

METHODOLOGY ARTICLE

Open Access

# Constructing stochastic models from deterministic process equations by propensity adjustment

Jialiang Wu<sup>1</sup>, Brani Vidakovic<sup>2</sup> and Eberhard O Voit<sup>2,3\*</sup>

## Abstract

**Background:** Gillespie's stochastic simulation algorithm (SSA) for chemical reactions admits three kinds of elementary processes, namely, mass action reactions of 0<sup>th</sup>, 1<sup>st</sup> or 2<sup>nd</sup> order. All other types of reaction processes, for instance those containing non-integer kinetic orders or following other types of kinetic laws, are assumed to be convertible to one of the three elementary kinds, so that SSA can validly be applied. However, the conversion to elementary reactions is often difficult, if not impossible. Within deterministic contexts, a strategy of model reduction is often used. Such a reduction simplifies the actual system of reactions by merging or approximating intermediate steps and omitting reactants such as transient complexes. It would be valuable to adopt a similar reduction strategy to stochastic modelling. Indeed, efforts have been devoted to manipulating the chemical master equation (CME) in order to achieve a proper propensity function for a reduced stochastic system. However, manipulations of CME are almost always complicated, and successes have been limited to relative simple cases.

**Results:** We propose a rather general strategy for converting a deterministic process model into a corresponding stochastic model and characterize the mathematical connections between the two. The deterministic framework is assumed to be a generalized mass action system and the stochastic analogue is in the format of the chemical master equation. The analysis identifies situations: where a direct conversion is valid; where internal noise affecting the system needs to be taken into account; and where the propensity function must be mathematically adjusted. The conversion from deterministic to stochastic models is illustrated with several representative examples, including reversible reactions with feedback controls, Michaelis-Menten enzyme kinetics, a genetic regulatory motif, and stochastic focusing.

**Conclusions:** The construction of a stochastic model for a biochemical network requires the utilization of information associated with an equation-based model. The conversion strategy proposed here guides a model design process that ensures a valid transition between deterministic and stochastic models.

## Background

Most stochastic models of biochemical reactions are based on the fundamental assumption that no more than one reaction can occur at the exact same time. A consequence of this assumption is that only elementary chemical reactions can be converted directly into stochastic analogues [1]. These include: 1) zero-order reactions, such as the generation of molecules at a constant rate; 2) first-order reactions, with examples including

elemental chemical reactions as well as transport and decay processes; and 3) second-order reactions, which include heterogeneous and homogeneous bimolecular reactions (dimerization). Reactions with integer kinetic orders other than 0, 1 and 2 are to be treated as combinations of sequential elementary reactions. The advantage of the premise of non-simultaneous reaction steps is that the stochastic reaction rate can be calculated from a deterministic, equation-based model with some degree of rigor, even though the derivation is usually not based on first physical principles but instead depends on other assumptions and on macroscopic information, such as a fixed rate constant in the

\* Correspondence: eberhard.voit@bme.gatech.edu

<sup>2</sup>The Wallace H. Coulter Department of Biomedical Engineering, Georgia Institute of Technology, Atlanta, GA30332, USA

Full list of author information is available at the end of the article

equation-based model. The severe disadvantage is that this rigorous treatment is not practical for modelling larger biochemical reaction systems. The reasons include the following. First, in many cases, elementary reaction rates are not available. Secondly, even in the case that all reaction parameters are available, the computational expense is very significant when the system involves many species and reactions, and this fact ultimately leads to a combinatorial explosion of required computations. Within a deterministic modelling framework, the common practice in this situation is to fit the transient and steady-state experimental data with a phenomenological, (differential) equation-based model, which explicitly or implicitly eliminates or merges some intermediate species and reactions. The best-known examples are probably Michaelis-Menten and Hill rate laws, which are ultimately explicit, but in truth approximate a multivariate system of underlying chemical processes.

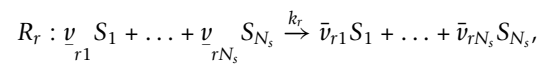
Similar model reduction efforts have been carried out for stochastic modelling. For instance, the use of a complex-order function (which corresponds to a reduced equation-based model) was shown to be justified for some types of stochastic simulations. A prominent example is again the Michaelis-Menten rate law, which can be reduced from a system of elementary reactions to an explicit function by means of the *quasi-steady-state assumption* (see Result section and [2,3]). However, model reduction within the stochastic framework has proven to be far more difficult than in the deterministic counterpart. The difficulties are mainly due to the fact that the reduction must be carried out on the chemical master equation (CME). This process is nontrivial and has succeeded only in simple cases.

In general, the construction of a stochastic model for a large biochemical network requires the use of information available from an equation-based model. In the past, several strategies have been proposed for this purpose and within the context of Gillespie's exact stochastic simulation algorithm (SSA; [1]) and its variants [4]. For example, Tian and Burrage [5] proposed that a stochastic model could be directly formulated from the deterministic model through a Poisson leaping procedure. However, a rigorous mathematical justification for such a conversion is lacking. Typical moment-based approaches [6-8] derive ODEs for the statistical moments of the stochastic model from an equation-based model where the 0<sup>th</sup>, 1<sup>st</sup> and 2<sup>nd</sup> order reactions follow mass action rate laws. More recently the moment method was extended to cover models consisting of rational rate laws [9]. Moreover, it was realized that the moment method is complementary to, but cannot fully replace, stochastic simulations, because it does not cover situations like genetic switches [6,10].

In this article, we explore the mathematical connection between deterministic and stochastic frameworks for the pertinent case of Generalized Mass Action (GMA) systems, which are frequently used in Biochemical Systems Theory (BST; [11-13]). Specifically, we address two questions: First, under what conditions can a deterministic, equation-based model be converted directly into a stochastic simulation model? And second, what is a proper way of implementing this conversion? We will develop a method to answer these questions and demonstrate it for functions in the canonical power-law format of GMA systems. However, the results are applicable to other functions and formats as well, as we will demonstrate with several examples.

### Representations of systems of biochemical reactions

Consider a well-stirred biochemical reaction system with constant volume and temperature, where  $N_s$  different chemical species  $\{S_s\}_{s=1}^{N_s}$ , interact through  $N_r$  unidirectional reaction channels  $\{R_r\}_{r=1}^{N_r}$ . Each reaction channel can be characterized as



where  $\underset{r_s}{\nu}$  and  $\bar{\nu}_{r_s}$  are the counts of molecular species  $S_s$  consumed and produced due to reaction  $R_r$ , respectively, and  $k_r$  is the rate constant. The changed amount of  $S_s$  is  $\bar{\nu}_{r_s} - \underset{r_s}{\nu}$ , which is due to the firing of reaction  $R_r$ , defines the stoichiometric coefficient of  $S_s$  in  $R_r$ . The stoichiometric coefficients of all species can be summarized according to each reaction  $R_r$  in the stoichiometric vector

$$\mathbf{v}_r \triangleq \begin{Bmatrix} \nu_{r1} \\ \vdots \\ \nu_{rN_s} \end{Bmatrix} \in \mathbb{Z}^{N_s}.$$

The stoichiometric vectors of all reactions can further be arranged as the stoichiometric matrix of the system

$$V \triangleq [\mathbf{v}_1, \dots, \mathbf{v}_{N_r}] \in \mathbb{Z}^{N_s \times N_r}.$$

The size of the system is defined as  $\Phi = AU$ , where  $A$  is the Avogadro number and  $U$  is the reaction volume.

The modelling of biochemical reaction networks typically uses one of two conceptual frameworks: deterministic or stochastic. In a deterministic framework, the state of the system is given by the a non-negative vector  $[\mathbf{X}(t)] = [[X_1(t)], \dots, [X_{N_s}(t)]]^T \in \mathbb{R}^{N_s}$ , where component  $[X_s(t)]$  represents the concentration of species  $S_s$ , measured in moles per unit volume. The temporal evolution of the state of the system is modelled by a set of

ordinary differential equations, which in our case are assumed to follow a generalized mass action (GMA) kinetic law. By contrast, in a stochastic framework, the state of the systems is characterized by a vector  $\mathbf{x}(t) = [x_1(t), \dots, x_{N_s}(t)]^T \in \mathbb{Z}^{N_s}$ , whose values are non-negative integers. Specifically,  $x_s(t) = \Phi[X_s(t)]$  is the count of  $S_s$  molecules, which is a sample value of the random variable  $X_s(t)$ . The system dynamics of this process is typically described with the chemical master equation (CME). Both GMA and CME will be discussed in detail in the following sections.

### Motivation for the power-law formalism: reactions in crowded media

Power-law functions with non-integer kinetics have proven very useful in biochemical systems analysis, and forty years of research have demonstrated their wide applicability (e.g., see [11-13]). Generically, this type of description of a biochemical reaction can be seen either as a Taylor approximation in logarithmic space or as a heuristic or phenomenological model that has been applied successfully hundreds of times and in different contexts, even though it is difficult or impossible in many situations to trace it back to first mechanistic principles. A particularly interesting line of support for the power-law format can be seen in the example of a bimolecular reaction occurring in a spatially restricted environment. Savageau demonstrated that the kinetics of such a reaction can be validly formulated as a generalization of the law of mass action, where non-integer kinetic orders are allowed [14,15]. Neff and colleagues [16-18] showed with careful experiments that this formulation is actually more accurate than alternative approaches.

Within the conceptual framework of power-law representations, the rate of the association reaction between molecules of species  $S_1$  and  $S_2$  is given as  $k[X_1(t)]^{f_1}[X_2(t)]^{f_2}$ . Here,  $k$  is the *rate constant* and  $f_1$  and  $f_2$  are real-valued *kinetic orders*, which are no longer necessarily positive integers as it is assumed in a mass action law. As an example, consider the reversible bimolecular reaction  $S_1 + S_2 \xrightleftharpoons[k_b]{k_f} S_3$ . Like Neff and colleagues [17], we begin by formulating a discrete update function for the population of  $S_3$  molecules as

$$x_3(t + \Delta t) - x_3(t) = f([X_1], [X_2])\Delta t x_1 x_2 - g([X_3])\Delta t x_3. \quad (1)$$

The first term on the right-hand side of this equation,  $f([X_1], [X_2])\Delta t x_1 x_2$ , describes the production of  $S_3$ : it depends on the totality of possible collisions  $x_1 x_2$  and also on some fraction  $f([X_1], [X_2])\Delta t$  that actually reacts and forms the product. In a dilute environment,  $f([X_1], [X_2])$  equals a traditional rate constant, and the reaction

obeys the law of mass action, while in a spatially restricted environment, such as the cytoplasm, one needs to take crowding effects into account. As shown in Savageau [14,15], the desired fraction of a reaction in a crowded environment becomes a *rate function* that depends on the current concentrations of  $S_1$  and  $S_2$ . The second term,  $g([X_3])\Delta t x_3$ , describes the fraction  $g([X_3])\Delta t$  of species  $S_3$  that dissociates back into  $S_1$  and  $S_2$ . This fraction may depend on some functional form of  $[X_3]$  because in a crowded environment the complex may not be able to dissociate effectively. Thus, rate constants in the generalized mass action setting become rate functions (cf. [17]).

By taking the limit  $\Delta t \rightarrow 0$ , one obtains the differential equation

$$\frac{dx_3}{dt} = f([X_1], [X_2])x_1 x_2 - g([X_3])x_3. \quad (2)$$

Savageau used Taylor series expansion to approximate the functions  $f$  and  $g$  in the logarithmic space ( $\log [X_1]$ ,  $\log [X_2]$ ) around some operating point  $(a, b)$ . The result for  $f$  is

$$\begin{aligned} \log f([X_1], [X_2]) &\triangleq F(\log [X_1], \log [X_2]) \\ &= F(a, b) + \left. \frac{\partial f([X_1], [X_2])}{\partial [X_1]} \right|_{(a,b)} (\log [X_1] - a) \\ &\quad + \left. \frac{\partial f([X_1], [X_2])}{\partial [X_2]} \right|_{(a,b)} (\log [X_2] - b) + \text{HOT} \\ &\approx k_f + \alpha \log [X_1] + \beta \log [X_2], \end{aligned} \quad (3)$$

where  $k_f$ ,  $\alpha$ , and  $\beta$  are constants related to the chosen operating point  $(a, b)$ . The final step is achieved by ignoring all higher order terms (HOT) beyond the constant and linear terms. Transformation back to the Cartesian space yields

$$f([X_1], [X_2]) \approx k_a [X_1]^\alpha [X_2]^\beta, \quad k_a = e^{k_f}. \quad (4)$$

The same procedure leads to the power-law expression for the degradation term:  $g([X_3]) \approx k_d [X_3]^\gamma$ . By combining constants we arrive at a power-law representation for the dynamics of species  $S_3$  as

$$\begin{aligned} \frac{d[X_3]}{dt} &= k_a [X_1]^\alpha [X_2]^\beta [X_1][X_2] - k_d [X_3]^\gamma [X_3] \\ &= k_a [X_1]^a [X_2]^b - k_d [X_3]^c, \end{aligned} \quad (5)$$

where  $a = \alpha + 1$ ,  $b = \beta + 1$ , and  $c = \gamma + 1$ . As long as  $k_f$ ,  $k_d$ ,  $a$ ,  $b$  and  $c$  remain more or less constant throughout a relevant range, the power-law model is mathematically well justified. In actual applications, the values of rate constants and kinetic orders can be estimated from experimental data [19]. When the functions  $f$  and  $g$  are

originally not in power-law format, they can be locally approximated by power-law functions with a procedure similar to the one shown above (Equations (3) to (5)). An illustration will be given in the example section.

### The Generalized Mass Action (GMA) format

In the GMA format within Biochemical Systems Theory, each process is represented as a univariate or multivariate power-law function. GMA models may be developed *de novo* or as an approximation of some other nonlinear rate laws. GMA models characterize the time evolution of the system state given that the system was in the state  $X(t_0)$  at some initial time  $t_0$ . Generically, the state of the system is changed within a sufficiently small time interval by one out of the  $N_r$  possible reactions that can occur in the system. The reaction velocity through reaction channel  $R_r$  is:

$$\left[ \frac{[X_1(t)]'}{v_{r1}} = \dots = \frac{[X_{N_s}(t)]'}{v_{rN_s}} \right] = k_r \prod_{s=1}^{N_s} [X_s(t)]^{f_{rs}} \quad (6)$$

for those  $v_{rs} = \bar{v}_{rs} - v_{rs} \neq 0$ ,  $s = 1, \dots, N_s$ . As shown in the example of a bimolecular reaction, the kinetic order  $f_{rs}$  associated with species  $S_s$  captures the effects of both reactant properties (such as the stoichiometric coefficient  $v_{rs}$ ) and environmental influences (such as temperature, pressure, molecular crowding effects, etc.). Therefore  $f_{rs}$  does not necessarily equal an integer  $v_{rs}$ , which is assumed to be the case in mass action kinetics, but is possibly real-valued and may be negative. Summing up the contributions of all reactions, one obtains a GMA model describing the dynamics of  $S_s$  as

$$\frac{d}{dt} [X_s(t)] = \sum_{r=1}^{N_r} v_{rs} k_r \prod_{s=1}^{N_s} [X_s(t)]^{f_{rs}} \quad (7)$$

for every  $s = 1, \dots, N_s$ . Each reaction contributes either a production flux or a degradation flux to the dynamics of a certain species. Positive terms ( $v_{rs} > 0$ ) represent the production of  $S_s$ , while negative terms ( $v_{rs} < 0$ ) describe degradation. If  $f_{rs}$  is positive, then  $S_s$  accelerates the reaction  $R_r$ ; a negative value represents that  $S_s$  inhibits the reaction, and  $f_{rs} = 0$  implies that  $S_s$  has no influence on the reaction. The rate constant  $k_r$  for reaction  $R_r$ , is either positive or zero. Both, the rate constant and the kinetic order, are to be estimated from data.

### Proper use of equation-based functions for stochastic simulations

The fundamental concept of a stochastic simulation is the propensity function  $\alpha(\mathbf{X})$ , and  $\alpha(\mathbf{X})dt$  describes the probability that a reaction will change the value of a

system variable within the next (infinitesimal) time interval  $(t, t + dt)$ . While a formal definition will be given later (Equation 18), it is easy to intuit that the propensity function is in some sense analogous to the rate in the corresponding deterministic model. In fact, the propensity function is traditionally assumed to be  $\alpha(\mathbf{X}) = f_s(\mathbf{X})$ , if the deterministic model is  $X_s' = f_s(\mathbf{X}, t)$ ,  $s = 1, \dots, N_s$ . However, a proper justification for this common practice is by and large missing. Indeed, we will show that the direct use of a rate function as the propensity function in a stochastic simulation algorithm requires that at least one of the following assumptions be true:

- 1)  $f$  is a linear function;
- 2) the reaction is monomolecular;
- 3) all  $X_i$  in the system are noise-free variables, *i.e.*, without (or with ignorable) fluctuations, which implies that the covariance of any two participating reactants is zero (or close to zero).

Each of these assumptions constitutes a sufficient condition for the direct use of a rate function as the propensity function and applies, in principle, to GMA as well as other systems. The validity of these conditions will be discussed later. Specifically, the first condition will be addressed in the Results section under the headings "0<sup>th</sup>-order reaction kinetics" and "1<sup>st</sup>-order reaction kinetics," while the second condition will be discussed under the heading "Real-valued order monomolecular reaction kinetics." The third condition will be the focus of Equations (29-36) and their associated explanations.

In reality, the rates of reactions in biochemical systems are commonly nonlinear functions of the reactant species, and fluctuations within each species are not necessarily ignorable. Therefore, to the valid use of an equation-based model in a stochastic simulation mandates that we know how to define a proper propensity function. The following section addresses this issue. It uses statistical techniques to characterize estimates for both the mean and variance of the propensity function, and these features will allow an assessment of the validity of the assumption  $\alpha(\mathbf{X}) = f_s(\mathbf{X})$  and prescribe adjustments if the assumption is not valid.

## Methods

### Deriving the mean and variance of a power-law function of random variables

Consider a generic power-law function of random variables  $X_s$  with the format  $PL(\mathbf{X}) = k \prod_{s=1}^{N_s} X_s^{f_s}$ . Estimates of its mean  $\mu_{PL}$  and variance  $\sigma_{PL}$  are given as

$$\mu_{PL} \approx k \prod_{s=1}^{N_s} \mu_s^{f_s} \exp \left( \sum_{i<j}^{N_s} f_i f_j \text{cov} [\log X_i, \log X_j] \right) \quad (8)$$

$$\sigma_{PL}^2 \approx \mu_{PL}^2 \Omega \quad (9)$$

(for details, see Additional file 1). Here,

$$\Omega = \sum_{s=1}^{N_s} f_s \mu_s^{-2} \sigma_s^2 + 2 \sum_{i<j}^{N_s} f_i f_j \text{cov} [\log X_i, \log X_j] \quad (10)$$

and  $\mu_s = E[X_s]$  and  $\sigma_s^2 = E[(X_s - \mu_s)^2]$  are the mean and variance of random variable  $X_s$ , respectively. If we choose to express  $\text{cov} [\log X_i, \log X_j]$  as a function of  $\mu_s$ ,  $\sigma_s^2$  and covariance  $\sigma_{ij} = \text{cov} [X_i, X_j]$ , using a Taylor approximation, we obtain

$$\mu_{PL} \approx k \prod_{s=1}^{N_s} \mu_s^{f_s} \exp \left( -\frac{1}{2} \sum_{s=1}^{N_s} f_s \sigma_s^2 / \mu_s^2 + \frac{1}{2} \Omega \right) \quad (11)$$

$$\sigma_{PL}^2 \approx \mu_{PL}^2 \Omega, \quad (12)$$

where

$$\begin{aligned} \Omega \approx & \sum_{s=1}^{N_s} f_s (\sigma_s / \mu_s)^2 + 2 \sum_{i<j}^{N_s} f_i f_j \left\{ \sigma_{ij} / (\mu_i \mu_j) \right. \\ & + \frac{1}{2} \log(\mu_i) (\sigma_j / \mu_j)^2 + \frac{1}{2} \log(\mu_j) (\sigma_i / \mu_i)^2 \\ & \left. - \frac{1}{4} (\sigma_i / \mu_i)^2 (\sigma_j / \mu_j)^2 \right\}. \end{aligned} \quad (13)$$

Since many biochemical variables approximately follow a log-normal distribution [20-22], it is valuable to consider the special situation where  $(X_1, \dots, X_s)$  is log-normally distributed (*i.e.*,  $(\log X_1, \dots, \log X_s)$  is normally distributed). In such a case, a simpler alternative way to calculate  $\text{cov} [\log X_i, \log X_j]$  is

$$\text{cov} [\log X_i, \log X_j] = \log \left( 1 + \frac{\sigma_{ij}}{\mu_i \mu_j} \right). \quad (14)$$

[23]. By substituting this result into (8)-(10), one obtains

$$\mu_{PL} \approx k \prod_{s=1}^{N_s} \mu_s^{f_s} \prod_{i<j}^{N_s} \left( 1 + \frac{\sigma_{ij}}{\mu_i \mu_j} \right)^{f_i f_j} \quad (15)$$

$$\sigma_{PL}^2 \approx \mu_{PL}^2 \Omega, \quad (16)$$

where

$$\Omega = \sum_{s=1}^{N_s} f_s \left( \frac{\sigma_s}{\mu_s} \right)^2 + 2 \sum_{i<j}^{N_s} f_i f_j \log \left( 1 + \frac{\sigma_{ij}}{\mu_i \mu_j} \right). \quad (17)$$

The approximation formulae for  $\mu_{PL}$  and  $\sigma_{PL}^2$  in eqns. (8)-(10) provide an easy numerical implementation if observation data are available to estimate  $\text{cov} [\log X_i, \log X_j]$ . Furthermore, Equations (11)-(13) demonstrate how  $\mu_{PL}$  and  $\sigma_{PL}^2$  are related to  $\mu_s$ ,  $\sigma_s^2$  and  $\sigma_{ij}$ ; however, the price of this insight is paid by the possible inaccuracy introduced through the Taylor approximation. Equations (15)-(17) also provide a functional dependence of  $\mu_{PL}$  and  $\sigma_{PL}^2$  on  $(\mu_s, \sigma_s^2, \sigma_{ij})$ , but it is only valid if the additional assumption of log-normality is acceptable.

### Deriving proper propensity functions for stochastic simulations from differential equation-based models

Assuming that the GMA model faithfully captures the average behaviour of a biochemical reaction system and recalling  $[\mathbf{X}(t)] = ([X_1(t)], \dots, [X_{N_s}(t)])^T$ , the expected metabolite numbers are defined as the expectation

$$E[\mathbf{X}] = [\mathbf{X}] \Phi, \quad (18)$$

where  $\Phi$  is the system size as defined above.

To describe the reaction channel  $R_r$  stochastically, one needs the state update vector  $\mathbf{v}_r$  and must characterize the quantity of molecules flowing through of reaction channel  $R_r$  during a small time interval. The key concept of this type of description is the propensity function  $\alpha_r(\mathbf{x})$ , which is defined as

$\alpha_r(\mathbf{x})dt$  = the probability that exactly one reaction  $R_r$  will occur some where inside  $U$  within infinitesimal(19) interval  $(t, t + dt)$ , given current state  $\mathbf{X}(t) = \mathbf{x}$ .

[1]. Because of the probabilistic nature of the propensity function,  $\mathbf{X}(t)$  is no longer deterministic, and the result is instead stochastic and based on the transition probability

$$P(\mathbf{x}, t | \mathbf{x}_0, t_0) = \text{Prob}\{\mathbf{X}(t) = \mathbf{x}, \text{ given } \mathbf{X}(t_0) = \mathbf{x}_0\}, \quad (20)$$

which follows the chemical master equation (CME)

$$\begin{aligned} \frac{\partial P(\mathbf{x}, t | \mathbf{x}_0, t_0)}{\partial t} = & \sum_{r=1}^{N_r} [\alpha_r(\mathbf{x} + \mathbf{v}_r) P(\mathbf{x} + \mathbf{v}_r, t | \mathbf{x}_0, t_0) \\ & - \alpha_r(\mathbf{x}) P(\mathbf{x}, t | \mathbf{x}_0, t_0)] \end{aligned} \quad (21)$$

Updating CME requires knowledge of every possible combination of all species counts within the population, which immediately implies that it can be solved analytically for only a few very simple systems and that

numerical solutions are usually prohibitively expensive [24]. To address the inherent intractability of CME, Gillespie developed an algorithm, called the *Stochastic Simulation Algorithm (SSA)*, to simulate CME models [1]. SSA is an exact procedure for numerically simulating the time evolution of a well-stirred reaction system. It is rigorously based on the same microphysical premise that underlies CME and gives a more realistic representation of a system's evolution than a deterministic reaction rate equation represented by ODEs. SSA requires knowledge of the propensity function, which however is truly available only for elementary reactions. These reactions include: 1) 0<sup>th</sup> order reactions, exemplified with the generation of a molecule at a constant rate; 2) 1<sup>st</sup> order monomolecular reactions, such as an elemental chemical conversion or decay of a single molecule; 3) 2<sup>nd</sup> order bimolecular reactions, including reactive collisions between two molecules of the same or different species. The reactive collision of more than two molecules at exactly the same time is considered highly unlikely and modelled as two or more sequential bimolecular reactions.

For elementary reactions, the propensity function of reaction  $R_r$  is computed as the product of a stochastic rate constant  $c_r$  and the number  $h_r$  of distinct combinations of reactant molecules, i.e.

$$\alpha_r(\mathbf{x}) = c_r h_r(\mathbf{x}), \quad r = 1, \dots, N_r. \quad (22)$$

$$\text{Here } h_r(\mathbf{x}) = \begin{cases} \prod_{s=1}^{N_s} \binom{x_s}{\nu_{rs}} \approx \frac{\prod_{s=1}^{N_s} x_s^{\nu_{rs}}}{\prod_{s=1}^{N_s} \nu_{rs}!}, & \text{for } x_s \geq \nu_{rs} > 0, \\ 0, & \text{otherwise} \end{cases}$$

where  $x_s$  is the sample value of random variable  $X_s$ . The approximation is invoked when  $x_s$  is large and  $(x_s - 1), \dots, (x_s - \nu_{rs} + 1)$  are approximately equal to  $x_s$ .

In Gillespie's original formulation [1]  $c_r$  is a constant that only depends on the physical properties of the reactant molecules and the temperature of the system, and  $c_r dt$  is the probability that a particular combination of reactant molecules will react within the next infinitesimally small time interval  $(t, t + dt)$ . The constant  $c_r$  can be calculated from the corresponding deterministic rate constants, if they are known.

Since the assumption of mass action kinetics is not valid generally, especially in spatially restricted environments and in situations dominated by macromolecular crowding, we address the broader scenario where  $c_r$  is not a constant but a function of the reactant concentrations. Thus, we denote  $c_r$  as a *stochastic rate function*, while retaining the definition of  $h_r$  as above. Knowing that any positive-valued differentiable function can be

approximated locally by a power-law function, we assume the functional form of the stochastic rate function as

$$c_r(\mathbf{x}) = \kappa_r \prod_{s=1}^{N_s} x_s(t)^{\varepsilon_{rs}}. \quad (23)$$

Here,  $\kappa_r$  and  $\varepsilon_{rs}$  are constants that will be specified in the next section, and  $r = 1, \dots, N_r$ . Note that  $\varepsilon_{rs}$  are now real-valued. Once the stochastic rate function is determined (see below), the propensity function can be calculated as

$$\alpha_r(\mathbf{x}) = c_r(\mathbf{x}) h_r(\mathbf{x}) = \frac{\kappa_r}{\prod_{s=1}^{N_s} \nu_{rs}!} \prod_{s=1}^{N_s} x_s^{\nu_{rs} + \varepsilon_{rs}}. \quad (24)$$

In order to identify the functional expression for a stochastic rate function, and thus the propensity function, we consider the connection between the stochastic and the deterministic equation models. By multiplying CME with  $\mathbf{x}$  and summing over all  $\mathbf{x}$ , we obtain

$$\frac{d}{dt} E[\mathbf{X}(t)] = \sum_{r=1}^{N_r} \mathbf{v}_r E[\alpha_r(\mathbf{X}(t))]. \quad (25)$$

Similarly, the expectation for any species  $X_s(t)$  is given as

$$\frac{d}{dt} E[X_s(t)] = \sum_{r=1}^{N_r} \nu_{rs} E[\alpha_r(\mathbf{X}(t))], \quad s = 1, \dots, N_s. \quad (26)$$

The details of these derivations are shown in Additional file 1.

We can use these results directly to compute the propensity function for a stochastic GMA model, assuming that its deterministic counterpart is well defined. Specifically, we start with the deterministic GMA equation for  $X_s$ ,

$$\frac{d}{dt} [X_s(t)] = \sum_{r=1}^{N_r} \nu_{rs} k_r \prod_{s'=1}^{N_s} [X_{s'}(t)]^{f_{rs'}}, \quad s = 1, \dots, N_s, \quad (27)$$

where  $\nu_{rs}$ ,  $k_r$  and  $f_{rs'}$  are again the stoichiometric coefficients, rate constants, and kinetic orders, respectively.

By substituting  $[X_s] = \frac{E[X_s]}{\Phi}$  from Equation (18) into this GMA model, we obtain a "particle-based" equation of the format

$$\frac{d}{dt} \left( \frac{E[X_s]}{\Phi} \right) = \sum_{r=1}^{N_r} \nu_{rs} k_r \prod_{s'=1}^{N_s} \left( \frac{E[X_{s'}]}{\Phi} \right)^{f_{rs'}}, \quad s = 1, \dots, N_s.$$

Elementary operations allow us to rewrite this equation as

$$\frac{d}{dt} (E[X_s]) = \sum_{r=1}^{N_r} v_{rs} k_r \Phi^{1-F_r} \prod_{s'=1}^{N_s} E[X_{s'}]^{f_{rs'}}, \quad s = 1, \dots, N_s, \quad (28)$$

where  $F_r = \sum_{s=1}^{N_s} f_{rs}$ . In this formulation, the differential operator is justified only when large numbers of molecules are involved. The assumption that the deterministic equations precisely capture the average behaviour of the biochemical reaction system directly equates the stochastic CME (25) to the deterministic equation based model (28)

$$E[\alpha_r(\mathbf{X}(t))] = k_r \Phi^{1-F_r} \prod_{s'=1}^{N_s} E[X_{s'}]^{f_{rs'}}. \quad (29)$$

Now we have two choices for approximating the expectation of the propensity function on left-hand side of equation (29):

1) adopt a zero-covariance assumption as was done in [25], which implies ignoring random fluctuations within every species as well as their correlations. This assumption is only justified for some special cases such as monomolecular and bimolecular reactions under the thermodynamic limit (cf. [4,6]), but is not necessary valid in generality. Here the thermodynamic limit is defined as a finite concentration limit which the system reaches when both population and volume approach infinity. Under this assumption, the left hand side of (29) becomes

$$E[\alpha_r(\mathbf{x})] = E \left[ \frac{\kappa_r}{\prod_{s=1}^{N_s} \frac{v_{rs}}{x_s^{v_{rs} + \epsilon_{rs}}}} \right] \quad (30)$$

$$= \frac{\kappa_r}{\prod_{s=1}^{N_s} v_{rs}!} \prod_{s=1}^{N_s} E[X_s]^{v_{rs} + \epsilon_{rs}}$$

for every  $r = 1, \dots, N_r$ , and Equation (24) yields

$$\epsilon_{rs} = f_{rs} - v_{rs}$$

$$\kappa_r = k_r \Phi^{1-F_r} \prod_{s=1}^{N_s} v_{rs}! \quad (31)$$

$$c_r(\mathbf{x}) = k_r \Phi^{1-F_r} \prod_{s=1}^{N_s} \frac{v_{rs}! x_s^{\epsilon_{rs}}}{v_{rs}!}$$

and

$$\alpha_{r_0}(\mathbf{x}) = k_r \Phi^{1-F_r} \prod_{s=1}^{N_s} x_s^{f_{rs}}. \quad (32)$$

Here, the index  $r_0$  is used to distinguish this 0-covariance propensity function from a second type of propensity in the next section.

With the zero-covariance assumption, one can substitute (32) back into the equation for the expectation for each species, which yields

$$\frac{d}{dt} E[X_s(t)] = \sum_{r=1}^{N_r} v_{rs} k_r \Phi^{1-F_r} \prod_{s=1}^{N_s} \mu_s^{f_{rs}} \quad (33)$$

for every  $s = 1, \dots, N_s$ . Note that this result is exactly equivalent to the equation-based model (27).

Equation (33) is based on assumption that both the fluctuations within species and their correlations are ignorable, which is not necessarily true in reality. If one uses it in simulations where the assumptions are not satisfied, it is possible that the means for the molecular species are significantly different from the corresponding equation-based model values. This discrepancy arises because the evolution of each species in the stochastic simulation is in truth affected by the covariance which is not necessarily zero, as it was assumed. This phenomenon was observed by Paulsson and collaborators [26] and further discussed in different moment-based approaches [6,7]. To assess the applicability limit of the propensity defined by (32), we can apply approximation techniques as shown in eqns. (8)-(10) on the functional expression of  $\alpha_{r_0}$  and obtain mean and variance as

$$\mu_{\alpha_{r_0}} = E[\alpha_{r_0}(\mathbf{X}(t))] \approx k_r \Phi^{1-F_r} \prod_{s'=1}^{N_s} E[X_{s'}]^{f_{rs'}} \exp \left( \sum_{i<j}^{N_s} f_{ri} f_{rj} \text{cov}[\log X_i, \log X_j] \right) \quad (34)$$

$$\sigma_{\alpha_{r_0}}^2 \approx \mu_{\alpha_{r_0}}^2 \Omega_r \quad (35)$$

where

$$\Omega_r = \sum_{s=1}^{N_s} f_{rs} \mu_s^{-2} \sigma_s^2 + 2 \sum_{i<j}^{N_s} f_{ri} f_{rj} \text{cov}[\log X_i, \log X_j], \quad (36)$$

for every  $s = 1, \dots, N_s$ . These expressions demonstrate that even with large numbers of molecules the mean of CME does not always converge to the GMA model. Indeed, the convergence is only guaranteed in one of the following special situations: 1) the reaction is of 0<sup>th</sup> order; 2) the reaction is a real value-order monomolecular reaction, with 1<sup>st</sup> order reaction as a special case; 3) the covariance contribution in (34) is sufficiently small

to be ignored for all participating reactant species of a particular reaction channel. Except for these three special situations, the covariance as shown in (34) significantly affects the mean dynamics. Therefore, stochastic simulations using zero-covariance propensity functions will in general yield means different from what the deterministic GMA model produces. How large these differences are cannot be said in generality. Under the assumption that the GMA model correctly captures the mean dynamics of every species, this conclusion means that  $\alpha_{r_0}$  is not necessarily an accurate propensity function for stochastic simulations, and the direct conversion of the equation-based model into a propensity function must be considered with caution.

Moreover, there is no theoretical basis to assume that there are no fluctuations in the molecular species or that these are independent. Therefore, we need to consider the second treatment of the expectation of the propensity function and study the possible effects of a non-zero covariance.

2) We again assume that the GMA model is well defined, which implies that information regarding the species correlations and fluctuations has been captured in the parameters of the GMA model on the left hand side of Equations (7) and (28). To gain information regarding correlations, we use Taylor expansion to approximate the propensity function (see Additional file 1 for details):

$$\begin{aligned}
 E[\alpha_r(\mathbf{X}(t))] &= E\left[\frac{\kappa_r}{\prod_{s=1}^{N_s} \frac{v_{rs}!}{X_s^{v_{rs}}}}\right] \\
 &\approx \frac{\kappa_r}{\prod_{s=1}^{N_s} v_{rs}!} E[X_s^{v_{rs}}] \\
 &\times \exp\left(\sum_{i<j}^{N_s} \left(\frac{v_{ri} + \varepsilon_{ri}}{v_{rj}} + \varepsilon_{rj}\right) \text{cov}[\log X_i, \log X_j]\right)
 \end{aligned} \tag{37}$$

After substitution of (37) in (29), one obtains

$$\begin{aligned}
 \kappa_r &= k_r \Phi^{1-F_r} \prod_{s=1}^{N_s} \frac{v_{rs}!}{X_s^{v_{rs}}} \exp\left(-\sum_{i<j}^{N_s} f_{ri}f_{rj} \text{cov}[\log X_i, \log X_j]\right) \\
 \varepsilon_{rs} &= f_{rs} - \frac{v_{rs}}{X_s}.
 \end{aligned}$$

Given the state  $\mathbf{x}$  of the system at time  $t$ , the stochastic rate function of reaction  $R_r$  is

$$\begin{aligned}
 c_r(\mathbf{x}) &= \kappa_r \prod_{s=1}^{N_s} X_s^{\varepsilon_{rs}} \\
 &= k_r \Phi^{1-F_r} \prod_{s=1}^{N_s} \frac{v_{rs}!}{X_s^{v_{rs}}} \\
 &\times \exp\left(-\sum_{i<j}^{N_s} f_{ri}f_{rj} \text{cov}[\log X_i(t), \log X_j(t)]\right) \prod_{s=1}^{N_s} X_s^{f_{rs}-v_{rs}}.
 \end{aligned} \tag{38}$$

Here it is important to understand that although the random variables  $\{X_s\}_{s \in S}$  appear in the expression  $c_r(\mathbf{x})$ ,  $c_r(\mathbf{x})$  is not a function of random variables but a deterministic function. The reason is that the  $\text{cov}[\log X_i(t), \log X_j(t)]$  in the composition of  $c_r(\mathbf{x})$ , which as the numerical characteristic of the random variables  $\{X_s\}_{s \in S}$ , is deterministic. Therefore, the stochastic rate function  $c_r(\mathbf{x})$  is a well-justified deterministic function that is affected by both the state of the system  $[x_1, \dots, x_{N_s}]$  and  $\text{cov}[\log X_i(t), \log X_j(t)]$ , the numerical characteristic of fluctuations in the random variables  $\{X_s\}_{s \in S}$ .

Given the expression  $c_r(\mathbf{x})$ , the propensity function is

$$\begin{aligned}
 \alpha_r(\mathbf{x}) &= c_r(\mathbf{x})h_r(\mathbf{x}) \\
 &= k_r \Phi^{1-F_r} \prod_{s=1}^{N_s} X_s^{f_{rs}} \\
 &\times \exp\left(-\sum_{i<j}^{N_s} f_{ri}f_{rj} \text{cov}[\log X_i(t), \log X_j(t)]\right).
 \end{aligned} \tag{39}$$

These results are based on the assumption that there are large numbers of molecules for all reactant species participating in reaction  $R_r$ . For simplicity of discussion, we define the *propensity adjustment factor* (paf) of reaction  $R_r$  as

$$\text{paf}(t) \triangleq \exp\left(-\sum_{i<j}^{N_s} f_{ri}f_{rj} \text{cov}[\log X_i(t), \log X_j(t)]\right). \tag{40}$$

*paf* is a function of time  $t$  and represents the contribution of the reactants to correlations among species in the calculation of the propensity function for reaction  $R_r$ . We denote the propensity function in (39), which accounts for the contribution of the covariance, as  $\alpha_{r\_cov}$ , in order to distinguish it from the propensity function  $\alpha_{r_0}$  (32), which is based on the assumption of zero-covariance, i.e.,

$$\alpha_{r\_cov}(\mathbf{x}) = \text{paf}(t)k_r \Phi^{1-F_r} \prod_{s=1}^{N_s} X_s^{f_{rs}}. \tag{41}$$

Remembering that  $\text{cov}[\log X_i(t), \log X_j(t)]$ , which is a component in both the stochastic rate function  $c_r(\mathbf{x})$  and now in the function *paf*( $t$ ), is a deterministic function rather than a function of random variables, *paf*( $t$ ) is



a deterministic correction to the kinetic constant  $k_r$  in the construction of  $\alpha_{r\_cov}$  in (41), which corrects the stochastic simulation toward the correct average.

In contrast to the propensity function  $\alpha_{r\_0}$ ,  $\alpha_{r\_cov}$  leads to accurate stochastic simulations. To illustrate this difference, we analyze  $\frac{d}{dt}E[X_s(t)]$  as follows: We apply the approximation techniques in eqns. (9)-(11) in order to obtain the mean and variance of the propensity function  $\alpha_{r\_cov}$ :

$$\mu_{\alpha_{r\_cov}} = E[\alpha_{r\_cov}(\mathbf{X}(t))] \approx k_r \Phi^{1-F_r} \prod_{s'=1}^{N_s} E[X_{s'}]^{f_{rs'}} \quad (42)$$

$$\sigma_{r\_cov}^2 \approx \mu_{\alpha_{r\_cov}}^2 \Omega_r. \quad (43)$$

Here

$$\Omega_r = \sum_{s=1}^{N_s} f_{rs} \mu_s^{-2} \sigma_s^2 + 2 \sum_{i<j}^{N_s} f_{ri} f_{rj} \text{cov}[1 \log X_i, \log X_j]. \quad (44)$$

By substituting (42) back into the derivation of CME (26), one obtains

$$\begin{aligned} & \frac{d}{dt} E[X_s(t)] \\ &= \sum_{r=1}^{N_r} v_{rs} E[\alpha_{r\_cov}(\mathbf{X}(t))] \\ &\approx \sum_{r=1}^{N_r} v_{rs} k_r \Phi^{1-F_r} \prod_{s'=1}^{N_s} \mu_{s'}^{f_{rs'}} \end{aligned} \quad (45)$$

for every  $s = 1, \dots, N_s$ , which is equivalent in approximation to the GMA model (28). In the other words, the mean of every molecular species obtained by using  $\alpha_{r\_cov}$  in the CME derived equation (27) is approximately identical to the corresponding macroscopic variable in the GMA model.

#### Calculation of $\text{cov}[\log X_i(t), \log X_j(t)]$

When data in the form of multiple time series for all the reactants are available, it is possible to compute  $\text{cov}[\log X_i(t), \log X_j(t)]$  directly from these data. Once this covariance is known, the function  $pdf$ ,  $\alpha_{r\_cov}$  and the mean dynamics can all be assessed. Alas, the availability of several time series data for all reactants under comparable conditions is rare, so that  $\text{cov}[\log X_i(t), \log X_j(t)]$  must be estimated in a different manner.

If one can validly assume that the covariance based on  $\alpha_{r\_0}$  does not differ significantly from the covariance based on  $\alpha_{r\_cov}$ , one may calculate  $\text{cov}[\log X_i(t), \log X_j(t)]$  by one of following methods.

Method 1:

One uses  $\alpha_{r\_0}$  to generate multiple sets of time series data of all reactants and then computes  $\text{cov}[\log X_i(t), \log X_j(t)]$ .

Method 2:

First,  $\text{cov}[\log X_i(t), \log X_j(t)]$  is expressed as a function of mean and covariance in one of the following ways; either as

$$\begin{aligned} \text{cov}[1 \log X_i, \log X_j] &\approx \sigma_{ij}/(\mu_i \mu_j) + \frac{1}{2} \log(\mu_i)(\sigma_j/\mu_j)^2 \\ &+ \frac{1}{2} \log(\mu_j)(\sigma_i/\mu_i)^2 - \frac{1}{4}(\sigma_i/\mu_i)^2(\sigma_j/\mu_j)^2 \end{aligned} \quad (46)$$

or as Equation (14):

$$\text{cov}[1 \log X_i, \log X_j] = \log \left( 1 + \frac{\sigma_{ij}}{\mu_i \mu_j} \right).$$

The first functional expression of  $\text{cov}[\log X_i(t), \log X_j(t)]$  is achieved by Taylor approximation, whereas the second expression is obtained by the additional assumption that the concentrations ( $X_1, \dots, X_s$ ) are log-normally distributed [8,23]. The consideration of a log-normal distribution is often justified by the fact that many biochemical data have indeed been observed to be log-normally distributed (e.g., [20-22]).

Second, one uses  $\alpha_{r\_0}$  to approximate the mean and covariance either by direct simulation, as shown in method 1, or by a moment-based approach, which is explained in Additional file 2, and which yields the differential equations

$$\begin{aligned} \frac{\partial \mu_s}{\partial t} &\approx \sum_{r=1}^{N_r} v_{r,s} \left\{ \alpha_{r\_0}(\boldsymbol{\mu}) + \frac{1}{2} \sum_{m,n=1}^{N_s} \frac{\partial^2 \alpha_{r\_0}(\boldsymbol{\mu})}{\partial X_m \partial X_n} \sigma_{mn} \right\} \\ \frac{\partial \sigma_{ij}}{\partial t} &\approx \sum_{r=1}^{N_r} \left\{ v_{r,i} \sum_{s=1}^{N_s} \frac{\partial \alpha_{r\_0}(\boldsymbol{\mu})}{\partial X_s} \sigma_{js} + v_{r,j} \sum_{s=1}^{N_s} \frac{\partial \alpha_{r\_0}(\boldsymbol{\mu})}{\partial X_s} \sigma_{is} \right. \\ &\quad \left. + v_{r,i} v_{r,j} \left[ \alpha_{r\_0}(\boldsymbol{\mu}) + \frac{1}{2} \sum_{m,n=1}^{N_s} \frac{\partial^2 \alpha_{r\_0}(\boldsymbol{\mu})}{\partial X_m \partial X_n} \sigma_{mn} \right] \right\} \end{aligned}$$

For convenience of computational implementation, the above equations can be written in matrix format

$$\begin{aligned} \frac{\partial \boldsymbol{\mu}}{\partial t} &\approx V^T \left( \boldsymbol{\alpha} + \frac{1}{2} \boldsymbol{\alpha}'' \odot \boldsymbol{\sigma} \right) \\ \frac{\partial \boldsymbol{\sigma}}{\partial t} &\approx (\boldsymbol{\sigma} (\boldsymbol{\alpha}' V))^T + \boldsymbol{\sigma} (\boldsymbol{\alpha}' V) + V^T \Lambda V. \end{aligned}$$

Here for  $r = 1, \dots, N_r$ , and  $s, m, n = 1, \dots, N_s$ ,  $\boldsymbol{\mu} = (\mu_1, \dots, \mu_{N_s})^T$ ,  $(V)_{rs} = v_{rs}$ ,  $\boldsymbol{\alpha} = (\alpha_1, \dots, \alpha_{N_r})^T$ ,

$$(\boldsymbol{\alpha}'')_{mn} = \frac{\partial^2 \alpha_r(\mathbf{X})}{\partial X_m \partial X_n}, \quad (\boldsymbol{\alpha}')_{mn} = \frac{\partial \alpha_r(\mathbf{X})}{\partial X_m \partial X_n},$$

$$\boldsymbol{\alpha}'' \odot \boldsymbol{\sigma} \triangleq (\boldsymbol{\alpha}''_1 \odot \boldsymbol{\sigma}, \dots, \boldsymbol{\alpha}''_{N_r} \odot \boldsymbol{\sigma})^T,$$

$$\alpha_r' = \left( \frac{\partial \alpha_r(\boldsymbol{\mu})}{\partial X_1}, \dots, \frac{\partial \alpha_r(\boldsymbol{\mu})}{\partial X_{N_s}} \right)^T, \quad \alpha' = (\alpha_1', \dots, \alpha_{N_r}'),$$

$$\alpha_r' = \left( \frac{\partial \alpha_r(\boldsymbol{\mu})}{\partial X_1}, \dots, \frac{\partial \alpha_r(\boldsymbol{\mu})}{\partial X_{N_s}} \right)^T, \text{ and } \Lambda \text{ is a diagonal}$$

matrix with  $(\Lambda)_{rr} = \alpha_r(\boldsymbol{\mu}) + \frac{1}{2} \sum_{m,n=1}^{N_s} \frac{\partial^2 \alpha_r(\boldsymbol{\mu})}{\partial X_m \partial X_n} \sigma_{mn}$ .

### Statistical criteria for propensity adjustment

Suppose an equation-based model captures the average behavior of a stochastic system and one intends to find the propensity function for a stochastic simulation that will reproduce that means. One can use the 95% confidence interval to evaluate the need for a propensity adjustment. Specifically, for stable systems that will reach a steady state, we use the reversible reaction model as an example. If the steady state of the ODE  $x_{st}$  is within the 95% confidence interval of  $n$  runs of stochastic simulations, i.e.

$$x_{st} \in \left[ \mu_{st} - 1.96 \frac{\delta_{st}}{\sqrt{n}}, \mu_{st} + 1.96 \frac{\delta_{st}}{\sqrt{n}} \right],$$

then the rate function in the original ODEs can be used as the propensity without adjustment; otherwise propensity adjustment is needed. Here  $\mu_{st}$  and  $\delta_{st}$  can be attained from either a moment-base method or from  $n$  independent runs of stochastic simulations using propensity without adjustment. An example discussing a reversible reaction with feedback controls can be found in the results section.

For other systems that do not reach a steady state, but where instead transient characteristics are of the highest interest, one can judge the need of propensity adjustment by whether the pertinent characteristics of the ODEs are within the 95% confidence interval of the corresponding characteristic, which is given by a prediction from the moment-based method or from  $n$  runs of stochastic simulations. The Repressilator example in the result section will serve as a demonstration.

## Results

### Generic special cases

It is generally not valid to translate a rate from a deterministic biochemical model into a propensity function of the corresponding stochastic simulation without adjustment (see Equations. (34)-(36)). However, in some situations, the propensity adjustment (e.g., Equations (40)-(44)) is not needed, and in some other cases it becomes relatively simple.

#### 1) 0<sup>th</sup>-order reaction kinetics

Consider a very simple equation-based model of the type

$$\frac{d[X_s(t)]}{dt} = k_r \text{ or } \frac{dE[X_s(t)]}{dt} = k_r \Phi, \quad (47)$$

for all  $s = 1, \dots, N_s, f_{rs} = 0$ . According to Equations (40)-(44), one obtains

$$\begin{aligned} \Omega_r &= 0 \\ \sigma_{\alpha_r}^2 &\approx 0 \\ \mu_{\alpha_r} &\approx \exp(\log(k_r \Phi)) = k_r \Phi \\ \text{i.e. } E[\alpha_r(\mathbf{X})] &\approx \alpha_r(E[\mathbf{X}]) \\ \alpha_{r\text{-cov}} &\approx k_r \Phi = \alpha_{r_0}. \end{aligned}$$

Thus, for a 0<sup>th</sup>-order reaction, its rate equation can be taken directly as the propensity function in stochastic simulations.

#### 2) 1<sup>st</sup>-order reaction kinetics

Direct application of Equations (40)-(44) yields

$$\begin{aligned} \frac{d[X_i(t)]}{dt} &= k_r [X_j(t)] \\ \text{or } \frac{dE[X_i(t)]}{dt} &= k_r E[X_j(t)], \end{aligned} \quad (48)$$

$f_{rs} = \delta_{sj}, i, j = 1, \dots, N_s$ . Therefore, according to Equations (40)-(44)

$$\begin{aligned} \Omega_r &= (\sigma_j / \mu_j)^2 \\ (\sigma_{\alpha_r} / \mu_{\alpha_r})^2 &= (\sigma_j / \mu_j)^2 \\ \mu_{\alpha_r} &\approx \exp(\log(k_r \mu_j)) = k_r \mu_j \\ \text{i.e. } E[\alpha_r(\mathbf{X})] &\approx \alpha_r(E[\mathbf{X}]) \\ \alpha_{r\text{-cov}}(\mathbf{X}) &\approx k_r X_j = \alpha_{r_0}(\mathbf{X}). \end{aligned}$$

Thus, for 1<sup>st</sup>-order reactions, the rate equation can again be taken directly as the propensity function in stochastic simulations.

#### 3) Real-valued order monomolecular reaction kinetics

Consider a reaction with kinetics of the type

$$\begin{aligned} \frac{d[X_i(t)]}{dt} &= k_r [X_j(t)]^{f_{rj}} \\ \text{or } \frac{dE[X_i(t)]}{dt} &= k_r \Phi^{1-f_{rj}} E[X_j(t)]^{f_{rj}}, \end{aligned} \quad (49)$$

$f_{rj} \neq 0, f_{rs} = 0$ , for any  $s \neq j, s = 1, \dots, N_s$ . Equations (40)-(44) lead to

$$\begin{aligned} \Omega_r &= (\sigma_j / \mu_j)^2 \\ (\sigma_{\alpha_r} / \mu_{\alpha_r})^2 &= (\sigma_j / \mu_j)^2 \\ \mu_{\alpha_r} &\approx k_r \Phi^{1-f_{ij}} \mu_j^{f_{ij}} \text{ i.e. } E[\alpha_r(\mathbf{X})] \approx \alpha_r(E[\mathbf{X}]) \\ \alpha_{r\_cov}(\mathbf{X}) &\approx k_r \Phi^{1-f_{ij}} X_j^{f_{ij}} = \alpha_{r\_0}(\mathbf{X}). \end{aligned}$$

Thus, for reaction kinetics involving a single variable and a real-valued order, the rate equation can again be taken as the propensity function in stochastic simulations.

#### 4) 2<sup>nd</sup>-order reaction kinetics

This type of reaction can be expressed as

$$\begin{aligned} \frac{d}{dt} [X_s(t)] &= k_r [X_i(t)] [X_j(t)] \\ \text{or } \frac{dE[X_s(t)]}{dt} &= k_r \Phi^{-1} E[X_i(t)] E[X_j(t)], \end{aligned} \quad (50)$$

$i, j \in \{1, \dots, N_s\}$ ,  $i \neq j$ ,  $f_{ri} = f_{rj} = 1$ , and  $f_{rs} = 0$ , for all  $s \neq i, j$ . Therefore, according to Equations (40)-(44)

$$\begin{aligned} \Omega_r &= (\sigma_i / \mu_i)^2 + (\sigma_j / \mu_j)^2 + 2\text{cov}[\log X_i, \log X_j] \\ &= (\sigma_i / \mu_i)^2 + (\sigma_s / \mu_s)^2 + 2\{\text{cov}[X_i / \mu_i, X_j / \mu_j] \\ &\quad + \frac{1}{2} \log(\mu_i)(\sigma_j / \mu_j)^2 + \frac{1}{2} \log(\mu_j)(\sigma_i / \mu_i)^2 - \frac{1}{4}(\sigma_i / \mu_i)^2 (\sigma_j / \mu_j)^2\} \end{aligned}$$

$$\begin{aligned} (\sigma_{\alpha_r} / \mu_{\alpha_r})^2 &= \Omega_r \\ \partial &= (\sigma_i / \mu_i)^2 + (\sigma_j / \mu_j)^2 + 2\text{cov}[\log X_i, \log X_j] \\ \mu_{\alpha_r} &\approx k_r (N_A V)^{-1} \mu_i \mu_j \\ \alpha_{r\_cov}(\mathbf{X}) &= k_r \Phi^{-1} X_i X_j \exp(-\text{cov}[\log X_i, \log X_j]) \neq \alpha_{r\_0}(\mathbf{X}). \end{aligned}$$

Thus, the proper propensity function for 2<sup>nd</sup>-order reactions is different from the rate equation. The difference can be ignored only if the contribution from the covariance is insignificant. In general, the rate equation yields only an approximate propensity function for stochastic simulations, and the approximation quality must be assessed on a case-by-case basis.

#### 5) Bimolecular reaction with real-valued order kinetics

This type of reaction can be formulated as

$$\begin{aligned} \frac{d[X_s(t)]}{dt} &= k_r [X_i(t)]^{f_i} [X_j(t)]^{f_j} \\ \text{or } \frac{dE[X_s(t)]}{dt} &= k_r \Phi^{1-f_i-f_j} E[X_i(t)]^{f_i} E[X_j(t)]^{f_j}, \end{aligned} \quad (51)$$

$i, j \in \{1, \dots, N_s\}$ ,  $i \neq j$ ,  $f_{ri}, f_{rj} \neq 0$ , and  $f_{rs} = 0$ , for all  $s \neq i, j$ . According to Equations (40)-(44) we obtain

$$\begin{aligned} \Omega_r &= (\sigma_i / \mu_i)^2 + (\sigma_j / \mu_j)^2 + 2f_i f_j \text{cov}[\log X_i, \log X_j] \\ &= (\sigma_i / \mu_i)^2 + (\sigma_j / \mu_j)^2 + 2f_i f_j \{\text{cov}[X_i / \mu_i, X_j / \mu_j] \\ &\quad + \frac{1}{2} \log(\mu_i)(\sigma_j / \mu_j)^2 + \frac{1}{2} \log(\mu_j)(\sigma_i / \mu_i)^2 - \frac{1}{4}(\sigma_i / \mu_i)^2 (\sigma_j / \mu_j)^2\} \end{aligned}$$

$$\begin{aligned} (\sigma_{\alpha_r} / \mu_{\alpha_r})^2 &= \Omega_r \\ &= (\sigma_i / \mu_i)^2 + (\sigma_j / \mu_j)^2 + 2f_i f_j \text{cov}[\log X_i, \log X_j] \\ \mu_{\alpha_r} &\approx k_r \Phi^{f_i+f_j-1} \mu_i^{f_i} \mu_j^{f_j} \\ \alpha_{r\_cov}(\mathbf{X}) &= k_r \Phi^{f_i+f_j-1} X_i^{f_i} X_j^{f_j} \exp(-f_i f_j \text{cov}[\log X_i, \log X_j]) \\ &\neq \alpha_{r\_0}(\mathbf{X}). \end{aligned}$$

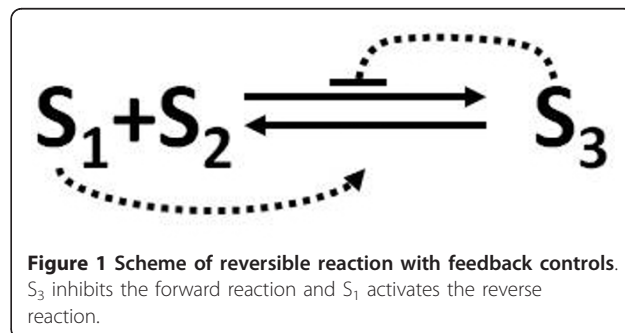
For bimolecular reactions of complex order, the propensity function is different from the rate equation. The difference can be ignored only if the contribution from the covariance is insignificant.

#### Power-law representation of a reversible reaction with feedback controls

We consider a reversible reaction with feedback controls (see Figure 1) whose average behaviour is accurately described by the following GMA model

$$\begin{aligned} \frac{dx_1}{dt} &= \frac{dx_2}{dt} = -\frac{dx_3}{dt} \\ &= -k_f \Phi^{1-f_1-f_2-f_3} x_1^{f_1} x_2^{f_2} x_3^{f_3} + k_b \Phi^{1-g_1-g_3} x_1^{g_1} x_3^{g_3}. \end{aligned} \quad (52)$$

Here  $S_3$  feeds back to inhibit the forward reaction and  $S_1$  feeds back on the reverse reaction and accelerates it. The task is to develop a stochastic model whose performance converges to that of the deterministic GMA model. We can see from equations (52) that three variables  $x_1$ ,  $x_2$  and  $x_3$  contribute to the forward flux  $k_f \Phi^{1-f_1-f_2-f_3} x_1^{f_1} x_2^{f_2} x_3^{f_3}$  and two variables  $x_1$  and  $x_3$  contribute to the backward flux  $k_b \Phi^{1-g_1-g_3} x_1^{g_1} x_3^{g_3}$ . Because several variables are involved, their covariance has the potential of affecting the forward and the backward propensity functions in a stochastic simulation. To obtain the covariance information, we formulate the moment equations (53) from the ODE model (52).



To simplify the calculation, as explained in detail in Additional file 2, we set the third central moment to zero and obtain a closed-form set of ODEs. Expressed differently, the rate of change in mean and covariance depends only on the functions of mean and covariance themselves, but not on higher-order moments. Thus,

$$\begin{aligned} \frac{\partial \mu}{\partial t} &\approx V^T \left( \alpha + \frac{1}{2} \alpha'' \odot \sigma \right) \\ \frac{\partial \sigma}{\partial t} &\approx (\sigma (\alpha' V))^T + \sigma (\alpha' V) + V^T \Lambda V. \end{aligned} \quad (53)$$

Here  $\mu = (\mu_1, \mu_2, \mu_3)^T$ ,  $V = \begin{bmatrix} -1 & -1 & 1 \\ 1 & 1 & -1 \end{bmatrix}$ ,

$$\alpha = (\alpha_1, \alpha_2)^T = \begin{bmatrix} k_f \Phi^{1-f_1-f_2-f_3} x_1^{f_1} x_2^{f_2} x_3^{f_3} \\ k_b \Phi^{1-g_1-g_2-g_3} x_1^{g_1} x_2^{g_2} x_3^{g_3} \end{bmatrix}.$$

Moreover, for  $r = 1, 2$  and  $m, n = 1, 2, 3$ ,

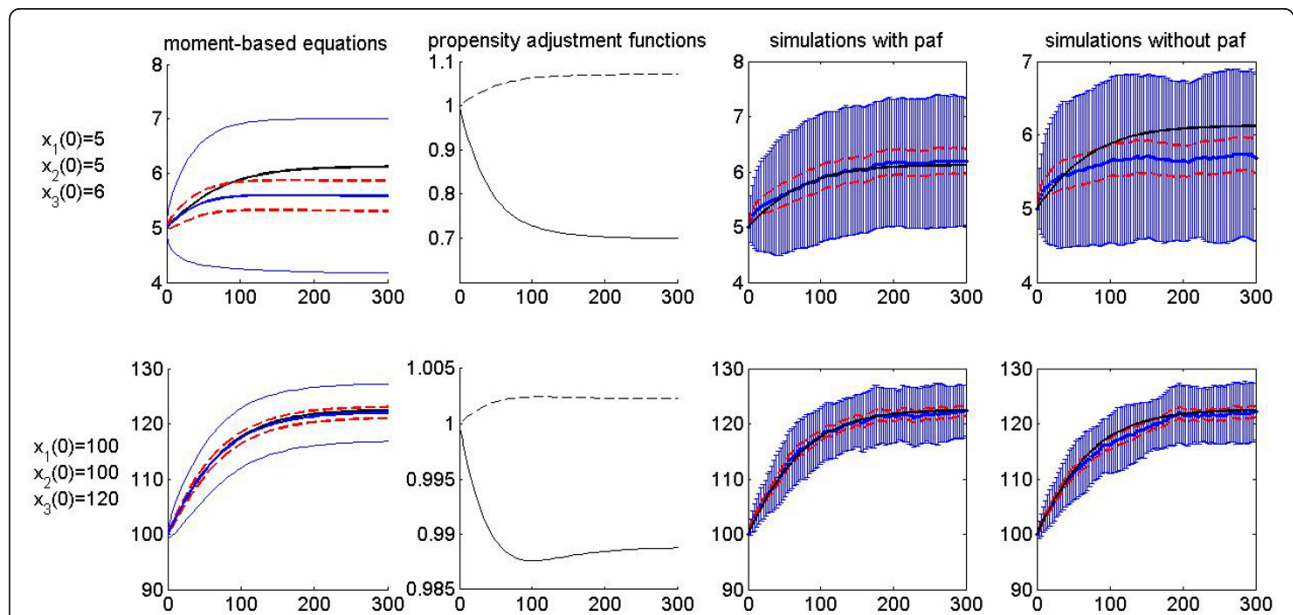
$$(\alpha''_{r'})_{mn} = \frac{\partial^2 \alpha_r(\mathbf{X})}{\partial x_m \partial x_n}, \alpha'' = (\alpha''_1, \alpha''_2)^T, \sigma = \begin{bmatrix} \sigma_{11} & \sigma_{12} & \sigma_{13} \\ \sigma_{21} & \sigma_{22} & \sigma_{23} \\ \sigma_{31} & \sigma_{32} & \sigma_{33} \end{bmatrix},$$

$$\alpha''_r \odot \sigma \triangleq \sum_{m,n=1}^3 \frac{\partial^2 \alpha_r(\mathbf{X})}{\partial x_m \partial x_n} \Big|_{\mathbf{X}=\mu} \sigma_{mn}, \alpha'' \odot \sigma \triangleq (\alpha''_1 \odot \sigma, \alpha''_2 \odot \sigma)$$

$$\sigma)^T, \alpha' = (\alpha'_1, \alpha'_2), \alpha'_r = \left( \frac{\partial \alpha_r(\mu)}{\partial x_1}, \frac{\partial \alpha_r(\mu)}{\partial x_2}, \frac{\partial \alpha_r(\mu)}{\partial x_3} \right)^T,$$

$$\text{and } \Lambda = \begin{bmatrix} \alpha_1(\mu) + \frac{1}{2} \sum_{m,n=1}^3 \frac{\partial^2 \alpha_1(\mu)}{\partial x_m \partial x_n} \sigma_{mn} & 0 \\ 0 & \alpha_2(\mu) + \frac{1}{2} \sum_{m,n=1}^3 \frac{\partial^2 \alpha_2(\mu)}{\partial x_m \partial x_n} \sigma_{mn} \end{bmatrix}.$$

Two initial conditions are chosen for representative simulations; they differ by a factor of 20 in species populations and reaction volume between the upper and lower panels of Figure 2. The purpose is to observe the thermodynamic limit of the systems: both scenarios have the same initial concentrations, but the system in the lower panel case has a larger species populations and reaction volume and can thus be regarded as the thermodynamic limit sample of system in the upper panel. As demonstrated by the figures in the first column, the moment approach predicts that for both population sizes the average trajectories of the stochastic model (without propensity adjustment) dynamics is lower than that of the equation-based model: the differences are about 10% of the steady-state value of the equation-based model in the upper figure and 1% in the lower figure; for 100 runs of the stochastic simulation,



**Figure 2 Comparative simulation results for a reversible reaction with feedback controls.** In all panels, the x-axis denotes time in seconds and the y-axis represents the number of molecules of species  $S_1$ . The upper and lower panels use two different sets of initial numbers of molecules, namely:  $(x_1(0), x_2(0), x_3(0), U) = (5, 5, 6, 1\mu\text{m}^3)$  and  $(x_1(0), x_2(0), x_3(0), U) = (100, 100, 120, 20\mu\text{m}^3)$ , respectively. Other simulation parameters are  $(f_1, f_2, f_3, g_1, g_2, g_3, k_f, k_b) = (1.3, 1.8, -1, 1, 1, 0.5, 0.5)$ . In both the upper and lower panels, the first column compares the time evolution of  $S_1$  molecules by different methods: the black line shows the ODE solution of Equation (52) for  $x_1$ ; the blue lines are the solutions of Equation (53) for  $\mu_1$  and for  $\mu_1 \pm \sigma_1$ , respectively. The red dotted lines framing the mean indicate the 95% confidence interval. The second column shows the propensity adjustment functions for the forward reaction (solid line) and the backward reaction (dashed line). The third column shows 100 independent stochastic simulations with propensity adjustment (blue means and error bars), in comparison with the ODE (Equation (52)) prediction (black line). The fourth column shows a second set of 100 independent stochastic simulations without propensity adjustment (blue means and error bars), in comparison with the ODE (Equation (52)) prediction (black line). The red dotted lines framing the mean in columns 3 and 4 again indicate the 95% confidence intervals.

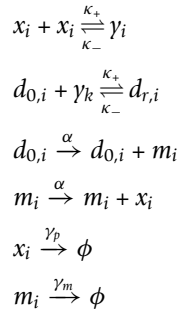
the steady-state value of the equation-based model lies outside the 95% confidence interval in the upper figure, while it is inside the interval in the lower figure. Therefore, we can expect that the propensity adjustment will significantly contribute to the stochastic simulation for the upper case while not for the lower case. This expectation is confirmed by the simulation results in the third and fourth columns. With the common assumption that the deterministic equations precisely capture the system's average behaviour, the case in the upper panel represents the situation where propensity adjustment is needed, while the lower panel represents the situation that a propensity without adjustment is sufficient when the system approaches its thermodynamics limit. This example furthermore demonstrates that either the moment approach or the stochastic simulations without propensity adjustment can be used to estimate whether there is a need to construct a propensity adjustment function for stochastic simulations.

### Repressilator

Interestingly, a propensity function may even be obtained through power-law approximation of some function that describes complex transient behaviours of a reaction network. As an example, consider the so-called *Repressilator* [27], which is a three-component genetic circuit where each component represses its downstream neighbour. More specifically (as shown in Figure 3), gene  $G_1$  codes for protein  $x_1$ , whose dimer  $y_1$  subsequently represses the transcription of the gene  $G_2$ . Similarly,  $y_2$ , the dimer of gene  $G_2$ 's protein product  $x_2$ , represses the transcription of gene  $G_3$ , and  $y_3$ , the dimer of gene  $G_3$ 's protein product  $x_3$ , represses the transcription of gene  $G_1$ . The corresponding differential equation model following mass action kinetics is given by [28]

$$\begin{aligned} x_i' &= -2\kappa_+x_i^2 + 2\kappa_-y_i + \sigma m_i - \gamma_p x_i \\ y_i' &= \kappa_+x_i^2 - \kappa_-y_i - k_+y_id_{0,j} + k_-d_{r,j} \\ d_{0,i}' &= -k_+\gamma_k d_{0,i} + k_-d_{r,i} \\ d_{r,i}' &= k_+\gamma_k d_{0,i} - k_-d_{r,i} \\ m_i' &= d_{0,i} - \gamma_m m_i \end{aligned} \tag{54}$$

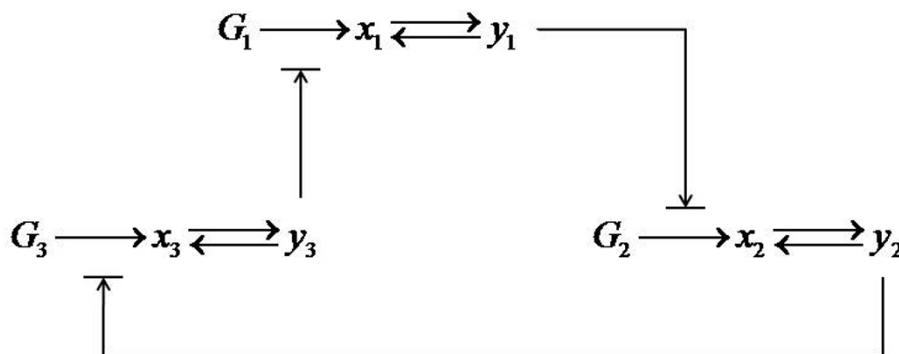
where  $i = 1, 2, 3; j = 2, 3, 1; k = 3, 1, 2$ ; the rate constants are explained in the diagram below



Assuming that the reversible dimerization and the dissociation/association of a protein dimer from/to the promoter are much faster than other processes, the full systems can be reduced to

$$\begin{aligned} x_i' &= \sigma p(x_i)^{-1} m_i - \gamma_p p(x_i)^{-1} x_i \\ m_i' &= \frac{\alpha d}{1 + c_d c_p x_i^2} - \gamma_m m_i \end{aligned} \tag{55}$$

[28]. Here  $\Phi = 1$ ,  $p(x_i) = 1 + 4c_p x_i + \frac{4c_d c_p d x_i}{(1 + c_d c_p x_i^2)^2}$ ,  $c_p = \kappa_+/\kappa_-$ ,  $c_d = k_+/k_-$  and  $d = d_{0,i} + d_{r,i}$  for  $i = 1, 2, 3$ . It has been shown that the simplified ODEs rather accurately approximate the transient dynamics of the full system by retaining the original oscillation period and amplitude.



**Figure 3 Reaction scheme of the Repressilator.** Gene  $G_1$  codes for protein  $x_1$ , whose dimer  $y_1$  represses the transcription of gene  $G_2$ . Similarly,  $y_2$ , the dimer of gene  $G_2$ 's protein product  $x_2$ , represses the transcription of gene  $G_3$ , and  $y_3$ , the dimer of gene  $G_3$ 's protein product  $x_3$ , represses the transcription of gene  $G_1$ .

In [28], the system (55) is further rescaled by setting  $\tilde{t} = \gamma_m t$ ,  $\tilde{x}_i = \sqrt{c_d c_p} x_i$  and  $\tilde{m}_i = (\sigma \sqrt{c_d c_p} m_i) / (\gamma_m \beta)$ , which yields

$$\begin{aligned} \frac{d\tilde{x}_i}{d\tilde{t}} &= \beta p(\tilde{x}_i)^{-1} \tilde{m}_i - \beta p(\tilde{x}_i)^{-1} \tilde{x}_i \\ \frac{d\tilde{m}_i}{d\tilde{t}} &= \frac{\kappa d'}{1 + \tilde{x}_k^2} - \tilde{m}_i. \end{aligned} \quad (56)$$

Intriguingly, one makes the following observation. The scaled ODE system (56) is consistent with the original system (55) in oscillation amplitude and period. However, its corresponding stochastic model produces results that deviate substantially from the average responses. To see the effects of the transition from a deterministic to a stochastic model, we apply SSA to the scaled system (56). The main result is that the oscillation periods of both  $x_i$  and  $m_i$  are reduced to half (Figure 4). The reason is that, in the stochastic simulation, the oscillation period is very sensitive to the ratio of  $x_i$  and  $m_i$ , which has been altered by the scaling operation. Therefore, in general one needs to pay attention to how scaling may affect the stochastic performance when the model is generated through the conversion of an ODE model.

We can see from equations (55) that two variables  $x_i$  and  $m_i$  contribute to the production of  $x_i$ ; hence, their covariance could affect the propensity function of  $x_i$  in the production reaction of a stochastic simulation. Similar to the example of a reversible reaction (Equation 52), it is therefore necessary to evaluate covariance effects and to judge whether the propensity function needs adjusting. Thus, we need to compare the difference between the dynamics of the phenomenological model (55) and the dynamics under the influence of covariance, which can be produced by either stochastic simulation or the moment approach.

The influence of the covariance on the dynamics of the stochastic simulation is relatively easy to assess: we simply use the terms on the right-hand side of the differential equations (54) as the propensity functions in SSA and obtain simulation results shown in the 2<sup>nd</sup> and the 4<sup>th</sup> panels of Figure 5. Obtaining the covariance-influenced dynamics with the moment-based approach is more complicated, and we need to discuss some implementation issues.

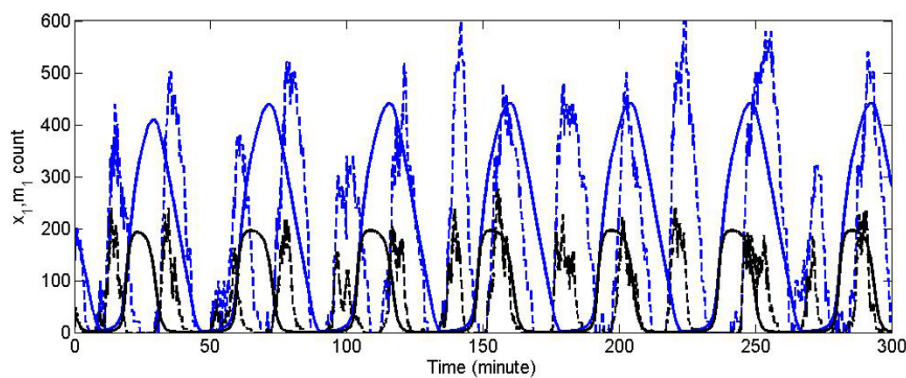
First, the moment-based approach requires information regarding the first and the second derivatives of  $p(x_i)^{-1}$ , which have rather complicated functional forms. To simplify the calculation, we replace the function  $p(x_i)^{-1}$  with an approximating power-law function. Specifically, suppose the original parameter values are  $\kappa_+ = k_+ = 5$ ,  $\kappa_- = k_- = 100$  and  $d = 20$ . Plotting the data  $(x_i, p(x_i)^{-1})$  in log-log space (Figure 5) indicates that the original function is represented well by a straight line:

$$\log y_i = \log 3.5188 - 0.9384 \log x_i.$$

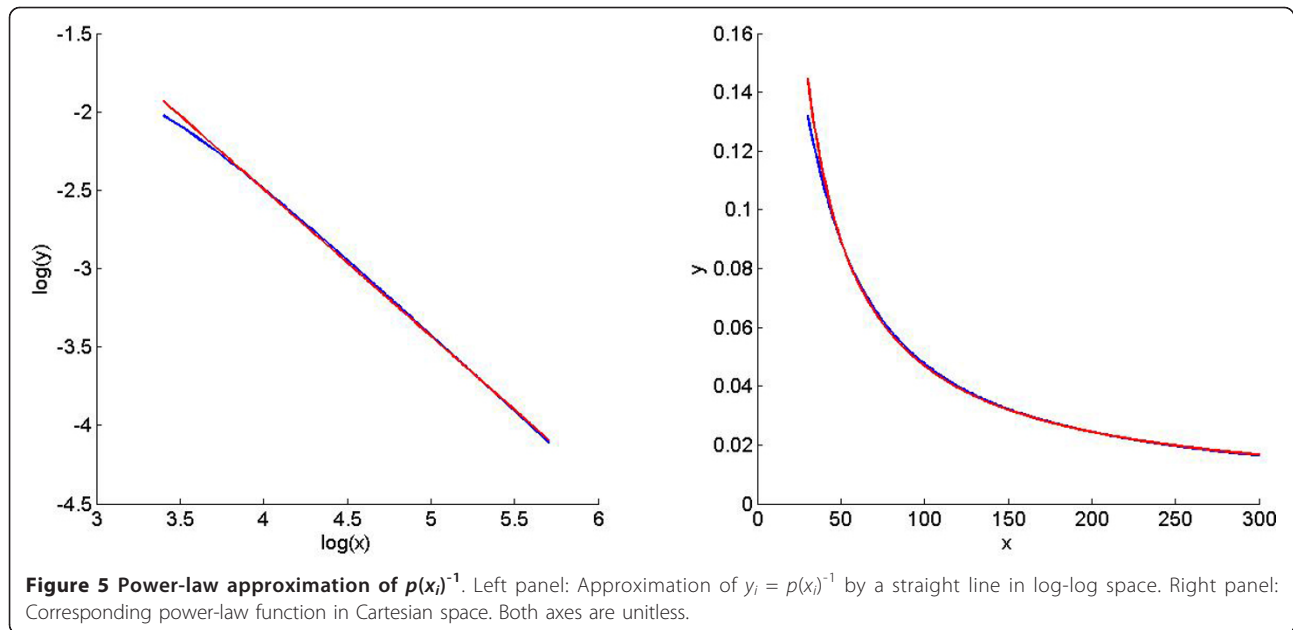
for  $x_i \in [30, 300]$ . In Cartesian space, this line corresponds to the power-law function

$$y_i = 3.5188 x_i^{-0.9384},$$

which models the original function very well (see Figure 5). For  $x_i \in [1, 30]$ , this power-law function does not fit the original function precisely; the effect of this imprecision can be evaluated later after we use this power-law function in the moment-based method. Moreover, using the truncated moment equations to estimate the mean and variance involves multiple approximations: First, the function  $p(x_i)^{-1}$  on the right-hand side of (55) is replaced by a power-law function (see Figure 5). Second, the result is approximated by Taylor expansion to the second order. Third, similar to the example of a reversible reaction, the central moment



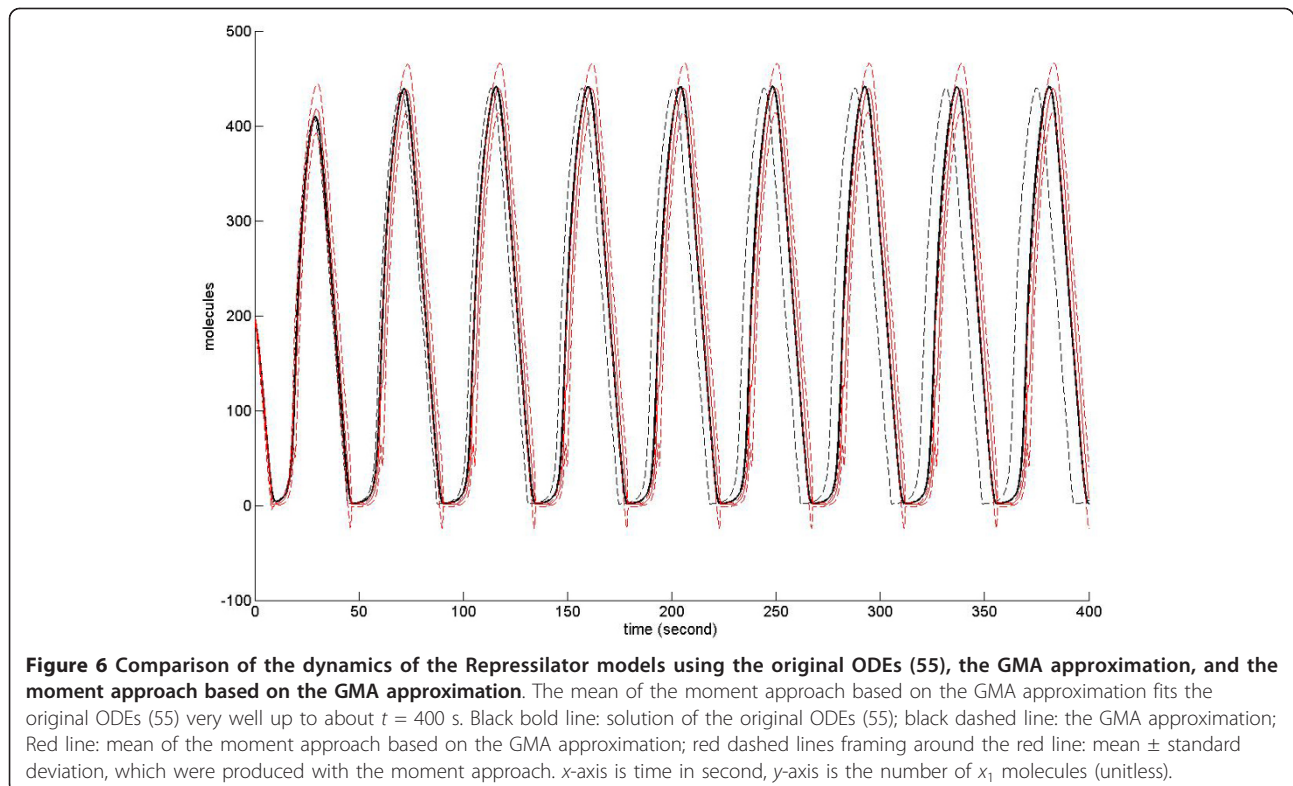
**Figure 4** Scaling of the Repressilator equations changes the oscillation period in the stochastic simulation. Solid lines represent solutions of ODEs (56), while dotted lines are trajectories of a stochastic simulation; blue lines represent  $x_i$  and black lines represent  $m_i$ .



of the third degree is assumed to be zero, which leads to a closed-form ODE for the first two moments.

Solving the technical issues as described, one obtains the corresponding moment-based model of (55) (not shown) with results shown in Figure 6. Suppose one is particularly interested in the period and amplitude of

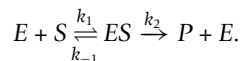
the oscillation within a time interval between 0 and 400 seconds. As shown in Figure 6, the GMA approximation (black dashed line) fits the original ODEs (55) (bold black solid line) very well at the beginning, but as time goes on, the approximation error accumulates. As seen in the time interval [350, 400], the GMA approximation



deviates from the original ODEs significantly. However, this does not mean that the GMA approximation cannot be used as a propensity function for stochastic simulations; the moment-based method with the GMA approximation shows that, when the GMA approximation is used as propensity function (without adjustment) for stochastic simulations, the resulting mean (red solid line) consistently fits the trajectory of the original ODEs (bold solid black line) very well up to about  $t = 400$  seconds. The oscillation period and amplitude in the stochastic simulation based on the GMA approximation (without adjustment) are almost identical to those of the original ODEs. Therefore, a propensity adjustment for the GMA approximation is not needed, and the GMA approximation can be used as a propensity function for stochastic simulations. In other words, a stochastic model for the Repressilator system can be generated by using the scheme in (32) without propensity adjustment. Moreover, the imprecision caused by the power-law approximation can be tolerated when its corresponding moment-based mean matches the original ODEs sufficiently well with respect to the features of highest interest.

#### Enzymatic reaction using a quasi-steady state assumption (QSSA)

We consider an enzymatic reaction following the Michaelis-Menten mechanism:



Here enzyme E reacts with substrate S through a reversible reaction to form complex ES, which can proceed to yield product P and to release the enzyme E. By assuming the law of mass action for the reaction kinetics we obtain a set of differential equations for the system dynamics:

$$\begin{aligned} \frac{d[S]}{dt} &= k_{-1}[ES] - k_1[S]([E]_0 - [ES]) \\ \frac{d[ES]}{dt} &= k_1[S]([E]_0 - [ES]) - (k_1 + k_2)[ES] \\ \frac{d[P]}{dt} &= k_2 [ES], \end{aligned} \quad (57)$$

where the total amount of enzyme in the form of free enzyme and complex  $[E]_0 \triangleq [E] + [ES]$  is assumed to be constant. In addition, by making the so-called *quasi steady state assumption* (QSSA) [29,30], assuming that the complex ES is essentially in steady state, we can assert  $\frac{d[ES]}{dt} \approx 0$ . As it has been discussed many times in the literature, QSSA reduces the system and leads to the approximate form

$$\begin{aligned} \frac{d[S]}{dt} &= -\frac{V_{\max}[S]}{K_m + [S]} \\ \frac{d[P]}{dt} &= \frac{V_{\max}[S]}{K_m + [S]}, \end{aligned} \quad (58)$$

which is known as *Michaelis-Menten kinetics* [30]. The characterizing parameters are  $V_{\max} = k_2[E]_0$  and  $K_m = (k_{-1} + k_2)/k_1$ .

Applying QSSA, Rao and Arkin [2] were able to reduce the CME of S and ES to a CME only containing S. For the reduced CME, the propensity function for the overall reaction  $S \rightarrow P$  is

$$\alpha(s) = \frac{V_{\max}s}{K_m + s}, \quad (59)$$

where the volume was scaled so that  $\Phi = 1$  and the lower-case letter  $s$  denotes the molecule count of species S. Instead of reviewing the relatively complicated manipulations with CME, we show in the following that the techniques described above lead directly from the equation-based model to the propensity function for the reduced system.

First, we recast the equation-based model into the GMA format [31], by introducing an auxiliary variable  $[T] \triangleq K_m + [S]$ . The result,

$$\begin{aligned} \frac{d[S]}{dt} &= -V_{\max}[S][T]^{-1} \\ \frac{d[T]}{dt} &= \frac{d[S]}{dt} = -V_{\max}[S][T]^{-1} \end{aligned} \quad (60)$$

is exactly equivalent to the reduced system in (58) with the initial condition  $[S]_0$  and  $[T]_0 = K_m + [S]_0$ . The corresponding stochastic model has only one reaction channel and the propensity function is

$$\alpha(s, t) = V_{\max}st^{-1}. \quad (61)$$

The propensity adjustment factor can be set to 1 because  $T$  is a function of  $s$  and its covariance with  $s$  is therefore 1. By applying  $t = K_m + s$ , the propensity function can be simplified as

$$\alpha(s, t) = V_{\max}st^{-1} = V_{\max}s(K_m + s)^{-1} = \alpha(s). \quad (62)$$

Thus, we arrive at the propensity function for the reduced system, which is identical to the result of Rao and Arkin obtained through manipulations of CME.

In the above derivation, we used the simplest type of recasting, where a new, auxiliary variable simply consists of an old variable plus a constant. This reformulation of the Michaelis-Menten process as a pair of GMA equations is a special case of a much more general recasting technique that permits the equivalent conversion of any

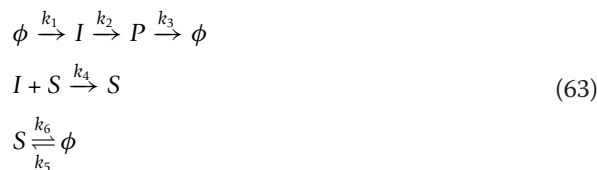


system of ordinary differential equations into a power-law format [31]. However, this equivalence transformation imposes constraints on the variables of the GMA equations, and it is at this point unclear whether there are mathematical warranties ensuring that the proposed transition from differential to stochastic equations in general preserves these constraints in all cases. This question will require further investigation.

### Stochastic Focusing

Stochastic focusing [26] describes the phenomenon that the fluctuations of a chemical species can drive the system to reach a different steady state than what a deterministic ODE model predicts. To demonstrate the utility of propensity adjustment, we derive a stochastic model which produces consistent results with those of the deterministic model.

Following [32], we consider the following reactions system



This system can be interpreted as follows: the intermediate species  $I$  is produced at constant rate  $k_1$  from some source  $\Phi$  and degrades with rate  $k_4$  through the catalysis with signalling molecule  $S$ ; the end product  $P$  is converted from species  $I$  at rate  $k_2$  and degrades at rate  $k_3$ ; the signalling molecule  $S$  is produced and degrades at rates  $k_5$  and  $k_6$ , respectively. Moreover, the value of  $k_5$  is reduced to half at a certain time point to achieve a significant divergence effect. In order to capture the average dynamics of the system accurately, we use a power-law model in GMA format instead of the mass action rate law in [32].

$$\begin{aligned}
 \frac{di}{dt} &= k_1 - k_2i - k_4i^f s^f \\
 \frac{dp}{dt} &= k_2i - k_3p \\
 \frac{ds}{dt} &= k_5 - k_6s
 \end{aligned}
 \tag{64}$$

The system size is set to 1. We can see from equations (64) that two variables  $i$  and  $s$  contribute to the degradation of  $I$  and that their covariance could therefore affect the propensity function of  $I$  in the degradation reaction of a stochastic simulation. To calculate the propensity adjustment function  $pdf_4(t) = \exp(-f_i f_s \text{cov}[\log I(t), \log S(t)])$  for reaction  $R_4 : I + S \xrightarrow{k_4} S$ , we formulate equations (cf. (60)) for the moments as

$$\begin{aligned}
 \frac{\partial \mu}{\partial t} &\approx V^T \left( \alpha + \frac{1}{2} \alpha'' \odot \sigma \right) \\
 \frac{\partial \sigma}{\partial t} &\approx (\sigma(\alpha'V))^T + \sigma(\alpha'V) + V^T \Lambda V.
 \end{aligned}
 \tag{65}$$

$$\text{Here } \mu = (\mu_I, \mu_P, \mu_S)^T, \quad V = \begin{bmatrix} 1 & 0 & 0 \\ -1 & 1 & 0 \\ 0 & -1 & 0 \\ -1 & 0 & 0 \\ 0 & 0 & 1 \\ 0 & 0 & -1 \end{bmatrix},$$

$$\alpha = \begin{bmatrix} \alpha_1 \\ \alpha_2 \\ \alpha_3 \\ \alpha_4 \\ \alpha_5 \\ \alpha_6 \end{bmatrix} = \begin{bmatrix} k_1 \\ k_2i \\ k_3p \\ k_4i^f s^f \\ k_5 \\ k_6s \end{bmatrix}.$$

Moreover, for  $r = 1, \dots, 6$  and  $m, n = i, p, s$ ,

$$(\alpha''_{r})_{mn} = \frac{\partial^2 \alpha_r}{\partial m \partial n}, \quad \alpha'' = (\alpha_1'', \dots, \alpha_6'')^T,$$

$$\alpha''_r \odot \sigma = \sum_{m,n=i,p,s} \frac{\partial^2 \alpha_r(\mu)}{\partial m \partial n} \sigma_{mn},$$

$$\alpha''_r \odot \sigma = \sum_{m,n=i,p,s} \frac{\partial^2 \alpha_r(\mu)}{\partial m \partial n} \sigma_{mn}, \quad \alpha'' \odot \sigma \triangleq (\alpha_1'' \odot \sigma, \dots,$$

$$\alpha_6'' \odot \sigma)^T, \quad \alpha' = (\alpha_1', \dots, \alpha_6'), \quad \alpha'_r = \left( \frac{\partial \alpha_r}{\partial i}, \frac{\partial \alpha_r}{\partial p}, \frac{\partial \alpha_r}{\partial s} \right)^T,$$

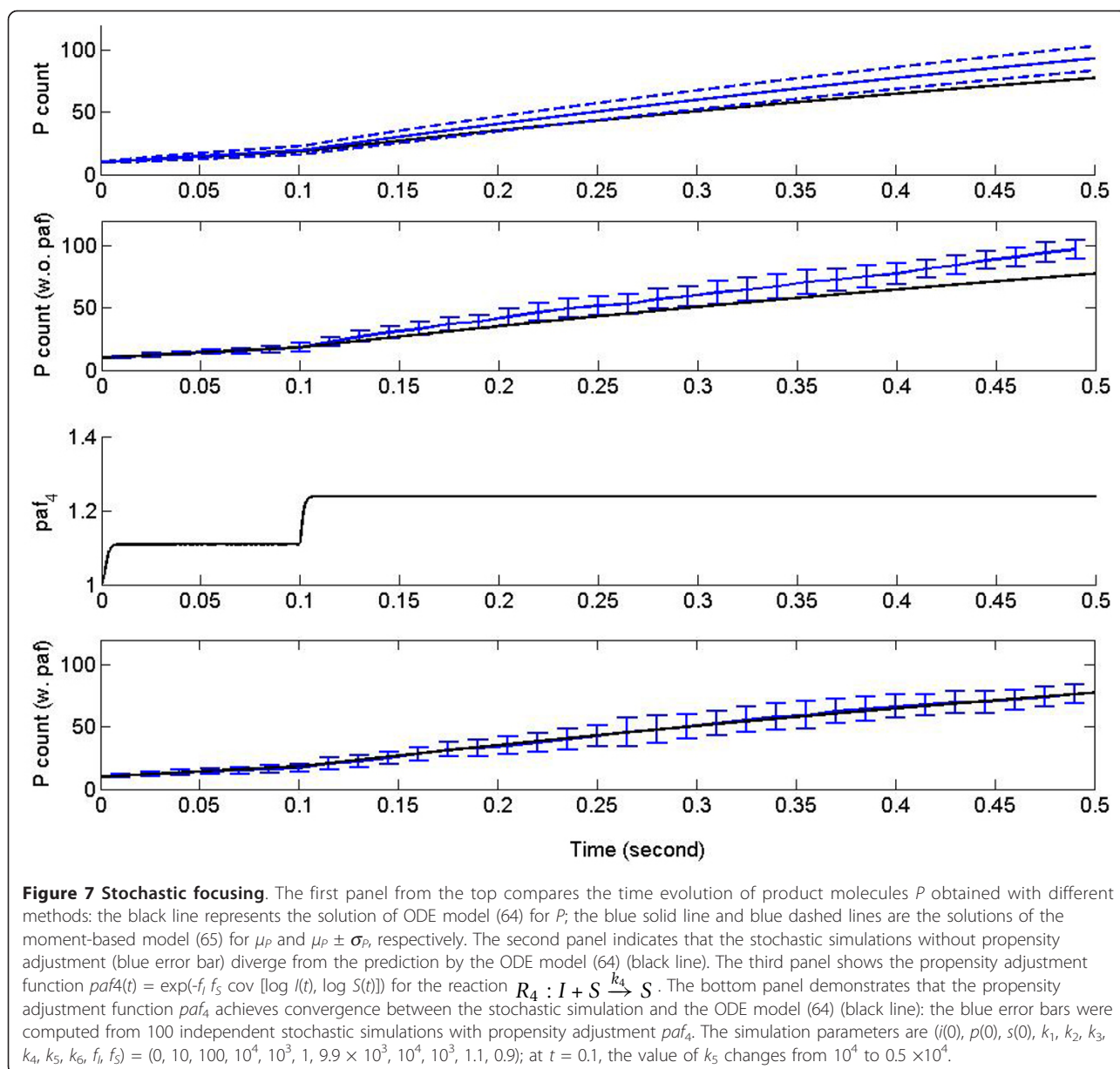
and the diagonal matrix  $\Lambda$  is defined by

$$(\Lambda)_{rr} = \alpha_r + \frac{1}{2} \sum_{m,n=i,p,s} \frac{\partial^2 \alpha_r}{\partial m \partial n} \sigma_{mn}.$$

The stochastic focusing model without propensity adjustment yields results quite different from those of the deterministic model, as is illustrated in Figure 7. In this figure, the blue lines in the 1<sup>st</sup> panel are predicted from the moment equations (65) and the blue error bars for  $\mu_P$  in the 2<sup>nd</sup> panel are obtained from ten independent stochastic simulations. Both diverge systematically from the black line predicted by ODE model (64). By contrast, the stochastic model with propensity adjustment produces results consistent with the deterministic model, as shown by the 4<sup>th</sup> panel.

### Discussion

It is often implicitly assumed that the rate of a dynamic process can be directly taken as the propensity for a corresponding stochastic process. We have shown here that this is sometimes, but not always, true. Our results fall into three categories. The first develops conditions for a valid conversion of a rate to a propensity, the second presents a general conversion procedure, and the third discusses computational issues of propensity adjustment.



### Conditions for the direct use of a rate constant (function) as propensity function

We have shown that the direct use of a rate constant or a rate function  $f$  as the propensity function in a stochastic simulation algorithm requires that at least one of the following assumptions be true:

- 1)  $f$  is a linear function; this assumption has been validated in the Results sections addressing 0<sup>th</sup>-order and 1<sup>st</sup>-order reaction kinetics.
- 2) the reaction is monomolecular; this assumption was evaluated in the Results section describing real-valued order monomolecular reaction kinetics.

3) all  $X_i$  in the system are noise-free variables, *i.e.*, without (or with ignorable) fluctuations; this assumption implies that the covariance of any two participating reactants is zero (or close to zero). This assumption is assessed in equations (29 - 36).

Each of these three conditions is a sufficient condition for the direct use of a rate function  $f$  as the propensity function. Moreover, these statements are valid for functions of a general format, not just for GMA. This is so because the functional formats in cases 1 and 2 above are special cases of the GMA format. For the third case, a formal proof is only given for functions in GMA format, because this structured format allows us to give an

explicit estimation on how the covariance can affect the average behavior of a stochastic simulation through equation (34). For functions not in GMA format, the conclusion is still holds, although an analogous explicit estimation is lacking. The argument is as follows. The bimolecular reaction  $E[\alpha_r(\mathbf{X}(t))]$  contains at least one quadratic moment of the form  $E[X_i(t)X_j(t)]$  (cf. [4] and page 38). Therefore, by definition of the covariance,  $E[X_i(t)X_j(t)] = E[X_i(t)]E[X_j(t)] + \text{cov}(X_i(t), X_j(t))$ , we obtain

$$E[X_i(t)X_j(t)] = E[X_i(t)]E[X_j(t)] \Leftrightarrow \text{cov}(X_i(t), X_j(t)) = 0.$$

This result implies the following: If the covariance between every pair of random variables is zero (or ignorable), we have  $E[X_i(t)X_j(t)] = E[X_i(t)]E[X_j(t)]$  and therefore  $E[\alpha_r(\mathbf{X}(t))] = \alpha_r(E[\mathbf{X}(t)])$ . Expressed in words, the expectation of the propensity function on left-hand side of equation (29) equals its rate function, and the rate function can be directly used as propensity function in stochastic simulations.

If at least one of the three assumptions is satisfied, the stochastic simulation algorithm (SSA) is applicable without changes.

#### A general procedure for converting an equation-based model into a stochastic analogue

In the past, efforts have been made to manipulate the chemical master equation (CME) in order to achieve a proper propensity function for a reduced system (e.g., see [2]). However, manipulations of CME are usually complicated, and successes have been modest and rare. Here we propose an alternative strategy for converting a reduced dynamical model into a stochastic analogue. To achieve this conversion, we addressed two fundamental issues: First, under what conditions can a deterministic, equation-based model be validly used in stochastic simulations? And second, what is a proper strategy to implement such a conversion?

To address the first question, we showed that the following steps are necessary:

- (1) A concentration-based model needs to be converted into a particle-based model by accounting for the size of the system; if the concentration-based model is scaled (as was illustrated with the repressilator example), it may first have to be un-scaled in order to render the conversion valid;
- (2) The difference between the mean of a stochastic model without propensity adjustment and the corresponding quantities of the equation-based model should be evaluated. The mean of the stochastic model is obtained either through stochastic simulations or through a moment-based approach. If the difference is significant, then an adjustment of the

propensity function for a non-elementary reaction is necessary.

To answer the second question, we need to execute the following steps

- (3) Compute a propensity adjustment function, either through simulated or experimental data or through a moment-based approach, in order to achieve the corrected propensity function (41);
- (4) Apply SSA or one of its variants using a propensity function with adjustment to obtain valid simulation trajectories.

#### Computational issues of propensity adjustments

When the propensity needs adjusting, an accurate propensity adjustment function (*paf*) is essential for obtaining the proper correction of the propensity. It is usually impossible to compute *paf* exactly, which necessitates a suitable approximation. The approximation error in *paf* originates from the following sources:

- 1) The expression of *paf* in Equation (40) is a function of the mean, variance, and covariance, which are computed with a 2<sup>nd</sup>-order Taylor expansion in log space.
- 2) The moment-based approach, from which the functions of mean, variance and covariance are usually derived, is an approximation method that yields a closed ODE system for the moments. In the method used here, the propensity function is approximated by a 2<sup>nd</sup>-order Taylor expansion, and the moments up to a certain degree (2 in our treatment) are retained, while all higher moments are assumed to be zero. One might expect that a higher-order Taylor expansion would improve the accuracy of *paf*, but it would come with a much higher computational cost. The error control of *paf* and the relative computational issues should be addressed in future studies.

Since computation cost is a major concern with the stochastic simulation of large biochemical reaction networks, another issue has yet to be addressed. Namely, how does the propensity function of a reduced system affect the accuracy and efficiency of various leaping methods that have been proposed to speed up SSA? Moreover, the question of molecular population sizes requires further analysis. Our derivation assumed large reactant populations, but simulations of a reversible pathway indicated that the method works rather well even for small populations. A more careful investigation of this issue of population size in different scenarios is

still needed and should be the subject of further research.

## Conclusions

Gillespie's stochastic simulation algorithm (SSA), as well as later variants, permits three kinds of elementary reactions to be modelled: 0<sup>th</sup>, 1<sup>st</sup> and 2<sup>nd</sup> order reactions that are assumed to follow the law of mass action. All other types of reactions, containing non-integer kinetic orders and/or following other types of kinetic law, are assumed to be convertible to one of these three kinds, so that SSA can validly be applied. However, the conversion to elementary reactions is often difficult, infeasible, or simply impossible. First, the kinetic parameters of the underlying elementary reactions are in many cases unknown for a complex-order reaction. Second, even when all elementary kinetic parameters are available, the multitude of reaction channels and participating species creates a combinatorial complexity that renders SSA simulations computationally impractical. Within a deterministic framework, model reduction is a possible and often-used strategy to address such challenges. For example, a reduced mechanistic model, such as the Michaelis-Menten rate law, is often proposed to fit the experimental data, at the cost of sacrificing the original mechanistic interpretation. The reduction in these cases simplifies the original formulation by approximating, merging, or omitting intermediate reaction steps and reactants.

In this article, we propose a rather general strategy for converting a deterministic process model into a corresponding stochastic model and characterize the mathematical connections between the two. The deterministic framework is assumed to be a generalized mass action system and the stochastic analogue is in the format of the chemical master equation. The analysis identifies situations: where a direct conversion is valid; where internal noise affecting the system needs to be taken into account; and where the propensity function must be mathematically adjusted. The conversion from deterministic to stochastic models is illustrated with several representative examples, including reversible reactions with feedback controls, Michaelis-Menten enzyme kinetics, a genetic regulatory motif, and stochastic focusing. The construction of a stochastic model for a biochemical network requires the utilization of information associated with an equation-based model. The conversion strategy proposed here guides a model design process that ensures a valid transition between deterministic and stochastic models.

## Additional material

**Additional file 1: Derivation of the mean and variance of a power-law function of random variables.**

**Additional file 2: Computation of approximate mean and covariance for a generic propensity function to be used in stochastic simulations.**

## Acknowledgements

The authors thank Dr. Yi Jiang for useful comments and for providing seminal references. The authors also appreciate Dr. Mukhtar Ullah's and Dr. Olaf Wolkenhauer's insightful comments and conceptual clarifications. This work was supported in part by a Molecular and Cellular Biosciences Grant (MCB-0946595; E.O. Voit, PI) from the National Science Foundation. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the sponsoring institutions.

## Author details

<sup>1</sup>Department of Mathematics, Bioinformatics Program, Georgia Institute of Technology, Atlanta, GA30332, USA. <sup>2</sup>The Wallace H. Coulter Department of Biomedical Engineering, Georgia Institute of Technology, Atlanta, GA30332, USA. <sup>3</sup>Integrative BioSystems Institute and The Wallace H. Coulter Department of Biomedical Engineering, Georgia Institute of Technology, Atlanta, GA30332, USA.

## Authors' contributions

JW developed the mathematical derivations, designed and performed the simulation, and drafted the manuscript. BV contributed to the statistical reasoning and revised the manuscript. EV supervised the research and revised the manuscript. All authors read and approved the final manuscript.

Received: 18 July 2011 Accepted: 8 November 2011

Published: 8 November 2011

## References

1. Gillespie D: Exact Stochastic Simulation of Coupled Chemical Reactions. *J Phys Chem* 1977, **81**(25):2340-2361.
2. Rao CV, Arkin AP: Stochastic chemical kinetics and the quasi-steady-state assumption: Application to the Gillespie algorithm. *The Journal of chemical physics* 2003, **118**(11):4999-5010.
3. Cao Y, Gillespie DT, Petzold LR: Multiscale stochastic simulation algorithm with stochastic partial equilibrium assumption for chemically reacting systems. *J Comput Phys* 2005, **206**:395.
4. Gillespie DT: Stochastic simulation of chemical kinetics. *Annual Review of Physical Chemistry* 2007, **58**:35-55.
5. Tian T, Burrage K: Stochastic models for regulatory networks of the genetic toggle switch. *Proc Natl Acad Sci USA* 2006, **103**(22):8372-7.
6. Gomez-Urbe CA, Verghese GC: Mass fluctuation kinetics: Capturing stochastic effects in systems of chemical reactions through coupled mean-variance computations. *The Journal of chemical physics* 2007, **126**(2):024109-12.
7. Lee CH, Kim K-H, Kim P: A moment closure method for stochastic reaction networks. *The Journal of chemical physics* 2009, **130**(13):134107-15.
8. Singh A, Hespanha J: LogNormal Moment Closures for Biochemical Reactions. In *Proc of the 45th Conf on Decision and Contr* 2006.
9. Milner P, Gillespie CS, Wilkinson DJ: Moment closure approximations for stochastic kinetic models with rational rate laws. *Mathematical Biosciences* 2011, **231**(2):99-104.
10. Chevalier MW, El-Samad H: A rigorous framework for multiscale simulation of stochastic cellular networks. *The Journal of chemical physics* 2009, **131**(5):054102-17.
11. Voit EO: In *Computational analysis of biochemical systems: a practical guide for biochemists and molecular biologists. Volume xii*. Cambridge University Press; 2000.
12. Savageau MA: Biochemical systems analysis. I. Some mathematical properties of the rate law for the component enzymatic reactions. *Journal of Theoretical Biology* 1969, **25**(3):365-9.

13. Savageau MA: In *Biochemical systems analysis. A study of function and design in molecular biology. Volume xvii*. Addison-Wesley; 1976.
14. Savageau M: **Michaelis-Menten mechanism reconsidered: implications of fractal kinetics.** *Journal of Theoretical Biology* 1995, **176**(1):115-124.
15. Savageau MA: **Influence of fractal kinetics on molecular recognition.** *Journal of Molecular Recognition* 1993, **6**(4):149-157.
16. Bajzer Z, Huzak M, Neff KL, Prendergast FG: **Mathematical analysis of models for reaction kinetics in intracellular environments.** *Mathematical Biosciences* 2008, **215**(1):35-47.
17. Neff KL: *Biochemical reaction kinetics in dilute and crowded solutions: Predictions of macroscopic and mesoscopic models and experimental observations* Mayo Clinic: Rochester, MN; 2010.
18. Neff , Kevin L, Offord Chetan P, Caride Ariel J, Strehler Emanuel E, Prendergast Franklyn G, Bajzer eljko: **Validation of Fractal-Like Kinetic Models by Time-Resolved Binding Kinetics of Dansylamide and Carbonic Anhydrase in Crowded Media.** *Biophysical journal* 2011, **100**(10):2495-2503.
19. Chou I-C, Voit EO: **Recent developments in parameter estimation and structure identification of biochemical and genomic systems.** *Math Biosc* 2009, **219**:57-83.
20. Walton RJ, Preston CJ, Bartlett M, Smith R, Russell RGG: **Biochemical measurements in Paget's disease of bone.** *European Journal of Clinical Investigation* 1977, **7**(1):37-39.
21. Koch AL: **The logarithm in biology 1. Mechanisms generating the log-normal distribution exactly.** *Journal of Theoretical Biology* 1966, **12**(2):276-290.
22. Limpert E, Stahel WA, Abbt M: **Log-normal Distributions across the Sciences: Keys and Clues.** *BioScience* 2001, **51**(5):341.
23. Law AM, Kelton WD: *Simulation Modeling and Analysis*. 3 edition. Boston: McGraw Hill; 2000.
24. Gillespie D, Petzold L: **Numerical Simulation for Biochemical Kinetics.** In *In System Modelling in Cellular Biology*. Edited by: Szallasi Z, Stelling J, Perival V. MIT Press; 2006.
25. Wolkenhauer O, Ullah M, Kolch W, Cho K: **Modelling and Simulation of IntraCellular Dynamics: Choosing an Appropriate Framework.** *IEEE Transactions on NanoBioscience* 2004, **3**:200-207.
26. Paulsson J, Berg OG, Ehrenberg M: **Stochastic focusing: Fluctuation-enhanced sensitivity of intracellular regulation.** *Proceedings of the National Academy of Sciences* 2000, **97**(13):7148-7153.
27. Elowitz MB, Leibler S: **A synthetic oscillatory network of transcriptional regulators.** *Nature* 2000, **403**(6767):335-338.
28. Bennett MR, Volfson D, Tsimring L, Hasty J: **Transient Dynamics of Genetic Regulatory Networks.** *Biophysical journal* 2007, **92**(10):3501-3512.
29. Segel LA: **On the validity of the steady state assumption of enzyme kinetics.** *Bull Math Biol* 1988, **50**:579-593.
30. Michaelis L, Menten ML: **Die Kinetik der Invertinwirkung.** *Biochem Zeitschrift* 1913, **49**:333-369.
31. Savageau MA, Voit EO: **Recasting Nonlinear Differential-Equations As S-Systems - A Canonical Nonlinear Form.** *Mathematical Bioscience* 1987, **87**:83-115.
32. Twomey A: *On the Stochastic Modelling of Reaction-Diffusion Processes* University of Oxford; 2007.

doi:10.1186/1752-0509-5-187

**Cite this article as:** Wu et al.: Constructing stochastic models from deterministic process equations by propensity adjustment. *BMC Systems Biology* 2011 5:187.

**Submit your next manuscript to BioMed Central and take full advantage of:**

- Convenient online submission
- Thorough peer review
- No space constraints or color figure charges
- Immediate publication on acceptance
- Inclusion in PubMed, CAS, Scopus and Google Scholar
- Research which is freely available for redistribution

Submit your manuscript at  
[www.biomedcentral.com/submit](http://www.biomedcentral.com/submit)

