

# Increasing the Time Resolution of Single-Molecule Experiments with Bayesian Inference

Colin D. Kinz-Thompson<sup>1,\*</sup> and Ruben L. Gonzalez, Jr.<sup>1,\*</sup>

<sup>1</sup>Department of Chemistry, Columbia University, New York, New York

**ABSTRACT** Many time-resolved single-molecule biophysics experiments seek to characterize the kinetics of biomolecular systems exhibiting dynamics that challenge the time resolution of the given technique. Here, we present a general, computational approach to this problem that employs Bayesian inference to learn the underlying dynamics of such systems, even when they are much faster than the time resolution of the experimental technique being used. By accurately and precisely inferring rate constants, our Bayesian inference for the analysis of subtemporal resolution dynamics approach effectively enables the experimenter to super-resolve the poorly resolved dynamics that are present in their data.

## INTRODUCTION

Given their inherent ability to eliminate ensemble averaging, time-resolved single-molecule biophysical methods have revolutionized the study of biological mechanisms by enabling distributions of molecular properties to be observed, stochastic fluctuations from equilibrium to be investigated, and transiently sampled reaction intermediates to be characterized (1). Generally, the majority of these methods involve making sequential measurements of an experimental signal that acts as a proxy for the underlying, time-dependent state of a biomolecule. As a result, this process yields a time-ordered series of discrete measurements from which the underlying dynamics of the corresponding biomolecule can be inferred (2). Unfortunately, the ability to resolve the continuously varying dynamics of the corresponding biomolecule from a series of discrete measurements is fundamentally limited. Indeed, whereas a biomolecule may exchange between multiple conformational states during a single measurement acquisition period, these states are collectively represented by a single, time-averaged measurement of the experimental signal. This effect is akin to chemical exchange effects in NMR experiments, in which distinct resonance peaks can coalesce into a single, averaged resonance peak when a nucleus rapidly exchanges between distinct magnetic environments (3). As a result of this effect, many time-resolved single-molecule biophysical methods often fail to detect or properly charac-

terize mechanistically critical biomolecular processes that occur on or faster than the time resolution of the technique, including early steps in ligand binding, local folding events, and rapid conformational fluctuations (4,5).

To push beyond the time-resolution limits of these single-molecule methods, we have developed a Bayesian inference-based computational approach, which we call Bayesian inference for the analysis of subtemporal resolution dynamics (BIASD), to infer the rate constants governing transitions between discrete states of a single molecule from the analysis of a time-resolved single-molecule experimental signal, even if those rate constants are substantially faster than the time resolution of the recorded experimental signal. Much like learning the point spread function describing the fluorescence signal from a single fluorophore in a super-resolution imaging experiment enables the spatial position of the fluorophore to be inferred beyond the spatial resolution of the experiment, learning the model describing the kinetic behavior of a single molecule in a time-resolved single-molecule experiment using BIASD enables the kinetic behavior of the single molecule to be inferred beyond the temporal resolution of the experiment. By using Bayesian inference, BIASD can also integrate information from other experiments to further enhance its resolving power, and it also employs a natural framework with which to describe the precision that the amount of data collected during the single-molecule experiment will lend to the determination of the parameters governing the single-molecule kinetics (2,6,7). It is worth noting that, in a close parallel to the approach we describe here, Bayesian inference has been previously employed to improve the time resolution of the time-dependent free induction

Submitted July 6, 2017, and accepted for publication November 21, 2017.

\*Correspondence: [cdk2119@columbia.edu](mailto:cdk2119@columbia.edu) or [rlg2118@columbia.edu](mailto:rlg2118@columbia.edu)

Editor: Joseph Puglisi.

<https://doi.org/10.1016/j.bpj.2017.11.3741>

© 2017 Biophysical Society.



decay in NMR spectroscopy experiments, resulting in an orders-of-magnitude improvement in spectral resolution (7,8).

Here, we first describe the Bayesian inference-based framework underlying BIASD. We then use BIASD to analyze computer-simulated signal versus time trajectories (signal trajectories) and investigate the accuracy and precision with which we can infer the known rate constants for transitions between states that were used to generate the signal trajectories. We next use BIASD to analyze experimentally recorded fluorescence resonance energy transfer efficiency ( $E_{\text{FRET}}$ ) versus time trajectories ( $E_{\text{FRET}}$  trajectories) to infer the unknown rate constants for transitions between states in the  $E_{\text{FRET}}$  trajectories. Notably, the  $E_{\text{FRET}}$  trajectories that we have analyzed here had previously eluded analysis due to the presence of transitions that are much faster than the time resolution of the electron-multiplying charge-coupled device camera that was used to record them (9). Finally, we describe and demonstrate an extension of the BIASD framework that can be used to infer rate constants for experimental systems consisting of static or interconverting subpopulations of molecular properties within an individual or ensemble of molecules. Remarkably, we find that BIASD permits accurate inference of rates constants from time-resolved single-molecule experiments, even when the rate constants are orders of magnitude larger than the time resolution of the signal trajectories.

### Bayesian inference-based framework underlying BIASD

In biomolecular systems, functional motions—such as those involved in ligand binding and dissociation processes, or large-scale conformational rearrangements—very often involve the simultaneous formation and/or disruption of numerous, noncovalent interactions. The relatively low probability of simultaneously forming and/or disrupting these numerous interactions can therefore result in large, entropically dominated, transition-state energy barriers for such functional motions (10,11). Consequently, individual biomolecules are generally expected to exhibit effectively discrete and instantaneous transitions between relatively long-lived states (5), an expectation that is consistent with the step-like transitions that are generally observed in time-resolved single-molecule experiments (12).

An important consideration when analyzing the signal trajectories from such single-molecule experiments is that whenever an individual molecule undergoes a transition from one state to another, the transition occurs stochastically during the time period,  $\tau$ , over which the detector collects and integrates the signal to record a data point in the signal trajectory. Thus, the probability that a transition will coincide exactly with the beginning or end of the  $\tau$  in which it takes place is essentially zero. As a result, when a transition takes place, the signal value that is recorded during that  $\tau$

does not solely represent either of the states involved in that transition. Instead, it represents the average of the signal values corresponding to the states that are sampled during  $\tau$ , weighted by the time spent in each of those states. This time averaging makes it imprudent to assign the signal value recorded during such a  $\tau$  to any one particular state, a process called idealization, because the molecule will have occupied multiple states during that  $\tau$ . Notably, when the rate constants for transitions between states become comparable to or greater than  $\tau^{-1}$ , there is a large probability that the  $\tau$ s of a signal trajectory will contain one or more transitions, and that, consequently, many of the signal values of the signal trajectory will exhibit this time averaging. Given such a scenario, analysis methods in which individual  $\tau$ s are assigned to particular states (e.g., the widely used strategy of idealizing signal trajectories using signal thresholds (13) or hidden Markov models (HMMs) (14,15)) will introduce significant errors into the calculated rate constants for transitions between states and into the signal values assigned to those states (2).

To overcome the potential errors associated with determining rate constants and signal values from the analysis of signal trajectories, BIASD instead analyzes a different parameter that depends upon the dynamics of the biomolecular system: the fraction of time that is spent in each state during the  $\tau$ s in a signal trajectory (16–21). To illustrate this approach, consider the case of an individual molecule that undergoes stochastic, uncorrelated (i.e., Markovian), and reversible transitions between two states, denoted 1 and 2 (i.e.,  $1 \rightleftharpoons 2$ , with forward and reverse rate constants of  $k_1$  and  $k_2$ , respectively), which have unique signal values of  $\epsilon_1$  and  $\epsilon_2$ . If the fraction of time that the molecule spends in state 1 during a particular  $\tau$  is  $f$ , then, because of the two-state nature of the system, the fraction of time that the molecule spends in state 2 during that  $\tau$  is  $1 - f$ . It is important to note that, although the molecule is in an equilibrium between states 1 and 2, the value of  $f$  for any particular  $\tau$  will not necessarily be the equilibrium value of  $f = (1 + k_1/k_2)^{-1}$ . This is because  $\tau$  might not be long enough for sufficient time averaging to occur (i.e., to invoke ergodicity). Instead, each  $\tau$  will exhibit a distinct, time-averaged value of  $f$ .

The exact value of  $f$  for a particular  $\tau$  will depend upon the molecule's stochastic path through state-space during  $\tau$ . As such, a probabilistic description of  $f$ , which accounts for all possible paths through state-space, is needed to describe the likelihood of observing a particular value of  $f$  during a  $\tau$  (22). In particular, for the reversible two-state system considered here, such a description, which has roots in the analysis of the NMR chemical exchange effects described above (23) and in sojourn-time probability distributions (24), was first given by Dobrushin (25). This particular expression (derived in the [Supporting Material](#)) is a function of  $k_1$ ,  $k_2$ , and  $\tau$ , and has been used in many single-molecule studies, though mostly in the context of

photon-counting experiments and without Bayesian inference-based implementations (16–26). Experimentally, if the exact values of  $f$ ,  $\epsilon_1$ , and  $\epsilon_2$  during each  $\tau$  were known, one would be able to calculate the expected value of the corresponding time-averaged signal,  $\mu$ , for each  $\tau$ , because it would be the linear combination  $\mu = f\epsilon_1 + (1 - f)\epsilon_2$ . However, the analysis of time-resolved single-molecule experiments deals with the opposite problem: observing a signal value,  $d$ , during each  $\tau$  and trying to infer  $f$ ,  $\epsilon_1$ , and  $\epsilon_2$ .

Generally, the values of  $d$  that are recorded during each  $\tau$  are random variables, which are distributed according to a probability distribution function (PDF). For any number of states, this PDF for the observed values of  $d$  is the convolution of the PDFs for the signal values associated with each individual state, weighted by the fraction of time spent in that state (i.e.,  $d \sim f_1 p(\epsilon_1) * \dots * f_n p(\epsilon_n)$ , where  $f_i$  is the fractional occupancy of the  $i^{\text{th}}$  state,  $p(\epsilon_i)$  is the PDF of the signal values associated with the  $i^{\text{th}}$  state, and  $*$  denotes a convolution). For many experimental techniques, the signal values associated with each state are, or are approximately, distributed according to a normal PDF (i.e., a Gaussian) with mean  $\epsilon_i$  and variance  $\sigma_i^2$  for the  $i^{\text{th}}$  state. Because the convolution of two normal PDFs is another normal PDF, in this case the PDF for the observed values of  $d$  is a normal distribution with mean  $\mu = \sum_i f_i \epsilon_i$  and variance  $\sigma^2 = \sum_i f_i \sigma_i^2$ . Furthermore, we can also account for noise from the detection process (e.g., a normal PDF with mean 0 and variance  $\sigma_{\text{noise}}^2$ ), as well as a time-dependent baseline (e.g., baseline drift at time  $t$ ,  $b_t$ , that is driven by white-noise is a normal PDF with mean  $b_{t-1}$  and variance  $\sigma_{\text{drift}}^2$ ) through additional convolutions; in these examples, the resulting PDF of  $d$  is again a normal PDF with mean  $\mu = b_{t-1} + \sum_i f_i \epsilon_i$  and variance  $\sigma^2 = \sigma_{\text{drift}}^2 + \sigma_{\text{noise}}^2 + \sum_i f_i \sigma_i^2$ . However, since  $\mu$  and  $\sigma$  depend upon the set of fractional occupancies,  $\{f\}$ , which are not experimental observables, we have no way of knowing the exact form of this PDF, information that is required to calculate the probability of observing a particular value of  $d$ .

To circumvent this experimental limitation, the dependence of the PDF upon  $\{f\}$  can be removed by marginalizing  $\{f\}$  out of the expression for the PDF that was described above. This marginalized probability distribution of  $d$  then describes the likelihood of experimentally observing a particular value of  $d$  during a  $\tau$  as a function of the set of rate constants for transitions between the states,  $\{k\}$ , the set of signal values corresponding to those states,  $\{\epsilon\}$ , and the set of the amounts of noise in those states,  $\{\sigma\}$ , regardless of the exact values of  $\{f\}$  (Fig. S1 A). As expected from the discussion in the previous section, this expression describes effects similar to those of chemical exchange in NMR experiments, in which rates with which nuclei exchange that are larger than the resonance frequency difference between exchanging nuclei cause distinct resonances to coalesce into a single, averaged resonance. As shown in Fig. 1 for a two-state system (see Eq. S10), the effect of increasing rate constants  $k_1$  and  $k_2$  results in distinct signal peaks

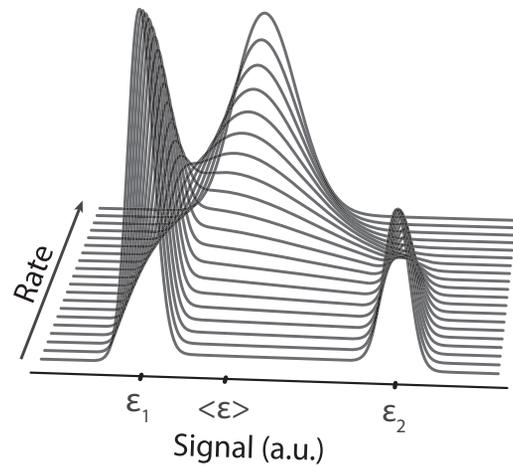


FIGURE 1 The marginalized, two-state likelihood function as a function of increasing rate of transitions. Arbitrary signal peaks at  $\epsilon_1$  and  $\epsilon_2$  coalesce into a single peak located at the equilibrium average  $\langle \epsilon \rangle$  as the rate constants for transitions between the two states increase relative to a fixed  $\tau$ . a.u., arbitrary units.

centered at  $\epsilon_1$  and  $\epsilon_2$  to coalesce into a single, averaged peak centered at  $\langle \epsilon \rangle$ .

With such an expression describing the marginalized probability distribution of  $d$ , we can then use Bayesian inference to estimate the parameters governing the single-molecule system (i.e.,  $\{\epsilon\}$ ,  $\{\sigma\}$ , and  $\{k\}$ ) from the series of the  $d$  that comprise each of the signal trajectories. Primarily due to recent developments in computational tractability, Bayesian inference has become a powerful method for the analysis of biophysical data, such as determining the phases of x-ray reflections in x-ray crystallographic studies (27), performing simultaneous phylogenetic analysis of nucleotide and protein data sets (28), elucidating the number of structural classes present in cryogenic electron microscopy images (29), and ascertaining the number of states and the rates of transitions between those states present in single-molecule signal trajectories (30,31). For an introduction to Bayesian inference, see (6,7,32) and the [Supporting Material](#).

Unfortunately, performing Bayesian inference on a multi-parameter system, such as the one described here, results in a multidimensional, joint-probability distribution of the model parameters, known as a posterior probability distribution, which is difficult to evaluate (32). To overcome this difficulty, we have chosen to evaluate the posterior probability distribution of the model parameters by numerically sampling it using a Markov chain Monte Carlo (MCMC) (6,33) method with affine-invariant ensemble sampling (34,35). Although alternative methods that approximate the posterior probability distribution of the model parameters, such as the Laplace approximation or variational inference, might be more computationally tractable, MCMC sampling is advantageous in that, unlike such approximation methods, it can provide an exact result that does not assume a particular structure of the posterior probability distribution (6). However, regardless of

the choice of method, the most important aspect of the approach described here is that we can evaluate the posterior probability distribution of the model parameters from the series of  $d$  that comprise a single-molecule signal trajectory in a manner that analytically accounts for the time resolution of the experimental technique.

To quantify the performance of BIASD and maximize its usefulness to the experimentalist community, we have also conducted a comprehensive analysis of how the posterior probability distribution behaves as a function of the parameters of the input signal trajectories (c.f., *Dependence of BIASD Performance on Parameter Values* in [Supporting Material](#)). This analysis reveals that collecting additional data points in a signal trajectory increases the performance of BIASD, as does optimizing the sensitivity of the instrumentation so as to increase the signal-to-noise ratio (SNR) of the signal trajectory. Although the results of these analyses can be used to determine the signal trajectory lengths, and/or SNRs that would be needed to accurately and precisely infer rate constants that are some arbitrary factor greater than  $\tau^{-1}$ , we find that the lengths and SNRs of the signal trajectories obtained using typical single-molecule instrumentation render BIASD useful for characterizing dynamics governed by rate constants that are up to three orders of magnitude greater than  $\tau^{-1}$ .

## METHODS

### Simulating signal trajectories

State trajectories were simulated with the stochastic simulation algorithm (36); briefly, sequential random lifetimes were drawn from exponential distributions with the specified rate constants, and subsequent states were chosen randomly according to the splitting probabilities. A random starting point for the initiation of the trajectory ( $t = 0$  s) was selected with a uniform distribution from the first lifetime. The fractional occupancies of each state during each sequential  $\tau$  were then calculated from the sequence of lifetimes. The resulting fractional occupation versus time trajectories were turned into signal trajectories by computing  $\mu$ , and then adding normally distributed noise with standard deviation (SD),  $\sigma$ . Simulations of the titration experiment were performed such that  $\epsilon_1 = 0$ ,  $\epsilon_2 = 1$ ,  $\sigma_1 = \sigma_2 = 0.04$ ,  $k_1^* = 10 \mu\text{M}^{-1}\text{s}^{-1}$ ,  $k_2 = 10 \text{s}^{-1}$ , and  $\tau = 0.1$  s, and each signal trajectory was 600 data points in length. Parameters for the simulations with the transitioning subpopulations are given in the [Supporting Material](#).

### Bayesian thresholding analysis

Signal trajectories were idealized by thresholding any measurement period with signal less than  $(\epsilon_1 + \epsilon_2)/2 = 0.5$  into state 1, and otherwise into state 2. Rate constants from the  $i^{\text{th}}$  state to  $j^{\text{th}}$  state were then calculated as  $k_{ij} = -\ln(1 - p_{ij})/\tau$ , where  $p_{ij}$  is the transition matrix from the idealized trajectory (2). Credible intervals for the transition probabilities and rate constants were calculated with uniform prior distributions (2). The joint posterior probability distributions of  $\epsilon_1$  and  $1/\sigma_1^2$ , and  $\epsilon_2$  and  $1/\sigma_2^2$ , were inferred using the analytical formulas for Bayesian inference with a joint normal-gamma prior probability distribution using those data points that were idealized into the respective states (6). The marginalized posterior distributions of the  $\epsilon_i$  and  $1/\sigma_i^2$  (T and gamma distributions, respectively) were used to calculate means and 95% credible intervals for each parameter (6).

An aggregate  $\sigma$  was then calculated by weighting  $\sigma_1^2$  and  $\sigma_2^2$  by the fraction of data points idealized into each state and taking the square root of their sum.

## Maximum-likelihood HMM analysis

Signal trajectories were analyzed using a two-state, discrete maximum likelihood HMM (ML HMM) with normal distribution emissions using the expectation maximization, and forward-backward algorithms (6,37). Each trajectory was analyzed with 20 randomized restarts, including one initialized at the simulated values, until the likelihood of each restart converged to a relative value of  $10^{-10}$ . From these, the point estimate with the greatest likelihood was used in subsequent analyses. Rate constants were calculated directly from the transition probability matrix point estimate, and an aggregate SD was calculated as described above for Bayesian thresholding.

## BIASD analysis

Adaptive Gauss-Kronrod (G10, K21) quadrature was used to numerically integrate the BIASD likelihood function on an Nvidia GeForce 1080 GTX graphics card; the likelihood of each data point took  $\sim 1 \mu\text{s}$  to compute. The posterior probability distribution was sampled using *emcee*, an ensemble, affine-invariant MCMC method (34,35). For each trajectory, 100 MCMC walkers were employed to draw 2000 samples each, and the first 1000 samples were discarded to burn in the chain. From the remaining samples, independent samples were chosen spaced apart by the maximum parameter autocorrelation time, and credible intervals and means were calculated from these samples. It took approximately 2 min to sample a single, 600-data point signal trajectory, such as those used in the computer-simulated titration.

## E<sub>FRET</sub> analysis of a ribosomal pretranslocation complex analog lacking a transfer RNA at the ribosomal aminoacyl-transfer RNA-binding site

Previously published Cy3 and Cy5 fluorescence intensity,  $I_{\text{Cy3}}$  and  $I_{\text{Cy5}}$ , versus time trajectories from a ribosomal pretranslocation (PRE) complex analog lacking a transfer RNA (tRNA) at the ribosomal aminoacyl-tRNA-binding (A) site (PRE<sup>-A</sup>) from the study by Wang and coworkers (9) were transformed into E<sub>FRET</sub> trajectories by calculating  $E_{\text{FRET}} = I_{\text{Cy5}} / (I_{\text{Cy3}} + I_{\text{Cy5}})$  at each measurement period. Outliers where  $E_{\text{FRET}} < -0.4$  or  $E_{\text{FRET}} > 1.4$  were clipped. The number of E<sub>FRET</sub> trajectories in the 22, 25, 28, 31, 34, and 37°C data sets were 490, 456, 435, 452, 270, and 459, respectively. Uniform distributions were used for the prior probability distributions. The first and second moments,  $E[k]$  and  $E[k^2]$ , of the marginalized posterior probability distributions for  $k_{\text{GS1}}$  or  $k_{\text{GS2}}$  were used to infer the values of  $\Delta H^\ddagger$ ,  $\Delta S^\ddagger$ , and a precision  $\lambda$  using Bayesian inference with the likelihood function  $p(\{E[k_i], E[k_i^2]\} | \Delta H^\ddagger, \Delta S^\ddagger, \lambda, \{T_i\}) = \prod \mathcal{N}(E[k_i] | \mu = k_{\text{TST}}(\Delta H^\ddagger, \Delta S^\ddagger, T_i), \sum = \lambda^{-1} + (E[k_i^2] - E[k_i]^2))$ , where  $\mathcal{N}$  is the normal distribution with mean  $\mu$  and variance  $\sum$ , where  $k_{\text{TST}}(\Delta H^\ddagger, \Delta S^\ddagger, T_i)$  is the rate constant calculated at temperature  $T_i$  with transition state theory, and where  $i$  indexes the set of temperatures. The resulting posterior probability distribution for  $\Delta H^\ddagger$ ,  $\Delta S^\ddagger$ , and  $\lambda$  was sampled using MCMC from which credible intervals and means were calculated.

## RESULTS AND DISCUSSION

### Analysis of computer-simulated single-molecule signal trajectories reporting on the kinetics of a ligand binding and dissociation process

To demonstrate the use of the analytical formulas underlying BIASD, we analyzed simulated single-molecule signal

trajectories that mimic the binding and dissociation of a ligand to its target biomolecule, a receptor, using the two-state, reversible kinetic scheme discussed in the previous section (36). In this example,  $\epsilon_1$  and  $\epsilon_2$  represent the signal values of the receptor in the ligand-free state and the ligand-bound state, respectively, and  $\sigma$  represents the SD of the signal values for both states. Correspondingly,  $k_1$  and  $k_2$  represent the pseudofirst-order rate constant of ligand binding to the receptor, and the first-order rate constant of ligand dissociation from the receptor, respectively. As such,  $k_1$  is dependent on  $[L]$  with a dependence that is given by  $k_1 = k_1^* \times [L]$ , where  $k_1^*$  is the second-order rate constant for binding of the ligand to the receptor,  $[L]$  is the ligand concentration, and  $k_2$  is not dependent on  $[L]$ . To emulate a titration experiment, we varied the  $[L]$  to alter the fraction of ligand-bound receptor from  $\sim 0.1$  to  $\sim 99.9\%$ , and simulated a series of individual signal trajectories where the  $[L]$  spanned six decades centered around the  $[L]$  corresponding to the equilibrium dissociation constant,  $[L] = K_D$ . Notably, as is always the case for experimentally recorded signal trajectories, the finite length of each simulated signal trajectory presents an intrinsic limit to the amount of kinetic information contained in each signal trajectory. Finally, estimates of the parameters that

were used to simulate the signal trajectories were obtained by analyzing the simulated signal trajectories using two idealization-based approaches: 1) half-amplitude signal thresholding (13) followed by Bayesian inference to infer the transition probability and quantify the kinetics (2) (Bayesian threshold), and 2) an HMM (6) that used the maximum-likelihood framework to estimate the transition probability and quantify the kinetics (ML HMM). In addition to these two idealization-based approaches, the simulated signal trajectories were also analyzed using the BIASD approach presented here.

As shown in Fig. 2 A, the values of  $k_1$  and  $k_2$  obtained using both idealization approaches are inaccurate (*green and red curves*). Interestingly, however, neither approach absolutely outperforms the other in accuracy, and both plateau at the acquisition rate (i.e.,  $\tau^{-1} = 10 \text{ s}^{-1}$ ). Notably, the use of Bayesian inference in the Bayesian threshold approach tempers the fluctuations that are seen in the rate constants obtained from the ML HMM at high  $[L]$ . These fluctuations originate from the maximum-likelihood framework used to estimate the transition probability in the ML HMM approach. Tempering of these fluctuations in the Bayesian threshold approach results from the use of a prior probability distribution, which describes the initial

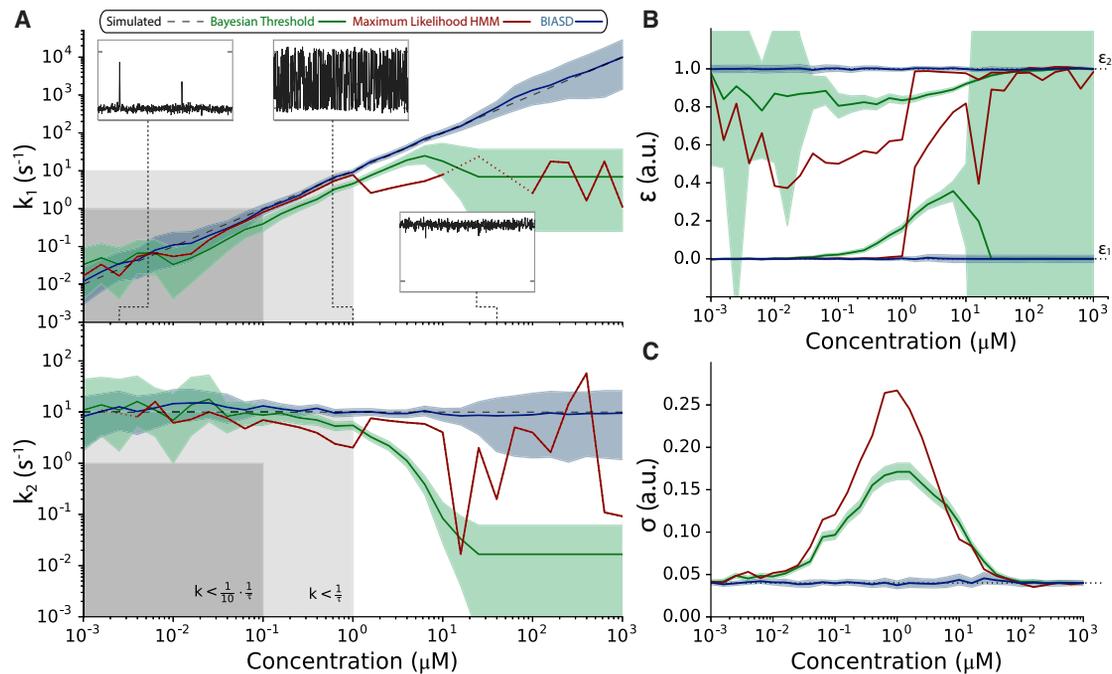


FIGURE 2 Analysis of  $k_1$  and  $k_2$  using BIASD (blue) and idealization-based (green and red) methods for a computer-simulated titration of a ligand to a receptor.  $[L]$  was varied three decades above and below the concentration where the equilibrium occupation probability of both states is equal to 0.5 (i.e.,  $K_D = 1$ ). (A) Analysis of estimated rate constants  $k_1$  and  $k_2$ . The regions where the rate constants are less than  $1/10$ th of the acquisition rate,  $\tau^{-1}$ , is shown in dark gray; the regions where the rate constants are less than the acquisition rate are shown in light gray. The simulated rate constants are plotted as the black dashed lines. The red line denotes the ML HMM estimate of the rate constants; dotted red lines indicate interpolated values due to transition probability estimates of unity. The green and blue areas denote the 95% credible intervals of the posterior probability distributions from analysis with half-amplitude thresholding-based Bayesian transition probability analysis (2,13) and BIASD, respectively. Insets show the simulated signal trajectory corresponding to the indicated concentration. (B) Analysis of estimated signal values  $\epsilon_1$  and  $\epsilon_2$ . Simulated values are plotted as black dashed lines. (C) Analysis of estimated signal noise  $\sigma$ . a.u., arbitrary units. To see this figure in color, go online.

knowledge of the model parameters, in Bayesian inference (2,6). Regardless, the rate constants were systematically underestimated across nearly the entire range of  $[L]$ s that were simulated, and this underestimation worsens with increasing  $[L]$ . It is striking that the values of  $k_1$  and  $k_2$  obtained using the Bayesian threshold idealization are also relatively precise, a misleading consequence of using idealization methods in general (2).

With regard to the values of  $\epsilon_1$  and  $\epsilon_2$  obtained using both idealization approaches, Fig. 2 B demonstrates that, whereas these methods can accurately determine the value of  $\epsilon_1$  if the receptor preferentially occupies the ligand-free state (low  $[L]$ ) or  $\epsilon_2$  if the receptor preferentially occupies the ligand-bound state (high  $[L]$ ), the time averaging caused by large values of  $k_1$  shift the inferred values of  $\epsilon_1$ , sometimes quite significantly, toward the simulated value of  $\epsilon_2$ , and vice versa. Here, the Bayesian threshold approach provided more accurate estimates of  $\epsilon_1$  and  $\epsilon_2$  than did the ML HMM. However, this was only because the signal threshold was set halfway between  $\epsilon_1$  and  $\epsilon_2$  using the known simulation parameters, thereby optimally minimizing the misclassification of states 1 and 2, and, consequently, maximizing the accuracy with which  $\epsilon_1$  and  $\epsilon_2$  are estimated; this is also why the estimates of  $\sigma$  from the Bayesian threshold approach are more accurate than those from the ML HMM (Fig. 2 C) (13).

In contrast to both idealization approaches, the values of  $k_1$  and  $k_2$  obtained using BIASD are highly accurate (Fig. 2 A). The simulated values of  $k_1$  and  $k_2$  are well encompassed by the 95% credible interval of the posterior probability distribution across the entire range of  $[L]$ s that were simulated, which includes rate constants that are up to three orders of magnitude larger than the simulated time resolution. In addition, the values of  $k_1$  and  $k_2$  are highly precise, as the 95% credible intervals of the posterior probability distribution are strikingly narrow over a range of  $[L]$ s corresponding to a value of  $k_1$  that is over an order of magnitude smaller than  $\tau^{-1}$  to one that is over an order of magnitude larger than  $\tau^{-1}$ . Importantly, as the amount of data that is analyzed increases, the contribution that the choice of prior probability distribution (i.e., the initial knowledge of  $k_1$ ,  $k_2$ ,  $\epsilon_1$ ,  $\epsilon_2$ , and  $\sigma$ ) makes to the posterior probability distribution diminishes. Consistent with the amount of data that are typically analyzed in single-molecule biophysical experiments, the results reported here are relatively insensitive to the prior probability distributions used for the analysis.

At the lower  $[L]$ s, the broadening of the posterior probability distribution that limits the precision for the estimates of  $k_1$  and  $k_2$  in both BIASD as well as the Bayesian threshold idealization arises from the finite amount of information regarding  $k_2$  and  $\epsilon_2$ , which is contained in signal trajectories that exhibit very low occupation of the ligand-bound state of the receptor. Likewise, at the higher  $[L]$ s, the broadening of the posterior probability distribution, and the implied limita-

tions to the precision for estimating  $k_1$  and  $k_2$  that is observed, arises from the finite amount of information regarding  $k_1$  and  $\epsilon_1$  that is contained in signal trajectories that exhibit very low occupation of the ligand-free state of the receptor. However, this broadening is somewhat attenuated, because the posterior probability distribution maintains consistency with the amount of previously known information about the underlying system contained within the prior probability distribution (c.f., *Analysis Using BIASD* in Supporting Material). Regardless, the uncertainty at the lower and higher  $[L]$ s is a consequence of the finite amount of information in a finite-length signal trajectory, as many reciprocal pairs of  $k_1$  and  $\epsilon_1$  values (i.e., a larger  $k_1$  and a smaller  $\epsilon_1$ , or a smaller  $k_1$  and a larger  $\epsilon_1$ ) are consistent with the data. In an experimental situation, this imprecision can be alleviated by employing prior probability distributions for the  $\{\epsilon\}$  values using the results of experiments performed under conditions in which one state is preferentially occupied, for instance, the values of  $\epsilon$  observed in the absence of ligand could be used to construct a prior probability distribution for the values of  $\epsilon$  associated with the ligand-free state of the receptor, whereas the values of  $\epsilon$  observed in the presence of saturating  $[L]$  could be used to construct a prior probability distribution for the values of  $\epsilon$  associated with the ligand-bound state of the receptor. In the case of large-scale conformational rearrangements, one could similarly use a buffer condition, ligand, temperature, or mutation that preferentially stabilizes one state, or, alternatively, one could use molecular structures or models to estimate prior probability distributions of  $\{\epsilon\}$  values. With regard to the values of  $\epsilon_1$  and  $\epsilon_2$  obtained using BIASD, Fig. 2 B demonstrates that these values were accurately inferred regardless of the value of  $[L]$ , even at  $[L]$ s at which the idealization approaches drastically misestimate them. Finally, unlike the idealization approaches, which were only able to successfully infer  $\sigma$  when the signal trajectories were almost entirely in the ligand-bound or ligand-free states, BIASD was also able to accurately and precisely infer  $\sigma$  from the simulated signal trajectories with intermediate values of  $[L]$  (Fig. 2 C).

In summary, we were able to use BIASD to obtain accurate and precise posterior probability distributions for  $k_1$ ,  $k_2$ ,  $\epsilon_1$ ,  $\epsilon_2$ , and  $\sigma$  across the entire range of  $[L]$ s that were simulated. Notably, BIASD was even successful when the rate constants in the simulated, single-molecule signal trajectories were much smaller than  $\tau^{-1}$ , although we note that, in this regime, the conventional analysis of idealizing the signal trajectories is much more computationally efficient. Most importantly, BIASD was able to accurately and precisely infer the rate constants and the signal values for simulated, single-molecule signal trajectories in which the rate constants were up to three orders of magnitude larger than  $\tau^{-1}$ , and up to about four orders of magnitude larger than the  $\tau^{-1}$ s where conventional idealization of signal trajectories begins to yield significant errors in the rate constants.

This finding is consistent with the expected performance of BIASD (c.f., *Dependence of BIASD Performance on Parameter Values* in [Supporting Material](#)).

### Analysis of experimentally observed single-molecule $E_{\text{FRET}}$ trajectories reporting on the kinetics of a large-scale conformational rearrangement

To demonstrate the use of BIASD in the analysis of experimental data, we chose to analyze experimentally observed, single-molecule  $E_{\text{FRET}}$  trajectories reporting on a large-scale conformational rearrangement of the ribosome. This essential, two-subunit, ribonucleoprotein-based biomolecular machine is universally responsible for the translation of messenger RNAs into proteins in living cells. The ribosome synthesizes proteins by repeatedly incorporating amino acids, delivered in the form of aminoacyl-tRNA substrates, into a nascent polypeptide chain in the order dictated by the messenger RNA being translated. During the elongation stage of protein synthesis (38), the ribosomal PRE complex undergoes stochastic, thermally driven fluctuations between two major, on-pathway conformational states that we refer to as global state 1 (GS1) and global state 2 (GS2), defining a dynamic equilibrium,  $\text{GS1} \rightleftharpoons \text{GS2}$  (39,40). These transitions between GS1 and GS2 constitute large-scale rearrangements of the PRE complex that involve relative rotations of the ribosomal subunits, reconfigurations of the ribosome-bound tRNAs, and repositioning of a ribosomal structural domain known as the L1 stalk (Fig. 3 A) (41).

Previously, we have conducted wide-field microscopy single-molecule fluorescence resonance energy transfer (smFRET) studies of the temperature dependence of the rate constants governing  $\text{GS1} \rightarrow \text{GS2}$  and  $\text{GS2} \rightarrow \text{GS1}$  transitions by imaging a Cy3 fluorescence resonance energy transfer (FRET) donor fluorophore- and Cy5 FRET acceptor fluorophore-labeled  $\text{PRE}^{-\text{A}}$  in a temperature-controlled, microfluidic observation flowcell (9). In these experiments, the increase in thermal energy that accompanied the increasing temperature caused the rate constants for the transitions between GS1 and GS2 to increase such that, at the highest temperatures, the  $E_{\text{FRET}}$  trajectories contained a significant number of time-averaged data points at the  $\tau = 50$  ms time resolution of the experiment (Fig. S2). Regrettably, the time averaging in these  $E_{\text{FRET}}$  trajectories precluded accurate determination of the rate constants and, correspondingly, an analysis of the thermodynamic properties of the transition-state energy barriers that control the  $\text{GS1} \rightarrow \text{GS2}$  and  $\text{GS2} \rightarrow \text{GS1}$  conformational rearrangements (9). To overcome these limitations, we have used BIASD to analyze the sets of  $E_{\text{FRET}}$  trajectories of  $\text{PRE}^{-\text{A}}$  complexes that we have previously collected at 22, 25, 28, 31, 34, and 37°C (9). Here, we assume that the  $\text{GS1} \rightleftharpoons \text{GS2}$  equilibrium can be represented by a single, reversible

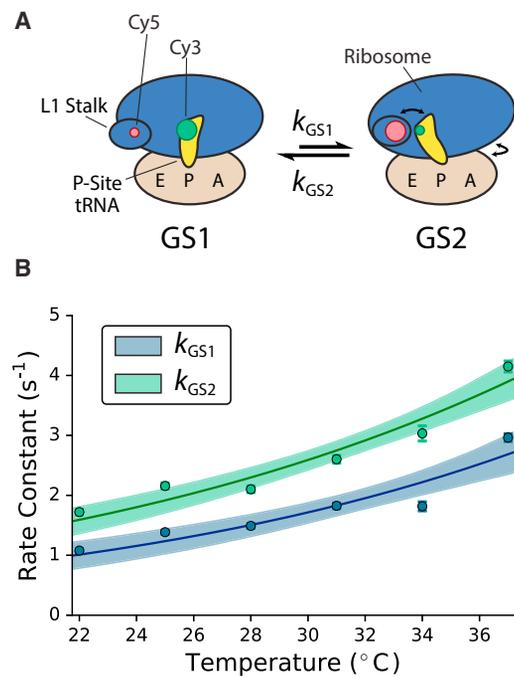


FIGURE 3 (A) Cartoon schematic of the  $\text{GS1} \rightleftharpoons \text{GS2}$  equilibrium on the  $\text{PRE}^{-\text{A}}$  complex previously studied by Wang et al. (9). Approximate positions of the Cy3 FRET donor and Cy5 FRET acceptor fluorophores of the “L1-tRNA” labeling scheme used by Wang et al. are shown as green and red circles, respectively. The size of the fluorophores denotes the relative fluorescence intensity of each fluorophore in each state due to FRET. A, P, and E denote the aminoacyl-tRNA-binding, peptidyl-tRNA-binding, and exit sites of the ribosome, respectively. (B) Temperature dependence of  $k_{\text{GS1}}$  and  $k_{\text{GS2}}$  for the  $\text{PRE}^{-\text{A}}$  complexes using BIASD. The scatter plots show the expectation value of the posterior probability distributions of  $k_{\text{GS1}}$  and  $k_{\text{GS2}}$ , and the error bars represent the 95% credible interval. The solid lines denote expectation values, and the shaded regions denote the 95% credible interval of the predictive posterior probability distribution from the transition-state theory analysis. To see this figure in color, go online.

two-state kinetic scheme (Fig. 3 A). In this kinetic scheme,  $k_{\text{GS1}}$  and  $k_{\text{GS2}}$  represent the unimolecular rate constants for the  $\text{GS1} \rightarrow \text{GS2}$  and  $\text{GS2} \rightarrow \text{GS1}$  conformational rearrangements, respectively. Correspondingly,  $\epsilon_{\text{GS1}}$  and  $\epsilon_{\text{GS2}}$  represent the  $E_{\text{FRET}}$  values of GS1 and GS2, respectively.

These six sets of  $E_{\text{FRET}}$  trajectories were analyzed using BIASD to provide estimates of  $k_{\text{GS1}}$ ,  $k_{\text{GS2}}$ ,  $\epsilon_{\text{GS1}}$ ,  $\epsilon_{\text{GS2}}$ , and  $\sigma$  that best describe the entire set of  $E_{\text{FRET}}$  trajectories observed at each temperature. The values of  $k_{\text{GS1}}$  and  $k_{\text{GS2}}$  that were inferred using BIASD increase with temperature (Fig. 3 B) and, at the highest temperatures, were greater than  $1/10$ th of  $\tau^{-1}$  (i.e.,  $2 \text{ s}^{-1}$ ), the regime where idealization approaches begin to systematically underestimate rate constants. We note that although the values of  $k_{\text{GS1}}$  and  $k_{\text{GS2}}$  inferred using BIASD are those that best describe the entire set of  $E_{\text{FRET}}$  trajectories observed at a particular temperature, inspection of individual  $E_{\text{FRET}}$  trajectories suggests the presence of kinetic heterogeneity, as some are consistent with rate constants  $> 45 \text{ s}^{-1}$ , whereas others are consistent with rate constants  $< 0.1 \text{ s}^{-1}$ . This broad

range of kinetic behaviors suggests the possibility that the  $\text{PRE}^{-\text{A}}$  complexes are compositionally heterogeneous (e.g., subpopulations of  $\text{PRE}^{-\text{A}}$  complexes that differ in the aminoacylation status of the tRNA at the ribosomal peptidyl-tRNA-binding (P) site, the presence or absence of a tRNA at the ribosomal tRNA exit site, and/or the presence or absence of a particular ribosomal protein) and/or are conformationally heterogeneous due to structural rearrangements that are slow on the timescale of the experiment and effect  $k_{\text{GS1}}$  and  $k_{\text{GS2}}$ , but not necessarily  $\epsilon_{\text{GS1}}$  and  $\epsilon_{\text{GS2}}$ . Additionally, we note that the posterior probability distributions of  $\epsilon_{\text{GS1}}$  and  $\epsilon_{\text{GS2}}$  that were inferred using BIASD have means of 0.13 and 0.78, respectively, which are values of  $\epsilon_{\text{GS1}}$  and  $\epsilon_{\text{GS2}}$  that very closely match the values of the mean  $E_{\text{FRET}}$  of GS1 and GS2 reported in previous, room-temperature studies of the analogous  $\text{PRE}^{-\text{A}}$  complex (0.16 and 0.76, respectively) (42). This correspondence strongly suggests that the values of  $\epsilon_{\text{GS1}}$  and  $\epsilon_{\text{GS2}}$  inferred using BIASD are accurate, regardless of the presence of time averaging in the  $E_{\text{FRET}}$  trajectories recorded at the highest temperatures.

With the inferred values of  $k_{\text{GS1}}$  and  $k_{\text{GS2}}$  as a function of temperature, we then used transition-state theory to quantify the apparent transition-state energy barriers along the apparent GS1  $\rightarrow$  GS2 and GS2  $\rightarrow$  GS1 reaction coordinates (43–46). Kramers' barrier-crossing theory, which was developed to analyze thermally activated, condensed-phase transitions of a Brownian particle (44–46) and is increasingly being used to analyze the conformational dynamics and folding of small, globular proteins (12,47), may ultimately provide a more exact analysis of the apparent transition-state energy barriers along the apparent GS1  $\rightarrow$  GS2 and GS2  $\rightarrow$  GS1 reaction coordinates. However, its application requires knowledge regarding the viscosity of the aqueous buffer in which the  $\text{PRE}^{-\text{A}}$  complex is dissolved and the "internal friction" of the  $\text{PRE}^{-\text{A}}$  complex, which are unavailable in the current study (12,48). As such, we have opted to use transition-state theory, and regard the results as upper limits on the apparent transition-state energy barriers along the apparent GS1  $\rightarrow$  GS2 and GS2  $\rightarrow$  GS1 reaction coordinates, which do not account for internal friction or transition-state recrossings. To apply transition-state theory, we used the marginalized posterior probability distributions of the rate constants at each temperature to infer  $\Delta H^\ddagger$  and  $\Delta S^\ddagger$  from the equation  $k_{\text{TST}} = \kappa k_{\text{B}}T/h \exp(-(\Delta H^\ddagger - T\Delta S^\ddagger)/(k_{\text{B}}T))$ , where  $\kappa$  is the transmission coefficient and is taken to be unity,  $k_{\text{B}}$  is the Boltzmann constant,  $h$  is Planck's constant, and  $\Delta H^\ddagger$  and  $\Delta S^\ddagger$  are the enthalpic and entropic differences between the transition and ground states (Fig. 3 B). The marginalized results for the GS1  $\rightarrow$  GS2 transition are  $\Delta H_{\text{GS1}}^\ddagger = 11.3$  (7.9, 15.0) kcal mol $^{-1}$  and  $\Delta S_{\text{GS1}}^\ddagger = -20.1$  (-31.3, -8.1) cal mol $^{-1}$  K $^{-1}$ , and for the GS2  $\rightarrow$  GS1 transition are  $\Delta H_{\text{GS2}}^\ddagger = 10.3$  (8.1, 12.7) kcal mol $^{-1}$  and

$\Delta S_{\text{GS2}}^\ddagger = -22.7$  (-29.9, -15.0) cal mol $^{-1}$  K $^{-1}$ , where the numbers in parentheses represent the lower and upper bounds for the 95% credible interval. Notably, the posterior probability distributions for  $\Delta H^\ddagger$  and  $\Delta S^\ddagger$  are highly correlated such that across all of the temperatures measured here,  $\Delta G_{\text{GS1}}^\ddagger$  and  $\Delta G_{\text{GS2}}^\ddagger$  were sufficiently resolved; for instance, at 37°C,  $\Delta G_{\text{GS1}}^\ddagger = 17.57$  (17.50, 17.65) kcal mol $^{-1}$ , and  $\Delta G_{\text{GS2}}^\ddagger = 17.34$  (17.29, 17.39) kcal mol $^{-1}$ . Structure-based interpretation of the absolute  $\Delta H^\ddagger$  and  $\Delta S^\ddagger$  values for the GS1  $\rightarrow$  GS2 and GS2  $\rightarrow$  GS1 transitions of a single  $\text{PRE}^{-\text{A}}$  complex is significantly complicated by the complexity of the enthalpic and entropic changes that are associated with conformational rearrangements of large macromolecular complexes, and by the inherent limitations of transition-state theory (5,43,49). Nonetheless, structure-based interpretations of the relative changes of the  $\Delta H^\ddagger$ s and  $\Delta S^\ddagger$ s ( $\Delta\Delta H^\ddagger$ s and  $\Delta\Delta S^\ddagger$ s) between different pairs of  $\text{PRE}^{-\text{A}}$  complexes (e.g., containing different tRNAs at the P site, containing wild-type or mutant P site tRNAs, or consisting of wildtype or mutant ribosomes, etc.) are much more straightforward and can reveal the thermodynamic contributions that particular structural features of tRNAs or ribosomes make to the apparent transition-state energy barriers along the apparent GS1  $\rightarrow$  GS2 and GS2  $\rightarrow$  GS1 reaction coordinates. Combined with the temperature-controlled, single-molecule microscopy platform that we have previously described (9), the analytical framework presented in this section now enables the collection, analysis, and interpretation of such data.

### Inferring rate constants and signal values from systems with subpopulations of molecular properties

BIASD can be extended to address the presence of multiple, time-averaged subpopulations of molecular properties. These subpopulations may be static or interconvert, and may be present in an individual molecule or found among an ensemble of molecules. In such a situation, we can classify each data point as belonging to one of  $K$  different types of time-averaged subpopulations, and then use a "1-of- $K$ " vector,  $\bar{z}_{ij}$ , to denote to which of the  $K$  subpopulations the  $i^{\text{th}}$  data point from the  $j^{\text{th}}$  molecule belongs. Given the one particular subpopulation specified by  $\bar{z}_{ij}$ , the likelihood of this data point being described by the parameters of this subpopulation is calculated as described above for the case of the time-averaged, single-population system. Unfortunately, in an experimental situation, there is no way of knowing which subpopulation a particular data point belongs to, thereby preventing the likelihood of this data point from being evaluated; this situation is similar to that of the unknown fractional occupancy,  $f$ , described above.

To address this shortcoming, we could try to infer the values of all the  $\vec{z}_{ij}$  along with all of the other BIASD model parameters, but this is an unreasonable number of variables for an inference procedure. Additionally, we are often not concerned with the exact values of  $\vec{z}_{ij}$ , so much as with the occupancies of the  $K$  states (e.g., the steady-state occupation probabilities) or with the rate constants that describe transitions between the  $K$  states. Fortunately, instead of performing inference to learn the model parameters and the set of  $\vec{z}_{ij}$ s,  $\{\vec{z}_{ij}\}$ , we can marginalize out all of the  $\{\vec{z}_{ij}\}$  with the expressions for the probability of each  $\vec{z}_{ij}$ . For instance, in the case of a mixture of static subpopulations of molecular properties among an ensemble of molecules (e.g., a mixture of posttranscriptionally or posttranslationally modified and unmodified molecules within an ensemble), these probabilities would be time-independent variables that specify the fraction of each subpopulation of the ensemble; this approach is called a mixture model. Marginalization would then involve summing the likelihoods for the different subpopulations, weighted by the probabilities of those subpopulations. Consequently, during the inference procedure, the probabilities of the subpopulation occupancies would then become model parameters that are also inferred using Bayes' rule.

Additionally, it is possible to have a time-dependent system with hierarchical transitions between the different subpopulations. In this case, the probabilities of each  $\vec{z}_{ij}$  in the  $\{\vec{z}_{ij}\}$  would not be constant for each subpopulation, as they would be for a mixture model, but would instead depend upon the subpopulation of the previous data point

$\vec{z}_{i-1,j}$  and a  $K \times K$  transition matrix,  $\mathbf{A}_{ij} = e^{\mathbf{Q}t_{ij}}$ , where  $\mathbf{Q}$  is the rate matrix that depends upon the set of rate constants for transitioning between the  $K$  different states, and  $t_{ij}$  is the time that has elapsed since the previous data point, which may not necessarily be equal to  $\tau$  (Fig. S1 B). Here, marginalization is efficiently performed with the forward-backward algorithm (37) and the steady-state probabilities, as calculated from the rate constants for the kinetic scheme under consideration, for instance by using the diagram method (50), are used to set the initial probability of each  $\vec{z}_{0j}$ . In total, this approach amounts to a hierarchical, continuous-time ensemble HMM for subtemporal resolution systems, where inference is performed directly upon the rate constants, instead of the transition probabilities. Consequently, this approach can handle shuttering of the laser light source in fluorescence microscopy experiments or other types of irregular spacing of data points, subtemporal resolution data, and population-level analyses with nonparametric posterior distributions, which can be used to ascertain the underlying thermodynamic landscape of the mesoscopic ensemble.

To highlight this hierarchical approach, consider a single-molecule fluorescence microscopy experiment in which a fluorophore-labeled biomolecule transitions between two states, 1 and 2, with forward and reverse rate constants  $k_{12}$  and  $k_{21}$ , respectively (Fig. 4 A). Such fluorescence microscopy experiments often suffer from photophysical phenomena such as fluorophore photoblinking, in which a fluorophore temporarily transitions into a long-lived, "dark," excited molecular electronic state and thus transiently stops fluorescing, or

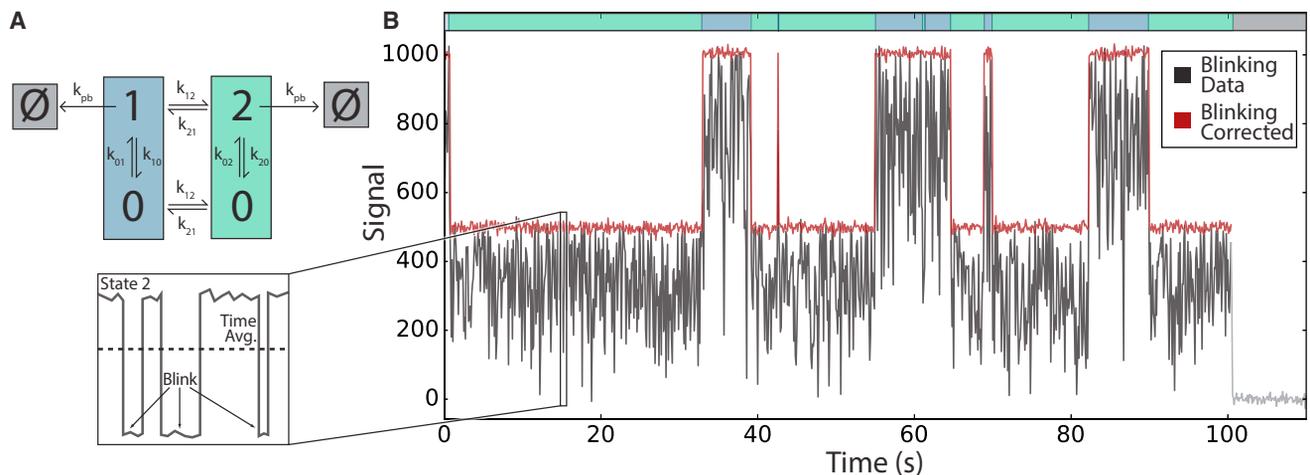


FIGURE 4 (A) Kinetic mechanism used to simulate an experimental system in which a biomolecule transitions between two conformational states that can both transition, with subtemporal resolution dynamics, into and out of a photoblinked state, until eventually photobleaching. Exact values of the rate constants used in the simulation are given in the Supporting Material. (B) Plot of a simulated signal trajectory and an estimated, corrected fluorescence intensity signal trajectory in the absence of photoblinking. The green- and blue-colored regions at the top of the plot denote the time spent in states 1 (blue) and 2 (green). The signal corrupted by subtemporal resolution photoblinking that was analyzed by BIASD is shown in black; the inset shows a cartoon of a single  $\tau$  where fast transitions are taking place between the fluorescent and photoblinked state 2. After analysis with BIASD, the maximum *a posteriori* (MAP) solution of the posterior probability distribution was used to generate a Viterbi-idealized path, which is plotted in red with noise added back from the MAP value of  $\sigma$ ; this is a plausible corrected fluorescence intensity signal trajectory in the absence of photoblinking. Although one excursion to state 2 is not present in this Viterbi path, the forward-backward algorithm used in the HMM analysis takes it into consideration. To see this figure in color, go online.

fluorophore photobleaching, in which a fluorophore that has transitioned into an excited molecular electronic state undergoes a photochemical reaction and permanently stops fluorescing (51). Often, the transition rates into and out of the dark states responsible for photoblinking are faster than the time resolution of techniques such as wide-field, fluorescence microscopy. As a result, instead of detecting a steady level of fluorescence intensity from the fluorophore, subtemporal resolution transitions between fluorescent and dark states of the fluorophore manifest as an extra, and often dominant, source of “noise” in the single-molecule fluorescence intensity signal trajectory (Fig. 4 B, inset). Intense experimental effort has gone into minimizing these photo-physical effects, including the use of fluorophores, such as Cy3B, that have been chemically altered so as to minimize transitions to dark states (52); elaborate excitation laser modulation schemes, such as triplet-state relaxation and dark-state relaxation schemes, which minimize transitions to higher-order dark states (53,54); photostabilizing additives, such as Trolox, that accelerate transitions out of dark states (55,56); and fluorophore-photostabilizer conjugates, such as Cy3- and Cy5-triplet-state quencher conjugates, that accelerate transitions out of dark states (57,58). Here, we show how extending BIASD with the hierarchical HMM described above allows us to computationally overcome these photo-physical effects.

To demonstrate this ability, we simulated the kinetic scheme shown in Fig. 4 A, where a fluorophore-labeled biomolecule transitions between conformational states 1 and 2 with signal values of  $\epsilon_1$  and  $\epsilon_2$ , respectively. However, in this simulation, both of these states can rapidly transition into and out of a photoblinded state, denoted 0 with signal value  $\epsilon_0 = 0$ , at rates much faster than the time resolution of the simulated data. These dynamics continue until the system eventually transitions into a photobleached state, denoted  $\emptyset$  with signal value  $\epsilon_{\emptyset} = 0$ . A similar situation has been recently investigated by Chung et al. (20) to analyze FRET photon trajectories reporting on the subtemporal resolution folding and unfolding dynamics, and photoblinking dynamics, of the villin subdomain protein. As expected, by analyzing this simulation using this hierarchical approach, the posterior probability distribution of the parameters describing the fluorescence emission from each subpopulation ( $\epsilon_1$ ,  $\epsilon_2$ ,  $\sigma$ ,  $k_{10}$ ,  $k_{01}$ ,  $k_{20}$ , and  $k_{02}$ ; see Fig. S1 B), as well as the rate constants describing the transitions between states 1 and 2 ( $k_{12}$ , and  $k_{21}$ ; see Fig. 4 A), were all found both accurately and precisely, as the parameter values used for the simulation fall within the inferred 95% credible intervals (Fig. S3). To provide visual intuition into this result, we also have shown the Viterbi-idealized path from the maximum *a posteriori* estimate of the model parameters to show the most likely fluorescence intensity signal trajectory in the absence of photoblinking (Fig. 4 B). Detection noise from the marginalized posterior distribution of  $\sigma$  was added to this path to show what the data might have

looked like in the absence of photoblinking. Regardless, we note that this particular path is essentially a point estimate of the  $\{\bar{z}_{ij}\}$ , whereas, by marginalizing out all of the  $\{z_{ij}\}$  during the inference procedure, we have actually considered all the other possible paths, regardless of whether or not a transition is missed in the Viterbi path. As such, the posterior probability distribution of the model parameters is a more encompassing result (Fig. S3). Finally, we note that the hierarchical HMM treatment that we present here is general and applicable to not just two, but to any number of K subpopulations.

## CONCLUSIONS

By analyzing the fraction of time that a single molecule spends in each state of a defined kinetic scheme during each  $\tau$  in a signal trajectory, BIASD adopts a fundamentally different approach to the analysis of time-resolved single-molecule experiments than that which has been traditionally employed by methods that idealize the trajectories (e.g., signal thresholding, HMMs, etc.). Using computer-simulated and experimentally observed data, we have demonstrated that this powerful approach enables BIASD to accurately and precisely infer the rate constants of a two-state kinetic scheme as well as the signal values corresponding to these two states, even when the rates of transitions between the states are orders of magnitude larger than the time resolution of the signal trajectories. When used to analyze experimental  $E_{\text{FRET}}$  trajectories reporting on the dynamics of single  $\text{PRE}^{-\text{A}}$  complexes recorded as a function of temperature (9), BIASD allowed us to infer the thermodynamic activation parameters characterizing the transition-state energy barriers along the  $\text{GS1} \rightarrow \text{GS2}$  and  $\text{GS2} \rightarrow \text{GS1}$  reaction coordinates, which had thus far remained inaccessible to traditional smFRET data analysis approaches. Moreover, we have demonstrated that a straightforward extension of the BIASD framework enables the kinetics of experimental systems exhibiting multiple subpopulations of molecular properties to be accurately and precisely inferred.

It is important to note that the BIASD framework is general and can be applied to any experimentally observed signal trajectory that exhibits stochastic transitions between distinct states, regardless of the nature or the origin of the signal. Thus, BIASD can be used to temporally resolve data collected using virtually any time-resolved single-molecule experimental method, including single-molecule fluorescence microscopy, force spectroscopy, conductance, and tethered particle motion methods. Moreover, although here we have developed BIASD to analyze single-molecule signal trajectories, we have not considered the temporal ordering of the data. Consequently, in addition to analyzing individual single-molecule signal trajectories, BIASD can also be used to analyze the distribution of fractional occupancies observed across an entire ensemble of individual

molecules during a given  $\tau$ . This could allow nonequilibrium phenomena to be monitored across an ensemble of single molecules, time period by time period (e.g., stopped-flow delivery of a ligand, substrate, cofactor, or inhibitor to an enzyme or other biomolecule). In addition, BIASD can be expanded to include the time evolution of the state occupation probabilities (c.f., Eq. S2), or to incorporate time dependence into the model parameters  $\{k\}$ ,  $\{\epsilon\}$ , and  $\{\sigma\}$  (e.g., the varying of  $\{\epsilon\}$  in single-molecule particle tracking experiments).

Regarding the performance of BIASD on experimental data, we note that the rate constants and signal values of a system can be more precisely inferred from experiments that collect higher SNR data, because then there is less uncertainty in the time-averaged fractional occupancies of the signal trajectories. Therefore, somewhat counterintuitively, subtemporal resolution dynamics can to some degree be more precisely inferred from signal trajectories recorded with lower time resolutions but higher SNRs (e.g., due to better photon conversion efficiencies on an electron-multiplying charge-coupled device), than those recorded with higher time resolutions but lower SNRs. Additionally, although we have focused the current work on the most widely applicable case of a Markovian, two-state system in which the noise of the signal can be modeled using a normal distribution, the Bayesian inference-based framework underlying BIASD can be readily extended to non-Markovian dynamics (21,59), N-state kinetic schemes (60,61), or systems in which the noise of the signal can be modeled using distributions other than a Normal distribution (18,62). However, it should be noted that such developments will come with added computational expenses. To facilitate the analysis of single-molecule data using BIASD, as well as to enable the future extension of BIASD along the lines described here, we have made the BIASD source code available at <http://github.com/ckinzthompson/biasd>. The source code is written in Python and integrated with computationally intensive functions provided in C as well as in CUDA (for GPU-based computation), to balance accessibility with high performance.

## SUPPORTING MATERIAL

Supporting Materials and Methods, five figures, and two tables are available at [http://www.biophysj.org/biophysj/supplemental/S0006-3495\(17\)34973-1](http://www.biophysj.org/biophysj/supplemental/S0006-3495(17)34973-1).

## AUTHOR CONTRIBUTIONS

C.D.K-T. conceived of and performed the research. C.D.K-T. and R.L.G. discussed results and wrote the manuscript.

## ACKNOWLEDGMENTS

The authors thank Prof. Jan-Willem van de Meent for his comments on an early version of this manuscript.

This work was supported by two National Institutes of Health (NIH)-National Institute of General Medical Sciences grants (R01 GM084288 and R01 GM 119386), an American Cancer Society Research Scholar grant (RSG GMC-117152), and a Camille Dreyfus Teacher-Scholar Award (DRFSCH CU11-0665) to R.L.G. C.D.K-T. was supported by the Department of Energy Office of Science Graduate Fellowship Program—made possible in part by the American Recovery and Reinvestment Act of 2009, administered by the Oak Ridge Institute for Science and Education-Oak Ridge Associated Universities under contract number DE-AC05-06OR23100—and by Columbia University's NIH Training Program in Molecular Biophysics (T32-GM008281).

## REFERENCES

1. Tinoco, I., Jr., and R. L. Gonzalez, Jr. 2011. Biological mechanisms, one molecule at a time. *Genes Dev.* 25:1205–1231.
2. Kinz-Thompson, C. D., N. A. Bailey, and R. L. Gonzalez, Jr. 2016. Precisely and accurately inferring single-molecule rate constants. *Methods Enzymol.* 581:187–225.
3. Cavanagh, J., W. J. Fairbrother, ..., N. J. Skelton. 1996. Protein NMR Spectroscopy: Principles and Practice. Elsevier Academic Press, San Diego, CA.
4. Boehr, D. D., H. J. Dyson, and P. E. Wright. 2006. An NMR perspective on enzyme dynamics. *Chem. Rev.* 106:3055–3079.
5. McCammon, J. A. 1984. Protein dynamics. *Rep. Prog. Phys.* 47:1–46.
6. Bishop, C. M. 2006. Pattern Recognition and Machine Learning. Springer, New York.
7. Jaynes, E. T. 2003. Probability Theory: The Logic of Science. Cambridge University Press, Cambridge, UK.
8. Bretthorst, G. L., C. C. Hung, ..., J. J. H. Ackerman. 1988. Bayesian analysis of time-domain magnetic resonance signals. *J. Magn. Reson.* 79:369–376.
9. Wang, B., J. Ho, ..., Q. Lin. 2011. A microfluidic approach for investigating the temperature dependence of biomolecular activity with single-molecule resolution. *Lab Chip.* 11:274–281.
10. Fenimore, P. W., H. Frauenfelder, ..., F. G. Parak. 2002. Slaving: solvent fluctuations dominate protein dynamics and functions. *Proc. Natl. Acad. Sci. USA.* 99:16047–16051.
11. Lubchenko, V., P. G. Wolynes, and H. Frauenfelder. 2005. Mosaic energy landscapes of liquids and the control of protein conformational dynamics by glass-forming solvents. *J. Phys. Chem. B.* 109:7488–7499.
12. Chung, H. S., and W. A. Eaton. 2013. Single-molecule fluorescence probes dynamics of barrier crossing. *Nature.* 502:685–688.
13. Colquhoun, D., and F. J. Sigworth. 1995. Fitting and statistical analysis of single channel records. In *Single Channel Recording*. Springer, Boston, MA, pp. 483–587.
14. Chung, S. H., J. B. Moore, ..., P. W. Gage. 1990. Characterization of single channel currents using digital signal processing techniques based on Hidden Markov Models. *Philos. Trans. R. Soc. Lond. B Biol. Sci.* 329:265–285.
15. Qin, F., A. Auerbach, and F. Sachs. 2000. A direct optimization approach to hidden Markov modeling for single channel kinetics. *Biophys. J.* 79:1915–1927.
16. Hanson, J. A., K. Duderstadt, ..., H. Yang. 2007. Illuminating the mechanistic roles of enzyme conformational dynamics. *Proc. Natl. Acad. Sci. USA.* 104:18055–18060.
17. Nir, E., X. Michalet, ..., S. Weiss. 2006. Shot-noise limited single-molecule FRET histograms: comparison between theory and experiments. *J. Phys. Chem. B.* 110:22103–22124.
18. Gopich, I. V., and A. Szabo. 2012. Theory of the energy transfer efficiency and fluorescence lifetime distribution in single-molecule FRET. *Proc. Natl. Acad. Sci. USA.* 109:7747–7752.

19. Kalinin, S., E. Sisamakias, ..., C. A. M. Seidel. 2010. On the origin of broadening of single-molecule FRET efficiency distributions beyond shot noise limits. *J. Phys. Chem. B.* 114:6197–6206.
20. Chung, H. S., T. Cellmer, ..., W. A. Eaton. 2013. Measuring ultrafast protein folding rates from photon-by-photon analysis of single molecule fluorescence trajectories. *Chem. Phys.* 422:229–237.
21. Berezhkovskii, A. M., A. Szabo, and G. H. Weiss. 1999. Theory of single-molecule fluorescence spectroscopy of two-state systems. *J. Chem. Phys.* 110:9145–9150.
22. Flomenbom, O., and R. J. Silbey. 2007. Properties of the generalized master equation: Green's functions and probability density functions in the path representation. *J. Chem. Phys.* 127:034103.
23. Anderson, P. W. 1954. A mathematical model for the narrowing of spectral lines by exchange or motion. *J. Phys. Soc. Jpn.* 9:316–339.
24. Good, I. 1961. The frequency count of a Markov chain and the transition to continuous time. *Ann. Math. Stat.* 32:41–48.
25. Dobrushin, R. 1953. Limit theorems for a Markov chain of two states. *Izv. Ross. Akad. Nauk. USSR Seriya Mat.* 17:291–330.
26. Weiss, G. H. 1976. The two-state random walk. *J. Stat. Phys.* 15:157–165.
27. Gilmore, C. J. 1996. Maximum entropy and Bayesian statistics in crystallography: a review of practical applications. *Acta Crystallogr. A.* 52:561–589.
28. Ronquist, F., and J. P. Huelsenbeck. 2003. MrBayes 3: Bayesian phylogenetic inference under mixed models. *Bioinformatics.* 19:1572–1574.
29. Scheres, S. H. W. 2012. RELION: implementation of a Bayesian approach to cryo-EM structure determination. *J. Struct. Biol.* 180:519–530.
30. Bronson, J. E., J. Fei, ..., C. H. Wiggins. 2009. Learning rates and states from biophysical time series: a Bayesian approach to model selection and single-molecule FRET data. *Biophys. J.* 97:3196–3205.
31. van de Meent, J.-W., J. E. Bronson, ..., R. L. Gonzalez, Jr. 2014. Empirical Bayes methods enable advanced population-level analyses of single-molecule FRET experiments. *Biophys. J.* 106:1327–1337.
32. Sivia, D. S., and J. Skilling. 2006. *Data Analysis: A Bayesian Tutorial.* Oxford University Press, Oxford, UK.
33. Metropolis, N., A. W. Rosenbluth, ..., E. Teller. 1953. Equation of state calculations by fast computing machines. *J. Chem. Phys.* 21:1087.
34. Goodman, J., and J. Weare. 2010. Ensemble samplers with affine invariance. *Comm. App. Math. Comp. Sci.* 5:65–80.
35. Foreman-Mackey, D., D. W. Hogg, ..., J. Goodman. 2013. emcee: the MCMC hammer. *Publ. Astron. Soc. Pac.* 125:306–312.
36. Gillespie, D. T. 1977. Exact stochastic simulation of coupled chemical reactions. *J. Phys. Chem.* 81:2340–2361.
37. Rabiner, L. R. 1989. A tutorial on hidden Markov models and selected applications in speech recognition. *Proc. IEEE.* 77:257–286.
38. Voorhees, R. M., and V. Ramakrishnan. 2013. Structural basis of the translational elongation cycle. *Annu. Rev. Biochem.* 82:203–236.
39. Fei, J., P. Kosuri, ..., R. L. Gonzalez, Jr. 2008. Coupling of ribosomal L1 stalk and tRNA dynamics during translation elongation. *Mol. Cell.* 30:348–359.
40. Fei, J., J. E. Bronson, ..., R. L. Gonzalez, Jr. 2009. Allosteric collaboration between elongation factor G and the ribosomal L1 stalk directs tRNA movements during translation. *Proc. Natl. Acad. Sci. USA.* 106:15702–15707.
41. Frank, J. 2012. Intermediate states during mRNA-tRNA translocation. *Curr. Opin. Struct. Biol.* 22:778–785.
42. Sternberg, S. H., J. Fei, ..., R. L. Gonzalez, Jr. 2009. Translation factors direct intrinsic ribosome dynamics during translation termination and ribosome recycling. *Nat. Struct. Mol. Biol.* 16:861–868.
43. Fersht, A. R. 1999. *Structure and Mechanism in Protein Science. A Guide to Enzyme Catalysis and Protein Folding.* W.H. Freeman and Co., Ltd., New York.
44. Zwanzig, R. 2001. *Nonequilibrium Statistical Mechanics.* Oxford University Press, New York.
45. Van Kampen, N. G. 2007. *Stochastic Processes in Physics and Chemistry, Third edition.* Elsevier, Amsterdam, The Netherlands.
46. Hänggi, P., P. Talkner, and M. Borkovec. 1990. Reaction-rate theory: fifty years after Kramers. *Rev. Mod. Phys.* 62:251–341.
47. Schuler, B., E. A. Lipman, and W. A. Eaton. 2002. Probing the free-energy surface for protein folding with single-molecule fluorescence spectroscopy. *Nature.* 419:743–747.
48. Soranno, A., B. Buchli, ..., B. Schuler. 2012. Quantifying internal friction in unfolded and intrinsically disordered proteins with single-molecule spectroscopy. *Proc. Natl. Acad. Sci. USA.* 109:17800–17806.
49. Chandler, D. 1978. Statistical mechanics of isomerization dynamics in liquids and the transition state approximation. *J. Chem. Phys.* 68:2959–2970.
50. Hill, T. L. 2005. *Free Energy Transduction and Biochemical Cycle Kinetics.* Dover Publications, Inc., Mineola, NY.
51. Ha, T., and P. Tinnefeld. 2012. Photophysics of fluorescent probes for single-molecule biophysics and super-resolution imaging. *Annu. Rev. Phys. Chem.* 63:595–617.
52. Cooper, M., A. Ebner, ..., R. West. 2004. Cy3B: improving the performance of cyanine dyes. *J. Fluoresc.* 14:145–150.
53. Donnert, G., C. Eggeling, and S. W. Hell. 2007. Major signal increase in fluorescence microscopy through dark-state relaxation. *Nat. Methods.* 4:81–86.
54. Donnert, G., C. Eggeling, and S. W. Hell. 2009. Triplet-relaxation microscopy with bunched pulsed excitation. *Photochem. Photobiol. Sci.* 8:481–485.
55. Rasnik, I., S. A. McKinney, and T. Ha. 2006. Nonblinking and long-lasting single-molecule fluorescence imaging. *Nat. Methods.* 3:891–893.
56. Cordes, T., J. Vogelsang, and P. Tinnefeld. 2009. On the mechanism of Trolox as antiblinking and antibleaching reagent. *J. Am. Chem. Soc.* 131:5018–5019.
57. Altman, R. B., D. S. Terry, ..., S. C. Blanchard. 2011. Cyanine fluorophore derivatives with enhanced photostability. *Nat. Methods.* 9:68–71.
58. van der Velde, J. H. M., J. Oelerich, ..., T. Cordes. 2016. A simple and versatile design concept for fluorophore derivatives with intramolecular photostabilization. *Nat. Commun.* 7:10144.
59. Abate, J., and W. Whitt. 2006. A unified framework for numerically inverting Laplace transforms. *INFORMS J. Comput.* 18:408–421.
60. Gibson, A., and B. Conolly. 1971. On a three-state sojourn time problem. *J. Appl. Probab.* 8:716–723.
61. Berezhkovskii, A. M., A. Szabo, and G. H. Weiss. 2000. Theory of the fluorescence of single molecules undergoing multistate conformational dynamics. *J. Phys. Chem. B.* 104:3776–3780.
62. Gopich, I. V., and A. Szabo. 2010. FRET efficiency distributions of multistate single molecules. *J. Phys. Chem. B.* 114:15221–15226.