

LA-UR-08-7864

Approved for public release;
distribution is unlimited.

Title: Coarse-Graining Stochastic Biochemical Networks:
Adiabaticity and Fast Simulations

Author(s): Ilya Nemenman, CCS-3, Z#202598
Nikolai Sinitsyn, CCS-3, Z#218021
Nick Hengartner, CCS-3, Z#186926

Intended for: Submission to Proceedings of the National Academy of
Sciences, USA



Los Alamos National Laboratory, an affirmative action/equal opportunity employer, is operated by the Los Alamos National Security, LLC for the National Nuclear Security Administration of the U.S. Department of Energy under contract DE-AC52-06NA25396. By acceptance of this article, the publisher recognizes that the U.S. Government retains a nonexclusive, royalty-free license to publish or reproduce the published form of this contribution, or to allow others to do so, for U.S. Government purposes. Los Alamos National Laboratory requests that the publisher identify this article as work performed under the auspices of the U.S. Department of Energy. Los Alamos National Laboratory strongly supports academic freedom and a researcher's right to publish; as an institution, however, the Laboratory does not endorse the viewpoint of a publication or guarantee its technical correctness.

Coarse-graining stochastic biochemical networks: adiabaticity and fast simulations

N.A. Sinitsyn^{*†}, Nicolas Hengartner[†], and Ilya Nemenman^{*†}

^{*}Center for Nonlinear Studies, and [†]Computer, Computational and Statistical Sciences Division, Los Alamos National Laboratory, Los Alamos, NM 87545 USA

Submitted to Proceedings of the National Academy of Sciences of the United States of America

We propose a universal approach for analysis and fast simulations of stiff stochastic biochemical kinetics networks, which rests on elimination of fast chemical species without a loss of information about mesoscopic, non-Poissonian fluctuations of the slow ones. Our approach, which is similar to the Born-Oppenheimer approximation in quantum mechanics, follows from the stochastic path integral representation of the cumulant generating function of reaction events. In applications with a small number of chemical reactions, it produces analytical expressions for cumulants of chemical fluxes between the slow variables. This allows for a low-dimensional, interpretable representation and can be used for coarse-grained numerical simulation schemes with a small computational complexity and yet high accuracy. As an example, we derive the coarse-grained description for a chain of biochemical reactions, and show that the coarse-grained and the microscopic simulations are in an agreement, but the coarse-grained simulations are three orders of magnitude faster.

stochastic processes | Monte-Carlo | Michaelis-Menten | Langevin

Abbreviations: MM, Michaelis-Menten; CGF, cumulant generating function; MGF, moment generating function; SPI, stochastic path integral

Introduction

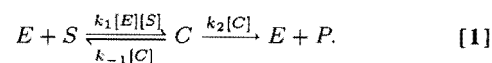
Computer simulations are often the method of choice to explore an agreement between a model and the observed experimental data in systems biology, especially in the context of single-molecule experiments [1, 2, 3]. Unfortunately, even the simplest biochemical simulations often face serious problems, both conceptual and practical. First, the networks usually involve *combinatorially many* chemical species and reaction processes: for example, a single molecule with n modification sites can exist in 2^n states, with an even larger number of reactions connecting them [4]. Second, while it is widely known that *some* molecules occur in the cell at very low copy numbers (e.g., a single copy of the DNA), which give rise to important stochastic effects [5, 6, 7, 8], it is less appreciated that the combinatorial complexity makes this true for *many* molecular species. Indeed, even for a large total number of molecules, typical abundances of microscopic species may be small if the number of the species is combinatorially large. Third, and perhaps the most profound difficulty of the simulations approach, is that only very few of the kinetic parameters underlying the networks are experimentally observed or even observable.

While some day computers may be able to tackle the formidable problem of modeling astronomically complex biochemical processes as a series of random reaction events, and then performing sweeps through parameter spaces in search of an agreement with experiments, these days are still far away. More importantly, even if the computing power were available, it would not help in building a comprehensible, tractable interpretation of the modeled system and in identifying connections between its microscopic and macroscopic features.

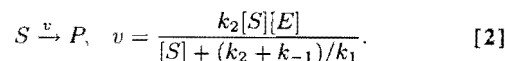
Clearly, such an interpretation can be aided by coarse-graining, that is, by merging or eliminating certain nodes and/or reaction processes (this would be called blocking or decimation in statistical physics, and is in the spirit of the real-space renormalization group). Ideally, one wants to substitute multiple elementary (that is, single-step, Poisson-distributed) biochemical reactions with a few complex processes linking the species that survive the coarse-graining in a way

that retains predictability of the system. Not incidentally, this would help with each of the three roadblocks mentioned above by reducing the number of interacting elements, increasing the copy numbers of agglomerated hyperspecies, and combining multiple microscopic rates into a smaller number of effective parameters.

The importance of coarse-graining in biochemistry is long understood [9], and the prime example is the MM kinetic scheme:



Here k_1 , k_2 , and k_{-1} are kinetic rates, S , P , E , and C denote the substrate, the product, the enzyme, and the enzyme-substrate complex molecules, respectively, and $[\dots]$ represent the abundances. The enzyme catalyzes the $S \rightarrow P$ transformation by merging with S to create an unstable complex C , which then dissociates either back into $E + S$ or forward into $E + P$, leaving E unmodified. If $[S] \gg [E]$, then the enzyme cycles many times before $[S]$ and $[P]$ change appreciably. Thus the enzymes equilibrate quickly, resulting in a coarse-grained, complex reaction with the decimated enzyme species:



However, this simple reduction is insufficient when stochasticity is important: each MM reaction consists of multiple elementary steps, thus the statistics of the number of the reactions per unit time is non-Poissonian in general. While some attempts have been made to extend deterministic coarse-graining to the stochastic domain [10, 11, 12], a systematic set of such tools for realistic biochemical networks has not been found yet. In this article, we make a step towards the goal.

We start by noting that, in addition to the three conceptual problems, a technical difficulty stands in the way of stochastic simulations in systems biology: molecular species have diverse dynamical time scales, making the systems stiff and difficult to simulate. We propose to use this property of multiple time scales to our advantage.

Many related approaches have been explored, differing largely by the definition of fast and slow variables. Commonly, *reaction rates* are used for this purpose [10]. However, if two species of different typical abundances are coupled by one reaction, then a relatively small change in the concentration of the high abundance species can have a dramatic effect on that of the low abundance one. This notion of *species* rather than *reaction* based adiabaticity is at the heart of the original MM derivation, as well as of our arguments.

Our method builds upon the SPI technique from mesoscopic physics [13, 14, 15] and provides three major improvements that make

Reserved for Publication Footnotes

the approach more applicable to biology. First, we extend the method, initially developed for large copy numbers, to discrete degrees of freedom, such as a single enzyme. Second, we explain how to use the SPI for a network of multiple reactions, reducing the entire network to a few complex reaction links. Finally, we show how the procedure can be turned into an efficient algorithm for coarse-grained simulations, preserving statistical characteristics of the original dynamics. The algorithm is akin to the widely used Langevin [16] or τ -leaping [17] schemes, but it simulates complex reactions in a single step. We believe that this development of a fast, yet precise numerical algorithm is the most important practical contribution of our work.

For pedagogical reasons, we develop the method using a model system that is simple enough for a detailed analysis, yet is complex enough to support our goals, and we provide a generalization later.

The model. Consider an enzyme attached to a cell membrane, Fig. 1. S_B substrate molecules are distributed over the bulk cell volume. Each molecule can either be adsorbed by the membrane, forming the species S_M , or dissociate from it, and the enzyme can interact only with the membrane-bound substrates. The enzyme-substrate complex C can split either into $E + S_M$ or into $E + P$. Let's suppose that the latter reaction is observable; for example, a GFP tag sparks each time a product molecule is created [3]. Finally, we assume that $C \rightarrow E + P$ is irreversible (e.g., the product leaves the membrane). This as a simple model of receptor signaling, such as in vision or immune system, or a model of a reaction-diffusion MM enzyme, where the membrane/bulk play the roles of the nearby/far away regions around the enzyme, and diffusion takes the molecules between them.

The full set of elementary reactions is

1. adsorption of the bulk substrate, $S_B \rightarrow S_M$ (rate $q_0 S_B$);
2. reemission of the substrate into the bulk, $S_M \rightarrow S_B$ (rate $q S_M$);
3. Michaelis-Menten conversion of S_M into P , consisting of
 - (a) complex formation, $S_M + E \rightarrow C$, (rate $k_1 S_M$);
 - (b) complex backward decay, $C \rightarrow S_M + E$ (rate k_{-1});
 - (c) product emission $C \rightarrow E + P$ (rate k_2).

Note that here and in the rest of the article we don't make a distinction between a species name and the number of its molecules.

In this setup, only emission of the product is directly observable. Our goal is to coarse-grain the above system of five reaction processes into a single complex reaction $S_B \rightarrow P$, as in Fig. 2(c). That is, we want to eliminate all intermediate species and processes, while preserving their effects on the statistics of the complex reaction $S_B \rightarrow P$ on time scales appropriate for its dynamics.

We stress again that this model is used for concreteness only, and the methodology we develop can be applied for other systems as well.

Results

There are three effective time scales in our model. One is the scale τ_B of the variation of the bulk substrate abundance. We assume that $S_B \gg S_M$. Therefore, this scale is the slowest, and we will be interested in studying the response of the system to changes in S_B on this scale. A faster time scale, τ_M is given by the dynamics of S_M . Finally, the fastest scale, τ_E , is set by single reaction events, that is, the characteristic time between enzyme-substrate binding/unbinding. Overall, $\tau_E \ll \tau_M \ll \tau_B$. We emphasize that all species in the problem are connected by reactions that happen with similar rates, and the separation of the time scales is a result of the particle abundances, rather than the reaction speeds: it takes longer to change a high-abundance species drastically compared to a low-abundance one. We expect this to be a generic property of many chemical networks.

The hierarchy of times allows us to coarse-grain the system in two steps, as in Fig. 2. First, we remove the variable with the fastest dynamics: the binary substrate-enzyme complex C . This replaces the

three steps of the MM mechanism with a single reaction $S_M \rightarrow P$, Fig. 2(b). Additionally, we represent the other reactions in the system in a more convenient form. In the second step, we eliminate S_M , which evolves on the scale τ_M . This results in the characterization of the average $S_B \rightarrow P$ flux and its fluctuations, treating S_B as a time-dependent input parameter, cf. Fig. 2(c).

Preliminaries. Let δQ_μ stand for the number of reaction events for the reaction type μ ($\mu = 1, 2, 3$ corresponds to adsorption, detachment, and the MM reaction, respectively). Then $P(\delta Q_\mu | T)$ is the probability distribution of the number of events of type μ during a time window of length T . Instead of considering these distributions directly, we rely on the corresponding MGFs¹:

$$\mathcal{Z}_\mu(\chi, T) = e^{\mathcal{S}_\mu(\chi, T)} = \sum_{\delta Q_\mu=0}^{\infty} P(\delta Q_\mu | T) e^{i \delta Q_\mu \chi}. \quad [3]$$

where $\mathcal{S}_\mu(\chi)$ is the CGF. Then the cumulants of $P(\delta Q_\mu | T)$ are

$$c_{\mu, \alpha} = (-i)^{\alpha} \left. \frac{\partial^{\alpha}}{\partial \chi^{\alpha}} \right|_{\chi=0} \mathcal{S}_\mu(\chi, T), \quad [4]$$

where α stands for the cumulant order. In particular, the average flux for the reaction is $c_{\mu, 1}$, and the corresponding variance is $c_{\mu, 2}$.

Step 1: The generating function representation. This step can be viewed as a generalization of the τ -leaping approximation [17], which simulates elementary reactions, for example attachments/detachments in Fig. 2(a), by choosing a time step δt , over which the number of the reactions is much larger than one, yet the slowly varying reaction rates can be considered stationary. In τ -leaping, one then approximates $P(\delta Q_\mu | \delta t)$ as Poissons. Similarly, in our case, for $\tau_E \ll \delta t \ll \tau_M$, we can approximate CGFs of membrane attachment/detachment as those of Poisson processes, $\mathcal{S}_\mu(\chi) = r_\mu(t)(e^{i\chi} - 1)\delta t$, and the rates are $r_1 = q_0 S_B(t)$ and $r_2 = q S_M(t)$, respectively.

Unfortunately, not all biochemical processes can be treated in this simple manner. For example, due to the single-copy nature of the MM enzyme in Fig. 1, the instantaneous rate of the product creation is a fast varying function of time, switching between zero and k_2 every time the complex forms. Therefore, one cannot treat the product creation, $P(\delta Q_3 | \delta t)$, as a homogeneous Poisson process and use τ -leaping or Langevin methods [16, 17]. Still, we would like to avoid resorting to the Gillespie [18] or similar techniques since they are based on Monte-Carlo simulations of individual reaction events and are slow.

As an alternative, we derive an approximation for the non-Poisson distribution of δQ_3 by characterizing its CGF, \mathcal{S}_3 . To this end, we eliminate the binary substrate-enzyme complex C and reduce the MM reaction triplet to a single process, whose dynamics can be considered stationary over times much longer than a single reaction event. The details are in *Methods: Coarse-graining the Michaelis-Menten reaction*, Eq. (21), and the obtained expression is valid for times δt , $\tau_E \ll \delta t \ll \tau_M$, so that many enzyme turnovers happen, but the effect on the abundance of S_M is still relatively small.

This completes Step 1 of the coarse-graining in which each reaction, or a small complex of reactions, is subsumed by a quasi-stationary CGF \mathcal{S}_μ of the distribution of the number of its events.

Importantly, in this Step, we remove the only species that exists, at most, in a single copy, thus simplifying analysis of the system. Additionally, while we don't focus on this in what follows, in the MM mechanism, the backward reaction is often a simple dissociation, whereas the forward one requires crossing an energy barrier and is exponentially suppressed. As a result, often $k_{-1} \gg k_2$ requiring multiple binding events (and with them the instantaneous rate changes) for each released product. Thus the effect of replacing the

¹More precisely, \mathcal{Z} is the *characteristic function*, and the usual definition of the MGF is without i in the exponent. Same is true for our CGF, \mathcal{S} . We choose this nomenclature to emphasize that we use the functions for calculations of moments and cumulants, respectively.

entire MM mechanism with a single complex reaction on the simulation efficiency may be quite dramatic.

To illustrate the simplification, using Eq. (21), we write the first few cumulants of the number of MM product releases in time δt :

$$c_{3,1} = \frac{k_1 k_2 S_M}{K} \delta t, \quad K = k_1 S_M + k_2 + k_{-1}, \quad [5]$$

$$c_{3,2} = c_{3,1} F, \quad F = 1 - 2Q/K, \quad Q = c_{3,1}/\delta t, \quad [6]$$

$$c_{3,3} = c_{3,1} [1 - 6Q(K - 2Q)/K^2], \quad [7]$$

$$c_{3,4} = c_{3,1} [1 - 2Q(7K^2 - 36KQ + 60Q^2)/K^3]. \quad [8]$$

The coefficient F is called the Fano factor (see below). To the extent that $F \neq 1$, this complex reaction is non-Poisson (cf. Fig. 1 in *Supporting Information*).

Knowing cumulants of $P(\delta Q_3|\delta t)$ allows for a numerical simulation procedure

$$\delta Q_3(t) = \eta_3(t, \delta t), \quad [9]$$

$$S_M(t + \delta t) = S_M(t) - \delta Q_3(t) + J(t)\delta t, \quad [10]$$

$$P(t + \delta t) = P(t) + \delta Q_3(t), \quad [11]$$

where $\eta_3(t)$ is a random variable with the cumulants given by Eqs. (5-8), and $J(t)$ represents currents exogenous to the MM reaction, such as changes in S_M due to membrane binding/unbinding. Notice that we are now treating the reaction in a quasi-stationary, τ -leaping or Langevin-like way by drawing a (random) number of reaction events over a time δt directly, assuming that all parameters defining the reaction are constants over this time. The price for the coarse graining is that, instead of a single-rate Poisson distribution, one is forced to characterize this reaction by a prescribed sequence of cumulants.

In principle, generation of such random variables is an ill-posed task since the moments do not define the distribution uniquely. Additionally, once we allow for a nonzero third or fourth cumulant, the remaining higher order cumulants cannot be all zero, and the generated random variable will depend on the assumptions made about them. Fortunately, in our case, the situation is simplified because all $c_{3,k} \sim \delta t$. Thus higher cumulants have a progressively smaller effect, $\sim (\delta t)^{1/k}$, on a number drawn from the distribution, and our random variables are almost Gaussian. Then the Gram-Charlier series expansion [19] aided either by the importance or rejection sampling [20, 21] reduces the simulation scheme, Eqs. (9-11), to a simple Gaussian, Langevin simulation with a small penalty, as described in *Methods: Simulations with near-Gaussian distributions*; see Fig. 5 in *Methods* for illustration of the precision provided by these tools.

Step 2: Coarse-graining membrane reactions. In Step 2 of the coarse-graining, we start with the CGFs S_μ , $\mu = 1, 2, 3$, of the slowly varying reactions. Using the SPI technique, we then express the CGF of δQ , the number of the entire coarse-grained reactions $S_B \rightarrow P$ in Fig. 2(c) over time T , in terms of the component CGFs, and then simplify the expression to account for the time scale separation between τ_B and τ_M , see *Methods: Coarse-graining all membrane reactions*, Eq. (31). This formally completes the coarse-graining. That is, we find the CGF of the $S_B \rightarrow P$ particle flux for times $T \lesssim \tau_B$, much longer than τ_E and τ_M , the other time scales in the problem.

The full expression for CGF is cumbersome and non-illuminating. Fortunately, we only look for the first few cumulants of $P(\delta Q|T)$, and these are obtained by differentiating the CGF as in Eq. (4). The expressions for the first three cumulants, c_1 , c_2 , and c_3 are in *Supporting Information*. Then, similar to the MM reaction, we can simulate the whole five-reaction network in one Langevin-like step:

$$\delta Q(t) = \eta(t, T), \quad [12]$$

$$S_B(t + T) = S_B(t) - \delta Q(t) + J(t)T, \quad [13]$$

$$P(t + T) = P(t) + \delta Q(t), \quad [14]$$

where η is a random variable with the cumulants as in Eqs. (1-4) in *Supporting Information*, and $J(t)$ is an external current, such as production or decay of the bulk substrate in other cellular processes.

Fano factor in a single molecule experiment. In analyses of single molecule experiments, one often measures the ratio of the variance of $P(\delta Q|T)$ to its mean—the Fano factor [3], $F = c_2/c_1$. The factor is zero for deterministic systems and one for a Poisson process, providing a quantification of the importance of stochastic effects.

Traditionally, to compare experimental data about F to a mathematical model, one would simulate the model using the Gillespie algorithm [18], which takes a long time to converge to the necessary accuracy. In contrast, our coarse-grained quasi-stationary approach yields an analytic expression for the Fano factor of the $S_B \rightarrow P$ transformation, see Eq. (3) in *Supporting Information*. Similar analytical shortcuts should be possible for other kinetic schemes. In Fig. 3, we compare the analytic expression to stochastic simulations for the full set of reactions in Fig. 2(a), seeing an excellent agreement.

Note that $F \neq 1$, which indicates a non-Poissonian nature of the complex reaction. The backwards decay of C adds extra randomization, thus larger values of k_{-1} increase F . At the other extreme, when $k_{-1} = 0$, the Fano factor may be as small as $1/2$, so that the entire $S_B \rightarrow P$ chain is equal to a sequence of two Poisson events with similar rates. Finally, when $q = 0$, i.e., the substrates are removed from the membrane only via $S_M \rightarrow P$, $F = 1$. This is because then the only stochasticity in the problem is from Poisson membrane binding, and all bound substrates will eventually get converted to P .

Computational complexity of coarse-grained simulations. We expect the coarse-graining approach to be particularly useful for simulations in systems biology. This is due to an essential speedup provided by the method over the traditional Gillespie algorithm [18], by which all approaches are benchmarked. Indeed, for our model, the computational complexity of a single Gillespie simulation run is $O(MT/\tau_E)$, where $M = 5$ is the number of reactions in the system, and T is the duration of the simulated dynamics. In contrast, the complexity of the coarse-grained approach is $O[M^0(T/\tau_E)^0]$ since we have removed the internal species and simulate the dynamics in steps of $\sim T$, instead of $\sim \tau_E$. However, the Gillespie algorithm is (statistically) exact, while our analysis relies on quasi-stationary assumptions.

To gauge the practical utility of our approach in reducing the simulation time while retaining a high accuracy, we benchmarked it against the Gillespie algorithm. All simulations were performed using Fortran 90, on a single CPU AMD Barton 2500 (1.83 GHz, Windows 2000). In *Supporting Information* we provide the benchmark results for the single MM enzyme (reaction 3), where the coarse-graining approach achieves factor of 40 speedup. Here we focus on the full five-reaction model system viewed at different coarse-graining levels. We use $k_1 = 0.02$, $k_{-1} = 2$, $k_2 = 1$, $q = 0.01$, and $q_0 S_B = 1.5$.

Coarse-grained, Step 1: Total time of the evolution is $T = 1000$, and the initial number of the substrates on the membrane is $S_M(t = 0) = 120$. Then the relaxation time of a typical fluctuation of S_M is $\tau_M \sim 1/[q + (\partial k_{MM}/\partial S_M)] \sim 80$, where k_{MM} is the rate of the Michaelis-Menten reaction for a given S_M , and this sets the scale $\delta t = 20 \ll \tau_M \sim 80$. We simulate all three reactions that survive Step 1 of the coarse-graining (membrane binding/unbinding and MM transformation) by approximating their distributions with the Gram-Charlier series with three known cumulants, and we perform 10^6 simulation runs, which is sufficient for convergence of the third cumulant of the $S_B \rightarrow P$ aggregate reaction. As shown in Tbl. 1, the coarse-grained approach speeds simulations 60-fold relative to the Gillespie one with little apparent accuracy loss.

Coarse-grained, Step 2: We do similar benchmarking for the system represented as a single coarse-grained reaction $S_B \rightarrow P$. Here we use the time step $T = 1000$, $\tau_M \ll T \ll \tau_B$. The results in Tbl. 1 show that simulating all five reactions in a single step results in a dramatic speedup of about 4000. This number relates to the ratio of the slow and the fast time scales in the problem, but also to leaping over the futile bindings-unbindings in the coarse-grained scheme.

For all cumulants, coarse grained simulations and analytic results differ from exact Gillespie values by, at most, a per cent. It is hard to imagine a practical situation in modern biology where the kinetic parameters are known well enough so that such discrepancy matters. Yet the reduction of the simulation time by the factor of $10^3 \dots 10^4$ is certainly a tangible improvement.

Generalizations to a network of reactions. As discussed in detail in the original literature [14], in the SPI formalism, a network of M reactions with N chemical species (cf. Fig. 4) is generally described by $2MN$ ordinary differential equations specifying the saddle point solution of the corresponding path integral. *Methods: Coarse-graining all membrane reactions* provides a particular example, and we refer the readers to the original literature for generalizations. Here, we build on the Ref. [14] and focus on developing a relatively simple, yet general coarse-graining procedure for more complex reaction networks.

At intermediate time scales, δt , many fast species connecting slow ones can be considered statistically independent. Therefore, in the SPI, every separate chain of such species simply adds to the effective Hamiltonian. Namely, we enumerate slow chemical species by μ, ν, \dots . Fast chains connecting them can be marked by pairs of indexes, e.g., $\mu\nu$ (cf. Fig. 4). An entire such chain will contribute a single effective Hamiltonian term, $H_{\mu\nu}(\{N\}, \{\chi\}, \{\chi_C\})$, to the full CGF of the slow fluxes, where $\{N\}$, $\{\chi\}$, and $\{\chi_C\}$ are the slow species abundances and the conjugate counting variables. If necessary, the geometric correction to the CGF, $S_{\text{geom}}^{\mu\nu}(\{N\}, \{\chi\}, \{\chi_C\})$, can be written out as well [15]. Overall,

$$S(\{\chi_C\}, T) = \sum_{\mu < \nu} S_{\text{geom}}^{\mu\nu}(\{N(t)\}, \{\chi(t)\}, \{\chi_C\}, T) + \int_0^T dt \left[\sum_{\mu} i\chi_{\mu} \dot{N}_{\mu} + \sum_{\mu < \nu} H_{\mu\nu}(\{N(t)\}, \{\chi(t)\}, \{\chi_C\}) \right]. \quad [15]$$

This provides for the following coarse-graining procedure. First, one finds a time scale δt , small enough for the slow species to be considered stationary, and yet fast enough for the fast ones to equilibrate. If the fast species consist only of a few degrees of freedom, like in the case of a single enzyme, one derives the CGF of the transformations mediated by these species similar to *Methods: Coarse-graining the Michaelis-Menten reaction*. If instead the fast species are mesoscopic, one uses the SPI technique to derive the CGF by analogy with Step 2.

At the next step, the CGFs of the fast species are incorporated into the SPI over the abundances of the slow ones. For this, one writes down the full effective Hamiltonian, Eq. (15), assumes adiabatic evolution, and solves the ensuing saddle point equations. The extremum of the effective Hamiltonian determines the CGF of the coarse-grained process. For hierarchies of time scales, this reduction procedure is then repeated.

Discussion

Rigorous mathematical techniques developed in physics, chemistry, and engineering are finding applications in the biological domain. This article represents one such example, where the *adiabatic approach*, paired with the SPI formalism of statistical physics, allows to coarse-grain stochastic biochemical kinetics systems. For a system with a separation of time scales, we eliminate fast variables and reduce the network to a handful of slow species coupled by complex interactions with properties that account for the decimated nodes. The simplified system is smaller, non-stiff, and easier to analyze, resulting in orders of magnitude improvement in the computational complexity of its simulations. This has a potential for a wide impact on simulations in systems biology, at least for systems with diverse time scales.

Fortunately, such systems occur more often in Nature than one would expect naively. Consider, for example, the system briefly mentioned in the *Introduction*: a molecule must be modified on n sites

in an arbitrary order to get activated. The kinetic diagram for this system is an n -dimensional hypercube, and the number of states of the molecule with m modified sites is $\binom{n}{m}$. Therefore, if the total number of molecules is N , then a typical state with m modifications will have $N_m = N/\binom{n}{m}$ molecules in it. This number may be small, ensuring the need for a stochastic analysis. More importantly, it is quite different from either N_{m-1} or N_{m+1} , e.g., $N_m/N_{m+1} = (m+1)/(n-m)$, and, as we discussed at length, the different abundances result in different time scales.

The adiabatic SPI coarse-graining simplifies interpretation of biological systems. For example, the Fano factor of the $S_B \rightarrow P$ reaction, Fig. 3, may approach unity, suggesting a simple, yet rigorous, replacement of the entire reaction by a simple Poisson step. Then the list of relevant parameters becomes smaller than suggested by the *ab initio* description, improving interpretability and decreasing the effective number of biochemical features that must be measured experimentally. Recent analysis suggests that this may be a common property of biochemical networks [22, 23], and our methods may prove helpful in determining the relevant kinetic features.

While orders of magnitude improvement in simulation speed is certainly impressive, we are still far from coarse-graining cellular-scale reaction networks. However, the following properties of our approach suggest that we may be on the right track:

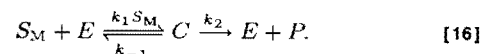
- We reduce a system of stochastic differential equations to a similar number of deterministic ones, which is a substantial simplification.
- We can operate with arbitrarily long series of cumulants of the number reaction events, keeping track of even rare fluctuations.
- Standard adiabatic approximations, well developed in classical and quantum physics, can be applied easily in the SPI context.
- Unlike some other coarse-graining techniques, the SPI approach can deal with copy numbers of order unity.
- Finally, the SPI is rigorous, mathematically justifiable, and allows for controlled approximations.

In the forthcoming publications, we expect to show how these advantageous properties of the adiabatic SPI technique allow to coarse-grain many standard biochemical network motifs.

Materials and Methods

Coarse-graining the Michaelis-Menten reaction.

Consider the $S_M \rightarrow P$ reaction, described mathematically as in Eq. (1):



The probabilities of transitions between bound, P_b , and unbound, $P_u = 1 - P_b$, states of the enzyme are given by a two state Markov process

$$\frac{d}{dt} \begin{bmatrix} P_u \\ P_b \end{bmatrix} = - \begin{bmatrix} k_1 S_M & -k_{-1} - k_2 \\ -k_1 S_M & k_{-1} + k_2 \end{bmatrix} \begin{bmatrix} P_u \\ P_b \end{bmatrix}. \quad [17]$$

Using Eq. (17) and the definition of Z_{μ} , Eq. (3), one can show that $Z_3(\chi, \delta t)$ satisfies a Schrödinger-like equation with a χ -dependent Hamiltonian, leading to a formal solution [2, 24, 11, 15]

$$Z_3(\chi, \delta t) = 1^+ \left(e^{-\hat{H}_{MM}(\chi, \delta t)} \right) p(t_0), \quad [18]$$

where $1^+ = (1, 1)$ is the unit vector, $p(t_0)$ is the probability vector of the initial enzyme states, and

$$\hat{H}_{MM}(\chi) = \begin{bmatrix} k_1 N_s & -k_{-1} - k_2 e^{i\chi} \\ -k_1 N_s & k_{-1} + k_2 \end{bmatrix}. \quad [19]$$

Similar Hamiltonians can be derived for a wide class of kinetic schemes [24, 11, 25, 15, 26], allowing for a straightforward extension of our methods.

The solution, Eq. (18), can be simplified if the MM reaction is considered in a quasi-steady state approximation, that is P_u is equilibrated at a current value

of the other parameters. This means that the time scale of interest, $\delta t \sim \tau_M$, is much larger than a characteristic time of a single enzyme turnover, τ_E , so we can consider $\delta t \rightarrow \infty$ in Eq. (18). Then only the eigenvalue $\lambda_0(\chi)$ of $\hat{H}_{MM}(\chi)$ with the smallest real part is relevant, and $Z_3(\chi, \delta t) = e^{-\lambda_0(\chi)\delta t}$.

It is possible to incorporate a slow time dependence of the parameters into this answer. By analogy with the quantum mechanical Berry phase [12, 15, 26], the lowest order non-adiabatic correction can be expressed as a geometric phase

$$Z_3(\chi) = e^{S_3(\chi)} = e^{\int_c A \cdot dk - \int dt \lambda_0(\chi, t)}, \quad [20]$$

where $A = \langle u_0(\chi) | \partial_k u_0(\chi) \rangle$, k is the vector in the parameter space that draws a contour c during the parameter evolution, and $\langle u_0(\chi) |$ and $|u_0(\chi)\rangle$ are the left and the right eigenvectors of $\hat{H}_{MM}(\chi, t)$ corresponding to the instantaneous eigenvalue $\lambda_0(\chi, t)$. The first term in Eq. (20) is the geometric phase, which is responsible for various ratchet-like fluxes [12, 27, 28].

The geometric phase gives rise to magnetic field-like corrections to the evolution of the slow variables. However, since these corrections depend on (small) time derivatives of these variables, they often are small and can be neglected, unless they break some important symmetry (such as the detailed balance [15, 28]), or the leading, non-geometric term is zero. In our model, the geometric effects are negligible when $\tau_E/\tau_M \sim 1/S_M \ll 1$, and we deemphasize them. However, we keep the geometric terms in several formal expressions for completeness, and the reader should be able to track their effects if desired.

Reading the value of $\lambda_0(\chi)$ from Ref. [12], we write the CGF of $P(\delta Q_3|\delta t)$, $\tau_E \ll \delta t \lesssim \tau_M$ in the adiabatic limit:

$$S_3(\chi, \delta t) = S_{\text{geom}}(\chi, S_M, \dot{S}_M) + \frac{\delta t}{2} [-(k_{-1} + k_2 + S_M k_1) + \sqrt{(k_{-1} + k_2 + S_M k_1)^2 + 4S_M k_1 k_2 (e^{\chi} - 1)}]. \quad [21]$$

Simulations with near-Gaussian distributions. A probability distribution $P(\delta Q)$ with known cumulants c_1, c_2, c_3, \dots , can be approximated as a limited Gram-Charlier expansion [19]

$$P(\delta Q) \approx \Psi(\delta Q, c_1, c_2) \left[1 + \frac{c_3(y^3 - y)}{6c_2^{3/2}} + \frac{c_4(y^4 - 6y^2 + 3)}{24c_2^2} + \frac{c_5(y^5 - 15y^3 + 45y - 15)}{72c_2^3} + \dots \right], \quad [22]$$

where $y = (\delta Q - c_1)/\sqrt{c_2}$ and $\Psi(\delta Q, c_1, c_2)$ is the Gaussian density with the mean c_1 and the variance c_2 . The leading term in Eq. (22) is a standard Gaussian approximation, and the subsequent terms account for skewness, kurtosis, etc. If all cumulants scale similarly (the near-Gaussian case), then the terms in the series become progressively smaller, ensuring rapid convergence.

Generation of random samples from the non-Gaussian Gram-Charlier series is still a difficult task. However, if, instead of the random numbers *per se*, the goal is to calculate the expectation of some function $f(\delta Q)$ over the distribution P , $\langle f(\delta Q) \rangle_P$, then the importance sampling [20] can be used. Specifically, we generate a Gaussian random number δQ from $\Psi(\delta Q, c_1, c_2)$ and define its importance factor according to its relative probability in the normal distribution and the considered Gram-Charlier series $\eta = P(\delta Q)/\Psi(\delta Q, c_1, c_2)$. After generating N such random numbers δQ_ν , $\nu = 1, \dots, N$, we get

$$\langle f(\delta Q) \rangle_P = \frac{\sum_{\nu=1}^N \eta_\nu f(\delta Q_\nu)}{\sum_{\nu=1}^N \eta_\nu}. \quad [23]$$

If a current random number draw represents just one reaction in a larger reaction network, then the overall importance factor of a Monte Carlo realization is a product of the factors for each of the random numbers drawn within it.

This reduces the complexity of simulations to that of a simple Gaussian, Langevin process with a small burden of (a) evaluating an algebraic expression for the Gram-Charlier series, and (b) keeping track of the importance factor. Yet this small computational investment allows to account for an arbitrary number of cumulants of the involved variables. To illustrate this, in Fig. 5, we compare

the Gram-Charlier, importance-sampling simulations of the MM reaction flux to the exact results in *Results: Step 1*. Introducing just the third and the fourth cumulant makes the simulations almost indistinguishable from the exact results.

Here we sound a note of caution: the Gram-Charlier series produces approximations that are not necessarily positive and hence are not, strictly speaking, probability distributions. However, the leading Gaussian term decreases so fast that this may not matter in practice. In fact, in our simulations, we simply rejected any random number that had a negative importance correction. However, this simplistic solution is inadequate for lengthy simulations, where the probability that one of random numbers in a long chain of events falls into a badly approximated region of the distribution approaches one. Then other means of generating random numbers, such as the well-known acceptance-rejection method [21] should be used. Since the true distributions of interest are near-Gaussian, a Gaussian with a slightly larger variance will be an envelope function for the Gram-Charlier approximation to the true distribution. Then the average random number acceptance probability will be similar to the ratio of the true and the envelope standard deviations, and it can be almost one. Then the rejection approach will require just a bit more than one normal and one uniform random numbers to generate a sample from the Gram-Charlier series. Importantly, in this case, the negativity of the series is not a problem since it will lead to an incorrect rejection of a single, highly improbable sample, rather than an entire sampling trajectory.

Coarse-graining all membrane reactions.

To perform the coarse-graining that connects Figs. 2(b) and 2(c), we look for the MGF of the total number of products Q_P produced over time $T \sim \tau_B$:

$$Z(\chi_C) = e^{S(\chi_C)} = \sum_{Q_P=0}^{\infty} P(Q_P|T) e^{iQ_P \chi_C}. \quad [24]$$

For this, we discretize the time into intervals t_k of duration δt , and we introduce random variables $\delta Q_\mu(t_k)$ ($\mu = 1, 2, 3$), which denote the numbers of each of the three different reactions in Fig. 2(b) (membrane binding, unbinding, and MM conversion) during each time interval. The probability distributions of $\delta Q_\mu(t_k)$ are given by inverse Fourier transforms of the corresponding MGFs:

$$P(\delta Q_\mu(t_k)) = \frac{1}{2\pi} \int d\chi_\mu(t_k) e^{-i\chi_\mu(t_k) \delta Q_\mu(t_k) + H_\mu(\chi_\mu(t_k), S_B(t_k)) \delta t}, \quad [25]$$

where the CGF are $S_\mu(\chi, S_B) = H_\mu(\chi, S_B) \delta t$.

Following [13, 14, 29, 15], and recalling that $Q_P = \sum_k \delta Q_3(t_k)$, we write the MGF of the total number of products created during time interval $(0, T)$ as the path integral over all possible trajectories of $\delta Q_\mu(t_k)$ and $S_M(t_k)$:

$$e^{S(\chi_C, T)} = \langle e^{i\chi_C Q_P} \rangle = \int DS_M(t_k) \prod_{k,\mu} \int D\delta Q_\mu(t_k) \times P[\delta Q_\mu(t_k)] e^{i\chi_C \sum_k \delta Q_3(t_k)} \times \delta[S_M(t_{k+1}) - S_M(t_k) - \delta Q_1(t_k) + \delta Q_2(t_k) + \delta Q_3(t_k)]. \quad [26]$$

The δ -function in Eq. (26) expresses the conservation law for the slowly changing number of substrate molecules S_M . We rewrite it as

$$\delta(\dots) = \frac{1}{2\pi} \int_{-\pi}^{+\pi} d\chi_M(t_k) \exp\{i\chi_M(t_k) \dots\}, \quad [27]$$

and we substitute the expression together with Eq. (25) into Eq. (26). Then the integration over $\delta Q_\mu(t_k)$ produces new δ -functions over χ_μ , which, in turn, are removed by integration over $\chi_\mu(t_k)$. This leads to an expression for the MGF:

$$e^{S(\chi_C, T)} = \int DS_M D\chi_M e^{\int_0^T dt [i\chi_M \dot{S}_M + H(S_M, \chi_M, \chi_C)]}, \quad [28]$$

$$H = H_1(-\chi_M, S_M, t) + H_2(\chi_M, S_M, t) + H_3(\chi_M + \chi_C, S_M, t) = q_0 S_B e^{-\chi_M} + S_M q e_{\chi_M} + \frac{1}{2} [-(k_{-1} + k_2 + S_M k_1) + \sqrt{(k_{-1} + k_2 + S_M k_1)^2 + 4S_M k_1 k_2 (e^{i\chi_M + \chi_C} - 1)}]. \quad [29]$$

where $e_{\pm\chi_M} = e^{\pm i\chi_M} - 1$. The original SPI work [13] assumed all component reactions to be Poisson. However, here H_3 is the CGF of the entire

complex, non-Poisson MM reaction, which we read as the coefficient in front of δt in Eq. (21). This ability to include subsystems with small number of degrees of freedom, such as the MM enzyme, opens doors to application of the method to a wide variety of coarse-graining problems.

Since $S_M \gg 1$, this path integral is dominated by the classical solution of the equations of motion (i. e., the saddle point), which, near the steady state, are

$$\dot{S}_M = \dot{\chi}_M = 0, \quad \frac{\partial H}{\partial \chi_M} = \frac{\partial H}{\partial S_M} = 0. \quad [30]$$

Let $\chi_{cl}(\chi_C)$ and $S_{M,cl}(\chi_C)$ solve Eq. (30). Then the cumulants generating function in the quasi-steady state approximation is

$$\mathcal{S}(\chi_C, T) = T H(S_{M,cl}(\chi_C), \chi_{cl}(\chi_C), \chi_C) \quad [31]$$

This completes the last step of the coarse-graining by deriving the CGF for the number of complex $S_B \rightarrow P$ transformation over long times.

ACKNOWLEDGMENTS. We thank F. Alexander, W. Hlavacek, B. Munsky, M. Wall for useful discussions and critical reading of the manuscript. This work was funded in part by DOE under Contract No. DE-AC52-06NA25396.

1. Orrit, M (2002) Photon statistics in single molecule experiments. *Single Mol.* 3:255.
2. Gopich, I, Szabo, A (2003) Statistics of transition in single molecule kinetics. *J. Chem. Phys.* 118:454.
3. English, B et al. (2006) Ever-fluctuating single enzyme molecules: Michaelis-Menten equation revisited. *Nat. Chem. Biol.* 2:87.
4. Hlavacek, W et al. (2006) Rules for modeling signal-transduction systems. *Sci. STKE* 344.
5. McAdams, H, Arkin, A (1997) Stochastic mechanisms in gene expression. *Proc. Natl. Acad. Sci. (USA)* 94:814.
6. Elowitz, M, Levine, A, Siggia, E, Swain, P (2002) Stochastic gene expression in a single cell. *Science* 297:1183.
7. Raser, J, O'Shea, E (2004) Noise in gene expression: Origins, consequences, and control. *Science* 304:1811.
8. Bialek, W, Setayeshgar, S (2005) Physical limits to biochemical signaling. *Proc. Natl. Acad. Sci. (USA)* 102:10040.
9. Michaelis, L, Menten, M (1913) The kinetics of invertase activity. *Biochem. Zeitschrift* 49:333 (in German).
10. Rao, C, Arkin, A (2003) Stochastic chemical kinetics and the quasi-steady-state assumption: application to the Gillespie algorithm. *J. Chem. Phys.* 118:4999.
11. Gopich, I, Szabo, A (2006) Theory of the statistics of kinetic transitions with application to single-molecule enzyme catalysis. *J. Chem. Phys.* 124:154712.
12. Sinitsyn, N, Nemenman, I (2007) Berry phase and pump effect in stochastic chemical kinetics. *EPL* 77:58001.
13. Pilgram, S et al. (2003) Stochastic path integral formulation of full counting statistics. *Phys. Rev. Lett.* 90:206801.
14. Jordan, A, Sukhorukov, E, Pilgram, S (2004) Fluctuation statistics in networks: A stochastic path integral approach. *J. Math. Phys.* 45:4386.
15. Sinitsyn, N, Nemenman, I (2007) Universal geometric theory of mesoscopic stochastic pumps and reversible ratchets. *Phys. Rev. Lett.* 99:220408.
16. Zinn-Justin, J (2002) *Quantum Field Theory and Critical Phenomena* (Oxford University Press, USA).
17. Gillespie, D (2001) Approximate accelerated stochastic simulation of chemically reacting systems. *J. Chem. Phys.* 115:1716–1733.
18. Gillespie, D (1977) Exact stochastic simulation of coupled chemical reactions. *J. Phys. Chem.* 81:2340.
19. Blinnikov, S, Moessner, R (1998) Expansions for nearly gaussian distributions. *Astron. Astrophys. Suppl. Ser.* 130:193.
20. Srinivasan, R (2002) *Importance sampling - Applications in communications and detection* (Springer-Verlag, Berlin).
21. von Neumann, J (1951) Various techniques used in connection with random digits. monte carlo methods. *Nat. Bureau Standards* 12:36.
22. Gutenkunst, R et al. (2007) Universally sloppy parameter sensitivities in systems biology. *PLoS Comput. Biol.* 3:e189.
23. Ziv, E, Nemenman, I, Wiggins, C (2007) Optimal information processing in small stochastic biochemical networks. *PLoS ONE* 2:e1077.
24. Bagrets, D, Nazarov, Y (2003) Full counting statistics of charge transfer in Coulomb blockade systems. *Phys. Rev. B* 67:085316.
25. Sukhorukov, E, Jordan, A (2007) Stochastic dynamics of a Josephson junction threshold detector. *Phys. Rev. Lett.* 98:136803.
26. Sinitsyn, N (2007) Reversible stochastic pump currents in interacting nanoscale conductors. *Phys. Rev. B* 76:153314.
27. Ohkubo, J (2008) The stochastic pump current and the non-adiabatic geometrical phase. *J. Stat. Mech.* p P02011.
28. Astumian, R (2007) Adiabatic operation of a molecular machine. *Proc. Natl. Acad. Sci. (USA)* 104:19715.
29. Jordan, A, Sukhorukov, E (2004) Transport statistics of bistable systems. *Phys. Rev. Lett.* 93:260604.

Table 1. Comparison of cumulants of the product flux for the full system calculated using the Gillespie simulations, the coarse-grained simulations at Step 1 and Step 2, and the analytical predictions; numbers in parentheses are the estimated errors in the last significant digits

cumulant	Gillespie	CG (step 1)	CG (step 2)	Analytics
c_1	418.7(1)	420.0(1)	418.9(1)	418.9
c_2/c_1	0.771(1)	0.764(2)	0.768(1)	0.767
c_3/c_1	0.50(3)	0.46(8)	0.48(3)	0.472
time	1h 14min	1min 17s	1s	N/A

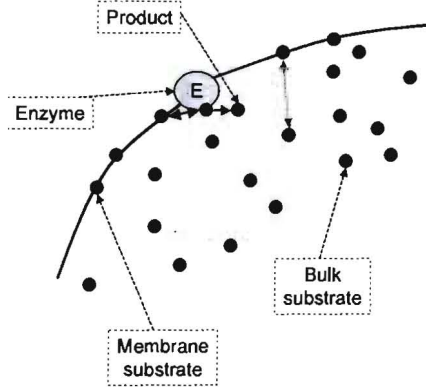


Fig. 1. The model system. Circles represent molecules and are labeled in the figure. Arrows stand for reactions: (1,2) adsorption and dissociation of S (orange); (3) multi-step MM conversion $S \rightarrow P$ (red).

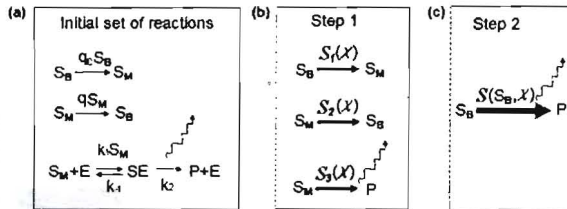


Fig. 2. Coarse-graining of the model system. Panel (a) shows the original set of reactions. Panel (b) represents the reactions after the first coarse-graining step: the MM mechanism has been replaced by a single complex reaction, and all the remaining reactions are now characterized by their slowly varying CGFs. Panel (c) shows the final reaction that describes the system at time scales $\delta t \gg \tau_M$. The wavy line corresponds to a spark of the tracer molecule [3], which counts the number of $S_B \rightarrow P$ transformations.

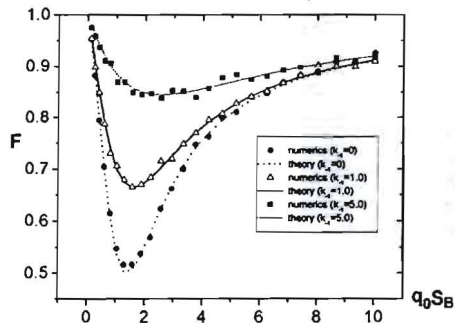


Fig. 3. Comparison of the analytically calculated Fano factor for the $S_B \rightarrow P$ reaction to Monte Carlo simulations with the Gillespie algorithm [18]. We use $q = 0.02$, $k_1 = 0.05$, $k_2 = 1$, and $T = 10000$. Each numerical data point averages 10000 simulation runs.

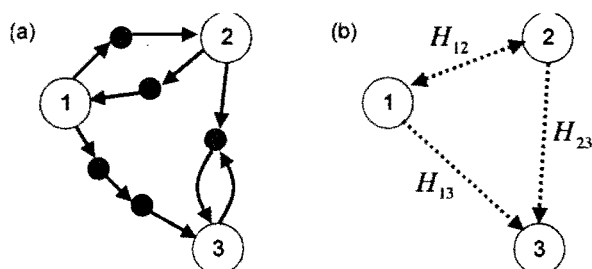


Fig. 4. Schematic coarse-graining of a network of reactions. (a) This network has $M = 10$ reactions (red arrows) and $N = 8$ species, of which three are slow (large circles), and five are fast (small circles). (b) Dynamics of each fast node can be integrated out, leaving effective, pairwise fluxes among the slow nodes (blue arrows), which are labeled by the corresponding effective Hamiltonians $H_{\mu\nu}$. Note that, for reversible pathways, the flux may be positive or negative (two-sided arrow), and it is strictly non-negative otherwise (one-sided arrows).

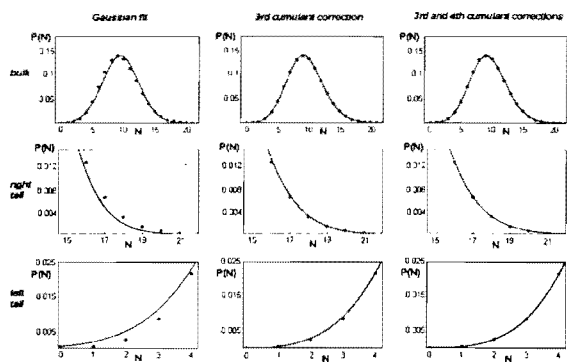


Fig. 5. Comparison of the exact discrete distribution of product molecules generated by the MM enzyme (points) with the continuous approximations by the Gram-Charlier series (lines). Left column compares the exact result to the Gaussian approximation. Central column shows improvements due to inclusion of the third cumulant correction. Including the fourth cumulant (right) makes the approximation and the exact result virtually indistinguishable. We used $S_M = 140 = \text{const}$, $k_1 = 0.02$, $k_{-1} = 2$, $k_2 = 1$, $q = 0.01$, and $\delta t = 35$.

Supporting Information: Coarse-graining stochastic biochemical networks: quasi-stationary approximation and fast simulations

N.A. Sinitzyn, Nicolas Hengartner, and Ilya Nemenman

September 17, 2008

1 Cumulants of the coarse-grained reaction

As described in the main text and *Methods*, the first three cumulants for the coarse-grained $S_B \rightarrow P$ reaction can be obtained by differentiating the corresponding CGF. This gives

$$c_1 = T \frac{1}{2k_1} \left[k_1(k_0 + k_2) + q(k_2 + k_{-1}) - \sqrt{k_1^2(k_0 - k_2)^2 + 2k_1q(k_0 + k_2)(k_2 + k_{-1}) + q^2(k_2 + k_{-1})^2} \right], \quad (1)$$

$$c_2 = F c_1, \quad (2)$$

$$F = 1 - \frac{q(2k_1k_0k_2 + k_1(k_0 + k_2)k_{-1} + qk_{-1}(k_2 + k_{-1}))}{k_1^2(k_0 - k_2)^2 + 2k_1q(k_0 + k_2)(k_2 + k_{-1}) + q^2(k_2 + k_{-1})^2} + \frac{qk_{-1}}{\sqrt{k_1^2(k_0 - k_2)^2 + 2k_1q(k_0 + k_2)(k_2 + k_{-1}) + q^2(k_2 + k_{-1})^2}}. \quad (3)$$

$$c_3 = -T \frac{\kappa}{\rho(-\kappa k_1 + \rho^2)^5} \left\{ \kappa^5 k_1^5 - \rho^{10} + \kappa \rho^7 [5k_1^2 k_2 + q(11k_1 + 6q)s] - \kappa^2 k_0^2 k_1^4 \rho^2 [5k_1^2 k_2^2 + 6k_2(k_1 - 2q)qs + 24q^2 s^2] + 2\kappa^2 k_0 k_1^2 \rho^3 [5k_1^3 k_2^2 + k_2 q(14k_1^2 - 9k_1 q - 6q^2)s + 6q^2(5k_1 + 3q)s^2] - 2\kappa k_0 k_1 \rho^4 [5k_1^4 k_2^3 + 19k_1^3 k_2^2 qs + 9k_1^2 k_2 q^2 s^2 + 6k_2 q^4 s^2 + 3k_1 q^3 s(-2k_2^2 + 8k_2 s + s^2)] \right\}, \quad (4)$$

where $s = k_1 \langle S_M \rangle + k_2 + k_{-1}$, $\langle S_M \rangle = \frac{1}{2k_1 q} \left\{ k_0 k_1 - k_1 k_2 - k_2 q - k_{-1} q + [4k_1 k_0 q(k_2 + k_{-1}) + (k_1 k_2 - k_1 k_0 + k_2 q + k_{-1} q)^2]^{1/2} \right\}$ is the average number of membrane-bound substrates, $k_0 = q_0 S_B$, $\kappa = k_0 k_1 k_2$, $\rho = k_1 k_2 + qs$, and, finally, T is the time step over which S_B changes by a relatively small amount, but many membrane reactions happen.

2 Simulating the Michaelis-Menten enzyme

We consider a MM enzyme with $S_M = 140 = \text{const}$, $k_1 = 0.01$, $k_{-1} = 2.0$, $k_2 = 1.0$. We analyze the number of product molecules produced by this enzyme over time $\delta t = 35$, with the enzyme initially in the (stochastic) steady state. To strain both Gillespie and our coarse-grained methods, we require a very high simulation accuracy, namely convergence of the fourth moment of the product flux distribution to two significant digits. For both methods, this means over 10 millions realizations of the same evolution.

In Tbl. 1 we report the results of our simulations. We see that the analytical coarse-grained results differ from the exact Gillespie simulations by, at most, two per cent, which is an expected deviation given the quality of the steady-state approximation. Further, the Langevin-like coarse-grained simulations, which accounted for the first four cumulants of the reaction events distribution, as in *Methods: Simulations with near-Gaussian distributions*, produce results nearly indistinguishable from the analytical expressions, and, at most two per cent different from the Gillespie runs. Yet coarse-grained simulations require only 1/40th the time of their Gillespie analogue since the time step is large, $\delta t = 35$.

Table 1: Comparison of the Gillespie and the coarse-grained simulation algorithms. The numbers are reported for 12 million realizations of the same evolution for each of the methods. To highlight deviations from the Poisson and the Gaussian statistics, we provide ratios of the higher order cumulants to the mean of the product flux distribution. In the last column, we report analytical predictions obtained from the quasi-steady state approximation to the CGF. Numbers in parentheses are the estimated errors of the last significant digits.

Cumulants	Gillespie	Coarse-grained	Analytics
c_1	11.24(1)	11.14(1)	11.14
c_2/c_1	0.843(1)	0.855(1)	0.855
c_3/c_1	0.613(4)	0.628(4)	0.628
c_4/c_1	0.32(2)	0.32(2)	0.319
time	8 min 45 s	12 s	N/A

3 Supporting figures

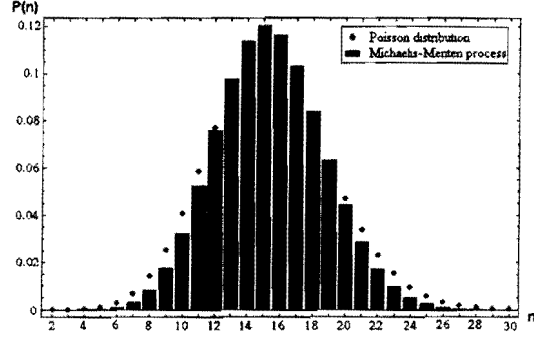


Figure 1: Distribution of the number of MM reactions over a time $\delta t = 35$ with $S_M = 140$, $k_1 = 0.01$, $k_{-1} = 1$, and $k_2 = 1$ vs. the Poisson distribution with the same mean. The distribution for the MM process is obtained using the Gram-Charlier expansion with four known cumulants, see *Methods* in the main article. The MM process is clearly non-Poissonian.