 13(11):e1005878. <https://doi.org/10.1371/journal.pcbi.1005878> PMID: 29186132
44. Eisenberg MC, Jain HV. A confidence building exercise in data and identifiability: Modeling cancer chemotherapy as a case study. *Journal of theoretical biology*. 2017; 431:63–78. <https://doi.org/10.1016/j.jtbi.2017.07.018> PMID: 28733187
45. Muñoz-Tamayo R, Puillet L, Daniel JB, Sauvant D, Martin O, Taghipoor M, et al. To be or not to be an identifiable model. Is this a relevant question in animal science modelling? *Animal*. 2018; 12(4):701–712. PMID: 29096725
46. Schmidt PJ, Emelko MB, Thompson ME. Recognizing Structural Nonidentifiability: When Experiments Do Not Provide Information About Important Parameters and Misleading Models Can Still Have Great Fit. *Risk Anal*. 2020; 40(2):352–369. <https://doi.org/10.1111/risa.13386> PMID: 31441953
47. Barreiro XR, Villaverde AF. Benchmarking tools for a priori identifiability analysis. *Bioinformatics*. 2023; 39:btad065. <https://doi.org/10.1093/bioinformatics/btad065>
48. Raue A, Kreutz C, Maiwald T, Bachmann J, Schilling M, Klingmüller U, et al. Structural and practical identifiability analysis of partially observed dynamical models by exploiting the profile likelihood. *Bioinformatics*. 2009; 25(15):1923–1929. <https://doi.org/10.1093/bioinformatics/btp358> PMID: 19505944
49. Stigter J, Joubert D. Computing measures of identifiability, observability, and controllability for a dynamic system model with the StruID App. *IFAC-PapersOnLine*. 2021; 54(7):138–143. <https://doi.org/10.1016/j.ifacol.2021.08.348>
50. Villaverde AF. Observability and structural identifiability of nonlinear biological systems. *Complexity*. 2019; Article ID 8497093.
51. Yates JW, Evans ND, Chappell MJ. Structural identifiability analysis via symmetries of differential equations. *Automatica*. 2009; 45(11):2585–2591. <https://doi.org/10.1016/j.automatica.2009.07.009>

52. Merkt B, Timmer J, Kaschek D. Higher-order Lie symmetries in identifiability and predictability analysis of dynamic models. *Physical Review E*. 2015; 92(1):012920. <https://doi.org/10.1103/PhysRevE.92.012920> PMID: 26274260
53. Villaverde AF. Symmetries in Dynamic Models of Biological Systems: Mathematical Foundations and Implications. *Symmetry*. 2022; 14(3):467. <https://doi.org/10.3390/sym14030467>
54. Massonis G, Banga JR, Villaverde AF. AutoRepar: a method to obtain identifiable and observable reparameterizations of dynamic models with mechanistic insights. *International Journal of Robust and Nonlinear Control*. 2023; 33:5039–5057. <https://doi.org/10.1002/mc.5887>
55. Villaverde AF, Tsiantis N, Banga JR. Full observability and estimation of unknown inputs, states and parameters of nonlinear biological models. *Journal of the Royal Society Interface*. 2019; 16(156):20190043. <https://doi.org/10.1098/rsif.2019.0043> PMID: 31266417
56. Díaz-Seoane S, Rey Barreiro X, Villaverde AF. STRIKE-GOLDD 4.0: user-friendly, efficient analysis of structural identifiability and observability. *Bioinformatics*. 2022; 39(1).
57. Lorenz EN. Deterministic Nonperiodic Flow. *Journal of the Atmospheric Sciences*. 1963; 20(2):130–141. [https://doi.org/10.1175/1520-0469\(1963\)020%3C0130:DNF%3E2.0.CO;2](https://doi.org/10.1175/1520-0469(1963)020%3C0130:DNF%3E2.0.CO;2)
58. Zhang Z. Mathematical Model of a Bacteria-Immunity System with the Influence of Quorum Sensing Signal Molecule. *Journal of Applied Mathematics and Physics*. 2016; 04(05):888–896. <https://doi.org/10.4236/jamp.2016.45097>
59. Süel GM, Garcia-Ojalvo J, Liberman LM, Elowitz MB. An excitable gene regulatory circuit induces transient cellular differentiation. *Nature*. 2006; 440(7083):545–550. <https://doi.org/10.1038/nature04588> PMID: 16554821
60. Evans ND, Chappell MJ. Extensions to a procedure for generating locally identifiable reparameterisations of unidentifiable systems. *Mathematical Biosciences*. 2000; 168(2):137–159. [https://doi.org/10.1016/S0025-5564\(98\)00004-2](https://doi.org/10.1016/S0025-5564(98)00004-2) PMID: 11121562
61. Johnston MD, Edwards CM, Bodmer WF, Maini PK, Chapman SJ. Examples of Mathematical Modeling: Tales from the Crypt. *Cell Cycle*. 2007; 6(17):2106–2112. <https://doi.org/10.4161/cc.6.17.4649> PMID: 17873520
62. Wolf J, Heinrich R. Effect of cellular interaction on glycolytic oscillations in yeast: a theoretical investigation. *Biochemical Journal*. 2000; 345(2):321–334. <https://doi.org/10.1042/bj3450321> PMID: 10702114
63. Castro M, de Boer RJ. Testing structural identifiability by a simple scaling method. *PLOS Computational Biology*. 2020; 16(11):e1008248. <https://doi.org/10.1371/journal.pcbi.1008248> PMID: 33141821
64. Holmberg A. On the practical identifiability of microbial growth models incorporating Michaelis-Menten type nonlinearities. *Mathematical Biosciences*. 1982; 62(1):23–43. [https://doi.org/10.1016/0025-5564\(82\)90061-X](https://doi.org/10.1016/0025-5564(82)90061-X)
65. La Cava W, Orzechowski P, Burlacu B, de Franca F, Virgolin M, Jin Y, et al. Contemporary Symbolic Regression Methods and their Relative Performance. In: Vanschoren J, Yeung S, editors. *Proceedings of the Neural Information Processing Systems Track on Datasets and Benchmarks*. vol. 1. Curran; 2021.
66. Karniadakis GE, Kevrekidis IG, Lu L, Perdikaris P, Wang S, Yang L. Physics-informed machine learning. *Nature Reviews Physics*. 2021; 3(6):422–440. <https://doi.org/10.1038/s42254-021-00314-5>
67. Gao TT, Yan G. Data-driven inference of complex system dynamics: A mini-review. *Europhysics Letters*. 2023; 142(1):11001. <https://doi.org/10.1209/0295-5075/acc3bf>
68. Omejc N, Gec B, Brence J, Todorovski L, Džeroski S. Probabilistic grammars for modeling dynamical systems from coarse, noisy, and partial data. *Research Square preprint*. 2023;.
69. Gelß P, Klus S, Eisert J, Schütte C. Multidimensional approximation of nonlinear dynamical systems. *Journal of Computational and Nonlinear Dynamics*. 2019; 14(6).
70. Kaheman K, Kaiser E, Strom B, Kutz JN, Brunton SL. Learning discrepancy models from experimental data. *arXiv preprint arXiv:190908574*. 2019;.
71. Reinbold PA, Kageorge LM, Schatz MF, Grigoriev RO. Robust learning from noisy, incomplete, high-dimensional experimental data via physically constrained symbolic regression. *Nature communications*. 2021; 12(1):3219. <https://doi.org/10.1038/s41467-021-23479-0> PMID: 34050155
72. Fasel U, Kutz JN, Brunton BW, Brunton SL. Ensemble-SINDy: Robust sparse model discovery in the low-data, high-noise limit, with active learning and control. *Proceedings of the Royal Society A*. 2022; 478(2260):20210904. <https://doi.org/10.1098/rspa.2021.0904> PMID: 35450025
73. Kaheman K, Brunton SL, Kutz JN. Automatic differentiation to simultaneously identify nonlinear dynamics and extract noise probability distributions from data. *Machine Learning: Science and Technology*. 2022; 3(1):015031.