

Bayesian inference and comparison of stochastic transcription elongation models

Jordan Douglas^{1,2}, Richard Kingston¹, Alexei J. Drummond^{2*}

1 School of Biological Sciences, University of Auckland, Auckland, New Zealand

2 Centre for Computational Evolution, School of Computer Science, University of Auckland, Auckland, New Zealand

* alexei@cs.auckland.ac.nz

Abstract

Transcription elongation can be modelled as a three step process, involving polymerase translocation, NTP binding, and nucleotide incorporation into the nascent mRNA. This cycle of events can be simulated at the single-molecule level as a continuous-time Markov process using parameters derived from single-molecule experiments. Previously developed models differ in the way they are parameterised, and in their incorporation of partial equilibrium approximations.

We have formulated a hierarchical network comprised of 12 sequence-dependent transcription elongation models. The simplest model has two parameters and assumes that both translocation and NTP binding can be modelled as equilibrium processes. The most complex model has six parameters makes no partial equilibrium assumptions. We systematically compared the ability of these models to explain published force-velocity data, using approximate Bayesian computation. This analysis was performed using data for the RNA polymerase complexes of *E. coli*, *S. cerevisiae* and Bacteriophage T7.

Our analysis indicates that the polymerases differ significantly in their translocation rates, with the rates in T7 pol being fast compared to *E. coli* RNAP and *S. cerevisiae* pol II. Different models are applicable in different cases. We also show that all three RNA polymerases have an energetic preference for the posttranslocated state over the pretranslocated state. A Bayesian inference and model selection framework, like the one presented in this publication, should be routinely applicable to the interrogation of single-molecule datasets.

Author summary

Transcription is a critical biological process which occurs in all living organisms. It involves copying the organism's genetic material into messenger RNA (mRNA) which directs protein synthesis on the ribosome. Transcription is performed by RNA polymerases which have been extensively studied using both ensemble and single-molecule techniques (see reviews: [1, 2]). Single-molecule data provides unique insights into the molecular behaviour of RNA polymerases. Transcription at the single-molecule level can be computationally simulated as a continuous-time Markov process and the model outputs compared with experimental data. In this study we use Bayesian techniques to perform a systematic comparison of 12 stochastic models of transcriptional elongation. We demonstrate how equilibrium approximations can strengthen or weaken the model, and show how Bayesian techniques can identify

necessary or unnecessary model parameters. We describe a framework to a) simulate, b) perform inference on, and c) compare models of transcription elongation.

Introduction

Transcription is carried out by RNA polymerases: RNAP in *Escherichia coli*, pol II in *Saccharomyces cerevisiae*, and T7 pol in Bacteriophage T7. It involves the copying of template double-stranded DNA (dsDNA) into single-stranded messenger RNA (mRNA). RNAP and pol II are comprised of multiple subunits, and their catalytic subunits are homologous [3,4]. In contrast, T7 pol exists as a monomer with a distinct sequence, and resembles the *E. coli* DNA polymerase I [5].

Optical trapping experiments have been performed on the transcription elongation complex (TEC) from a variety of organisms [6–12]. In a typical experimental setup, two polystyrene beads (around 600 nm in diameter) are tethered to the system; one attached to the RNA polymerase and the other to the DNA [6]. As transcription elongation progresses, the distance between the two beads increases and the velocity of a single TEC can be computed. Optical tweezers can be used to apply a force F to the system (Fig 1).

Fig 1. Effect of an applied force on elongation velocity. (A) Optical trapping setup showing dsDNA being transcribed by RNA polymerase (grey ellipse) into mRNA. Two polystyrene beads are tethered to the system allowing the application of force using optical tweezers. An assisting load $F > 0$ acts in the same direction as transcription (top) while a hindering load $F < 0$ acts in the opposing direction (bottom). Figure not to scale. (B) Schematic depiction of the effect of applying a force on RNA polymerase. Due to the stochastic nature of transcription at the single-molecule level, each experiment yields a different distance-time trajectory, even under the same applied force.

Single-molecule studies of the TEC have revealed that RNA polymerases progress in a discontinuous fashion [6,13–16] with step sizes that correspond to the dimensions of a single nucleotide (3.4 Å [17]). Consequently, at the single molecule level, transcription is best modelled as a discrete process rather than a continuous one.

A single cycle in the main transcription elongation pathway (Fig 2) requires (1) Forward translocation of the RNA polymerase, making the active site accessible; (2) Binding of the complementary nucleoside triphosphate (NTP); (3) Addition of the nucleotide onto the 3' end of the mRNA. This third step involves NTP hydrolysis. Nucleoside monophosphate is added onto the chain and pyrophosphate is released from the enzyme.

Our study aimed to identify the best model to describe this reaction cycle for RNAP, pol II, and T7 pol, based on analysis of published force-velocity data. As there are three reactions, up to six rate constants may be necessary for a kinetic model of a single nucleotide addition. These describe forward and backwards translocation (k_{fwd} and k_{bck}), binding and release of NTP (k_{bind} and k_{rel}), and NTP catalysis and reverse-catalysis (k_{cat} and k_{rev}), also known as pyrophosphorolysis [20]. However fewer than six parameters may be required in practice.

First, it is reasonable to assume that polymerisation is effectively irreversible [19,21–23], as pyrophosphorolysis is a highly exergonic reaction, reducing the number of rate constants to five. Second, translocation between the pretranslocated and posttranslocated states, and/or NTP binding, may occur on timescales significantly more rapid than the other steps, in which case they may be modelled as equilibrium processes. These assumptions simplify the model, as the respective forward and reverse

Fig 2. State diagrams of RNA polymerase. (A) The model of the main transcription elongation pathway, which shows the postulated states; the pathways for interconversion; and the rate constants that govern each part of the reaction. The transcription bubble is the set of $\beta_1 + h + \beta_2$ bases (see main text for definitions) in the double-stranded DNA which are unpaired. States are denoted by $S(l, t)$ where l is the length of the mRNA and t is the position of the polymerase active site (small grey rectangle) with respect to the 3' end of the mRNA. Polymerase translocation displaces the polymerase by a distance of $\delta = 1$ bp = 3.4 Å. During polymerisation the chain is extended by one nucleotide. (B) Instantiated posttranslocated state of RNA polymerase transcribing the *rpoB* gene sequence, with $\beta_1 = 2, h = 9, \beta_2 = 1$. Forward translocation requires melting two T/A basepairs (right arrows). Backward translocation requires melting two C/G basepairs (left arrows). The mRNA secondary structure would also require reconfiguration [18, 19].

reaction rate constants are subsumed by a single equilibrium constant. Third, thermodynamic models of nucleic acid structure can be used to estimate sequence-dependent translocation rates $k_{fwd}(l)$ and $k_{bck}(l)$, by invoking transition state theory, and this can sometimes result in parameter reduction [18, 19, 23].

Irrespective of equilibrium assumptions and parameterisation, transcription elongation under applied force can be modelled in two fundamentally distinct ways. First, there are the **deterministic** equations which can be used to calculate the mean pause-free elongation velocity $v(F, [\text{NTP}])$ as a function of force F and NTP concentration $[\text{NTP}]$. This kind of model can be derived from the differential equations describing the time evolution of all species, by application of the steady state approximation. Force effects on the translocation step are incorporated using transition state theory [24, 25].

An example is the following 3-parameter model [6].

$$v(F, [\text{NTP}]) = \frac{k_{cat}}{1 + \frac{K_D}{[\text{NTP}]} (1 + K_\tau e^{-F\delta/k_B T})} \quad (1)$$

where δ is the distance between adjacent basepairs (3.4 Å, [17]), $K_D = \frac{k_{rel}}{k_{bind}}$ is the equilibrium constant of NTP binding, $K_\tau = \frac{k_{bck}}{k_{fwd}}$ is the equilibrium constant of translocation, k_B is the Boltzmann constant, and T is the absolute temperature. Increasingly complex equations may be used as more parameters or states are added to the model [6, 8, 19]. Such equations describe the velocity averaged across an ensemble of molecules. Parameter inference applied to velocity-force-[NTP] experimental data is straightforward and computationally fast when using these equations. However these equations do not describe the distribution of velocity nor do they account for site heterogeneity across the nucleic acid sequence and therefore cannot predict local sequence effects.

Second, there are the **stochastic** models, which can be implemented via simulation of single-molecule behaviour using the Gillespie algorithm [26]. The mean velocity can be calculated by averaging velocities over a number of simulations for a given F and $[\text{NTP}]$. This offers not just the mean but a full distribution of velocities and could potentially explain emergent properties unavailable from a deterministic model. Unfortunately, simulating can be very slow and therefore parameter inference can be a problem.

In this study we used a Markov-chain-Monte-Carlo approximate-Bayesian-computation (MCMC-ABC) algorithm [27] to estimate transcription elongation parameters for **stochastic** models via simulation. The

observed pause-free velocities we are fitting to were measured at varying applied force and NTP concentration. For each RNA polymerase under study – *E. coli* RNAP, *S. cerevisiae* pol II, and T7 pol – we fit to one respective dataset from the single-molecule literature [6, 28, 29].

Models

Notation and state space

Suppose the TEC is transcribing a gene of length L . Then let $S(l, t)$ denote a TEC state, where the mRNA is currently of length $l \leq L$, and $t \in \mathbb{Z}$ describes the position of the active site with respect to the 3' end of the mRNA. When $t = 0$ the polymerase is pretranslocated and cannot bind NTP, and when $t = 1$ the polymerase is posttranslocated and *can* bind NTP (Fig. 2). This study is focused on the main elongation pathway and the observed velocities being fitted have pauses filtered out. Therefore, although additional backtracked states ($t < 0$) [6, 30, 31] and hypertranslocated states ($t > 1$) [32, 33] exist, these are not incorporated in the model.

Let β_1 and β_2 be the number of unpaired template nucleotides upstream and downstream of RNA polymerase, respectively, and let h be the number of basepairs in the DNA/mRNA hybrid (Fig. 2A). Although there are uncertainties in these parameters, they are held constant at $h = 9$, $\beta_1 = 2$, and $\beta_2 = 1$ [19, 34].

Transcription of the gene begins at state $S(l_0, 0)$ and ends upon reaching $S(L, 0)$, where $l_0 = \beta_1 + h + 2$.

Parameterisation of the NTP binding step

// New section. Content relocated from Model Space

NTP binding has been modelled as both a kinetic and equilibrium process in the literature [6, 19, 23].

In a kinetic binding model, NTP binding occurs at pseudo-first order rate $k_{bind}[\text{NTP}]$, while NTP release occurs at rate k_{rel} . In this case, k_{bind} and $\frac{k_{rel}}{k_{bind}}$ must be estimated.

Under a partial equilibrium approximation NTP binding and release are assumed to be rapid enough that equilibrium is achieved. In this case, the rate constants k_{bind} and k_{rel} are subsumed by the NTP dissociation constant $K_D = \frac{k_{rel}}{k_{bind}}$ which becomes the sole binding-related parameter to estimate.

Parameterisation of the translocation step

While inferences about the rate constants associated with NTP binding and catalysis (k_{bind} , $\frac{k_{rel}}{k_{bind}}$, and k_{cat}) can be made directly from the data, the translocation step is more complex. Transition state theory is invoked in order to estimate k_{fwd} and k_{bck} . Recasting the problem in this way (1) provides a way of accommodating the effects of applied force on the elongation process, and (2) allows the sequence-dependence of translocation to be incorporated by considering the energetics of basepairing. When allowing for sequence dependence, the total number of translocation rates required to model translocation of the full gene is $2(L - l_0)$.

Thermodynamic models of base pairing energies

// New section adapted from Nucleic Acid Thermodynamics (previously in Materials and Methods) and two paragraphs from the Introduction

The standard Gibbs free energies $\Delta_r G^0 (= \Delta G)$ involved in duplex formation are calculated using nearest neighbour models. The standard Gibbs energy of state S – arising from nucleotide basepairing and dangling ends – is calculated as

$$\Delta G_S^{(bp)} = \Delta G_{gene}^{(bp)} + \Delta G_{hybrid}^{(bp)} \quad (2)$$

where SantaLucia’s DNA/DNA basepairing parameters [35] are used to calculate $\Delta G_{gene}^{(bp)}$ and Sugimoto’s DNA/RNA parameters [36] are used for $\Delta G_{hybrid}^{(bp)}$. For the latter, dangling end energies are estimated as described by Bai et al. 2004 [23]. Here, and elsewhere, the $^{(bp)}$ superscript is used to denote a model parameter that can be evaluated from the sequence alone. Gibbs energies are expressed relative to the thermal energy of the system, in units of $k_B T$, where $k_B T = 0.6156$ kcal/mol at $T = 310$ K.

In order for RNA polymerase to translocate forward (backward), up to two basepairs must be disrupted: (1) the basepair at the downstream (upstream) edge of the transcription bubble, and (2) the basepair at the upstream (downstream) end of the DNA/mRNA hybrid (Fig 2B). Differences in the basepairing energies in these regions confer sequence-dependence on the rate of translocation.

Calculation of translocation rates or translocation equilibrium constant

// New section. The critical information regarding translocation kinetics has been brought to the front so that minor details can be skipped by the reader

The standard Gibbs energies of the pre and posttranslocated states, $\Delta G_{S(l,0)}^{(bp)}$ and $\Delta G_{S(l,1)}^{(bp)}$, respectively, are used with up to four additional terms – $\Delta G_{\tau 1}$, δ_1 , $\Delta G_{\tau}^{\ddagger}$, and $\Delta G_{T(l,t)}^{(bp)}$ – to calculate the translocation rates. The first three are model parameters which must be estimated while the latter is directly evaluated from the sequence.

Let $T(l, t)$ be the translocation transition state between $S(l, t)$ and $S(l, t + 1)$. Then $\Delta G_{T(l,t)}^{\ddagger} = \Delta G_{\tau}^{\ddagger} + \Delta G_{T(l,t)}^{(bp)}$ is the sequence-dependent standard Gibbs energy of activation which must be overcome in order to translocate (Fig 3).

Fig 3. Parameterisation of the translocation step. (A) Effects of model parameters on state energies. The figure displays a schematic Gibbs energy landscape of translocation, with backtracked states included for visualisation purposes. The solid red lines represent translocation states ($t = 0$: pretranslocated, $t = 1$: posttranslocated, and $t < 0$: backtracked), while the dashed red lines represent transition states. Applying an assisting force $F > 0$ tilts the landscape in favour of higher values of t . The effect of $\Delta G_{\tau 1}$ is observed at the posttranslocated state $t = 1$. In a translocation equilibrium model, the barrier height is assumed to be so small, = translocation is so rapid, that the transition states are disregarded. (B) A model for the sequence-dependent transition state between translocation states $S(l, 0)$ and $S(l, 1)$. This is required for estimating the Gibbs energy of basepairing $\Delta G_{T(l,t)}^{(bp)}$ in the transition state. The basepairing energy, added to a baseline term $\Delta G_{\tau}^{\ddagger}$, together specify the height of the activation barrier (Equation 10).

Given an applied force F , the translocation rates governing transition between the pre and posttranslocated states ($k_{fwd}(l)$ and $k_{bck}(l)$) are calculated from barrier height $\Delta G_{T(l,0)}^{\ddagger}$ using an Arrhenius type relation:

$$k_{fwd}(l) = A e^{-(\Delta G_{T(l,0)}^{\ddagger} - \Delta G_{S(l,0)}^{(bp)} - F\delta_1/k_B T)} \quad (3)$$

$$k_{bck}(l) = A e^{-(\Delta G_{T(l,0)}^{\ddagger} - (\Delta G_{S(l,1)}^{(bp)} + \Delta G_{\tau 1}) + F(\delta - \delta_1)/k_B T)} \quad (4)$$

The derived rates $k_{fd}(l)$ and $k_{bck}(l)$ are therefore dependent on the local sequence. The pre-exponential factor A is held constant at 10^6 s^{-1} . This term has been arbitrarily set to a variety of values in previous studies ($10^6 - 10^9 \text{ s}^{-1}$ [18, 19, 23]). This has little consequence for model fitting, however the value of $\Delta G_{T(l,t)}^\ddagger$ is entangled with the value of the pre-exponential factor A and can only be meaningfully interpreted in light of its value.

If the system has time to reach equilibrium, the probabilities of observing the pretranslocated state $S(l, 0)$ and posttranslocated state $S(l, 1)$ are

$$p(S(l, 0)) \propto e^{-(\Delta G_{S(l,0)}^{(bp)})} \quad (5)$$

$$p(S(l, 1)) \propto e^{-(\Delta G_{\tau 1} + \Delta G_{S(l,1)}^{(bp)})} \quad (6)$$

This is described by equilibrium constant K_τ .

$$K_\tau(l) = \frac{p(S(l, 0))}{p(S(l, 1))} \quad (7)$$

$$= \exp\{-(\Delta G_{S(l,0)} - \Delta G_{S(l,1)})\} \quad (8)$$

$$= \exp\{-(\Delta G_{S(l,0)}^{(bp)} - \Delta G_{S(l,1)}^{(bp)} - \Delta G_{\tau 1})\} \quad (9)$$

The physical meanings of the terms $\Delta G_{\tau 1}$, δ_1 , ΔG_τ^\ddagger , and $\Delta G_{T(l,t)}^{(bp)}$, and the way they are used in the model, are detailed below.

Energetic bias for the posttranslocated states

$\Delta G_{\tau 1}$ (units $k_B T$) is a parameter added to the standard Gibbs energy of the posttranslocated state. If $\Delta G_{\tau 1} = 0$, then the sequence alone determines the Gibbs energy difference between pre and posttranslocated states. In this case, pretranslocated states are usually favoured over posttranslocated states due to the loss of a single basepair in the hybrid of the latter.

$\Delta G_{\tau 1}$ has frequently been estimated for T7 pol [37–39] and there has been discussion around whether such a term is necessary for RNAP [8].

Polymerase displacement and formation of the transition state

δ_1 (units \AA) is the distance that the polymerase must translocate forward to facilitate the formation of the transition state. The distance between adjacent basepairs is held constant at an experimentally measured value $\delta = 3.4 \text{ \AA}$ [17], and $0 < \delta_1 < \delta$. The response of the system to an applied force F depends on this term. In general, the application of force F tilts the Gibbs energy landscape – the Gibbs energy difference between adjacent translocation states being augmented by a factor $\frac{F\delta}{k_B T}$ (Fig 3A, [1, 40]).

It may be necessary to estimate δ_1 to model the data adequately [19], or it may be sufficient to simply set $\delta_1 = \delta/2$ [40].

Energy barrier of translocation

ΔG_τ^\ddagger and $\Delta G_{T(l,t)}^{(bp)}$ (units $k_B T$) together determine the activation barrier height in the translocation step. It is assumed that the sequence-dependent standard Gibbs energy of activation $\Delta G_{T(l,t)}^\ddagger$ can be written as

$$\Delta G_{T(l,t)}^\ddagger = \Delta G_\tau^\ddagger + \Delta G_{T(l,t)}^{(bp)} \quad (10)$$

$\Delta G_{\tau}^{\ddagger}$ is therefore a sequence-independent baseline term used to compute the translocation barrier heights. The parameter $\Delta G_{\tau}^{\ddagger}$ must be estimated in order to evaluate translocation rates.

// The paragraph below was relocated and adapted from SI Appendix 3. SI Appendix 3 has since been deleted.

In contrast $\Delta G_{T(l,t)}^{(bp)}$ is a term that is evaluated directly from the sequence derived from a model of the transition state (Fig 3B). The term is evaluated as the standard Gibbs energy of a TEC containing all hybrid and gene basepairs found in both $S(l, t)$ and $S(l, t + 1)$, ie. the intersection between the two sets of basepairs.

Model space

The full transcription elongation model makes use of the following 6 parameters:

- k_{cat} (units s^{-1}).
- $K_D = \frac{k_{rel}}{k_{bind}}$ (units μM).
- k_{bind} (units $\mu M^{-1} s^{-1}$).
- $\Delta G_{\tau 1}$ (units $k_B T$).
- δ_1 (units \AA).
- $\Delta G_{\tau}^{\ddagger}$ (units $k_B T$).

However fewer than 6 parameters may be needed to adequately describe the data. If it is assumed that the energy differences between pre and posttranslocated states are determined by basepairing energies alone, the parameter $\Delta G_{\tau 1}$ does not need to be estimated. This is equivalent to holding $\Delta G_{\tau 1}$ constant at 0. If it is assumed that the displacement required for formation of the translocation intermediate state is half the distance between adjacent basepairs, the parameter δ_1 does not need to be estimated. This is equivalent to holding δ_1 constant at $\delta/2$.

Partial equilibrium approximations may also simplify the model, as detailed above. If binding is approximated as an equilibrium process, k_{bind} does not need to be estimated. If translocation is approximated as an equilibrium process, $\Delta G_{\tau}^{\ddagger}$ and δ_1 do not need to be estimated. One, both, or neither of these two steps (binding and translocation) could be assumed to achieve equilibrium, thus yielding four equilibrium model variants (Fig 4A). The introduction of partial equilibrium approximations for both the NTP binding and translocation steps has implications when specifying the prior distributions for the Bayesian analysis (S4 Appendix.) The chemical master equations for single nucleotide addition cycles of these models are presented in S2 Appendix.

Fig 4. // B: Arrows in figure now have labels **The space of models to be compared.** (A) The four equilibrium model variants. NTP binding, translocation, both, or neither, could be assumed to achieve equilibrium prior to catalysis. (B) The 12 transcription elongation models. An arrow connects model i to j if augmentation of model i with a single parameter generates model j . The number of parameters to estimate k is shown for each level in the network. Equilibrium approximation colour scheme is the same as in A. $\Delta G_{\tau 1}$ and δ_1 can each be estimated or set to a constant.

Incorporating these simplifications to the model in a combinatorial fashion results in a total of 12 related models, which together constitute the model space. Our objective was to determine which of these 12 models provides the best description of the

experimental data. The simplest model (Model 1) contains 2 parameters (k_{cat} and K_D). The most complex model (Model 12) contains all 6 parameters. The full model space is displayed in Fig 4B.

Stochastic modelling

For each model we performed stochastic simulations, appropriate for the modelling of single-molecule force-velocity data. The simulations, performed using the Gillespie algorithm [26, 41], can be used to estimate the mean elongation velocity under a model.

The estimation of mean velocity can be broken down into three steps. First, the system is initialised by placing the RNA polymerase at the 3' end of the template – state $S(l_0, 0)$ – with the transcription bubble open and a DNA/RNA hybrid formed. The force and NTP concentrations are assigned their experimentally set values. Second, a chemical reaction is randomly sampled. The probability that reaction $S \xrightarrow{k} S'$ is selected is proportional to its rate constant k (Fig 2). The amount of time taken for the reaction to occur is sampled from the exponential distribution. States which are subject to a partial equilibrium approximation are coalesced into a single state, which augments the outbound rate constants. The second step is repeated until the RNA polymerase has copied the entire template. Third, the previous two steps are repeated c times. The mean elongation velocity is evaluated as the mean of each mean elongation velocity across c simulations. For further information, see S1 Appendix.

Relation to previous models and stochastic simulations

// Section relocated and adapted from “Discussion”

There is an extensive literature concerned with the kinetic modelling of transcription elongation. Such models may incorporate backtracking, hypertranslocation, and other reactions. Here we are concerned only with the central elongation pathway.

A stochastic and sequence-dependent model was proposed by Bai et al. 2004 [23] for RNAP, with both NTP binding and translocation treated as equilibrium processes. The translocation equilibrium constant was calculated entirely from basepairing energies. Therefore this model is equivalent to Model 1, and the parameters were estimated as $k_{cat} = 24.7 \text{ s}^{-1}$ and $K_D = 15.6 \mu\text{M}$ from fit to experimental data. Maoiléidigh et al. 2011 also presented stochastic simulations of RNAP. The elongation component of their model is equivalent to Model 6 [19]. We build on this work by providing a systematic Bayesian framework for model comparison and parameter estimation.

While our analysis employed sequence-dependent stochastic models, comparisons can also be made with some deterministic models.

Abbondanzieri et al. 2005 [6], Larson et al. 2012 [42], Schweikhard et al. 2014 [28], and Thomen et al. 2008. [29, 39] described a deterministic model (for RNAP, pol II, pol II, and T7 pol respectively) which estimated k_{cat} , K_D and translocation equilibrium constant $K_\tau = \frac{k_{bck}}{k_{fwd}}$. These are most similar to Model 4.

Maoiléidigh et al. 2011 for RNAP, and Dangkulwanich et al. 2013 for pol II, however found that the translocation and catalysis were occurring on similar timescales, and modelled only NTP binding as an equilibrium process [19, 43]. They also estimated the distance of translocation. These deterministic models are most similar to Model 11.

Finally, Mejia et al. 2015 [44] used a model that is quite different to all the above models, as it does not explicitly treat translocation. Instead elongation is modelled with a two step kinetic scheme, the first step involving NTP binding and conformational change, and the second step involving nucleotide incorporation and product release. This model is most similar to a special case of Model 5 where ΔG_{τ_1} becomes extremely negative, driving the polymerase into the posttranslocated position.

Results and Discussion

Model selection with MCMC-ABC

Our aim was to 1) use Bayesian inference to select the best of 12 transcription elongation models for each RNA polymerase; and 2) estimate the parameters for those of the models appearing in the 95% credible set of the posterior distribution. Selecting prior distributions behind each parameter is a critical process in Bayesian inference. A prior distribution should reflect what is known about the parameter before observing the new data. We have explicated our prior assumptions, with justifications, in Table 1.

We performed MCMC-ABC experiments which estimated the parameters and model indicator M_i for $i \in \mathbb{Z}, 1 \leq i \leq 12$. Models which appear more often in this posterior distribution are better choices, given the data.

// The content below was relocated from the now-removed Materials and Methods

The datasets we fit our models to are all from the single-molecule literature and are presented in: Figures 5a and 5b of Abbondanzieri et al. 2005 [6] for *E. coli* RNAP, Figure 2a of Schweikhard et al. 2014 [28] for *S. cerevisiae* pol II, and Table 2 of Thomen et al. 2008 [29] for T7 pol. To computationally replicate these experiments as faithfully as we could with the available information and computational limitations, simulations in this study were run on the 4 kb *E. coli rpoB* gene for RNAP (GenBank: EU274658), the first 4.75 kb of the human *rpb1* gene for pol II (NCBI: NG_027747) the first 10 kb of the Enterobacteria phage λ genome for T7 pol (NCBI: NC_001416). The mean velocities from 32 (for RNAP), 10 (for pol II) and 3 (for T7 pol) simulations of the full respective sequences were used to estimate the mean elongation velocity during MCMC-ABC, given F and $[NTP]$.

For further information about the MCMC-ABC algorithm [27, 45], see S3 Appendix.

The posterior distributions

// Section relocated and adapted from "Model selection with MCMC-ABC"

The posterior distributions from our MCMC-ABC experiments are presented in Table 2, Fig 5, and Fig 6.

Fig 5. Posterior and prior distribution plots. Posterior distributions for all models which appear in the 95% credible set are displayed (two models for RNAP, two models for pol II, and one model for T7 pol). Plots show the prior probability density $P(\theta)$ of each parameter and posterior probability density of each parameter conditional on the model $P(\theta|D, M_i)$. The geometric median point-estimates and highest posterior density (HPD) intervals (calculated with Tracer 1.6 [54]) are displayed above each plot (3 sf).

Fig 6. // Figure relocated from "Supporting Information" Posterior distributions of simulated velocities. Black open circles represent experimentally measured mean velocities reported in the original publication for (A) RNAP, (B) pol II, and (C) T7 pol [6, 28, 29]. Each coloured dot represents a single sample simulated from the posterior distribution of parameters/models for the respective polymerase. 30 samples were generated from each of the three posterior distributions. For RNAP, $[NTP]_{eq}$ is defined as $[ATP] = 5 \mu M$, $[CTP] = 2.5 \mu M$, $[GTP] = 10 \mu M$, and $[UTP] = 10 \mu M$.

A large effective sample size (> 100 [54]) and a small \hat{R} (< 1.1 , as defined by Gelman et al. 1992 [55–57]) are essential for making reliable parameter estimates. Table

Table 1. Prior distributions used during Bayesian inference.

Parameter	Prior distribution(s)	Justification of prior distribution(s)
Model M	$P(M_i) = 2/16$ for $i \in \{1, 2, 4, 5\}$ $P(M_i) = 1/16$ for $i \in \{3, 6, 7, 8, 9, 10, 11, 12\}$	Each model should each have uniformly distributed values. Models with translocation at equilibrium have double the prior probability since these models do not use δ_1 .
k_{cat} (s^{-1})	Lognormal($\mu = 3.454$, $\sigma = 0.587$) for RNAP/pol II Lognormal($\mu = 4.585$, $\sigma = 0.457$) for T7 pol	k_{cat} and elongation velocity estimates for <i>E. coli</i> RNAP and <i>S. cerevisiae</i> pol II range from 18 to 50 s^{-1} for optical trapping experiments [8–10, 23, 44], but as much as 100 bp/s <i>in vivo</i> [46–49]. Distribution selected such that (10, 100) is central 95% interval. For T7 pol k_{cat} and elongation velocity estimates range from 43 - 240 bp/s [11, 50–52]. Distribution selected such that (40, 240) is central 95% interval.
K_D (μM)	Lognormal($\mu = 1.844$, $\sigma = 1.762$)	Estimates for K_D under binding equilibrium models range from 20-140 μM [8, 22, 39, 42, 53]. In models where binding is kinetic and slow, $K_D \equiv \frac{k_{rel}}{k_{bind}}$ could be much lower (S4 Appendix). To accommodate for both binding models, the prior distribution was selected such that the central 95% interval is (0.2, 200).
k_{bind} ($\mu M^{-1}s^{-1}$)	Lognormal($\mu = -1.498$, $S\sigma = 1.585$)	Central 95% interval set so that NTP binding is a slow kinetic step (S4 Appendix). Centered around (0.01, 5).
$\Delta G_{\tau 1}$ ($k_B T$)	Normal($\mu = 0$, $\sigma = 1.55$) for RNAP/pol II Normal($\mu = -3.3$, $\sigma = 1.55$) for T7 pol	For RNAP and pol II, centered around 0 with a standard deviation comparable to the free energy of a single nucleotide basepair doublet, and such that the 95% central interval is (-4, 4). For T7 pol $\Delta G_{\tau 1}$ has been estimated as -4.3 [39] and -4.87 $k_B T$ [37]. However these estimates are likely resulting partially from dangling ends. Thus, we subtracted the mean dangling end contribution of ~ -1 $k_B T$ [35] and centered the prior around this interval with a standard deviation the same as above.
$\Delta G_{\tau}^{\ddagger}$ ($k_B T$)	Normal($\mu = 5.5$, $\sigma = 0.97$) for RNAP/pol II Normal($\mu = 2.5$, $\sigma = 1.36$) for T7 pol	Central 95% interval set so that translocation is a slow kinetic step (S4 Appendix). Selected so that 99% central interval is (3, 8) for RNAP and pol II, and (-1, 6) for T7 pol.
δ_1 (\AA)	Uniform($l = 0$, $u = 3.4$)	Uniformly distributed across all possible values.

Prior distributions behind all estimated parameters and the model indicator. Unless specified otherwise, the prior distribution is used for all three RNA polymerases. Lognormal priors (parameterised in log space) are used for rates and equilibrium constants while normal priors are used for Gibbs energy terms. **To maintain statistical integrity of the Bayesian analysis, prior distributions were not derived from the data presented by** Abbondanzieri et al. 2005 [6] for RNAP, by Schweikhard et al. 2014 [28] for pol II, or by Thomen et al. 2008 [29] for T7 pol.

2 suggests that the parameters in the 95 % credible set of models are sufficiently estimated by these criteria.

These results indicate that the best models for the datasets examined are Models 11

Table 2. Summary of MCMC-ABC experiments.

Enzyme		<i>E. coli</i> RNAP		<i>S. cerevisiae</i> pol II		Bacteriophage T7 pol
ϵ		2.39		0.705		4.63
Combined chain length		3.5×10^7		6.2×10^7		1.2×10^8
i		11	12	11	12	5
Model	Description	Binding equilibrium, Translocation kinetic	Binding kinetic, Translocation kinetic	Binding equilibrium, Translocation kinetic	Binding kinetic, Translocation kinetic	Binding kinetic, Translocation equilibrium
ESS / \hat{R}	\hat{k}_{cat}	257 / 1.04	1441 / 1.03	549 / 1.02	1203 / 1.01	2110 / 1.00
	$\frac{\hat{k}_{cat}}{\hat{k}_{bind}}$	328 / 1.01	101 / 1.01	536 / 1.01	133 / 1.05	106 / 1.00
	\hat{k}_{bind}	—	705 / 1.09	—	516 / 1.02	154 / 1.00
	$\Delta\hat{G}_{\tau 1}$	466 / 1.02	1844 / 1.00	1145 / 1.00	2769 / 1.00	1626 / 1.02
	$\hat{\delta}_1$	300 / 1.04	2290 / 1.03	658 / 1.01	1469 / 1.00	—
	$\Delta\hat{G}_{\tau}^{\ddagger}$	340 / 1.02	1680 / 1.03	589 / 1.02	1179 / 1.00	—
Posterior	$P(M_i D)$	0.81	0.19	0.29	0.71	0.96

// Rhat values and each individual ESS is now shown in this plot. The parameter estimates were moved to Fig 5 Each column summarises the posterior distribution for the respective RNA polymerase, which arises from multiple independent MCMC chains. Approximate Bayesian computation threshold ϵ is shown for each enzyme; state Θ is accepted into the posterior distribution only if $X^2(\Theta) \leq \epsilon$ (S3 Appendix). Models which appear in an RNA polymerase's 95% credible set and their posterior probabilities $P(M_i|D)$ are shown. The effective sample size (ESS, calculated with Tracer 1.6 [54]) and R-hat (\hat{R} [55–57]) of each parameter, conditional on M_i , are displayed. A large ESS (> 100) and a small \hat{R} (< 1.1) imply that the MCMC experiment has converged. Where a parameter is not incorporated in the kinetic model, a ‘—’ is left in its place.

and 12 for both RNAP and pol II, and Model 5 for T7 pol (Fig 4B).

For pol II, Model 12 has the highest posterior probability $P(M_{12}|D) = 0.71$. This is the most complex model considered, with 6 estimated parameters. In Model 12 translocation, NTP binding and catalysis are all kinetic processes; the displacement required to facilitate formation of the translocation transition state, $\delta_1 < \delta$, is estimated ($\hat{\delta}_1 = 3.1$ Å); and the standard Gibbs energy of the posttranslocated state is influenced by parameter $\Delta G_{\tau 1} \neq 0$.

The posterior distribution for RNAP consists of the same set models as that of pol II. For RNAP, Model 11 has the highest probability $P(M_{11}|D) = 0.81$. This model is a submodel of Model 12 with one fewer parameter: in Model 11 NTP binding is treated as an equilibrium process while in Model 12 it is not.

The only model in the 95 % credible set for T7 pol is Model 5 $P(M_5|D) = 0.96$. In Model 5 (4 parameters) translocation, but not binding, is treated as an equilibrium process, and $\Delta G_{\tau 1}$ is estimated. This positions T7 pol in a quite different area of the model space to the other two polymerases.

Translocation rates differ among RNA polymerases

For RNAP and pol II, we estimate that a partial equilibrium approximation for the translocation step is inadequate. The posterior probability that such models are inadequate is 1.00 (see Table 2). For T7 pol, however, translocation is significantly faster than catalysis and is best modelled with a partial equilibrium approximation. Using estimates for $\Delta G_{\tau}^{\ddagger}$ and $\Delta G_{\tau 1}$ under the maximum posterior models (Model 11 for RNAP and Model 12 for pol II) we estimate the mean forward \bar{k}_{fwd} and backward \bar{k}_{bck} translocation rates averaged across the *rpoB* sequence as: 230 s^{-1} and 112 s^{-1} for RNAP, and 350 s^{-1} and 12.7 s^{-1} for pol II, respectively (3 sf). These estimates are within one order of magnitude of the respective estimate for the rate of catalysis (Fig 5) suggesting that translocation and catalysis indeed occur on similar timescales.

For RNAP and pol II, translocation has frequently been modelled as an equilibrium process [6, 23, 28, 42, 44], however in some recent analyses this assumption has been rejected [18, 19, 43, 58, 59]. Our Bayesian analysis supports this. In contrast, there is general agreement that translocation in T7 pol is adequately modelled as an equilibrium process [29, 60, 61].

The data does not determine the kinetics of the NTP binding step

// Content relocated and adapted from “Model selection with MCMC-ABC”

It remains unclear how to best model the NTP binding step. Models that describe NTP binding as a kinetic process have posterior probabilities of 0.19 for RNAP, 0.71 for pol II and 0.96 for T7 pol (Table 2). However, in our sensitivity analysis, where we used different a prior distribution for $\frac{k_{rel}}{k_{bind}}$, these probabilities were 0.65, 0.22, and 0.19, respectively (results not shown).

Furthermore, $\frac{k_{rel}}{k_{bind}}$ and k_{bind} (Models 5 and 12) are unable to be estimated simultaneously. For pol II and for T7 pol, k_{bind} is estimated at around 0.48 and $1.4 \mu\text{M}^{-1} \text{ s}^{-1}$ respectively with fairly narrow 95% highest posterior density (HPD) intervals (Fig 5). However, the HPD interval of $\frac{k_{rel}}{k_{bind}}$ spans three orders of magnitude and the value of this parameter was therefore poorly informed by the data. For RNAP, in contrast, neither k_{bind} nor $\frac{k_{rel}}{k_{bind}}$ were well-informed by the data and both have HPD intervals spanning 1-2 orders of magnitude. This non-identifiability – where two or more parameters are unable to be estimated simultaneously (S4 Appendix) – highlights the appeal in an NTP binding equilibrium model where only one parameter $\frac{k_{rel}}{k_{bind}}$ needs to be estimated, despite the unrealistic assumptions it may invoke. In the case of each enzyme, the data has taught us nothing about one or two of the binding parameters.

The pause-free mean velocities measured during transcription elongation follow Michaelis-Menten kinetics even though the reaction cycle is more complicated than that of a simple enzyme [62]. As such, the inability to resolve the timescale of the substrate binding step is unsurprising [63–65].

Overall, these three factors, i) the sensitivity of the posterior probabilities to the choice of prior, ii) the intermediate magnitude of said probabilities, iii) the inability to estimate both $\frac{k_{rel}}{k_{bind}}$ and k_{bind} simultaneously, collectively imply the data carries very little information about the rates of NTP binding and release.

In the transcription literature, NTP binding is almost always assumed to achieve equilibrium for RNAP, pol II, and T7 pol [6, 18, 19, 23, 28, 29, 39, 42, 43, 61]. However Mejia et al. 2015 [44] have shown that NTP binding is indeed rate-limiting, and that mutations in the RNAP trigger loop impair the binding rate thus suggesting that the trigger loop is coupled with NTP binding.

RNAP has an energetic preference for the posttranslocated state

In previous stochastic sequence-dependent models [18,23] the standard Gibbs energies of the pre and posttranslocated states have been based solely on the nucleic acid basepairing energies. Our models include an additional term, $\Delta G_{\tau 1}$, to account for potential interactions between the protein and the nucleic acid. The marginal posterior probability of a model in which an additional term $\Delta G_{\tau 1}$ is required is 1.00 in all three polymerases. In each case $\Delta G_{\tau 1}$ was estimated to be less than 0 $k_B T$ and 0 $k_B T$ is not included in the 95 % HPD interval (Fig 5). We find that $\Delta \hat{G}_{\tau 1}$ is the most significant in pol II and T7 pol: $-4.6 k_B T$ and $-4.0 k_B T$ respectively, while $\Delta \hat{G}_{\tau 1} = -2.0 k_B T$ for RNAP (2 sf).

These results suggest that structural elements within RNA polymerases can energetically favour posttranslocated states over pretranslocated states. We note that the sequence-dependent contribution of the dangling end of the DNA/RNA hybrid is included in the thermodynamic model. The energetic bias for the posttranslocated state is separable from this effect.

To facilitate comparison with previous deterministic models, using our estimates of $\Delta G_{\tau 1}$ we calculated the equilibrium constant between the pre and posttranslocated states. Geometrically averaged across the *rpoB* gene, these are

$$\bar{K}_{\tau} = \frac{1}{L - l_0} \exp\left\{\sum_{l=l_0}^{L-1} \ln(k_{bck}(l)/k_{fwd}(l))\right\} = \begin{cases} 0.77 \text{ for RNAP} \\ 0.057 \text{ for pol II} \\ 0.10 \text{ for T7 pol.} \end{cases} \quad (11)$$

Thus, for all three polymerases, $K_{\tau} < 1$, indicating that the small energetic preference that the protein has for the posttranslocated state is sufficient to override the loss of basepairing energy, thereby biasing the system towards population of the posttranslocated positions. This is in agreement with estimates made for pol II and T7 pol [28,29,37,38,42] and Kireeva et al. 2018 [59] for RNAP: “*forward translocation occurs in milliseconds and is poorly reversible*”. However these estimates are inconsistent with some RNAP and pol II studies which place this ratio above 1 [6,19,43,53].

Kinetic modelling can itself suggest no physical mechanism for the stabilisation. Yu et al. 2012 [38] have identified a conserved tyrosine residue near the active site of T7 pol that pushes against the 3' end of the mRNA, and thus stabilises the posttranslocated state. They propose a similar mechanism for the multi-subunit RNA polymerases.

δ_1 may be an important parameter but its physical meaning is unclear

Our results suggest that δ_1 , the distance that RNA polymerase must translocate forward by to reach the translocation transition state, is a necessary parameter to estimate for RNAP and pol II. Setting $\delta_1 = \delta/2$ is not sufficient. The marginal posterior probability of models which estimate this term is 1.00. δ_1 is irrelevant to the modeling of the T7 pol data because the best models invoke a partial equilibrium approximation for the translocation step.

While our prior distribution restricted δ_1 to lie in the range $(0, \delta)$, the upper end our 95% HPD intervals of δ_1 for RNAP and pol II are very close to $\delta = 3.4 \text{ \AA}$. If it was not for this prior distribution, δ_1 estimates would have included values higher than δ . Similar results have been observed by Maoiléidigh et al. 2011 [19] for RNAP.

Our interpretation of δ_1 implies it should never be greater than δ nor should δ be more than the width of one basepair. The physical meaning of δ_1 with values greater than δ is thus unclear. It is noted that δ_1 is only used when $F \neq 0$.

Comparing the kinetics of RNA polymerases

The *in vivo* rate of transcription elongation varies considerably across RNAP, pol II and T7 pol. The prokaryotic and eukaryotic RNA polymerases have a mean rate ranging from 20-120 bp/s [46, 47, 49, 50, 66–68], which may be slowed down in histone-wrapped regions of eukaryotic genomes [9]. In contrast, Bacteriophage T7 pol operates up to an order of magnitude faster (around 200-240 bp/s [50, 69]) and is known to be quite insensitive to transcriptional pause sites [11, 29].

In addition to these differences, we have shown that translocation is very rapid in T7 pol, relative to the rate of NTP incorporation, while the disparity is much less significant in RNAP and pol II. Furthermore, the model does not fit the data for T7 pol as closely it does for RNAP and pol II (Fig 6). T7 pol therefore seems to operate under quite a different kinetic scheme than that of the cellular polymerases, which is not unexpected given their distant evolutionary relationship [5].

In general, the elongation velocity of RNA polymerase is significantly slower in an optical trap (with estimates ranging from 9.7-22 bp/s for RNAP [13–15, 44, 70]) compared with that of the untethered enzyme (with estimates *in vitro* or *in vivo* ranging from 25-118 bp/s for RNAP [46, 50, 71, 72]). This relationship holds for multiple RNA polymerases including *E. coli* RNAP, *S. cerevisiae* pol II [42, 43, 53, 73], Bacteriophage T7 pol [11, 29, 50, 52], and Bacteriophage $\Phi 6$ P2 [12, 74]. This suggests that optical trapping perturbs the system to a significant extent. Additionally, varying degrees of heterogeneity in elongation rate have been observed across different polymerase complexes even under the same conditions [13, 15, 29].

The velocity perturbations resulting from the optical trapping apparatus will be propagated into the model parameters, especially k_{cat} , and $\Delta G_{\tau}^{\ddagger}$, and some caution is needed when extrapolating these results to untethered systems.

Bayesian inference of transcription elongation

To our knowledge we are the first to perform Bayesian inference on single-molecule models of transcription elongation. This was achieved by simulation which necessitated the use of approximate Bayesian computation. An alternative would be to build and use a likelihood function (ie. the probability of taking exactly t units of time for RNA polymerase to copy the sequence n times). The latter approach can be achieved using chemical master equations, as opposed to (Gillespie) sampling from the distribution. Finding analytical, stable numerical, or approximate solutions to the chemical master equations could provide a similar insight in less computational time, however is susceptible to a multitude of analytical and numerical issues associated with the exponentiation of an arbitrary transition rate matrix that grows with the length of the sequence (S2 Appendix) [75]. This problem would be amplified by the introduction of backtracking, hypertranslocation, or NTP misincorporation reactions into the model, for instance. The Bayesian framework we have presented, although computationally intensive due to its simulation requirement, is general and will work on any model of transcription without the need to resolve these issues. The path has been paved for modelling transcriptional pausing, for instance [18, 23, 76]. Nevertheless, likelihood-based Bayesian inference is an approach that should be explored in the future.

We have demonstrated that single-molecule data can be usefully analysed using a Bayesian inference and model selection framework. This analysis would have even greater statistical power if applied to the progression of individual RNA polymerase complexes instead of mean velocities averaged across multiple experiments.

Conclusion

In this article we evaluated some simple Brownian ratchet models of transcription elongation (Fig 2). By varying the parameterisation of the translocation step (Fig 3) and incorporating partial equilibrium approximations commonly invoked in the literature (Fig 4A) we enumerated a total of 12 related models (Fig 4B). Using stochastic simulations and approximate Bayesian computation, we then assessed which of these models were capable of describing the force-velocity data previously measured for several RNA polymerases (Table 2 and Fig 5) using single-molecule optical trapping experiments [6, 28, 29].

Our analysis suggests that 1) different partial equilibrium approximations of the translocation step are appropriate for the multisubunit RNA polymerases versus the single subunit T7 RNA polymerase. 2) Treatment of the NTP binding step remains a point of ambiguity. The existing data does not place strong constraints on the modelling of this step. 3) There is an energetic bias for posttranslocated state. 4) The model of the force-dependent translocation, which invokes transition state theory, is not physically realistic.

Supporting information

S1 Appendix. Stochastic simulation.

Reactions are simulated using the Gillespie algorithm [26]. Given the current state s and a set of possible reactions $s \rightarrow s_1, s \rightarrow s_2, \dots, s \rightarrow s_n$ with rate constants k_1, k_2, \dots, k_n , the next reaction to perform is sampled proportional to its rate:

$$p(s \rightarrow s_i) = \frac{k_i}{\sum_{j=1}^n k_j}. \quad (12)$$

The amount of time the reaction takes to occur is sampled from the exponential distribution with rate $\sum_{j=1}^n k_j$.

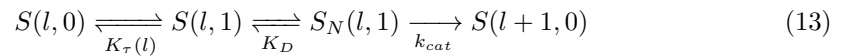
S2 Appendix. Chemical master equations.

// Added chemical master equations

The chemical master equations for the four equilibrium variants (Fig 4A) of single nucleotide addition cycles are provided in this section. Transcription of a full gene involves chaining multiple of these single-cycle models together.

Translocation and binding equilibrium model

The state pathway of a single-cycle of the translocation and binding equilibrium model is



where S_N denotes a state where NTP is bound. As states $S(l, 0)$, $S(l, 1)$, and $S_N(l, 1)$ are in mutual equilibrium they can be coalesced into one state. Let $S(l)$ be the coalesced state that exists in equilibrium between the three, then the pathway may be rewritten as



where $\omega_{cat}(l)$ is the *effective* rate of catalysis from this state. This term is equal to k_{cat} multiplied by the proportion of time that the coalesced state has NTP bound, ie. in state $S_N(l, 1)$.

$$\omega_{cat}(l) = p(S_N(l, 1)) k_{cat} \quad (15)$$

$$= \frac{\frac{[NTP]}{K_D}}{1 + \frac{[NTP]}{K_D} + \exp\{-(\Delta G_{S(l,0)}^{(bp)} - \Delta G_{S(l,1)}^{(bp)} - \Delta G_{\tau 1})\}} k_{cat} \quad (16)$$

Let $p(S, \mathbb{T})$ be the probability of the system existing in state S at time \mathbb{T} . Under this model, the chemical master equation of a single-cycle is:

$$\begin{pmatrix} \frac{dp(S(l), \mathbb{T})}{d\mathbb{T}} \\ \frac{dp(S(l+1, 0), \mathbb{T})}{d\mathbb{T}} \end{pmatrix} = \begin{pmatrix} -\omega_{cat}(l) & 0 \\ \omega_{cat}(l) & 0 \end{pmatrix} \begin{pmatrix} p(S(l), \mathbb{T}) \\ p(S(l+1, 0), \mathbb{T}) \end{pmatrix} \quad (17)$$

This system is simple enough to solve analytically. Let $p(\mathbb{T})$ be the Markov transition matrix after time \mathbb{T} and let Q be the Markov process transition rate matrix. In these matrices, entry i, j is the probability (p) or rate (Q) of transition from i to j .

$$p(\mathbb{T}) = \exp\{Q\mathbb{T}\} \quad (18)$$

$$= \exp\left\{ \begin{pmatrix} -\omega_{cat}(l) & \omega_{cat}(l) \\ 0 & 0 \end{pmatrix} \mathbb{T} \right\} \quad (19)$$

$$= \begin{pmatrix} e^{-\omega_{cat}(l)\mathbb{T}} & 1 - e^{-\omega_{cat}(l)\mathbb{T}} \\ 0 & 1 \end{pmatrix} \quad (20)$$

It is noted that as the time \mathbb{T} approaches infinity, the probability of the system existing in state $S(l+1, 0)$ approaches 1, because it is an absorbing state.

Let $f(\mathbb{T})$ be the probability density of taking *exactly* \mathbb{T} units of time to arrive at state $S(l+1, 0)$, starting from $S(l)$. Because $S(l+1, 0)$ is an absorbing state, computing $f(\mathbb{T})$ is trivial.

$$f(\mathbb{T}) = \frac{d}{d\mathbb{T}} (1 - e^{-\omega_{cat}(l)\mathbb{T}}) \quad (21)$$

$$= \omega_{cat}(l) e^{-\omega_{cat}(l)\mathbb{T}} \quad (22)$$

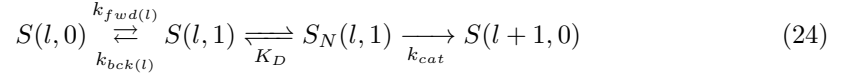
Transcribing the full gene requires the traversal of all $L - l_0 + 1$ states. The size of the transition rate matrix for transcribing the entire sequence therefore grows with the length of the DNA sequence.

$$Q = \begin{pmatrix} -\omega_{cat}(l_0) & \omega_{cat}(l_0) & 0 & 0 & \dots & 0 \\ 0 & -\omega_{cat}(l_0 + 1) & \omega_{cat}(l_0 + 1) & 0 & \dots & 0 \\ \vdots & \ddots & \ddots & \ddots & \ddots & \vdots \\ 0 & 0 & 0 & 0 & -\omega_{cat}(L - 1) & \omega_{cat}(L - 1) \\ 0 & 0 & 0 & 0 & 0 & 0 \end{pmatrix} \quad (23)$$

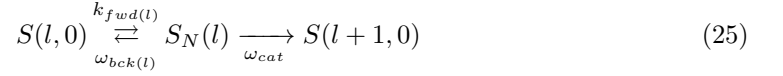
Applying the matrix exponential function to a linear pathway, such as the one described above, has an analytical solution [77]. However, for non-linear systems with an arbitrary number of states (eg. the other 3 equilibrium model variants), analytical solutions may not exist and numerical solutions likely exhibit numerical instabilities [75]. This option was not further investigated in this project (simulation was used instead). However, analytical or stable numerical solutions for $f(\mathbb{T})$ for an arbitrary model would facilitate the use of likelihood functions in Bayesian inference, thereby rendering simulation obsolete.

Binding equilibrium model

The state pathway of a single-cycle of the binding equilibrium model is



The binding equilibrium assumption permits the coalescence of $S(l, 1)$ and $S_N(l, 1)$ into a single state $S_N(l)$. Thus, the pathway can be rewritten as



where $\omega_{bck(l)}$ is the effective rate of backwards translocation from $S_N(l)$, and ω_{cat} is the effective rate of catalysis. These rates are derived by multiplying their composite rates – $k_{bck(l)}$ and k_{cat} – by the probability of the system existing in the required state to apply the reaction – $p(S(l, 1))$ and $p(S_N(l, 1))$.

$$\omega_{bck(l)} = p(S(l, 1)) k_{bck(l)} \quad (26)$$

$$= \frac{1}{1 + \frac{[NTP]}{K_D}} k_{bck(l)} \quad (27)$$

$$\omega_{cat}(l) = p(S_N(l, 1)) k_{cat} \quad (28)$$

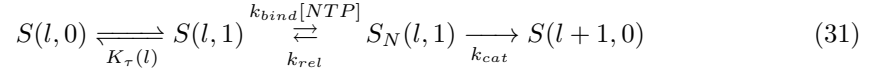
$$= \frac{\frac{[NTP]}{K_D}}{1 + \frac{[NTP]}{K_D}} k_{cat} \quad (29)$$

Let $p(S, \mathbb{T})$ be the probability of the system being at state S at time \mathbb{T} . Under this model, the chemical master equation of a single-cycle is:

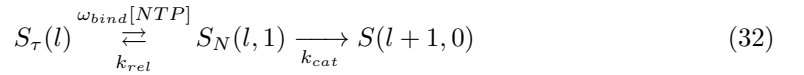
$$\begin{pmatrix} \frac{dp(S(l,0),\mathbb{T})}{d\mathbb{T}} \\ \frac{dp(S_N(l),\mathbb{T})}{d\mathbb{T}} \\ \frac{dp(S(l+1,0),\mathbb{T})}{d\mathbb{T}} \end{pmatrix} = \begin{pmatrix} -k_{fwd}(l) & \omega_{bck}(l) & 0 \\ k_{fwd}(l) & -\omega_{bck}(l) - \omega_{cat}(l) & 0 \\ 0 & \omega_{cat}(l) & 0 \end{pmatrix} \begin{pmatrix} p(S(l,0),\mathbb{T}) \\ p(S_N(l),\mathbb{T}) \\ p(S(l+1,0),\mathbb{T}) \end{pmatrix} \quad (30)$$

Translocation equilibrium model

The state pathway of a single-cycle of the translocation equilibrium model is



The translocation equilibrium assumption permits the coalescence of $S(l,0)$ and $S(l,1)$ into a single state $S_\tau(l)$. Thus, the pathway can be rewritten as



where ω_{bind} is the effective rate of NTP binding.

$$\omega_{bind} = p(S(l,1)) k_{bind} \quad (33)$$

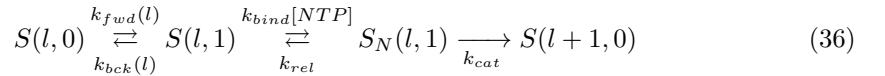
$$= \frac{\exp\{-(\Delta G_{S(l,1)}^{(bp)} + \Delta G_{\tau 1})\}}{\exp\{-(\Delta G_{S(l,1)}^{(bp)} + \Delta G_{\tau 1})\} + \exp\{-\Delta G_{S(l,0)}^{(bp)}\}} k_{bind} \quad (34)$$

Let $p(S, \mathbb{T})$ be the probability of the system being at state S at time \mathbb{T} . Under this model, the chemical master equation of a single-cycle is:

$$\begin{pmatrix} \frac{dp(S_\tau(l),\mathbb{T})}{d\mathbb{T}} \\ \frac{dp(S_N(l,1),\mathbb{T})}{d\mathbb{T}} \\ \frac{dp(S(l+1,0),\mathbb{T})}{d\mathbb{T}} \end{pmatrix} = \begin{pmatrix} -\omega_{bind}[NTP] & k_{rel} & 0 \\ \omega_{bind}[NTP] & -k_{rel} - k_{cat} & 0 \\ 0 & k_{cat} & 0 \end{pmatrix} \begin{pmatrix} p(S_\tau(l),\mathbb{T}) \\ p(S_N(l,1),\mathbb{T}) \\ p(S(l+1,0),\mathbb{T}) \end{pmatrix} \quad (35)$$

Full kinetic model

The state pathway of a single-cycle of the full kinetic model is



Let $p(S, \mathbb{T})$ be the probability of the system being at state S at time \mathbb{T} . Under this model, the chemical master equation of a single-cycle is:

$$\begin{pmatrix} \frac{dp(S(l,0),\mathbb{T})}{d\mathbb{T}} \\ \frac{dp(S(l,1),\mathbb{T})}{d\mathbb{T}} \\ \frac{dp(S_N(l,1),\mathbb{T})}{d\mathbb{T}} \\ \frac{dp(S(l+1,0),\mathbb{T})}{d\mathbb{T}} \end{pmatrix} = \begin{pmatrix} -k_{fwd}(l) & k_{bck}(l) & 0 & 0 \\ k_{fwd}(l) & -k_{bck}(l) - k_{bind}[NTP] & k_{rel} & 0 \\ 0 & k_{bind}[NTP] & -k_{rel} - k_{cat} & 0 \\ 0 & 0 & k_{cat} & 0 \end{pmatrix} \begin{pmatrix} p(S(l,0),\mathbb{T}) \\ p(S(l,1),\mathbb{T}) \\ p(S_N(l,1),\mathbb{T}) \\ p(S(l+1,0),\mathbb{T}) \end{pmatrix}$$

S3 Appendix. MCMC-ABC. Given model/parameters Θ and observed data $D = (D_1, D_2, \dots, D_n)$, Bayesian inference conventionally involves approximating the posterior probability distribution $P(\Theta|D)$ using the likelihood $P(D|\Theta)$ and the prior $P(\Theta)$.

$$P(\Theta|D) \propto P(D|\Theta)P(\Theta). \quad (38)$$

As there is no easily computed likelihood function, simulation is used. The chi-squared test statistic X^2 evaluates how well a given set of parameters fits the data.

$$X^2 = \sum_i \frac{(S_i - D_i)^2}{S_i} \quad (39)$$

where S_i is the mean velocity simulated under the same [NTP] and applied force F that D_i was measured under. The probability that $X^2 = 0$ is equal to the likelihood $P(D|\Theta)$, however it is computationally impractical to only accept parameters into the posterior distribution when the simulation yields $X^2 = 0$. Therefore a threshold ϵ is used and sample Θ_i is accepted into the posterior only if $X^2(\Theta_i) \leq \epsilon$. This method is called approximate Bayesian computation [27, 45]. This is coupled with Markov chain Monte Carlo (MCMC) to give the MCMC-ABC algorithm which is becoming increasingly popular among computational biologists [27, 78].

// Content below relocated and adapted from the since-removed “Materials and Methods”

Each MCMC chain estimated six parameters and the model indicator M . This means the 12 models share the same parameter objects in the MCMC. There are therefore seven terms to estimate: M , k_{cat} , $\frac{k_{rel}}{k_{bind}}$, k_{bind} , $\Delta G_{\tau 1}$, ΔG_{τ}^\ddagger , and δ_1 .

When the current model M_i does not use a parameter (eg. in Model 5 δ_1 is not used), the parameter is still estimated even though it is not being used. When the model requires a parameter to be held constant (eg. in Model 1 $\Delta G_{\tau 1} = 0$), the parameter is set to its constant during the simulation. This is done without affecting or being affected by its current estimate, which is used by other models.

To achieve convergence, we used an exponential cooling scheme on ϵ [78] where $\epsilon_{i+1} = \max(\epsilon_{min}, \epsilon_i \gamma)$ for manually tuned values of $0 < \gamma < 1$ and ϵ_0 . Chains which failed to converge were discarded. A heavy-tailed distribution [79] is used as a proposal function, and the parameter to change at each step in the MCMC is selected uniformly at random.

We ran one or more independent MCMC-ABC chains for each selected ϵ_{min} / RNA polymerase combination. Selecting the threshold ϵ_{min} is a critical process in approximate Bayesian computation. Threshold ϵ_{min} must be large enough to achieve convergence within finite computational resources, but small enough that the resulting posterior distribution is still an accurate approximation of the true posterior distribution. For each RNA polymerase we set ϵ_{min} to some initial guess. Then we ran the MCMC chain until the ESS for X^2 was large (> 300) and lowered ϵ_{min} to the bottom 0.05 quantile of the posterior distribution of X^2 . This step was repeated until either: a) the distribution of model indicators M converged (model posterior probabilities have changed by less than 0.01, on average). Or, b) the acceptance rate was less than 5%. The values of ϵ_{min} used in the final posterior distributions were 2.39 for RNAP, 0.705 for pol II, and 4.63 for T7 pol (Table 2).

Parameter point estimates (Fig 5) are the geometric median: that is the value which minimises the total Euclidean distance from the other posterior samples. Parameters were normalised into z-scores first. Our code is open source and available at <http://www.polymerase.nz>. Textfiles containing the posterior distributions and simulation settings are available to download or visualise with the software.

S4 Appendix. Prior distributions.

Prior for $\Delta G_{\tau}^{\ddagger}$, which governs the rates of translocation

RNAP/pol II: to select a prior for $\Delta G_{\tau}^{\ddagger}$ we simulated transcription on the *rpoB* gene under Model 3 – the simplest binding equilibrium model. $\Delta G_{\tau}^{\ddagger}$ and k_{cat} were sampled uniformly from a relevant range, with K_D held constant at 100 μM and $[NTP] = 1000 \mu M$. For each simulation, the mean elongation velocity was calculated. The results are displayed in S1 Fig.

This plot shows that as the energy barrier of translocation ($\Delta G_{\tau}^{\ddagger}$) increases, the velocity decreases. If $\Delta G_{\tau}^{\ddagger} \gtrsim 8 k_B T$ then it becomes impossible to achieve a realistic mean velocity, providing a relatively clear upper bound on this parameter. If $\Delta G_{\tau}^{\ddagger} \lesssim 3 k_B T$ then translocation becomes very rapid and the same distribution of velocities is obtained in simulations, irrespective of the exact value of $\Delta G_{\tau}^{\ddagger}$. In this case catalysis becomes strongly rate-limiting, and it would be appropriate to apply a partial equilibrium approximation to the translocation step. This provides an effective lower bound for parameter $\Delta G_{\tau}^{\ddagger}$. Therefore we centered our prior distribution for $\Delta G_{\tau}^{\ddagger}$ in this interval (a normal distribution with a mean of 5.5 and a standard deviation of 0.97, so that the central 99% interval is (3, 8)). We performed the same analysis with different values of K_D , as well as varying $\Delta G_{\tau 1}$, and arrived at the same interval for $\Delta G_{\tau}^{\ddagger}$ (results not shown).

T7 pol: the same analysis was performed, however with $\Delta G_{\tau}^{\ddagger}$ at its prior mean of $-3.3 k_B T$ (S1 Fig).

Prior for k_{bind} , which governs the rate of NTP binding

To select a prior for k_{bind} we performed similar simulations, but instead used Model 2 – the simplest kinetic binding model. k_{bind} and k_{cat} were sampled uniformly from relevant ranges, K_D was set to 100 μM and $[NTP] = 1000 \mu M$. (S1 Fig).

Depending on the exact value of k_{cat} , if $k_{bind} \lesssim 0.1 \mu M^{-1} s^{-1}$, then it is impossible to achieve a realistic velocity, providing a relatively clear lower bound on this parameter. If $k_{bind} \gtrsim 5 \mu M^{-1} s^{-1}$ then binding becomes very rapid and the same distribution of velocities is obtained in simulations, irrespective of the exact value of k_{bind} . Again this is because catalysis becomes strongly rate limiting in this region, and it would be appropriate to apply a partial equilibrium approximation to the binding step. Hence we centered our (lognormal) prior around the interval (0.01, 5) – the conservatively selected bounds reflecting that the experimental data has been collected at differing NTP concentrations, altering the rate. Performing the same analysis with different parameters gave us a similar prior.

Prior distribution related to rate of NTP release

A model is non-identifiable if two or more parameterisations can produce the same output. Our preliminary results suggested non-identifiability between $\frac{k_{rel}}{k_{bind}}$ and k_{bind} (S1 Fig). When k_{bind} is low (and hence binding is rate-limiting), there is an approximately linear relationship between $\frac{k_{rel}}{k_{bind}}$ and k_{bind} . As k_{bind} increases from 0, the dissociation constant $\frac{k_{rel}}{k_{bind}}$ must also increase in order for the system to achieve the same velocity. However, as binding comes closer to achieving equilibrium, $\frac{k_{rel}}{k_{bind}}$ converges. Most previous estimates of K_D have assumed binding to be at equilibrium. This assumption restrains the values which K_D may take, and subsequently estimates for K_D are typically in the order of $10^1 - 10^2 \mu M$. However for a model in which binding is slow it is expected that estimates of $\frac{k_{rel}}{k_{bind}}$ can be lower. This has indeed been demonstrated by Mejia et al. 2015 [44] who estimated $\frac{k_{rel}}{k_{bind}}$ to be 0.6 μM . Therefore the prior distribution for $\frac{k_{rel}}{k_{bind}}$ must permit both of these binding models to be tested fairly during Bayesian inference. We centered our lognormal prior for $\frac{k_{rel}}{k_{bind}}$ around a very broad range, with a central 95% interval of (0.2, 200).

It is noted that selecting a prior distribution which does not discriminate between

the kinetic and equilibrium binding models *a priori* may not be plausible.

S1 Fig. Simulations of the elongation pathway. Each point is a single simulation of the full *rpoB* gene (4029 nt). For (A-C), Parameters on the x- and z-axis are sampled uniformly at random from the displayed range at the beginning of each trial. The y-axis of each plot (mean elongation velocity) is then measured from the respective simulation. [NTP] and F held constant at 1000 μM and 0 pN respectively. (A) and (B): Relationship between $\Delta G_{\tau}^{\ddagger}$ and k_{cat} for the melting model with binding at equilibrium (Model 8). $\Delta G_{\tau-1}$ set to its prior mean (0 for RNAP and pol II, and -3.3 for T7 pol). (C) Relationship between k_{bind} and k_{cat} for the kinetic binding model with translocation at equilibrium (Model 2). (D) Relationship between K_D and k_{bind} with translocation held at equilibrium (Model 2). K_D and k_{bind} sampled uniformly from specified range and velocity is measured. Samples with simulated velocities outside of the range 1-2 bp/s were discarded. [NTP] = 10 μM and $k_{cat} = 100 \text{ s}^{-1}$.

Acknowledgments

We wish to acknowledge the contribution of NeSI high-performance computing facilities to the results of this research. NZ's national facilities are provided by the NZ eScience Infrastructure and funded jointly by NeSI's collaborator institutions and through the Ministry of Business, Innovation & Employment's Research Infrastructure programme. URL <https://www.nesi.org.nz>.

References

1. Herbert KM, Greenleaf WJ, Block SM. Single-molecule studies of RNA polymerase: motoring along. *Annu Rev Biochem.* 2008;77:149–176.
2. Dangkulwanich M, Ishibashi T, Bintu L, Bustamante C. Molecular mechanisms of transcription through single-molecule experiments. *Chemical reviews.* 2014;114(6):3203–3223.
3. Sweetser D, Nonet M, Young RA. Prokaryotic and eukaryotic RNA polymerases have homologous core subunits. *Proceedings of the National Academy of Sciences.* 1987;84(5):1192–1196.
4. Sosunov V, Sosunova E, Mustaev A, Bass I, Nikiforov V, Goldfarb A. Unified two-metal mechanism of RNA synthesis and degradation by RNA polymerase. *The EMBO journal.* 2003;22(9):2234–2244.
5. Sousa R, Chung YJ, Rose JP, Wang BC. Crystal structure of bacteriophage T7 RNA polymerase at 3.3 Å resolution. *Nature.* 1993;364(6438):593.
6. Abbondanzieri EA, Greenleaf WJ, Shaevitz JW, Landick R, Block SM. Direct observation of base-pair stepping by RNA polymerase. *Nature.* 2005;438(7067):460–465.
7. Adelman K, La Porta A, Santangelo TJ, Lis JT, Roberts JW, Wang MD. Single molecule analysis of RNA polymerase elongation reveals uniform kinetic behavior. *Proceedings of the National Academy of Sciences.* 2002;99(21):13538–13543.
8. Bai L, Fulbright RM, Wang MD. Mechanochemical kinetics of transcription elongation. *Physical review letters.* 2007;98(6):068103.

9. Hodges C, Bintu L, Lubkowska L, Kashlev M, Bustamante C. Nucleosomal fluctuations govern the transcription dynamics of RNA polymerase II. *Science*. 2009;325(5940):626–628.
10. Galburt EA, Grill SW, Wiedmann A, Lubkowska L, Choy J, Nogales E, et al. Backtracking determines the force sensitivity of RNAP II in a factor-dependent manner. *Nature*. 2007;446(7137):820–823.
11. Skinner GM, Baumann CG, Quinn DM, Molloy JE, Hoggett JG. Promoter binding, initiation, and elongation by bacteriophage T7 RNA polymerase a single-molecule view of the transcription cycle. *Journal of Biological Chemistry*. 2004;279(5):3239–3244.
12. Dulin D, Vilfan ID, Berghuis BA, Hage S, Bamford DH, Poranen MM, et al. Elongation-competent pauses govern the fidelity of a viral RNA-dependent RNA polymerase. *Cell reports*. 2015;10(6):983–992.
13. Neuman KC, Abbondanzieri EA, Landick R, Gelles J, Block SM. Ubiquitous transcriptional pausing is independent of RNA polymerase backtracking. *Cell*. 2003;115(4):437–447.
14. Davenport RJ, Wuite GJ, Landick R, Bustamante C. Single-molecule study of transcriptional pausing and arrest by *E. coli* RNA polymerase. *Science*. 2000;287(5462):2497.
15. Tolić-Nørrelykke SF, Engh AM, Landick R, Gelles J. Diversity in the rates of transcript elongation by single RNA polymerase molecules. *Journal of Biological Chemistry*. 2004;279(5):3292–3299.
16. Abbondanzieri EA, Shaevitz JW, Block SM. Picocalorimetry of transcription by RNA polymerase. *Biophysical journal*. 2005;89(6):L61–L63.
17. Watson JD, Crick FH, et al. Molecular structure of nucleic acids. *Nature*. 1953;171(4356):737–738.
18. Tadigotla VR, Maoiléidigh DÓ, Sengupta AM, Epshtein V, Ebright RH, Nudler E, et al. Thermodynamic and kinetic modeling of transcriptional pausing. *Proceedings of the National Academy of Sciences of the United States of America*. 2006;103(12):4439–4444.
19. Maoiléidigh DÓ, Tadigotla VR, Nudler E, Ruckenstein AE. A unified model of transcription elongation: what have we learned from single-molecule experiments? *Biophysical journal*. 2011;100(5):1157–1166.
20. Maitra U, Nakata Y, Hurwitz J. The Role of Deoxyribonucleic Acid in Ribonucleic Acid Synthesis XIV. A Study of the Initiation of Ribonucleic Acid Synthesis. *Journal of Biological Chemistry*. 1967;242(21):4908–4918.
21. Erie DA, Yager TD, Von Hippel PH. The single-nucleotide addition cycle in transcription: a biophysical and biochemical perspective. *Annual review of biophysics and biomolecular structure*. 1992;21(1):379–415.
22. Rhodes G, Chamberlin MJ. Ribonucleic acid chain elongation by *Escherichia coli* ribonucleic acid polymerase I. Isolation of ternary complexes and the kinetics of elongation. *Journal of Biological Chemistry*. 1974;249(20):6675–6683.

23. Bai L, Shundrovsky A, Wang MD. Sequence-dependent kinetic model for transcription elongation by RNA polymerase. *Journal of molecular biology*. 2004;344(2):335–349.
24. Bustamante C, Chemla YR, Forde NR, Izhaky D. Mechanical processes in biochemistry. *Annual review of biochemistry*. 2004;73(1):705–748.
25. Cleland W. Partition analysis and concept of net rate constants as tools in enzyme kinetics. *Biochemistry*. 1975;14(14):3220–3224.
26. Gillespie DT. Exact stochastic simulation of coupled chemical reactions. *The journal of physical chemistry*. 1977;81(25):2340–2361.
27. Beaumont MA. Approximate Bayesian computation in evolution and ecology. *Annual review of ecology, evolution, and systematics*. 2010;41:379–406.
28. Schweikhard V, Meng C, Murakami K, Kaplan CD, Kornberg RD, Block SM. Transcription factors TFIIF and TFIIS promote transcript elongation by RNA polymerase II by synergistic and independent mechanisms. *Proceedings of the National Academy of Sciences*. 2014;111(18):6642–6647.
29. Thomen P, Lopez P, Bockelmann U, Guillerez J, Dreyfus M, Heslot F. T7 RNA polymerase studied by force measurements varying cofactor concentration. *Biophysical journal*. 2008;95(5):2423–2433.
30. Wang MD, Schnitzer MJ, Yin H, Landick R, Gelles J, Block SM. Force and velocity measured for single molecules of RNA polymerase. *Science*. 1998;282(5390):902–907.
31. Shaevitz JW, Abbondanzieri EA, Landick R, Block SM. Backtracking by single RNA polymerase molecules observed at near-base-pair resolution. *Nature*. 2003;426(6967):684–687.
32. Artsimovitch I, Landick R. Pausing by bacterial RNA polymerase is mediated by mechanistically distinct classes of signals. *Proceedings of the National Academy of Sciences*. 2000;97(13):7090–7095.
33. Zhou Y, Navaroli DM, Enuameh MS, Martin CT. Dissociation of halted T7 RNA polymerase elongation complexes proceeds via a forward-translocation mechanism. *Proceedings of the National Academy of Sciences*. 2007;104(25):10352–10357.
34. Greive SJ, Von Hippel PH. Thinking quantitatively about transcriptional regulation. *Nature Reviews Molecular Cell Biology*. 2005;6(3):221–232.
35. SantaLucia J. A unified view of polymer, dumbbell, and oligonucleotide DNA nearest-neighbor thermodynamics. *Proceedings of the National Academy of Sciences*. 1998;95(4):1460–1465.
36. Wu P, Nakano Si, Sugimoto N. Temperature dependence of thermodynamic properties for DNA/DNA and RNA/DNA duplex formation. *The FEBS Journal*. 2002;269(12):2821–2830.
37. Yin YW, Steitz TA. The structural mechanism of translocation and helicase activity in T7 RNA polymerase. *Cell*. 2004;116(3):393–404.
38. Yu J, Oster G. A small post-translocation energy bias aids nucleotide selection in T7 RNA polymerase transcription. *Biophysical journal*. 2012;102(3):532–541.

39. Thomen P, Lopez PJ, Heslot F. Unravelling the mechanism of RNA-polymerase forward motion by using mechanical force. *Physical Review Letters*. 2005;94(12):128102.
40. Depken M, Galburt EA, Grill SW. The origin of short transcriptional pauses. *Biophysical journal*. 2009;96(6):2189–2193.
41. Lecca P. Stochastic chemical kinetics. *Biophysical reviews*. 2013;5(4):323–345.
42. Larson MH, Zhou J, Kaplan CD, Palangat M, Kornberg RD, Landick R, et al. Trigger loop dynamics mediate the balance between the transcriptional fidelity and speed of RNA polymerase II. *Proceedings of the National Academy of Sciences*. 2012;109(17):6555–6560.
43. Dangkulwanich M, Ishibashi T, Liu S, Kireeva ML, Lubkowska L, Kashlev M, et al. Complete dissection of transcription elongation reveals slow translocation of RNA polymerase II in a linear ratchet mechanism. *Elife*. 2013;2:e00971.
44. Mejia YX, Nudler E, Bustamante C. Trigger loop folding determines transcription rate of *Escherichia coli*'s RNA polymerase. *Proceedings of the National Academy of Sciences*. 2015;112(3):743–748.
45. Csilléry K, Blum MG, Gaggiotti OE, François O. Approximate Bayesian computation (ABC) in practice. *Trends in ecology & evolution*. 2010;25(7):410–418.
46. Vogel U, Jensen KF. The RNA chain elongation rate in *Escherichia coli* depends on the growth rate. *Journal of bacteriology*. 1994;176(10):2807–2813.
47. Ryals J, Little R, Bremer H. Temperature dependence of RNA synthesis parameters in *Escherichia coli*. *Journal of bacteriology*. 1982;151(2):879–887.
48. Richardson JP, Greenblatt J. Control of RNA chain elongation and termination. *Escherichia coli and Salmonella: cellular and molecular biology*. 1996;1:822–848.
49. Mason PB, Struhl K. Distinction and relationship between elongation rate and processivity of RNA polymerase II in vivo. *Molecular cell*. 2005;17(6):831–840.
50. Iost I, Guillerez J, Dreyfus M. Bacteriophage T7 RNA polymerase travels far ahead of ribosomes in vivo. *Journal of bacteriology*. 1992;174(2):619–622.
51. Bonner G, Lafer EM, Sousa R. Characterization of a set of T7 RNA polymerase active site mutants. *Journal of Biological Chemistry*. 1994;269(40):25120–25128.
52. Anand VS, Patel SS. Transient state kinetics of transcription elongation by T7 RNA polymerase. *Journal of Biological Chemistry*. 2006;281(47):35677–35685.
53. Kireeva ML, Nedialkov YA, Cremona GH, Purtov YA, Lubkowska L, Malagon F, et al. Transient reversal of RNA polymerase II active site closing controls fidelity of transcription elongation. *Molecular cell*. 2008;30(5):557–566.
54. Rambaut A, Drummond A. Tracer 1.6. University of Edinburgh, Edinburgh. UK. Technical report; 2013.
55. Gelman A, Rubin DB, et al. Inference from iterative simulation using multiple sequences. *Statistical science*. 1992;7(4):457–472.
56. Brooks SP, Gelman A. General methods for monitoring convergence of iterative simulations. *Journal of computational and graphical statistics*. 1998;7(4):434–455.

57. Brooks S, Gelman A, Jones G, Meng XL. Handbook of markov chain monte carlo. CRC press; 2011.
58. Nedialkov YA, Nudler E, Burton ZF. RNA polymerase stalls in a post-translocated register and can hyper-translocate. *Transcription*. 2012;3(5):260–269.
59. Kireeva M, Trang C, Matevosyan G, Turek-Herman J, Chasov V, Lubkowska L, et al. RNA–DNA and DNA–DNA base-pairing at the upstream edge of the transcription bubble regulate translocation of RNA polymerase and transcription rate. *Nucleic acids research*. 2018;46(11):5764–5775.
60. Guajardo R, Lopez P, Dreyfus M, Sousa R. NTP concentration effects on initial transcription by T7 RNAP indicate that translocation occurs through passive sliding and reveal that divergent promoters have distinct NTP concentration requirements for productive initiation. *Journal of molecular biology*. 1998;281(5):777–792.
61. Arnold S, Siemann M, Scharnweber K, Werner M, Baumann S, Reuss M, et al. Kinetic modeling and simulation of in vitro transcription by phage T 7 RNA polymerase. *Biotechnology and bioengineering*. 2001;72(5):548–561.
62. Wong F, Dutta A, Chowdhury D, Gunawardena J. Structural conditions on complex networks for the Michaelis–Menten input–output response. *Proceedings of the National Academy of Sciences*. 2018;115(39):9738–9743.
63. Briggs GE, Haldane JBS. A note on the kinetics of enzyme action. *Biochemical journal*. 1925;19(2):338.
64. English BP, Min W, Van Oijen AM, Lee KT, Luo G, Sun H, et al. Ever-fluctuating single enzyme molecules: Michaelis-Menten equation revisited. *Nature chemical biology*. 2006;2(2):87–94.
65. Schnell S. Validity of the Michaelis–Menten equation—steady-state or reactant stationary assumption: that is the question. *The FEBS journal*. 2014;281(2):464–472.
66. Tennyson CN, Klamut HJ, Worton RG. The human dystrophin gene requires 16 hours to be transcribed and is cotranscriptionally spliced. *Nature genetics*. 1995;9(2):184–190.
67. Darzacq X, Shav-Tal Y, De Turrís V, Brody Y, Shenoy SM, Phair RD, et al. In vivo dynamics of RNA polymerase II transcription. *Nature structural & molecular biology*. 2007;14(9):796–806.
68. Kainov DE, Lísál J, Bamford DH, Tuma R. Packaging motor from double-stranded RNA bacteriophage ϕ 12 acts as an obligatory passive conduit during transcription. *Nucleic acids research*. 2004;32(12):3515–3521.
69. Makarova OV, Makarov EM, Sousa R, Dreyfus M. Transcribing of *Escherichia coli* genes with mutant T7 RNA polymerases: stability of lacZ mRNA inversely correlates with polymerase speed. *Proceedings of the National Academy of Sciences*. 1995;92(26):12250–12254.
70. Mejia YX, Mao H, Forde NR, Bustamante C. Thermal probing of *E. coli* RNA polymerase off-pathway mechanisms. *Journal of molecular biology*. 2008;382(3):628–637.

71. Burns CM, Richardson LV, Richardson JP. Combinatorial effects of NusA and NusG on transcription elongation and rho-dependent termination in *Escherichia coli*. *Journal of molecular biology*. 1998;278(2):307–316.
72. Kingston R, Nierman W, Chamberlin M. A direct effect of guanosine tetraphosphate on pausing of *Escherichia coli* RNA polymerase during RNA chain elongation. *Journal of Biological Chemistry*. 1981;256(6):2787–2797.
73. Galburt EA, Grill SW, Bustamante C. Single molecule transcription elongation. *Methods*. 2009;48(4):323–332.
74. Usala SJ, Brownstein BH, Haselkorn R. Displacement of parental RNA strands during in vitro transcription by bacteriophage $\varphi 6$ nucleocapsids. *Cell*. 1980;19(4):855–862.
75. Moler C, Van Loan C. Nineteen dubious ways to compute the exponential of a matrix, twenty-five years later. *SIAM review*. 2003;45(1):3–49.
76. Bai L, Wang MD. Comparison of pause predictions of two sequence-dependent transcription models. *Journal of Statistical Mechanics: Theory and Experiment*. 2010;2010(12):P12007.
77. Jahnke T, Huisinga W. Solving the chemical master equation for monomolecular reaction systems analytically. *Journal of mathematical biology*. 2007;54(1):1–26.
78. Ratmann O, Jørgensen O, Hinkley T, Stumpf M, Richardson S, Wiuf C. Using likelihood-free inference to compare evolutionary dynamics of the protein networks of *H. pylori* and *P. falciparum*. *PLoS Computational Biology*. 2007;3(11):e230.
79. Brewer BJ, Foreman-Mackey D. DNest4: Diffusive Nested Sampling in C++ and Python. *arXiv preprint arXiv:160603757*. 2016;.