

Hidden Markov Model for Competitive Binding and Chain Elongation

R. M. Roberts, T. J. Cleland,* P. C. Gray, and J. J. Ambrosiano

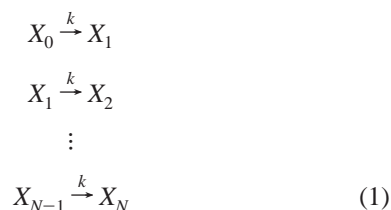
Los Alamos National Laboratory, Los Alamos, New Mexico 87545

Received: October 1, 2003; In Final Form: March 11, 2004

Many chemical systems of interest consist of sets of reactions that contain iterated sequences that elongate a molecular chain. For example, such sets of reactions are commonly found in the transcription of DNA or the translation of RNA. However, there are competitive reactions that can prematurely terminate the chain-elongation process. A hidden Markov method appropriate for modeling chains of reactions with competitive processes (i.e., premature chain termination) is developed. The method is an extension of a hidden Markov model suggested by Gibson and Bruck (*J. Phys. Chem. A* **2000**, *104*, 1876). The equivalence between results using this method and results from simulation of a system employing the full set of reactions will be demonstrated (Gillespie, D. T. *Markov Processes: An Introduction for Physical Scientists*; Academic: San Diego, CA, 1992). Examples of its use are shown for several test problems. As an example of a practical application, a comparison among results for a model of terminal modification of ligand-aggregated receptors (Hlavacek, W.; Redondo, A.; Wofsy, C.; Goldstein, B. *Bull. Math. Biol.* **2002**, *64*, 887) simulated by ordinary differential equations, the full set of reactions, and the hidden Markov method are shown.

1. Introduction

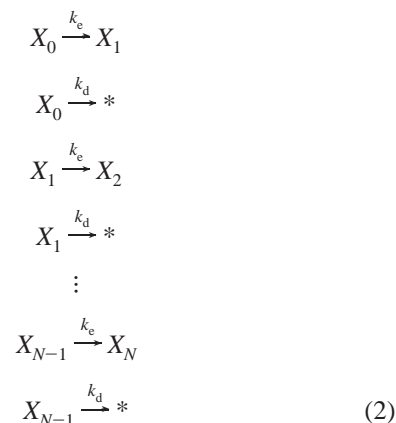
Often in modeling large systems of reactions one is faced with including chains of reactions in which similar subunits polymerize in an orderly iterative fashion. Examples of such reactions are DNA transcription or RNA translation (see ref 1 as a modeling example) where the reactions involve “ratcheting” along a very long string of nucleotides. For the dynamics of these systems, the inner details of these reactions are ancillary, though not entirely ignorable. Gibson and Bruck² presented a hidden Markov method for simulating such reactions, which obviates the need for explicitly simulating the transitions between intermediate states but is stochastically equivalent to the simulation of the full system of reactions. Their method works for systems such as eq 1, which is a simplified set of N reactions representing the elongation of a molecule to a chain of length N .



In such systems, N could be of the order of several hundred to several thousand or more. Monte Carlo simulation of so many reactions can become quite run time and/or storage intensive; moreover, the intermediate steps in the chain may be of little interest in the system being studied. The Gibson and Bruck² method dramatically reduces the number of chemical equations needed to simulate an elongation system by representing the N separate equations of eq 1 by the single *non-Markov* chemical reaction: $X_0 \xrightarrow{N:k} X_N$ (where the $N:k$ over the arrow indicates there are N steps each with reaction rate k).

One weakness of their model is that it assumes that once the system begins the chain of reactions, it will continue through

all N steps undisturbed until completion. However, this is not necessarily a good assumption for real chemical systems. Often, competitive degradation reactions can occur at each step in the chain (e.g., attenuation of DNA transcription or RNA translation). These are represented in eq 2 by the chemical reactions with rate constant k_d and product “*”, where “*” represents any other product(s) that end(s) the chain-elongation process.



To accommodate the more realistic scenario including competitive reactions of this sort, we have developed an extension to the Gibson and Bruck elongation method. In our method, the chemical equations of eq 2, both the elongation steps and competitive degradation steps, can be represented as two reactions



In this equation, the $N:k_e$ represents an elongation reaction of N steps, each with rate constant k_e . The second reaction of eq 3 represents the competitive degradation for *each step* of

the elongation. In both the elongation method of Gibson and Bruck and our method, the assumption for the reaction constant is that k_e of eq 1 and eq 2 must be the same for each elongation reaction along the chain. In addition, our method requires that the rate constant, k_d , for any competitive reaction be the same everywhere along the chain. (It should be noted here that although eq 2 includes only one type of competitive reaction (with rate constant k_d), the method described below is perfectly generalizable to the case where there is more than one type of competitive reaction. In that case, one would just have competitive reaction rates, k_{d1}, k_{d2}, \dots . The same condition must hold true in that each k_{di} must remain the same constant everywhere along the chain.)

2. Gibson and Bruck Elongation Method

In the implementation of Gibson and Bruck's "Next Reaction" method for simulating systems of standard chemical kinetic reactions, a time is randomly chosen for each reaction from an exponential distribution.² The reactions and their generated times are then stored in an *external* priority queue. The reaction with the earliest time (i.e., at the top of priority queue) is taken as the firing reaction.

In their elongation method, a chain of length N and reaction rate k can be thought of as one reaction whose time to completion, t , is sampled from a Gamma distribution parametrized on k and N .

$$\Pr(t \leq \tau < t + dt) = \Gamma_N(k, t) dt = \frac{k(k t)^{N-1}}{(N-1)!} \exp(-k t) dt \quad (4)$$

Instead of generating one t for each reaction of this kind, Gibson and Bruck generate one t for each *molecule* of X_0 . These times are stored in an *internal* priority queue for each elongation reaction. Whenever the smallest t in the *internal* priority queue is smaller than the smallest t in the system's *external* priority queue, the elongation reaction is taken as the firing reaction. That t is then removed from the *internal* priority queue, and the number of X_0 molecules is decremented while the number of X_N molecules is incremented.

As previously mentioned, the assumption by Gibson and Bruck is that once a molecule of X_0 is produced, no reactions other than the elongation steps can affect it.

3. Elongation with Competitive Reactions Method

Our competitive elongation method uses the same structure as described in Section 2 but with the addition of competitive reactions as in eq 2. First, we show what the proper reaction probability distribution is for such a system. To do so, we shall go stepwise through the reactions of eq 2. The probability of starting with a single X_0 at time $t = 0$ and having a reaction, $X_0 \rightarrow X_1$, occur at time τ , where $t_1 < \tau < t_1 + dt_1$, (the reason for the subscript on the time variable will become apparent below) is

$$P_{X_0 \rightarrow X_1}(t_1) dt_1 = \exp[-(k_e + k_d)t_1] k_e dt_1 \quad (5)$$

where the exponential factor is the probability that no reaction occurs in the interval $[0, t_1)$ (see ref 3 pages 413–414) and $k_e dt_1$ is the probability that $X_0 \rightarrow X_1$ occurs in the time interval $[t_1, t_1 + dt_1)$. Similarly, the probability of the next reaction in the chain, $X_1 \rightarrow X_2$, occurring in the time interval $[t_1, t_2 + dt_2)$ is

$$P_{X_1 \rightarrow X_2}(t_2 - t_1) dt_2 = \exp[-(k_e + k_d)(t_2 - t_1)] k_e dt_2 \quad (6)$$

Multiplying eqs 5 and 6 and integrating t_1 from 0 to t_2 gives the probability of starting with an X_0 at $t = 0$ and producing an X_2 at time τ , where $t_2 < \tau < t_2 + dt_2$:

$$P_{X_0 \rightarrow X_2}(t_2) dt_2 = k_e^2 t_2 \exp[-(k_e + k_d)t_2] dt_2 \quad (7)$$

Following the same procedure for the next step in the chain, we find the equation analogous to eq 6 for the reaction $X_2 \rightarrow X_3$

$$P_{X_2 \rightarrow X_3}(t_3 - t_2) dt_3 = \exp[-(k_e + k_d)(t_3 - t_2)] k_e dt_3 \quad (8)$$

and multiply it with eq 7 followed by integration over t_2 to get

$$P_{X_0 \rightarrow X_3}(t_3) dt_3 = k_e^3 \frac{t_3^2}{2} \exp[-(k_e + k_d)t_3] dt_3 \quad (9)$$

as the probability of starting with an X_0 at $t = 0$ and producing an X_3 at time τ where $t_3 < \tau < t_3 + dt_3$. This procedure taken through all N steps of the elongation gives

$$P_{X_0 \rightarrow X_N}(t) dt = \Gamma_N(k_e, t) \exp(-k_d t) dt \quad (10)$$

where we have now dropped the subscripts on t and $\Gamma_N(k_e, t)$ is the Gamma distribution defined in eq 4. Equation 10 is the probability of starting with a single X_0 at $t = 0$ and having it follow all N elongation steps resulting in a single X_N at t .

We must now find the probability distribution for starting with a single X_0 and having a competitive reaction occur somewhere in the elongation process. Each possible path which ends the elongation (e.g., $X_0 \rightarrow *$, $X_0 \rightarrow X_1 \rightarrow *$, $X_0 \rightarrow X_1 \rightarrow X_2 \rightarrow *$, etc.) is a separate alternative, and thus we should expect a solution which is a sum of N probability terms. We start by noting that the reactions $X_i \rightarrow *$ ($i = 0, 1, \dots, N-1$) are simple one-step reactions. Therefore, the appropriate probability distribution is an exponential distribution with parameter k_d :

$$P_{X_i \rightarrow *}(t_{i+1} - t_i) dt_{i+1} = k_d \exp[-k_d(t_{i+1} - t_i)] dt_{i+1} \quad (i = 1, 2, \dots, N-1) \quad (11)$$

where we have again subscripted the time variables as a convenience in anticipation of integrating over the various t_i variables. For each competitive pathway, note that for the $X_i \rightarrow *$ reaction to occur, we must first have completed the elongation to X_i . Following a similar stepwise procedure as that which gave us eq 10, while using eq 11 as the final step in each pathway, we get

$$P_{X_0 \text{all} \rightarrow *}(t) = k_d \exp(-k_d t) \sum_{i=0}^{N-1} \frac{(k_e t)^i}{i!} \exp(-k_e t) dt \quad (12)$$

where the form of the right-hand side was chosen to emphasize the Poisson-like structure of the terms in the sum.

The total probability distribution for the competitive elongation reactions is the sum of eq 10 and eq 12:

$$P_{\text{total}}(t) dt = \Gamma_N(k_e, t) \exp(-k_d t) dt + k_d \exp(-k_d t) \sum_{i=0}^{N-1} \frac{(k_e t)^i}{i!} \exp(-k_e t) dt \quad (13)$$

for which it can be shown that

$$\int_0^\infty P_{\text{total}}(t) dt = 1 \quad (14)$$

indicating it is properly normalized.

Equation 13 is quite complex, and randomly sampling a time from it would be both difficult and computationally expensive. However, there is a shortcut method which is entirely equivalent. The chain-elongation reaction t is still chosen from the same Gamma distribution (eq 4), but now the competing reaction is included by also choosing a different t from a standard exponential distribution parametrized by k_d and the current number of X_0 molecules in play, N_{X_0} ,

$$\Pr(t \leq \tau < t + dt) = N_{X_0} k_d \exp(-N_{X_0} k_d t) dt \quad (15)$$

When a competing reaction occurs (i.e., its t is the smallest in the *external* priority queue), we decrement the number of X_0 molecules by randomly deleting one entry from the *internal* priority queue. A new competing reaction time is then generated from eq 15 (with the newly updated value of N_{X_0}) and put into the *external* priority queue. If the deleted entry happens to be the one at the top of the *internal* priority queue (i.e., the one with the earliest firing time), the value of the firing time for the elongation reaction is then replaced by the next smallest value in the elongation reaction's *internal* priority queue. Thus, each elongation has a finite probability of premature termination. (In the more general case where there are multiple competitive reactions having rate constants k_{d1} , k_{d2} , ..., one would simply have a distribution analogous to eq 15 for each competitive reaction using the appropriate rate constant in each. Then a separate time t would be generated for each competitive reaction.)

In the next section we show mathematically that the above algorithm results in the same distribution as eq 13.

4. Mathematical Validation

Suppose that for each molecule of X_0 currently in play, we randomly generate a time, t_e , from the Gamma distribution (i.e., eq 4) with $k = k_e$ (corresponding to the time for completion of the elongation reaction for that molecule, $X_0 \rightarrow X_N$) and a time, t_d , from the exponential distribution, $P(t) dt = k_d \exp(-k_d t) dt$ (corresponding to the time that the elongation reaction would be disrupted at some point in the chain) and compare these times to determine if a complete elongation reaction occurred before being disrupted. The resulting probability for $X_0 \rightarrow X_N$ completion would have the form

$$P_{\text{elong}}(t_e) dt_e = \Gamma_N(k_e, t_e) dt_e \int_{t_e}^\infty k_d \exp(-k_d t) dt \quad (16)$$

where the Gamma distribution is the probability that t_e was generated for the elongation reaction and the integral is the probability that $t_d > t_e$. Similarly, the probability for competitive disruption would be

$$P_{\text{comp}}(t_d) dt_d = k_d \exp(-k_d t_d) dt_d \int_{t_d}^\infty \Gamma_N(k_e, t) dt \quad (17)$$

The total probability is the sum of eqs 16 and 17. Performing

the integrations and then dropping the subscripts on t gives

$$P_{\text{total}}(t) dt = \Gamma_N(k_e, t) \exp(-k_d t) dt + k_d \exp(-k_d t) \sum_{i=0}^{N-1} \frac{(k_e t)^i}{i!} \exp(-k_e t) dt \quad (18)$$

which matches eq 13. (In the case in which we have multiple competitive reactions having rate constants k_{d1} , k_{d2} , ..., one could follow a similar procedure and get an analogous distribution where the first term would be $\Gamma_N(k_e, t) \exp(-\sum_i k_{di} t) dt$ and the second term has a factor of $\sum_i k_{di} \exp(-k_{di} t)$, replacing the $k_d \exp(-k_d t)$ factor above. This is the same distribution one would get by doing the stepwise derivation outlined in Section 3 for the case of multiple competitive reactions.)

Of course the above proof assumes we are generating a separate t_d for each X_0 . We now show that choosing one time, t_d , from the distribution of eq 15 as described in the algorithm is equivalent to the above.

Suppose we had an ensemble of groups of particles. Each group in the ensemble has M identical particles, each of which is assigned a real-valued time, t_i ($i = 1, 2, \dots, M$), randomly generated from a probability distribution $f(t)$. We are interested only in the smallest-valued t_i for each of the groups in the ensemble. What is the probability distribution, $P_M^{(1)}(t)$, formed by collecting all the smallest t_i values (one from each group) together into one group? We can find $P_M^{(1)}(t)$ by multiplying $f(t)$ by the $(M - 1)$ -fold probability of a time generated from $f(t)$ being greater than or equal to t . That is,

$$P_M^{(1)}(t) dt = M f(t) \left[\int_t^\infty f(t') dt' \right]^{M-1} dt \quad (19)$$

where the leading factor of M is there because there are M possible t_i values and they all have a priori equal probability of being the smallest; therefore, each possible outcome contributes to the probability equally. It can be shown that $P_M^{(1)}$ is properly normalized by integrating from $t = 0$ to $t \rightarrow \infty$ (to integrate $P_M^{(1)}$ call the rhs integral $F(t)$ and note that $dF(t) = -f(t) dt$).

We can return to the problem at hand by substituting $k_d \exp(-k_d t)$ in for $f(t)$ and substituting N_{X_0} in for M . After performing the interior integration, we get

$$P_{N_{X_0}}^{(1)}(t) dt = N_{X_0} k_d \exp(-N_{X_0} k_d t) dt \quad (20)$$

which matches eq 15 and validates the use of eq 15 as our method of generation of t_d .

Although this shows that the distribution used by our algorithm is correct, we still have to show that the distribution is preserved after a random deletion of an X_0 when a competitive reaction fires. To prove this, it is necessary and sufficient to show that a random deletion of an X_0 transforms the distribution from $N_{X_0} k_d \exp(-N_{X_0} k_d t) dt$ to $(N_{X_0} - 1) k_d \exp(-(N_{X_0} - 1) k_d t) dt$. Following the logic that led to eq 19, suppose we take the same ensemble of groups of M particles with assigned times, except instead of the smallest-valued t_i , we take the second smallest t_i from every group and place them in their own group. The resulting distribution of these second-smallest times would be:

$$P_M^{(2)}(t) dt = M(M - 1) f(t) \left[\int_0^t f(t') dt' \right] \left[\int_t^\infty f(t') dt' \right]^{M-2} dt \quad (21)$$

where the integral from 0 to t is the probability that one of the t_i values is less than t and the other integral factor is the

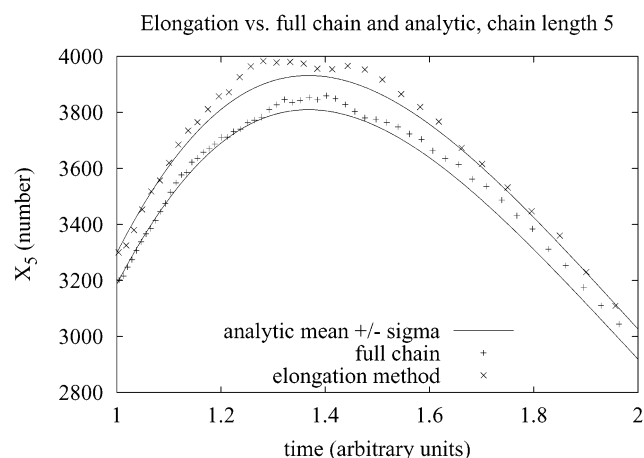
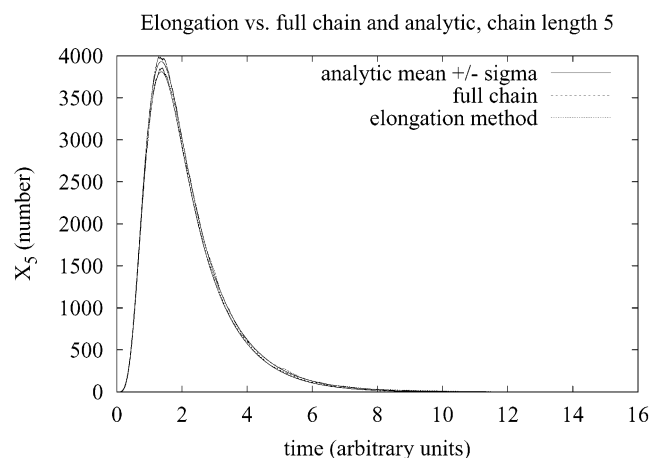


Figure 1. Top: Comparison of results for chain length of five using the full-chain simulation, our elongation method, and the analytical solution. (A feedback reaction $X_5 \xrightarrow{k_{fb}} X_0$ was included in all cases.) $k_c = 3.0$, $k_d = 2.0$, $k_{fb} = 1.0$, $X_0(0) = 100\,000$. Bottom: Rescaled to emphasize peak. Points are taken at every 500th reaction for clarity.

probability that $M - 2$ of the t_i values are greater than or equal to t . The leading factor of $M(M - 1)$ is there because there are $M(M - 1)$ alternative ways of permuting M objects into groups of size 1 (i.e., which t_i is smaller than t) and size $M - 2$ (i.e., which t_i values are larger than t), and each alternative contributes equally to the probability. It can be shown that by integrating eq 21 from $t = 0$ to $t \rightarrow \infty$, $\binom{2}{M}(t)$ is properly normalized.

Substituting our distribution and performing the internal integration gives

$$P_{N_{X_0}}^{(2)}(t) dt = N_{X_0}(N_{X_0} - 1)k_d \exp[-(N_{X_0} - 1)k_d t] [1 - \exp(-k_d t)] dt \quad (22)$$

A random deletion of one of the N_{X_0} particles would mean there would be a probability of $(N_{X_0} - 1)/N_{X_0}$ that we delete a particle other than the one with the smallest time (meaning the smallest time of the group does not change), and a $1/N_{X_0}$ probability that we delete the particle *with* the smallest time (meaning the smallest time shifts to be what was once the second-smallest time). In other words, the resulting distribution, $P_{N_{X_0}-1}^{(1)}(t)$, will be the weighted sum

$$P_{N_{X_0}-1}^{(1)}(t) dt = \left(\frac{1}{N_{X_0}}\right) P_{N_{X_0}}^{(2)}(t) + \left(\frac{N_{X_0} - 1}{N_{X_0}}\right) P_{N_{X_0}}^{(1)}(t) \quad (23)$$

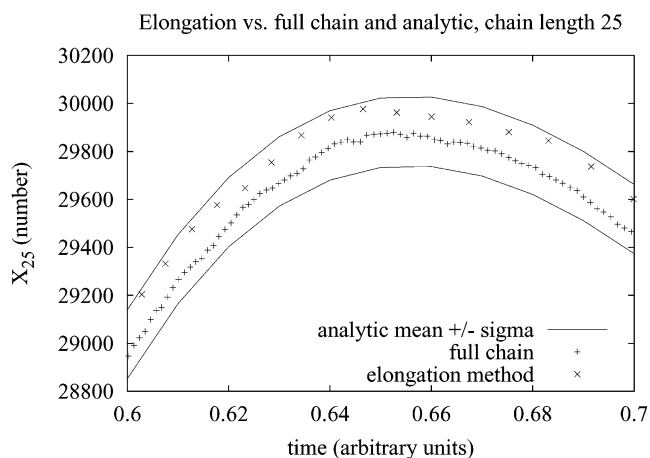
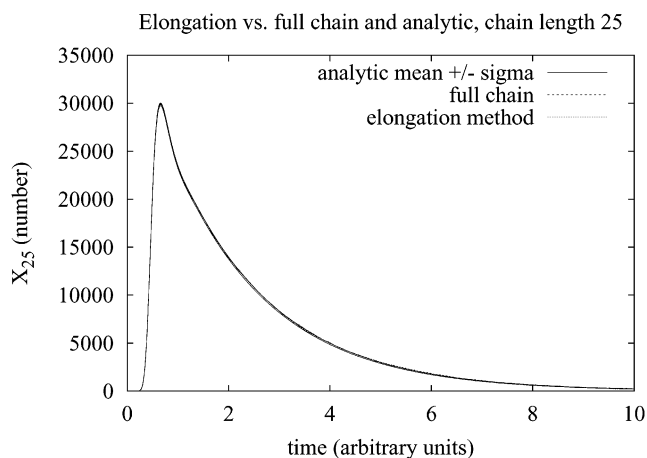


Figure 2. Top: Comparison of results for chain length of 25 using the full-chain simulation, our elongation method, and the analytical solution. (A feedback reaction $X_{25} \xrightarrow{k_{fb}} X_0$ was included in all cases.) $k_c = 50.0$, $k_d = 2.0$, $k_{fb} = 1.0$, $X_0(0) = 100\,000$. Bottom: Rescaled to emphasize peak. Points are taken at every 500th reaction for clarity.

After substitution of eq 20 and eq 22 and some algebra, we get

$$P_{N_{X_0}-1}^{(1)}(t) dt = (N_{X_0} - 1)k_d \exp[-(N_{X_0} - 1)k_d t] dt \quad (24)$$

which is the correct distribution and completes the proof.

5. Results and Discussion

We present three sets of results comparing the competitive elongation method with analytic calculations. All of the simulations were performed using *BioReactor*, an open-source reaction kinetics tool developed by us. *BioReactor* includes our competitive elongation method in the Monte Carlo solver. The first two sets compare the competitive elongation method, using the chemical equations of eq 3, with *BioReactor* runs using an explicit full chain description of the chemical equations, as per eq 2. In each case, we have included a feedback reaction, $X_N \rightarrow X_0$, to have an analytic solution for validation.

Figures 1 and 2 demonstrate typical runs using the full chain and competitive elongation method. These are plotted against an analytic calculation of the mean and standard deviation of X_N as a function of time, using a multivariate extension of the moment evolution equations for a Jump Markov process with discrete states (see ref 4, Section 5.1.E). We have also determined the statistical agreement of the two probability distributions of peak values for X_N from 1000 realization runs of *BioReactor* for each stochastic treatment.

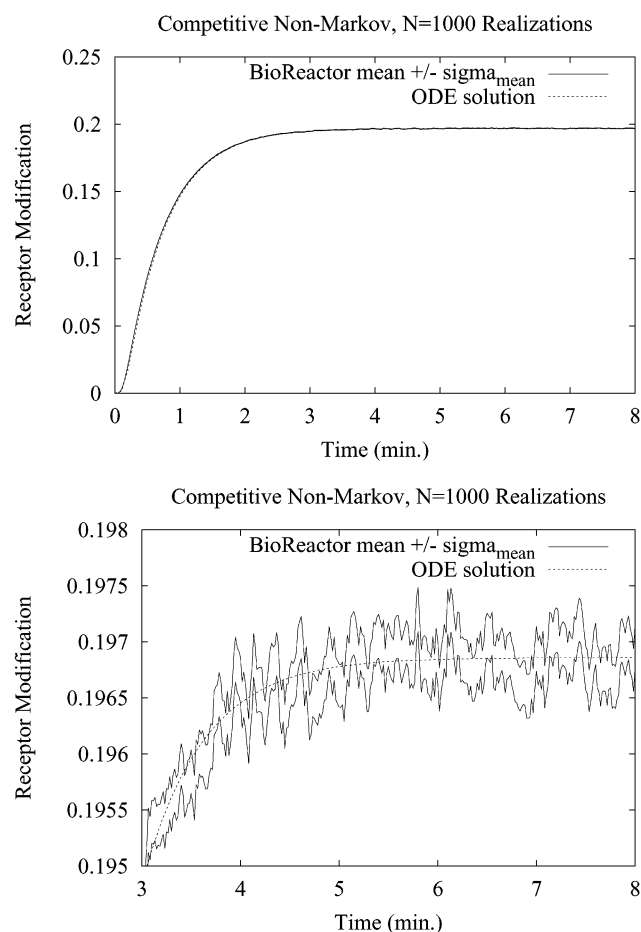


Figure 3. Top: Comparison of our competitive elongation model with ODE solution for terminal modification of ligand-aggregated receptors (see the high-affinity ligand result of Figure 5a of ref 5). Bottom: Rescaled to show stochastic fluctuations.

As a further example of where this stochastic method might be applied, we have performed comparison with a deterministic receptor-modification model. For details of the model, see Hlavacek et al.⁵ In Figure 3 the sample mean, \bar{X} and sample variance, s^2 , are calculated with the results from 1000 realizations of *BioReactor* on a model for terminal modification of ligand-aggregated receptors. The standard deviation of the sample mean, $\sigma_{\bar{X}}$, is calculated from the population standard deviation, σ , by $\sigma_{\bar{X}} = (\sigma/\sqrt{n})$, where n is the total number of realizations, in this case $n = 1000$ (see ref 6, Section 7.3). We use the square root of the sample variance as an approximation for σ .

The major motivation for developing this method was the need to simulate attenuated DNA transcription and RNA

translation. In the case of DNA transcription, we have an RNA-polymerase molecule attached to DNA (the resulting initial complex would play the role of X_0 above) while “ratcheting” along one of the DNA strands adding a nucleotide at each step and ultimately forming messenger RNA. The attenuation is the result of the fact that at any point along the DNA strand, the RNA-polymerase can “fall off” rather than ratchet forward.

Similarly for RNA translation, a ribosome attaches to messenger RNA, ratcheting along the chain and adding an amino acid at each step until a terminal codon is reached (ultimately producing a complete protein). The competing reaction is one in which the ribosome “falls off” midstream.

It should be noted that as the number of elongation steps, N , gets larger, the ratio k_e/k_d must get larger, otherwise no complete elongations will occur. For example, for the case of an elongation of length $N = 50$, if $k_e/k_d = 1$, then having a complete elongation occur would have the exact same probability as flipping a coin $N = 50$ times and getting all 50 heads. If anywhere along the way a coin flip results in a tail, one has to start over again. In other words, all it takes is one successful competing reaction to end the elongation, whereas N successful elongation reactions must occur for completion. Since most elongations of interest deal with potentially very large values of N , the competing reaction must have a much smaller probability of occurring (i.e., the ratio k_e/k_d must be large).

6. Conclusions

We have presented a method which extends the elongation reaction method of Gibson and Bruck² to allow for competitive disruption (attenuation) of the elongation as is generally possible in real chemical systems. Our method is very efficient in that it reduces the need to simulate $2N$ reactions down to only two reactions (one elongation, one disruption). We have shown the mathematical validity of our method and presented results which show the statistical equivalence of our method with the full $2N$ -fold simulation.

Acknowledgment. Los Alamos National Laboratory is operated by the University of California for the National Nuclear Security Administration of the U.S. Department of Energy under contract W-7405-ENG-36. We would like to thank W. Hlavacek for providing the test model and the ODE calculation results.

References and Notes

- (1) Arkin, A. P.; Ross, J.; McAdams, H. H. *Genetics* **1998**, *149*, 1633.
- (2) Gibson, M. A.; Bruck, J. *J. Phys. Chem. A* **2000**, *104*, 1876.
- (3) Gillespie, D. T. *J. Comput. Phys.* **1976**, *22*, 403.
- (4) Gillespie, D. T. *Markov Processes: An Introduction for Physical Scientists*; Academic: San Diego, CA, 1992.
- (5) Hlavacek, W.; Redondo, A.; Wofsy, C.; Goldstein, B. *Bull. Math. Biol.* **2002**, *64*, 887.
- (6) Rice, J. A. *Mathematical Statistics and Data Analysis*, 2nd ed.; Duxbury Press: Belmont, CA, 1995.