

**A peer-reviewed version of this preprint was published in PeerJ on 20 September 2016.**

[View the peer-reviewed version](https://doi.org/10.7717/peerj.2417) (peerj.com/articles/2417), which is the preferred citable publication unless you specifically need to cite this preprint.

Astola L, Stigter H, Gomez Roldan MV, van Eeuwijk F, Hall RD, Groenenboom M, Molenaar JJ. 2016. Parameter estimation in tree graph metabolic networks. PeerJ 4:e2417 <https://doi.org/10.7717/peerj.2417>

# Parameter estimation in tree graph metabolic networks

Laura Astola, Hans Stigter, Maria Victoria Gomez Roldan, Fred van Eeuwijk, Robert D Hall, Marian Groenenboom, Jaap J. Molenaar

We study the glycosylation processes that convert initially toxic substrates to nutritionally valuable metabolites in the flavonoid biosynthesis pathway of tomato (*Solanum lycopersicum*) seedlings. To estimate the reaction rates we use ordinary differential equations (ODEs) to model the enzyme kinetics. A popular choice is to use a system of linear ODEs with constant kinetic rates or to use Michaelis-Menten kinetics. In reality, the catalytic rates, which are affected among other factors by kinetic constants and enzyme concentrations, are changing in time and with the approaches just mentioned, this phenomenon cannot be described. Another problem is that, in general these kinetic coefficients are not always identifiable. A third problem is that, it is not precisely known, which enzymes are catalyzing the observed glycosylation processes. With several hundred potential gene candidates, experimental validation using purified target proteins is expensive and time consuming. We aim at reducing this task via mathematical modeling to allow for the pre-selection of most potential gene candidates. In this article we discuss a fast and relatively simple approach to estimate time varying kinetic rates, with three favorable properties: Firstly, it allows for identifiable estimation of time dependent parameters in networks with a tree-like structure. Secondly, it is very fast compared to the usually applied methods, since it is not based on an iterative scheme. Thirdly, by combining the metabolite concentration data with a corresponding microarray data, it can help in detecting the genes related to the enzymatic processes. By comparing the estimated time dynamics of the catalytic rates with time series gene expression data we may assess potential candidate genes behind enzymatic reactions. As an example, we show how to apply this method to select prominent glycosyltransferase genes in tomato seedlings.

# Parameter Estimation in Tree Graph Metabolic Networks

Laura Astola<sup>1</sup>, Hans Stigter<sup>2</sup>, Victoria Gomez Roldan<sup>3</sup>, Fred van Eeuwijk<sup>4</sup>, Robert D. Hall<sup>5</sup>, Marian Groenenboom<sup>6</sup>, and Jaap Molenaar<sup>7</sup>

<sup>1</sup>l.j.astola@tue.nl

<sup>2</sup>hans.stigter@wur.nl

<sup>3</sup>gomez@evry.inra.fr

<sup>4</sup>fred.vaneeuwijk@wur.nl

<sup>5</sup>robert.hall@wur.nl

<sup>6</sup>marian.groenenboom@gmail.com

<sup>7</sup>jaap.molenaar@wur.nl

## ABSTRACT

We study the glycosylation processes that convert initially toxic substrates to nutritionally valuable metabolites in the flavonoid biosynthesis pathway of tomato (*Solanum lycopersicum*) seedlings. To estimate the reaction rates we use ordinary differential equations (ODEs) to model the enzyme kinetics. A popular choice is to use a system of linear ODEs with constant kinetic rates or to use Michaelis-Menten kinetics. In reality, the catalytic rates, which are affected among other factors by kinetic constants and enzyme concentrations, are changing in time and with the approaches just mentioned, this phenomenon cannot be described. Another problem is that, in general these kinetic coefficients are not always identifiable. A third problem is that, it is not precisely known, which enzymes are catalyzing the observed glycosylation processes. With several hundred potential gene candidates, experimental validation using purified target proteins is expensive and time consuming. We aim at reducing this task via mathematical modeling to allow for the pre-selection of most potential gene candidates. In this article we discuss a fast and relatively simple approach to estimate time varying kinetic rates, with three favorable properties: Firstly, it allows for identifiable estimation of time dependent parameters in networks with a tree-like structure. Secondly, it is very fast compared to the usually applied methods, since it is not based on an iterative scheme. Thirdly, by combining the metabolite concentration data with a corresponding microarray data, it can help in detecting the genes related to the enzymatic processes. By comparing the estimated time dynamics of the catalytic rates with time series gene expression data we may assess potential candidate genes behind enzymatic reactions. As an example, we show how to apply this method to select prominent glycosyltransferase genes in tomato seedlings.

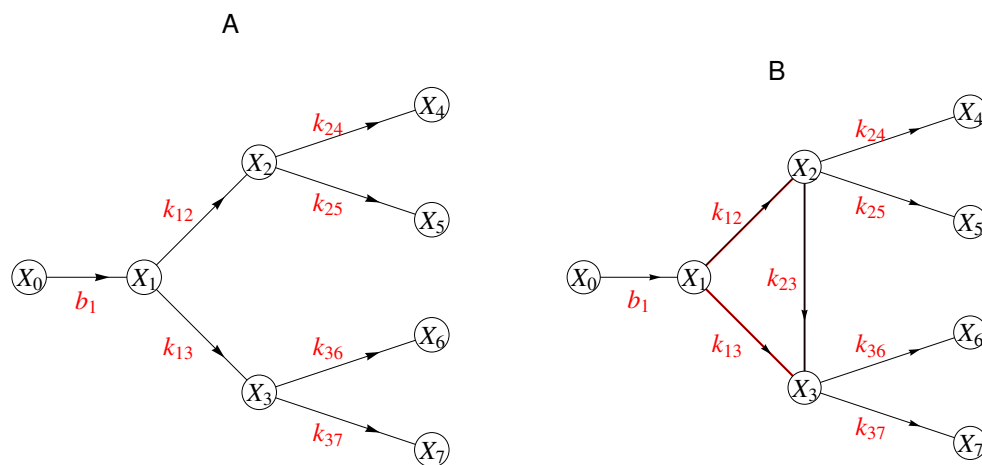
Keywords: Metabolic networks, Systems Biology, Kinetic models, Glycosylation, Network inference

## INTRODUCTION

In this paper we study metabolic network inference from given biological time-series data. The two main ingredients in general metabolic pathway inference are the reconstruction of the network topology and the estimation of the parameters involved. When the network is large and the concentrations of intermediates are unknown, or when there are no time series data available, one may still study the fluxes by setting up stoichiometric models for flux balance analysis Varma and Palsson (1995); Stelling et al. (2002); Orth et al. (2010). If time-series data of metabolites are available ordinary differential equations (ODEs) can often provide a suitable model Chen et al. (2010); Chou and Voit (2009); Srinath and Gunawan (2010); Hatzimanikatis et al. (1996). If also the enzymes involved are known, it is customary to use enzyme-kinetic models Steuer and Junker (2009); Schallau and Junker (2010); Liebermeister and Klipp (2006) with Michaelis-Menten kinetics, although the reliability of this approach has been questioned, especially when applied to *in vivo* measurements Savageau (1995); Hill et al. (1977). When (part of) the catalytic rates are not known, linear ODEs Astola et al. (2011) and general biochemical systems theory Voit et al. (2005) can be used. When the network topology is completely unknown, the situation

is more complicated, although some recent studies attempt to tackle this problem using methods based on genetic algorithms Schmidt et al. (2011). Still, the uniqueness of the reconstructed network is often compromised and the identifiability of the system remains an issue that needs to be investigated Craciun and Pantea (2008); Srinath and Gunawan (2010). Model identifiability is an essential prerequisite in making any conclusions from (by default limited number of) observations. The foremost categories of identifiability are the structural and the practical identifiability, the former related to the symbolic expression of the model itself and the latter related to the amount and nature of the available data. We will test our models and data on both conditions.

Here we discuss a special and relevant class of network topologies, which are so-called tree networks and show that in such networks linear models yield parameter estimates that are unique in the structural sense. As the name suggests, a tree graph looks like a branching tree where the edges (arrows) are directed so that the nutrients flow from root to leaf (cf. Fig. 1).



**Figure 1.** A: a graph with a tree structure. B: this graph contains a cycle and is thus not a tree graph. The catalytic rate corresponding to reaction between node  $i$  and  $j$  is indicated as  $k_{ij}$ . Here the node  $X_0$  represents a boundary node connecting this network to the surrounding larger network

As in real trees the branches do not form cycles. By a cycle we mean any closed chain of edges regardless of the directions of the edges. In many biological pathways, such as in the flavone and the flavonol biosynthesis KEGG (2010), a tree graph captures the network of the enzymatic reactions. Indeed metabolic networks with tree structures constitute a relevant class, including for example large parts of the biosynthetic pathways of, e.g.,  $\gamma$ -carotene, limonene, ansamycin and puromycin etc. KEGG (2010).

The paper is organized as follows: in section 1 to set the stage we review our earlier work in modeling metabolic pathways using time-invariant systems of linear differential equations and discuss the particular properties of tree-graph networks. In section 2 we consider the essential problem of model identifiability and show that our candidate networks satisfy the criteria for structural and practical identifiability. In section 3 we propose a novel application for our time-variant estimation scheme by showing how it can be employed in finding the most likely catalysts from a large set of enzymes.

# 1 PARAMETER ESTIMATION

## 1.1 Parameter Estimation in general networks

In this section we consider the parameter estimation problem in general linear time-invariant ordinary differential equation (LTI-ODE) systems. For convenience, we first briefly sketch the approach when the catalytic rates are constants over time as in our previous work Astola et al. (2011).

We recall that any network can be represented as a graph, where nodes are connected by edges when there is some interaction between these nodes. In a metabolic network a node represents a substrate or a product, and a directed edge from node  $i$  to node  $j$  means that  $i$  can be converted to  $j$  by enzymatic activity. To an edge from node  $i$  to  $j$ , we assign a weight, i.e., the catalytic rate  $k_{ij} \geq 0$ , which represents the rate of product formation. In parameter inference one estimates the  $k_{ij}$  from data.

Denoting the concentration of substrate  $i$  at time  $t$  as  $X_i(t)$ , a general time-invariant linear ODE model with a constant nonhomogeneous term, satisfying the mass conservation law, can be written as

$$\dot{X}_i(t) = -\sum_{j \neq i} k_{ij} X_j(t) + \sum_{j \neq i} k_{ji} X_j(t) + b_i, \quad (1)$$

for  $i = 1, \dots, n$ , with

$$b_i = \begin{cases} \text{constant} & \text{if } i = 1 \\ 0 & \text{otherwise} \end{cases}. \quad (2)$$

The first summation in (1) stands for the edges leaving  $X_i$ , the second for the incoming edges and  $b_i$  for the possible in or outflow to the system. To simplify the notation, we introduce a matrix  $A$  with components given by

$$\begin{cases} A_{ij} = k_{ji}, & i \neq j \\ A_{ii} = -\sum_{j \neq i} k_{ij}, \end{cases} \quad (3)$$

Then, (1) becomes

$$\dot{X}_i(t) = \sum_{j=1}^n A_{ij} X_j(t) + b_i, i = 1, \dots, n. \quad (4)$$

Equation (4) can be rewritten in a compact matrix form as

$$\dot{X}(t) = \begin{pmatrix} \dot{X}_1(t) \\ \dot{X}_2(t) \\ \vdots \\ \dot{X}_n(t) \end{pmatrix} = \begin{pmatrix} -\sum_{j \neq 1} k_{1j} & k_{n1} & b_1 \\ k_{12} & k_{n2} & 0 \\ \vdots & \vdots & \vdots \\ k_{1n} & -\sum_{j \neq n} k_{nj} & 0 \end{pmatrix} \begin{pmatrix} X_1(t) \\ X_2(t) \\ \vdots \\ X_n(t) \\ 1 \end{pmatrix} = \tilde{A} \cdot \tilde{X}(t), \quad (5)$$

where  $\tilde{X}(t)$  is obtained from  $X(t)$  by appending an extra 1 and matrix  $\tilde{A}$  is obtained from  $A$  by extending it with an extra column containing the constant  $b_1$ .

To reconstruct a metabolic network from time-series measurements, we have to estimate the reaction rates  $k_{ij}$ , i.e., the weights of the edges in the network and the flow terms  $b_i$ . In view of (5), it is sufficient to estimate the  $(n+1) \times (n+1)$  matrix  $\tilde{A}$ . We denote the data, i.e., measured concentrations of substrate  $i$  at time points  $t_j$ ,  $j = 1, \dots, m$ , as an  $(n \times m)$  matrix  $\mathbb{X}$ . Estimates of the derivatives of the data curves we will store in a matrix  $\dot{\mathbb{X}}$ . To compute these estimates we may proceed in two ways. First, construct two  $n \times m$  data matrices  $\mathbb{X}_0, \mathbb{X}_1$  as follows

$$\mathbb{X}_0 = \begin{pmatrix} \mathbb{X}_{1,m-1} & \mathbb{X}_{1,m-2} & \dots & \mathbb{X}_{1,0} \\ \mathbb{X}_{2,m-1} & \mathbb{X}_{2,m-2} & \dots & \mathbb{X}_{2,0} \\ \vdots & & \dots & \vdots \\ \mathbb{X}_{n,m-1} & \mathbb{X}_{n,m-2} & \dots & \mathbb{X}_{n,0} \end{pmatrix}, \mathbb{X}_1 = \begin{pmatrix} \mathbb{X}_{1,m} & \mathbb{X}_{1,m-1} & \dots & \mathbb{X}_{1,1} \\ \mathbb{X}_{2,m} & \mathbb{X}_{2,m-1} & \dots & \mathbb{X}_{2,1} \\ \vdots & & \dots & \vdots \\ \mathbb{X}_{n,m} & \mathbb{X}_{n,m-1} & \dots & \mathbb{X}_{n,1} \end{pmatrix}, \quad (6)$$

where  $m$  is the number of measurements. The matrix

$$\dot{\mathbb{X}} \equiv \frac{1}{\Delta t} (\mathbb{X}_1 - \mathbb{X}_0), \quad (7)$$

could then be used as an approximation for  $\dot{X}$ . For simplicity we assume the time grid to be equidistant with time step  $\Delta t$ . If this is not the case, the necessary modifications are easily implemented.

Secondly, we may use an alternative and often better approach to obtain approximations for  $\dot{X}_i$  by fitting splines to the time series data  $X_i$  Zhan and Yeung (2011). To obtain curves that interpolate the data faithfully, we require that the distances between the curves and the measurements are minimal and that at the same time the curves are smooth. To achieve this we fit P-splines, which are B-splines with a penalization for non-smoothness Eilers and Marx (1996).

From these splines, we evaluate the derivative estimates at time points  $t_j$ . These estimates are then used as entries in the matrix  $\dot{X}$ . Having at hand an estimate for matrix  $\dot{X}$ , the problem of network inference comes down to finding the the matrix  $\tilde{A}$  from the equation

$$\dot{X} = \tilde{A}\tilde{X}, \quad (8)$$

in which  $\dot{X}$  is known and  $\tilde{X}$  is obtained from the data matrix  $X$  by extending this with an extra row of ones. However, solving  $\tilde{A}$  directly from (8) often results in over-fitting, since all possible edges are included in the modeled network. Another serious shortcoming of such a matrix (pseudo-) inversion approach is the fact that we cannot control the positivity of the reaction rates. Although in Schmidt et al. (2005) negative coefficients were interpreted as inhibition of the compounds, in many biological pathways, negative coefficients are not permitted. Thus we take a more general approach in which one can exclude all edges that are biologically not acceptable, and in which one can constrain the reaction rates to be positive, without substantially compromising computation time.

To this end, we reformulate the equation as a minimization problem:

$$\arg \min_{\tilde{A}} (||\dot{X} - \tilde{A}\tilde{X}||) . \quad (9)$$

The matrix norm used here is the Frobenius norm:

$$||\tilde{A}|| = \sqrt{\sum_{i=1}^n \sum_{j=1}^m \tilde{A}_{ij}^2} . \quad (10)$$

This alternative formulation allows inclusion of expert knowledge in a simple way. We put  $\tilde{A}_{ij} = 0$ , when an edge from node  $i$  to node  $j$  cannot exist. Nearly all mathematical software packages (Mathematica, Matlab, Maple, etc. ) can numerically find the minimizer  $\tilde{A}$  (and thus the reaction rates  $k_{ij}$  and the flow term  $b_1$ ) with the constraint that  $k_{ij} \geq 0$ .

## 1.2 Parameter estimation in tree networks

As described in the introduction, tree networks are networks, whose graphs resemble trees in that they branch away from the root and the directions of the edges always point from the root towards the leaves. In Fig. 1 we presented, using an example, the difference between a tree and a non-tree graph. In a kinetic reaction system with a tree network, the parameters can be uniquely estimated even when they are time dependent. We could write this down in general. However, the proof is based on one central idea. We feel that the reader gains more insight if we simply show this idea through an example. To that end we use as example the network in the left hand side of Fig. 1. The extension to the general is straightforward.

For the network on the left in Fig. 1, we have the following kinetic mass balance model:

$$\begin{pmatrix} \dot{X}_1 \\ \dot{X}_2 \\ \dot{X}_3 \\ \dot{X}_4 \\ \dot{X}_5 \\ \dot{X}_6 \\ \dot{X}_7 \end{pmatrix} = \begin{pmatrix} -(k_{1,2} + k_{1,3}) & 0 & 0 & 0 & 0 & 0 & 0 & b_1 \\ k_{1,2} & -(k_{2,4} + k_{2,5}) & 0 & 0 & 0 & 0 & 0 & 0 \\ k_{1,3} & 0 & -(k_{3,6} + k_{3,7}) & 0 & 0 & 0 & 0 & 0 \\ 0 & k_{2,4} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & k_{2,5} & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & k_{3,6} & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & k_{3,7} & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \begin{pmatrix} X_1 \\ X_2 \\ X_3 \\ X_4 \\ X_5 \\ X_6 \\ X_7 \\ 1 \end{pmatrix}, \quad (11)$$

where the constant  $b_1$  represents the influx into the system and the  $k_{i,j}$  are the catalytic rates. Note that there are as many unknown parameters ( $k_{i,j}$ ,  $b_1$ ) as there are measured variables  $X_i(t_j)$ . Therefore, as can be directly verified, we can rewrite the previous matrix equation by exchanging the  $X_i$  and  $k_{i,j}$  as follows:

$$\begin{pmatrix} \dot{X}_1 \\ \dot{X}_2 \\ \dot{X}_3 \\ \dot{X}_4 \\ \dot{X}_5 \\ \dot{X}_6 \\ \dot{X}_7 \end{pmatrix} = \underbrace{\begin{pmatrix} 1 & -X_1 & -X_1 & 0 & 0 & 0 & 0 \\ 0 & X_1 & 0 & -X_2 & -X_2 & 0 & 0 \\ 0 & 0 & X_1 & 0 & 0 & -X_3 & -X_3 \\ 0 & 0 & 0 & X_2 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & X_2 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & X_3 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & X_3 \end{pmatrix}}_{\text{matrix } B} \begin{pmatrix} b_1 \\ k_{1,2} \\ k_{1,3} \\ k_{2,4} \\ k_{2,5} \\ k_{3,6} \\ k_{3,7} \end{pmatrix}. \quad (12)$$

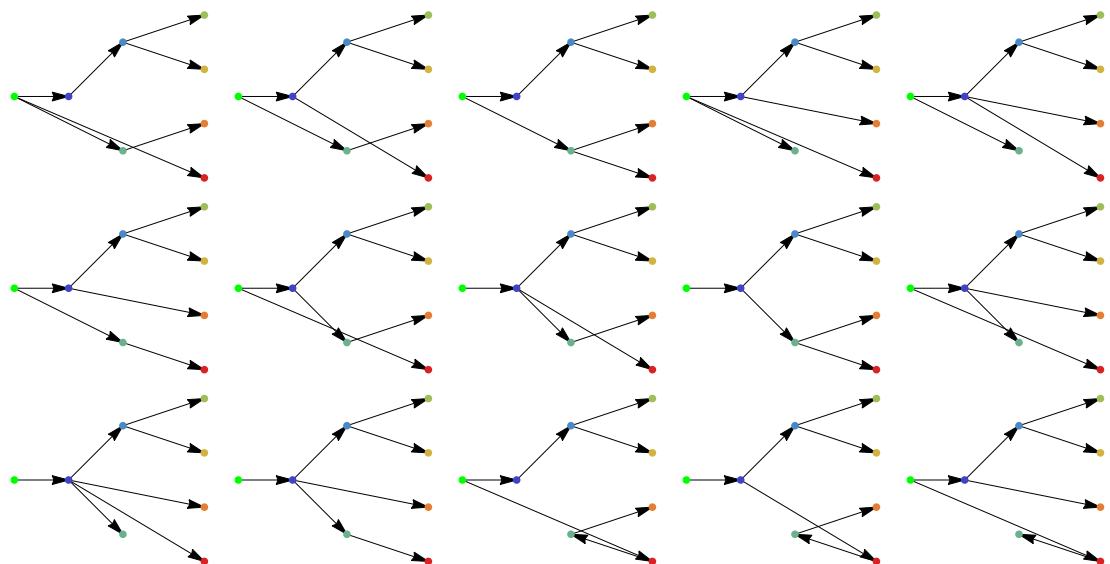
We immediately see that  $B$  is an upper triangular matrix since the entries below the diagonal are zero. This implies that the determinant of the matrix  $B$  in (12) is the product of the entries on the diagonal:  $X_1^2 \cdot X_2^2 \cdot X_3^2$ , and thus unequal to 0 since  $X_i \neq 0, \forall i = 1, \dots, n$ . So,  $B$  is invertible and the system of equations has the unique solution.

$$\begin{pmatrix} b_1 \\ k_{1,2} \\ k_{1,3} \\ k_{2,4} \\ k_{2,5} \\ k_{3,6} \\ k_{3,7} \end{pmatrix} = \underbrace{\begin{pmatrix} X_0^{-1} & X_0^{-1} & X_0^{-1} & X_0^{-1} & X_0^{-1} & X_0^{-1} & X_0^{-1} \\ 0 & X_1^{-1} & 0 & X_1^{-1} & X_1^{-1} & 0 & 0 \\ 0 & 0 & X_1^{-1} & 0 & 0 & X_1^{-1} & X_1^{-1} \\ 0 & 0 & 0 & X_2^{-1} & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & X_2^{-1} & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & X_3^{-1} & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & X_3^{-1} \end{pmatrix}}_{\text{matrix } B^{-1}} \begin{pmatrix} \dot{X}_1 \\ \dot{X}_2 \\ \dot{X}_3 \\ \dot{X}_4 \\ \dot{X}_5 \\ \dot{X}_6 \\ \dot{X}_7 \end{pmatrix} \quad (13)$$

## 2 MODEL IDENTIFIABILITY

### 2.1 Structural identifiability

A general problem in parameter estimation is that it is difficult and sometimes even impossible to be sure that the estimated parameters are unique. If the model is structurally unidentifiable, there is an infinite number of parameter sets that give equal results. This is a substantial challenge, especially when the network structure is not known, since an overly complex network can result in over-fitting. This problem is not present in any of the (biologically) potential networks as sketched in Fig. 2, since as tree graphs these all turn out to be locally structurally identifiable.

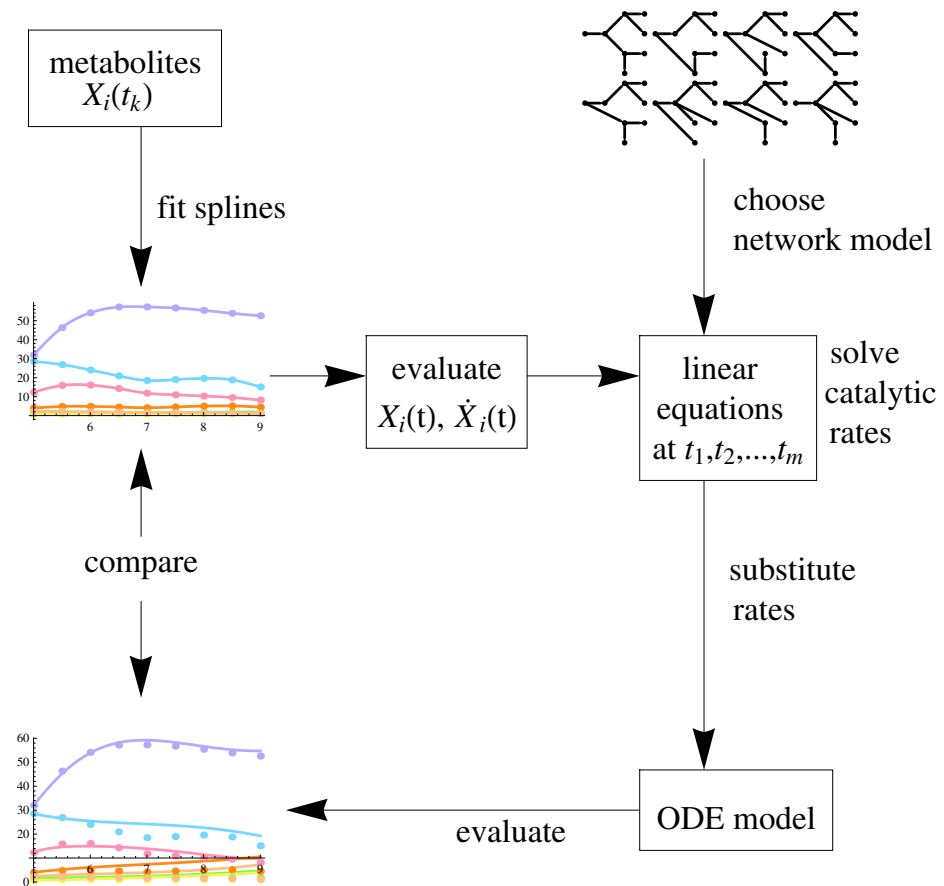


**Figure 2.** Here we depict all biologically feasible networks of the quercetin glycosylation pathway.

## 2.2 Practical identifiability

Structural identifiability does not imply practical identifiability and therefore we have studied the practical identifiability of the parameters in our system by means of profile likelihood Raue et al. (2009). We learned that all the kinetic parameters connecting substrates and products with concentrations above detection limit show also practical identifiability. Another observation is that if we allow a product to decay without constraints, the practical identifiability as well as the tree structure of the graph is lost.

The inference method proposed here is by no means restricted to tree networks, but in case the network has a tree structure, the parameters can be estimated in an unambiguous way. We summarize the general work flow for the proposed parameter inference in the schematic diagram in Fig. 3.



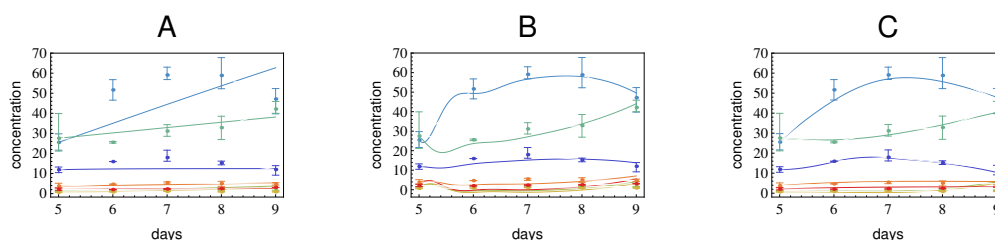
**Figure 3.** A schematic view of the inference procedure.

## 3 EXTENSIONS AND NOVELTIES

### 3.1 Time varying kinetic rates

In earlier work we developed a fast method to reconstruct metabolic networks Astola et al. (2011). The idea in this approach was to substitute the measurements directly into the model equations and not only in the objective function. This approach had as a limitation that all parameters were assumed to be constant in time. Here we extend our previous approach by allowing the catalytic rates to be time dependent, to better reflect the real situation, since in practice the enzyme concentrations are fluctuating in time. This has also immediately resulted in reconstructions that better fit the observed data as can be seen in Fig. 4. While the standard practice in enzyme kinetics is to either use constant catalytic rates in mass balance





**Figure 4.** In this figure we have used three different models to reconstruct a flavonol concentration data indicated as dots. The compounds shown here belong to a pathway with putative structure as on the left in Fig. 1. The colors of the reconstructed curves correspond to those of the dots. On the left hand side, a reconstruction with a tree network (as in Fig. 1) and constant catalytic rates. Middle, a reconstruction with the full network (all nodes are connected to each other) and constant catalytic rates. Note that the fit is still poor, although the number of parameters is much higher than in the case on the left. On the right, a reconstruction with the same tree structure as in the leftmost case but with time dependent catalytic rates

equation or to model product formation through a Hill function Goutelle et al. (2008) such as in the Michaelis-Menten equation Savageau (1995), none of these take into account the fact that the enzyme concentration is also changing in time. Since we also want to study the relation of gene expression and enzyme concentration in time, we need to capture their dynamics.

As the catalytic rate is now modeled as a function in time, and not as a constant, it is no longer possible to infer this with the standard procedure of solving for those parameters that fit the ordinary differential equations to data in the sense of maximum likelihood. We cannot clearly separate the substrate/product, enzyme concentrations and noise, for we have no measurements of the enzyme concentrations. To solve them, we would have to impose a model on them, which we don't have a priori. A reasonable approach in this situation is to first estimate a model for the metabolite concentrations on which we have several measurements. By fixing the concentrations first using spline approximations, we may then estimate the trends in the enzyme concentrations. This method assumes that the solutions are rather smooth. If this is not the case and the sampling frequency is low, the derivatives obtained by fitting splines can introduce errors that distort the reconstruction.

### 3.2 Time dependent parameter estimation

In this section we present three different schemes to estimate the  $k_{ij}(t)$  in model (4). In (9) we used the data at all time points simultaneously to estimate the time independent parameters. However, a remarkable feature of tree structured networks is that the data at one time point is already enough to calculate unique estimates for the parameter values at that particular time point. This is immediately clear from (13): as soon as we have estimates for the time derivatives  $\dot{X}(t_k)$  available, we may calculate estimates for the  $k_{ij}(t_k)$ .

**Scheme 1.** To estimate the derivatives at some time point one still needs the data of neighboring time points. So, the first step in this scheme is to fit, e.g., P-splines to the data time series O' Sullivan (1986); Eilers and Marx (1996). From these splines we calculate estimates for the time derivatives  $\dot{X}_i(t_k)$ . Then by substituting these estimates as well as the measurements into equation (4), we are left with a set of linear equations to solve  $k_{ij}(t_k)$  and  $b_1$  at all times  $t_k$ . Finally, for smooth and continuous catalytic rates, one may fit, e.g., a second order polynomial through these estimates.

**Scheme 2.** An alternative approach in which the number of parameters is smaller than in scheme 1, is to assume that the functions  $k_{ij}(t)$  can be adequately represented as polynomials in time of some order. In practice order 2 is often sufficient. With this choice we have then:

$$k_{ij}(t) = \alpha_{ij}t^2 + \beta_{ij}t + \gamma_{ij} . \quad (14)$$

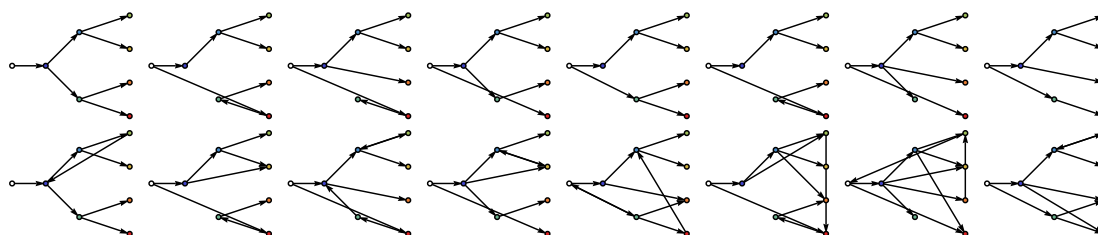
This implies that per  $k_{ij}$  we have 3 parameters to be estimated using the whole time series data. By substituting (14) into matrix  $\tilde{A}$  in (9) we then obtain estimates for  $\alpha_{ij}$ ,  $\beta_{ij}$  and  $\gamma_{ij}$ , and thus for  $k_{ij}(t)$ .

Scheme 3. As in the previous scheme, we assume (14). We construct an objective function like the following:

$$\sum_k ||X(t_k) - \mathbb{X}(t_k)||, \quad (15)$$

which is the sum of the distances between  $X(t_k)$  and the measurements. We look for a matrix  $\tilde{A}$ , such that the solutions  $X_i(t)$  to (4) minimize this objective function. Using suitable optimization algorithm we simultaneously estimate  $X_i$ ,  $k_{ij}$ , and  $b_1$ .

To compare the accuracy and speed of these three schemes we applied them using as test networks random tree networks (see top row in Fig. 5) that have equal numbers of nodes and edges as the network on the left in Fig. ??.



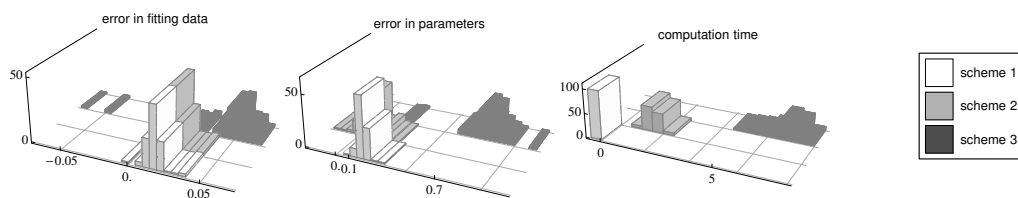
**Figure 5.** These networks are examples of the networks that were used in the simulations, in which we compared the three different parameter estimation schemes. Top row: tree networks. Bottom row: non-tree networks.

In these networks, we simulated time series data with time varying catalytic rates. To generate artificial data, we assigned random values to  $\alpha_{ij}$ ,  $\beta_{ij}$  and  $\gamma_{ij}$  in a range, such that the resulting solutions have approximately the same range as the metabolite concentration data for quercetin glycosides measured in tomato seedlings (cf. Fig. 4). To assess the reconstruction power of the three schemes, we also tested them on networks that are not trees (see bottom row in Fig. 5). The corresponding data generation process is the same but the network models have more edges and thus cycles. In the first set of simulations we used noiseless data and in the second set we added  $\pm 10\%$  uniformly distributed noise.

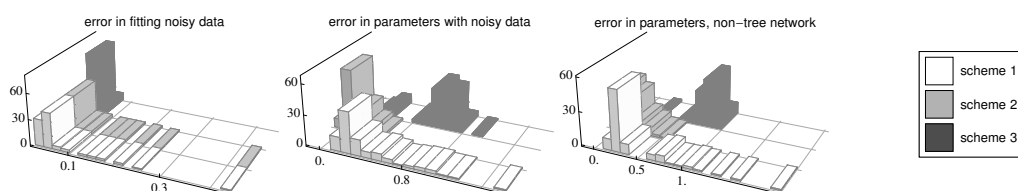
As can be seen from Fig. 6 and 7 scheme 1 and scheme 2 give almost equal results concerning the residuals. In principle they are solving the same optimization problem, only scheme 1 first solves the point wise rate values and then fits a polynomial, whereas scheme 2 searches for a polynomial-valued rates that fit to the data at each point. We measured the accuracy of the parameter estimation by computing the Frobenius norm (10) of the difference between the original parameter matrix used in simulation and the reconstructed matrix. Although scheme 3 can result in quite good fit to data, this does not imply that the parameters are estimated accurately. In terms of accuracy, scheme 2. performed best in our simulations. This is also intuitive, since scheme 2 is optimizing the polynomial-valued parameters to fit the whole time series and thus likely to be more robust against errors in point wise derivatives that are due to the initial spline approximation. Besides the actual estimation accuracy, also computation times are relevant. In terms of computation time, scheme 1 is the fastest and scheme 3 is around 1000 times slower. Notice that the comparisons in Fig. 6 and 7 were done in a setting where equal parameter constraints were given to the solvers and the parameters were estimated using constrained non-linear optimization (NMinimize in Mathematica) choosing for the fast Nelder-Mead algorithm.

This result is more or less to be expected, since when the data is reasonably accurate, it does not always make sense to re-estimate the data by using it as an unknown variable in the equations of the system. Rather, it pays off to substitute the data directly into the equations reducing the number of unknown elements. Since our method is based on initial fitting of splines, the major sensitivity is indeed with respect to data. This was also confirmed by the sensitivity analysis we conducted.

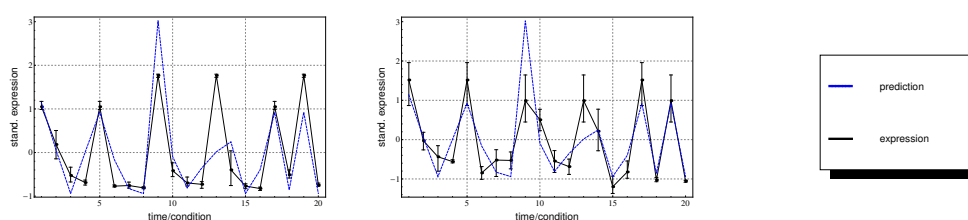
Our network models, although relatively small, belong to the general group of the so-called sloppy biochemical models Gutenkunst (2008). The range of eigenvalues of the Hessian of the residual (between predicted and measured values) varies from  $10^{-4}$  to  $10^5$ . We remark that since we are considering time varying parameters, to compute the derivatives at the point yielding minimal residual, we have taken time-averages of point-wise derivatives w.r.t. a parameter at each data point. Eigenvectors corresponding to very small eigenvalues, implying sloppiness in sensitivity, all point towards those parameters that are



**Figure 6.** We have compared three different reconstruction schemes in 150 simulations, when the underlying network has a tree structure. In each sub-figure the percentual histograms of simulation results are plotted. On the left: the errors in the estimated catalytic rates, using noiseless data. Middle: similarly but with 10% uniformly distributed random noise added to data. Right: the computation times in noiseless case. We observe that although scheme 3 occasionally yields a very precise solution, this happens relatively rarely. In terms of network inference, scheme 1 (corresponding to the method proposed here) and 2 give consistently lower error rates. The computation times for scheme 1 and scheme 2 overlap here, and for this reason only the latter is visible. The error rates and timings in all three figures are on logarithmic scales.



**Figure 7.** We have compared three different reconstruction schemes in 100 simulations, when the underlying network has a non-tree structure. Left: the errors in the estimated catalytic rates. Middle: similarly with 10% uniform random noise added to data. Right: computation times when data are noiseless. We observe a substantial growth in computation times used up by scheme 3, compared to the situation where the networks had tree structures. The computation times for scheme 1 and scheme 2 overlap and thus only the latter is visible. The error rates and timings in all three figures are on logarithmic scales.



**Figure 8.** The mean expression levels of three different putative glucosyl transferase (GT) genes and the estimated catalytic rate for reaction that converts quercetin to quercetin-3-O-glucoside. Assuming that the expression levels correlate with the predicted enzymatic activity, the leftmost gene looks promising in that the two have almost identical slopes and curvatures, whereas the rightmost gene seems to have a trend opposite to that of the catalytic rate.

206 associated with network nodes where the measured metabolite concentrations are very low. This is logical  
 207 since the parameters associated with concentration values close to zero have little effect on the residual,  
 208 because our objective function does not contain the standard deviation term in the denominator. By this  
 209 choice we explicitly wanted to avoid that those measurements that are close to noise level shall have equal  
 210 weight with the more abundant ones.

#### 4 PARAMETER INFERENCE AS A MEAN TO SELECT ACTIVE GENES

In addition, as a potentially powerful application, we show how we may infer the gene candidates likely to be involved in the enzymatic reactions. This can be done by comparing estimated time dependent catalytic rates with simultaneously measured gene expression data. If, according to the model, the formation of a metabolite necessitates higher/lower enzyme concentration, this should be also observable in the expression level of the gene that codes for this enzyme. Using this heuristics we were able to select from a large set of potential genes the most likely candidate genes for further experimental validation of their functioning in particular reactions. In view of this application small inaccuracies in parameters are not detrimental, since here we are mainly interested the dynamic trends of the catalytic rates instead of their precise numeric values.

As an example we take the quercetin glycosylation pathway in cotyledons, occurring during the development of tomato seedlings Koes et al. (1994). Quercetin glycosides are a subset of flavonoids, which are plant secondary metabolites naturally produced by plants. Flavonoids are being intensively studied for their proposed beneficial effects on prevention of chronic diseases Bovy et al. (2007); Rein et al. (2006); Moon et al. (2006).

We have measured the concentrations of several quercetin derivative compounds accumulating in cotyledon- and hypocotyl tissues. We have daily measurements from day 5 after sowing up to day 9. The same sample used for the metabolite analysis with liquid chromatography mass spectrometry were used for gene expression analysis. The expression levels of genes, putatively involved in the glycosylation of quercetin, were quantified using microarray analysis. Glycosyltransferases (GTs) are members of the multigene superfamily in plants that can transfer single or multiple sugars to various plant molecules, resulting in the glycosylation of these compounds Wang (2009). To date, it is not known exactly which GTs catalyze each glycosylation reaction. With more than 200 GT candidates an experimental validation of every single GT is costly. Therefore we wanted to make a pre-selection of the potentially strongest gene candidates, using mathematical modeling and simulations. We use the heuristics that if the kinetic ODE model describes the system of enzymatic reactions reasonably well, the estimated catalytic rates should reflect the real enzymatic activity. This in turn should correlate with the expression trends of the GTs observed using the time series microarray analysis.

Our procedure for the GT inference is as follows:

1. Given the time series metabolite concentration data, estimate the time dependent parameters using all biologically relevant networks. Select the network that gives the best fit to measurements with respect to residual or goodness of fit etc. Save the estimated catalytic rates corresponding to the best network as reference.
2. Compute correlations between the time series of expression levels of each GT and the previously saved series of catalytic rates.
3. Select those GTs whose dynamics correlate best with catalytic dynamics for further experimental validation.

For illustration see Figure 7, where we show the expression levels of two GTs together with the estimated catalytic rates for a reaction that glycosylates quercetin to quercetin-3-O-glucoside. We have standardized, i.e. subtracted the mean and divided by standard deviation both predicted and measured expressions for visual comparison. As can be seen from the Figure 7, the deviation of the expression levels between samples can vary from gene to gene. One could also weight the correlation according to this variation so that more precise observations are favored. To test experimentally whether the inferred genes are actually related to the enzymes that glycosylate the flavonols, a set of selected genes are currently being cloned.

As a computational validation of the selection procedure, we tested whether substituting the (scaled) expression levels of the selected genes into the model will result in a decreased residual (better likelihood of observing the measurements). The reason we want to do this post-analysis is two-fold. First of all, our GT candidates are ranked according to their correlation with the predicted enzymatic trends, but it may happen that several candidates have almost equal correlation coefficients. This makes it difficult to distinguish between the candidates, especially because the initial GT-population is already a result of an ontology-based selection. Another point is that, the selection of the most likely GT's is based on individual matchings with single dynamic parameters whose magnitudes are unknown. It is not absolutely clear, say, whether the combination of the very best candidates will always give better results than when

for example one candidate is actually the second best one (in terms of correlation). In each network combination, at most seven GT's are considered, but still the number of all possible combinations is very large. Also the expression levels need to be scaled to match the metabolic model.

To ensure a rich set of gene combinations, we ran a Markov Chain Monte Carlo-algorithm (MCMC) Calvetti and Somersalo (2007). To address the question, of whether the differences in correlations are significant enough, we first ordered the genes into a sequence according to their correlation with the predicted enzyme concentration levels and took two sets of genes according to their order number in the sequence: 1, 2, ..., 10 and 11, 12, ..., 20. We tested whether the residuals, obtained after 200 iterations of 1000 samples with MCMC algorithm using the data of these two sets, have equal means and variances. For the mean test we obtained a P-value less than 0.00001 and for the variance test a P-value of less than 0.006. We may conclude that in the context of a dynamic kinetic reaction model, those genes with expression levels highly correlating to the predicted enzyme dynamics, are significantly more likely to be responsible for the observations.

## 5 CONCLUSIONS

In this article, we consider the time dependence and unique estimability of kinetic rates in metabolic networks. Firstly, we show that when the underlying network has a structure of a tree graph, these rates can be unambiguously estimated. Secondly we propose a fast approach for the estimation of time dependent kinetic rates and demonstrate its performance on simulated data. Finally we also propose an application, where we utilize the estimation method to detect the genes that are potentially involved in particular enzymatic reactions using microarray data.

## REFERENCES

- Astola, L., Groenenboom, M., Gomes Roldan, V., Hall, R. D., Molenaar, J., Bovy, A., and Eeuwijk, F. (2011). Metabolic pathway inference from time series data: a non iterative approach. In *6th IAPR International Conference, Pattern Recognition in Bioinformatics, Lecture Notes in Bioinformatics*, volume 7036, pages 97–108. Springer.
- Bovy, A., Schijlen, E., and Hall, R. (2007). Metabolic engineering of flavonoids in tomato (*Solanum lycopersicum*): the potential for metabolomics. *Metabolomics*, 3(3):399–412.
- Calvetti, D. and Somersalo, E. (2007). Introduction to Bayesian Scientific Computing: Ten lectures on subjective computing. In Antman, S. S., Marsden, J. E., and Sirovich, L., editors, *Surveys and Tutorials in the Applied Mathematical Sciences*. Springer.
- Chen, W., Niepel, M., and Sorger, P. (2010). Classic and contemporary approaches to modeling biochemical reactions. *Genes & Development*, 24(17):1861–1875.
- Chou, I.-C. and Voit, E. (2009). Recent developments in parameter estimation and structure identification of biochemical and genomic systems. *Math Biosci.*, 219(2):57–83.
- Craciun, G. and Pantea, C. (2008). Identifiability of chemical reaction networks. *J Math Chem*, 44:244–259.
- Eilers, P. and Marx, B. (1996). Flexible smoothing with B-splines and penalties. *Statistical Science*, 11(2):89–121.
- Goutelle, S., Maurin, M., Rougier, F., Barbaut, X., Bourguignon, L., Ducher, M., and Maire, P. (2008). The hill equation: a review of its capabilities in pharmacological modelling. *Fundam Clin Pharmacol.*, 22(6):633–648.
- Gutenkunst, R. (2008). *Sloppiness, modeling and evolution in biochemical networks*. Ph.d thesis, Cornell University, New York.
- Hatzimanikatis, V., Floudas, C., and Bailey, J. (1996). Analysis and design of metabolic reaction networks via mixed-integer linear optimization. *AIChE Journal*, 42(5):1277–1292.
- Hill, C., Waight, R., and Bardsley, W. (1977). Does any enzyme follow the Michaelis-Menten equation? *Mol. Cell. Biochem.*, 15:173–178.
- KEGG (2010). Flavone and Flavonol Biosynthesis <http://www.genome.jp/kegg/pathway/map/map00944>
- Koes, R., Quattrocchio, F., and Mol, J. (1994). The flavonoid biosynthetic pathway in plants: Function and evolution. *BioEssays*, 16(2):123–132.
- Liebermeister, W. and Klipp, E. (2006). Bringing metabolic networks to life: convenience rate law and thermodynamic constraints. *Theor Biol Med Model.*, 3(41).



- 316 Moon, Y., Wang, X., and Morris, M. (2006). Dietary flavonoids: effects on xenobiotic and carcinogen  
317 metabolism. *Toxicol. in Vitro*, 20(2):187–210.
- 318 O' Sullivan, F. (1986). A statistical perspective on ill-posed inverse problems (with discussion). *Statistical*  
319 *Science*, 1:505–527.
- 320 Orth, J., Thiele, I., and Palsson, B. (2010). What is flux balance analysis? *Nature Biotechnology*,  
321 28:245–248.
- 322 Raue, A., Kreutz, C., Maiwald, T., Bachmann, J., Schilling, M., Klingmüller, U., and Timmer, J. (2009).  
323 Structural and practical identifiability analysis of partially observed dynamical models by exploiting  
324 the profile likelihood. *Bioinformatics*, 25(15):1923–1929.
- 325 Rein, D., Schijlen, E., Kooistra, T., Herbers, K., Verschuren, L. Hall, R., Sonnewald, U., Bovy, A., and  
326 Kleemann, R. (2006). Transgenic flavonoid tomato intake reduces c-reactive protein in human c-reactive  
327 protein transgenic mice more than wild-type tomato. *The Journal of Nutrition*, 136:2331–2337.
- 328 Savageau, M. (1995). Enzyme kinetics *in vitro* and *in vivo*: Michaelis-Menten revisited. In Bittar, E.,  
329 editor, *Principles of Medical Biology*, volume 4, chapter 5, pages 93–146.
- 330 Schallau, K. and Junker, B. (2010). Simulating plant metabolic pathways with enzyme-kinetic models.  
331 *Plant Physiology*, 152(4):1763–1771.
- 332 Schmidt, H., Cho, K.-H., and Jacobsen, E. (2005). Identification of small scale biochemical networks  
333 based on general type system perturbations. *The FEBS Journal*, 272:2141–2151.
- 334 Schmidt, M., Vallabhajosyula, R., Jenkins, J., Hood, J., Soni, A., and Wikswo, J. (2011). Automated  
335 refinement and inference of analytical models for metabolic networks. *Phys. Biol.*, 8.
- 336 Srinath, S. and Gunawan, R. (2010). Parameter identifiability of power-law biochemical system models.  
337 *Journal of Biotechnology*, 149(3):132–140.
- 338 Stelling, J., Klamt, S., Bettenbrock, K., Schustert, S., and Gilles, E. (2002). Metabolic network structure  
339 determines key aspects of functionality and regulation. *Nature, letters*, 420(14).
- 340 Steuer, R. and Junker, B. (2009). Computational models of metabolism: Stability and regulation in  
341 metabolic networks. *Advances in Chemical Physics*, 142.
- 342 Varma, A. and Palsson, B. (1995). Parametric sensitivity of stoichiometric flux balance models applied to  
343 wild-type escherichia coli metabolism. *Biotechnol Bioeng*, 45:69–79.
- 344 Voit, E., Marino, S., and Lall, R. (2005). Challenges for the identification of biological systems from in  
345 vivo time series data. *In Silico Biol.*, 5:83–92.
- 346 Wang, J. (2009). Glycosyltransferases: key players involved in the modification of plant secondary  
347 metabolites. *Front. Biol. China*, 4(1):39–46.
- 348 Zhan, C. and Yeung, L. (2011). Parameter estimation in systems biology models using spline approxima-  
349 tion. *BMC Systems Biology*, 5(14).