



THE UNIVERSITY *of* EDINBURGH

Edinburgh Research Explorer

## Model reduction for the chemical master equation

**Citation for published version:**

Öcal, K, Sanguinetti, G & Grima, R 2023, 'Model reduction for the chemical master equation: An information-theoretic approach', *Journal of Chemical Physics*, vol. 158, no. 11, 114113.  
<https://doi.org/10.1063/5.0131445>

**Digital Object Identifier (DOI):**

[10.1063/5.0131445](https://doi.org/10.1063/5.0131445)

**Link:**

[Link to publication record in Edinburgh Research Explorer](#)

**Document Version:**

Peer reviewed version

**Published In:**

Journal of Chemical Physics

**General rights**

Copyright for the publications made accessible via the Edinburgh Research Explorer is retained by the author(s) and / or other copyright owners and it is a condition of accessing these publications that users recognise and abide by the legal requirements associated with these rights.

**Take down policy**

The University of Edinburgh has made every reasonable effort to ensure that Edinburgh Research Explorer content complies with UK legislation. If you believe that the public display of this file breaches copyright please contact [openaccess@ed.ac.uk](mailto:openaccess@ed.ac.uk) providing details, and we will remove access to the work immediately and investigate your claim.



# Model Reduction for the Chemical Master Equation: an Information-Theoretic Approach

Kaan Öcal<sup>1,2</sup>, Guido Sanguinetti<sup>3</sup>, and Ramon Grima<sup>2,†</sup>

<sup>1</sup>School of Informatics, University of Edinburgh, Edinburgh EH8 9AB, UK

<sup>2</sup>School of Biological Sciences, University of Edinburgh, Edinburgh EH9 3JH, UK

<sup>3</sup>Scuola Internazionale Superiore di Studi Avanzati, 34136 Trieste, Italy

<sup>†</sup>Corresponding author: [ramon.grima@ed.ac.uk](mailto:ramon.grima@ed.ac.uk)

## Abstract

The complexity of mathematical models in biology has rendered model reduction an essential tool in the quantitative biologist's toolkit. For stochastic reaction networks described using the Chemical Master Equation, commonly used methods include time-scale separation, the Linear Mapping Approximation and state-space lumping. Despite the success of these techniques, they appear to be rather disparate and at present no general-purpose approach to model reduction for stochastic reaction networks is known. In this paper we show that most common model reduction approaches for the Chemical Master Equation can be seen as minimising a well-known information-theoretic quantity between the full model and its reduction, the Kullback-Leibler divergence defined on the space of trajectories. This allows us to recast the task of model reduction as a variational problem that can be tackled using standard numerical optimisation approaches. In addition we derive general expressions for the propensities of a reduced system that generalise those found using classical methods. We show that the Kullback-Leibler divergence is a useful metric to assess model discrepancy and to compare different model reduction techniques using three examples from the literature: an autoregulatory feedback loop, the Michaelis-Menten enzyme system and a genetic oscillator.

**Keywords:** Chemical Master Equation · Model Reduction · Systems Biology

## 1 Introduction

Stochastic biochemical reaction networks such as those involved in gene expression, immune response or cellular signalling [1–4] are often described using the Chemical Master Equation (CME). The CME describes the dynamics of biochemical processes on a mesoscopic level, viewing them as a discrete collection of molecules interacting and undergoing reactions stochastically; as such it is generally considered more accurate than continuum approximations such as rate equations and the Chemical Langevin Equation [4]. Despite its explanatory power, the CME poses significant analytical and computational difficulties to modellers that have limited its use in practice. Closed-form solutions to the CME are difficult to obtain and are only known for a small number of biologically relevant systems, and solving the CME numerically requires using approximations such as the Finite State Projection (FSP) [5]. Numerical approaches tend to scale poorly with the number of species and reactions present in a system, and as a result there is significant interest in finding ways to simplify a description of a stochastic reaction network that make it easier to analyse and study - this is the goal of model reduction.

Model reduction for deterministic and continuum-limit models in biology is an active research topic [6, 7], but very few existing methods can be applied to the discrete, stochastic setting of the CME. The Quasi-Steady State Approximation (QSSA) is perhaps the best known technique, first considered in the stochastic case in [8]. Here the system is partitioned into ‘slow’ and ‘fast’ species such that the fast species evolve very quickly on the timescale of the slow species. On the slow timescale the states of the fast species can therefore be approximated by their steady-state value (conditioned on the slow species), effectively allowing a description of the system in terms of the

slow species only. By its nature the QSSA is only applicable to systems with a clear separation of timescales between species, the existence of which cannot always be established. The QSSA for stochastic systems is generally believed to require more stringent conditions than in the deterministic case, but the exact validity conditions are not well-understood [9–14]. Recent work [15] has analysed a modification of the QSSA, the total QSSA (tQSSA), in systems involving reversible binding of molecules and shown that it has a wide range of applicability in the stochastic case.

Similar to the QSSA is the Quasiequilibrium Approximation (QEA), which was first considered in [16, 17] for stochastic reaction networks. Here the reaction network is decomposed into ‘slow’ and ‘fast’ *reactions*, and the fast reactions are assumed to equilibrate rapidly on the timescale of the slow reactions. Similar to the QSSA, the QEA can be used to reduce the number of species and reactions in a system, but it relies on the existence of a clear timescale separation between reactions, which is not always present for large systems with many distinct reactions. Much like the QSSA, the validity of the QEA for systems without the appropriate timescale separation has not been generally established, and from the asymptotic nature of the descriptions it is not usually possible to quantify the approximation error. Despite this, both the QSSA and the QEA are by far the most commonly used model-reduction technique for chemical reaction networks owing to their physical interpretability and analytical tractability, most famously in the Michaelis-Menten model of enzyme kinetics.

A distinct approach for model reduction with the Chemical Master Equation is state-space lumping, which originates from the theory of finite Markov chains, see e.g. [18]. Here different states in a system are lumped together such that the coarse-grained system is still Markovian and can be modelled using the CME. For a typical biochemical reaction network it may not be possible to perform any nontrivial lumping while preserving Markovianity, whence approximate lumping methods have been considered e.g. in [19–21]. Here the coarse-grained system is approximated by a Markov process, and the approximation error quantified using the KL divergence between the original model and a lift of the approximation to the state space of the original model. State-space lumping for the CME often occurs in the context of the QEA, as states related by fast reactions are effectively lumped together, or averaged [22–24]. For this reason we will not consider this approach separately, although many of our considerations, such as the optimal form of the lumped propensity functions, extend to state-space lumping.

Finally, a more recent model reduction technique specifically for gene expression systems is the Linear Mapping Approximation (LMA) [25]. The LMA replaces bimolecular reactions of the form  $G + P \rightarrow (\dots)$ , where  $G$  is a binary species such as a gene, by a linear approximation where  $P$  is replaced by its mean conditioned on  $[G] = 1$ . While the LMA does not reduce the number of species or reactions, reaction networks with linear reactions are generally easier to analyse: their moments can be computed exactly, and closed-form solutions are known for many cases [26–29].

At a first glance these approaches - timescale separation, state space lumping and the Linear Mapping Approximation - seem rather disparate, and it is unclear what, if any, relationships exist between them. In this paper we show that all of these methods can be viewed as minimising a natural information theoretic quantity, the Kullback-Leibler (KL) divergence, between the full and reduced models. In particular they can be seen as maximal likelihood approximations of the full model, and one can assess the quality of the approximation in terms of the resulting KL divergence. Based on these results we show how the KL divergence can be estimated and minimised numerically, therefore providing an automated method to choose between different candidate reductions of a given model in situations where none of the above model reduction techniques are classically applicable. In its full generality, the Chemical Master Equation describing a reduced model can be obtained by marginalisation of the original CME, and hence our approach recovers the method based on condi-

tional expectations presented in [30].

The KL divergences we consider in this paper are computed on the space of trajectories, and as such include both static information and dynamical information, in contrast to purely distribution-matching approaches. The KL divergence and similar information-theoretic measures between continuous-time Markov chains have previously been considered in [31, 32] in the context of variational inference (with the true model and the approximation reversed compared to our approach), in [33, 34] to obtain approximate non-Markovian reductions, in [35] to analyse information flow for stochastic reaction networks and in [36] in quantifying model discrepancy for Markovian agent-based models.

In Section 2 we introduce the mathematical framework in which we consider model reduction for the Chemical Master Equation, based on KL divergences between continuous-time Markov chains. We show how the KL divergence can be minimised analytically in some important cases, recovering standard results in the literature and providing a mathematical justification for commonly used mean-field arguments as in the QSSA, the QEA and the LMA. We furthermore provide numerical algorithms for estimating as well as minimising the KL divergence in cases where analytical solutions are not available. In Section 3 we illustrate the use of KL divergences as a metric for approximation quality using three biologically relevant examples: an autoregulatory feedback loop exhibiting critical behaviour, Michaelis-Menten enzyme kinematics, where we reanalyse the QSSA and the QEA, and an oscillatory genetic feedback system taken from [11], for which we compare different reductions using our approach. Finally in Section 4 we discuss our observations and how our approach could be used as a stepping-stone towards automated reduction of complex biochemical reaction pathways.

## 2 Methods

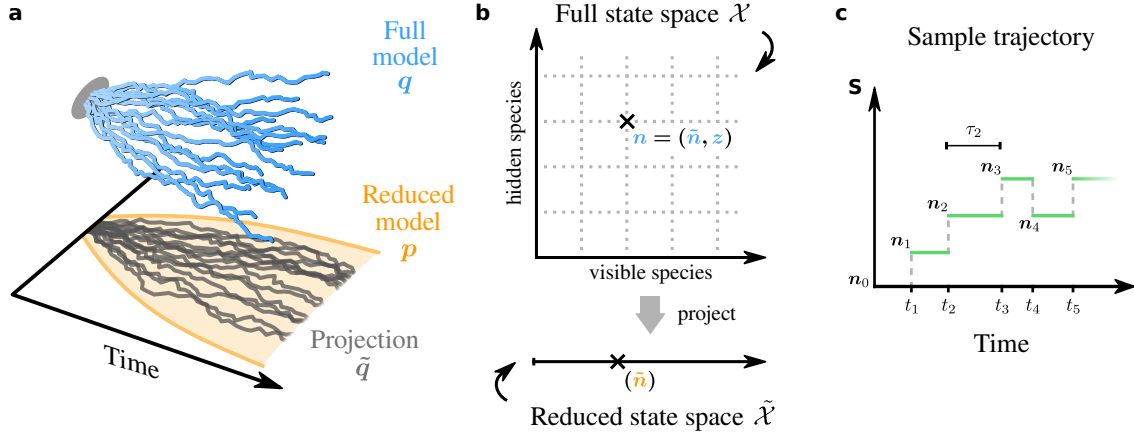
### 2.1 Stochastic Reaction Networks

The Chemical Master Equation describes a biochemical reaction network as a stochastic process on a discrete state space  $\mathcal{X}$ . We will use the letter  $q$  to denote such a stochastic process, which for the purposes of this paper can be seen as a probability distribution over *trajectories* on  $\mathcal{X}$ . For a biochemical reaction network the state space will be  $\mathcal{X} = \mathbb{N}^s$ , where  $s$  is the number of species in the system: every state consist of a tuple  $\mathbf{n} = (n_1, \dots, n_s)$  of  $s$  integers describing the abundances of each species.

Since the model  $q$  can consist of many species interacting in complicated ways, we often want to find a reduction that is more tractable, yet approximates  $q$  as closely as possible. The reduced model, which we will call  $p$ , should be of the same form as  $q$ , i.e. described by the CME, but will typically involve fewer species and simpler reactions. In particular  $p$  can be defined on a lower-dimensional state space  $\tilde{\mathcal{X}}$ . A state  $\mathbf{n}$  in the original model can then be described by its projection  $\tilde{\mathbf{n}}$  onto this lower-dimensional space, together with some unobserved components, which we will denote  $\mathbf{z}$ . See Fig. 1a and b for an illustration.

In this paper we assume that the basic structure of  $p$  is known *a priori*, i.e. the species and reactions we wish to retain are fixed. Our approach to model reduction therefore consists in finding the optimal propensity functions for the reduced model, and we shall see how this can give rise to various known approximations such as the QSSA, the QEA or the LMA depending on what reductions are performed. We will return to the related problem of choosing the structure of the reduced model  $p$  in the discussion.

The original model  $q$  defines a probability distribution over trajectories in  $\mathcal{X}$ , and projecting each trajectory onto the chosen reduced state space we get the (exact) projection of  $q$  onto this space, which we denote  $\tilde{q}$  (Fig. 1a). This is a stochastic process on  $\tilde{\mathcal{X}}$  that is generally not Markovian and



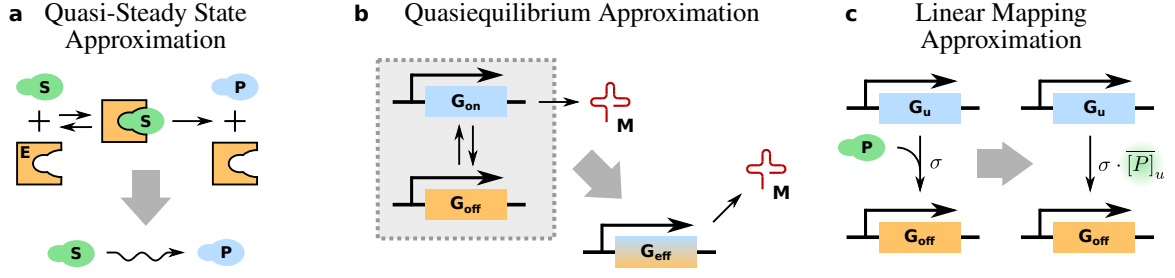
**Figure 1:** Model reduction for the Chemical Master Equation. **(a)** Model reduction approximates a high-dimensional model  $q$  by a lower-dimensional version. Since the direct projection  $\tilde{q}$  of the full model is not easy to describe, we approximate it using a family of tractable candidate models: in this paper, the approximation  $p$  is described by the CME. **(b)** Comparison of the full state space of a system, consisting of two species  $S_1$  and  $S_2$ , and a reduced state space containing  $S_1$  only. Species which are not deemed essential can be discarded in the reduction and become unobserved variables. The dynamics of the original system involves all species, whereas the reduced model aims at an effective description only in terms of the reduced species. **(c)** Sample trajectory for a one-dimensional system, defined by the sequence  $\mathbf{n}_0, \mathbf{n}_1, \dots$  of states visited and the corresponding jump times  $t_1, t_2, \dots$  (or alternatively the waiting times  $\tau_0, \tau_1, \dots$ ).

thus cannot be modelled using the CME. We aim to find a tractable approximation  $p$  to  $\tilde{q}$  that can be described using the CME, and we will do this by minimising the KL divergence  $\text{KL}(\tilde{q} \| p)$  between the two models on the space of trajectories. Several well-known examples of model reduction for the CME are illustrated in Fig. 2.

Jumps in  $q$  come in two kinds: those that affect the observed species  $\tilde{n}$ , which we will call visible jumps, and those that only change  $z$ , which we call hidden. The jumps in  $\tilde{q}$  correspond to visible jumps in  $q$ . In the context of the CME, jumps are typically grouped into reactions with fixed stoichiometry, often also called reaction channels, and we can similarly distinguish visible and hidden reactions in  $q$ . We will always assume that different reactions have different stoichiometries, so that every jump in  $q$  and  $p$  corresponds to a unique reaction. Reactions with the same stoichiometry can always be combined by summing their propensities.

We introduce some more notation at this point, which is summarised in Table 1 and illustrated in Fig. 1c. A single realisation, or trajectory, of  $q$  is defined by the sequence of states  $\mathbf{n}_0, \mathbf{n}_1, \mathbf{n}_2, \dots \in \mathcal{X}$  visited and jump times  $0 < t_1 < t_2 < \dots$ . We will write  $\mathbf{n}_{[0,T]} = \{\mathbf{n}(t)\}_{0 \leq t \leq T}$  for a trajectory, where  $\mathbf{n}(0) = \mathbf{n}_0$  and  $\mathbf{n}(t) = \mathbf{n}_i$  for  $t_i \leq t < t_{i+1}$ , and denote by  $\tau_i = t_{i+1} - t_i$  the waiting times between jumps.

For a continuous-time Markov process  $p$ , e.g. one defined using the CME, we denote the transition rate from state  $\mathbf{n}$  to  $\mathbf{m} \neq \mathbf{n}$  by  $p_{\mathbf{m} \leftarrow \mathbf{n}}$ . We let  $p_{\leftarrow \mathbf{n}} = \sum_{\mathbf{m} \neq \mathbf{n}} p_{\mathbf{m} \leftarrow \mathbf{n}}$  be the total transition rate out of  $\mathbf{n}$ . The transition probabilities at  $\mathbf{n}$  are then given by  $p_{\mathbf{m} | \mathbf{n}} = p_{\mathbf{m} \leftarrow \mathbf{n}} / p_{\leftarrow \mathbf{n}}$ . For completeness we define  $p_{\mathbf{n} | \mathbf{n}} := 0$ .



**Figure 2:** Common model reduction techniques for the Chemical Master Equation. **(a)** The QSSA eliminates intermediate species which evolve on a faster timescale than others, replacing them by their steady-state values. In the case of the pictured Michaelis-Menten enzyme system this is often applied to the enzyme  $E$  and the substrate-enzyme complex  $ES$ . **(b)** The QEA is an analogue of the QSSA that can be applied when a reaction and its reverse equilibrate rapidly on timescales of interest. This can occur e.g. when a gene switches rapidly between states. **(c)** The LMA replaces protein-gene binding reactions by effective unimolecular reactions, where the concentration of proteins is approximated by its mean conditioned on the gene state. While this does not remove any species from the system, it considerably simplifies the propensities of the reactions.

## 2.2 Minimising the KL Divergence

The Kullback-Leibler Divergence between two distributions  $\tilde{q}$  and  $p$  is defined as

$$\text{KL}(\tilde{q} \parallel p) = \int \tilde{q}(x) \log \frac{\tilde{q}(x)}{p(x)} dx = \mathbb{E}_{\tilde{q}} [\log \tilde{q}(x)] - \mathbb{E}_{\tilde{q}} [\log p(x)] \quad (1)$$

$$= \mathbb{E}_{\tilde{q}} [\log \tilde{q}] + H(\tilde{q}; p) \quad (2)$$

where the quantity  $H(\tilde{q}; p)$  is known as the cross-entropy:

$$H(\tilde{q}; p) = -\mathbb{E}_{\tilde{q}} [\log p(x)] \quad (3)$$

Minimising the KL divergence with respect to  $p$  is therefore equivalent to minimising the cross-entropy. Alternatively, it is equivalent to maximising the average log-likelihood under  $p$  of samples drawn from  $\tilde{q}$ , which is a statistical inference problem.

In this paper samples from  $\tilde{q}$  and  $p$  are trajectories. We therefore need to compute the log-likelihood of an arbitrary trajectory  $\mathbf{n}_{[0,T]}$  under  $p$ , which by assumption is a stochastic reaction network described by the CME. In Appendix A we show that the log-likelihood for a trajectory

Symbol	Explanation	Symbol	Explanation
$q$	Full model	$q_{m \leftarrow n}$	Transition rate
$\tilde{q}$	Projected model	$q_{\leftarrow n}$	Total transition rate
$p$	Reduced model	$q_{m n}$	Transition probability
$\mathbf{n}$	Full state vector	$q_0(\mathbf{n})$	Initial distribution
$\tilde{\mathbf{n}}$	Reduced state vector	$q_t(\mathbf{n})$	Single-time marginal
$\mathbf{z}$	Unobserved state vector		

**Table 1:** Notation used in this paper.

visiting the states  $\mathbf{n}_0, \mathbf{n}_1, \dots, \mathbf{n}_k$  with waiting times  $\tau_0, \tau_1, \dots, \tau_k$  is given by

$$\log p(\mathbf{n}_{[0,T]}) = \log p_0(\mathbf{n}_0) - \sum_{i=0}^k \tau_i p_{\leftarrow \mathbf{n}_i} + \sum_{i=1}^k \log p_{\mathbf{n}_i \leftarrow \mathbf{n}_{i-1}}. \quad (4)$$

Note that computing the likelihood of a fully observed trajectory is easier than computing the likelihood of a trajectory sampled at discrete time points, which is the form of data typically observed in biological experiments. The difference is that computing the likelihood in the latter case requires integrating over all possible courses of the trajectory between observations; this integral is computed implicitly by the Chemical Master Equation, which is generally hard to solve. For a fully observed trajectory there are no hidden variables to integrate out, which considerably simplifies the task of computing likelihoods.

Returning to the problem of minimizing the KL divergence (1), or equivalently the cross entropy, we need to compute expectations over  $\tilde{q}$  that are not in general available in closed form. It is easier to approximate these expectations by simulating  $N$  trajectories  $\tilde{\mathbf{n}}_{[0,T]}^{(i)}$ ,  $i = 1, \dots, N$ , from  $\tilde{q}$  and computing the estimate

$$\widehat{H}(\tilde{q}; p_\theta) = -\frac{1}{N} \sum_{i=1}^N \log p(\tilde{\mathbf{n}}_{[0,T]}^{(i)}). \quad (5)$$

We can then minimise this with respect to the reduced model parameters by gradient descent, noting that the cross-entropy is a convex function of  $p$ . Formulæ for the gradients can be computed by differentiating Eq. (4) by hand, or by using automatic differentiation software. Minimising the estimated cross-entropy via gradient descent yields a probabilistic estimate for the optimal parameters, which will generally converge to the true solution as  $N \rightarrow \infty$ . Summarising the contents of this section we arrive at Algorithm 1 for automatically fitting the optimally reduced model.

---

**Algorithm 1** Model reduction via gradient descent.

---

**Inputs:** number of simulations  $N$ , simulation length  $T$ , full model  $q$ , model family  $\{p_\theta\}$ , learning rate  $\eta$

**Output:** reduced model parameters  $\widehat{\theta}$

---

```

for all  $i = 1, \dots, N$  do
  sample  $\mathbf{n}_{[0,T]}^{(i)}$  from  $q$ 
  project  $\mathbf{n}_{[0,T]}^{(i)}$  to  $\tilde{\mathbf{n}}_{[0,T]}^{(i)}$ 
end
initialize parameters  $\hat{\theta}$ 
while not converged do
  compute  $\widehat{H}(\tilde{q}; p_{\hat{\theta}})$  using (4), (5)
  compute  $\nabla_{\theta} \widehat{H}(\tilde{q}; p_{\hat{\theta}})$ 
   $\hat{\theta} \leftarrow \hat{\theta} - \eta \nabla_{\theta} \widehat{H}(\tilde{q}; p_{\hat{\theta}})$ 
return  $\hat{\theta}$ 

```

---

This algorithm has the advantage of being completely general, but unlike standard model reduction algorithms it involves numerical optimisation. In the next section we will analyse our loss function more closely to obtain alternative expressions for the optimum, allowing us to bypass Alg. 1 for a wide range of problems.



Instead of minimising  $\text{KL}(\tilde{q} \| p)$ , an alternative approach would minimise the opposite KL divergence  $\text{KL}(p \| \tilde{q})$ . The latter has been analysed e.g. in [31, 32] for continuous-time Markov processes. Despite the almost symmetric definition, it is well-known that these two quantities behave rather differently in an optimisation context and lead to very different results, cf. [37]. From a pragmatic point of view, minimising  $\text{KL}(\tilde{q} \| p)$  involves computing likelihoods under  $p$ , while minimising  $\text{KL}(p \| \tilde{q})$  involves computing probabilities of trajectories under  $\tilde{q}$ , which is significantly more difficult.

### Example: Telegraph Model

We will illustrate the above by means of a well-known example, the telegraph model of gene transcription [38], simplified to neglect degradation and depicted in Fig. 3. The full model  $q$  consists of three species: a gene in the on state ( $G_{\text{on}}$ ) and the off state ( $G_{\text{off}}$ ), together with mRNA ( $M$ ). The gene switches between both states with constant activation rate  $\sigma_{\text{on}}$  and inactivation rate  $\sigma_{\text{off}}$ , and in the active state produces mRNA with rate  $\rho_{\text{on}}$ . The reduced model  $p$  consists only of mRNA ( $M$ ), which is produced at constant rate  $\rho_{\text{eff}}$ . This instructive example has previously been treated in [33] in a similar context.

A state in the full model can be represented as  $\mathbf{n} = (g, m)$ , where  $g$  is the gene state (on or off) and the  $m$  is the number of mRNA present, while a reduced state  $\tilde{\mathbf{n}} = (m)$  is given by the number of mRNA molecules only, with the unobserved component  $\mathbf{z} = (g)$  being the gene state. The only reaction in  $q$  that descends to  $\tilde{q}$  is mRNA production, whereas the two gene switching reactions are not observed. Note that the projection  $\tilde{q}$  of the telegraph model is not Markovian, as the instantaneous mRNA production rate depends on the current gene state  $g$ , but we can sample from  $\tilde{q}$  by simulating  $q$  and discarding the information about the gene state.

The reduced model  $p$  is a Poisson process whose only parameter is the production rate  $\rho_{\text{eff}}$ . The only trajectories possible under  $p$  are those starting at  $m = 0$  and increasing by one at every jump. The log-likelihood of such a trajectory  $\tilde{\mathbf{n}}_{[0,T]}$  under  $p$  is given by

$$\log p(\tilde{\mathbf{n}}_{[0,T]}) = \tilde{\mathbf{n}}(T) \log \rho_{\text{eff}} - \rho_{\text{eff}} T, \quad (6)$$

where  $\tilde{\mathbf{n}}(T)$  is the total number of mRNA produced up to time  $T$ , as can be verified using Eq. (4). The log-likelihood of any other trajectory is  $-\infty$ .

Note that Eq. (6) describes the likelihood for an entire trajectory, as opposed to the likelihood of observing  $\mathbf{n}_T$  molecules at time  $T$ . The latter follows a Poisson distribution with rate  $\rho_{\text{eff}} T$  and can be obtained by integrating Eq. (4) over all trajectories that end at the same  $\mathbf{n}_T$ . Indeed, conditioned on observing  $\mathbf{n}_T$  mRNA molecules at time  $T$  their individual production times are uniformly and independently distributed on  $[0, T]$  by the properties of the Poisson process, and integrating Eq. (6) over all possible combinations of production times we recover the usual Poisson likelihood, keeping in mind that any permutation of the production times yields the same trajectory.

The reduced model  $p$  can model every trajectory obtained from  $\tilde{q}$ , so the log-likelihood of any sample from  $\tilde{q}$  is finite. The mRNA transcription reactions in both models correspond. If  $q$  were to include e.g. mRNA degradation as is usual in the literature, some trajectories from the full model would feature decreases in  $\tilde{\mathbf{n}}$  and be impossible under  $p$ ; in this case the KL divergence would be infinite, which signals that the reduced model  $p$  is not appropriate and needs to be extended to include degradation.

For this simple example the likelihood of a reduced trajectory under  $p$  only depends on the total mRNA produced. We can therefore compute the cross-entropy explicitly:

$$H(\tilde{q}; p) = -\mathbb{E}_{\tilde{q}}[\tilde{\mathbf{n}}(T)] \log \rho_{\text{eff}} + \rho_{\text{eff}} T. \quad (7)$$



Minimising this with respect to  $\rho_{\text{eff}}$  yields the optimum

$$\rho_{\text{eff}}^* = \frac{1}{T} \mathbb{E}_{\tilde{q}} [\bar{\mathbf{n}}(T)], \quad (8)$$

which is the average mRNA production rate on the interval  $[0, T]$ . Thus the optimal approximation to the telegraph model is obtained by setting the mRNA production rate to its mean value. This is the result we obtain by applying the QEA to this system (see 2.5.2), which suggests that the approximation will be better if  $\sigma_{\text{on}}$  and  $\sigma_{\text{off}}$  are large compared to  $\rho_{\text{on}}$ . We will see how to evaluate the approximation quality in 2.4.

Even if the analytical minimisation above were not possible, the above minimisation can be performed using Alg. 1. To do so we start by sampling a fixed number  $N$  of mRNA trajectories from the telegraph model; the more trajectories we use, the better our estimate of  $\rho_{\text{eff}}$  will be. The next step is constructing the cross-entropy loss function for the reduced model, which only depends on one parameter  $\rho_{\text{eff}}$  and is given by Eq. (7). We then perform gradient descent on the cross-entropy, averaged over all trajectories, to arrive at our estimate of  $\rho_{\text{eff}}$ . Here the learning rate  $\eta$  in Alg. 1 has to be determined by experimentation. The code accompanying this paper contains an example implementation of this procedure in Julia.

### 2.3 Analysing the KL Divergence

In the last section we have derived a numerical procedure to minimise the KL divergence (1), or equivalently the cross-entropy (3). While this yields a direct computational procedure to fit reduced models, we are equally interested in the theoretical insights to be gained from considering model reduction as a variational problem. Understanding how our approach works in more detail will help us establish links with existing model reduction techniques and derive general principles that can help us generalise these to new systems.

As shown in Appendix B, for a Markovian reduced model  $p$  the cross-entropy can be written as

$$H(\tilde{q}; p)_{[0, T]} = H(\tilde{q}_0; p_0) + \int_0^T \sum_{\tilde{\mathbf{n}}} \tilde{q}_t(\tilde{\mathbf{n}}) \left( p_{\leftarrow \tilde{\mathbf{n}}}(t) - \sum_{\tilde{\mathbf{m}} \neq \tilde{\mathbf{n}}} \tilde{q}_{\tilde{\mathbf{m}} \leftarrow \tilde{\mathbf{n}}}(t) \log p_{\tilde{\mathbf{m}} \leftarrow \tilde{\mathbf{n}}}(t) \right) dt \quad (9)$$

$$= H(\tilde{q}_0; p_0) + \int_0^T H(\tilde{q}; p)_t dt, \quad (10)$$

with the instantaneous cross-entropy rate at time  $t$  defined as

$$H(\tilde{q}; p)_t = - \sum_{\tilde{\mathbf{n}}} \tilde{q}_t(\tilde{\mathbf{n}}) \left( \sum_{\tilde{\mathbf{m}} \neq \tilde{\mathbf{n}}} \tilde{q}_{\tilde{\mathbf{m}} \leftarrow \tilde{\mathbf{n}}}(t) \log p_{\tilde{\mathbf{m}} \leftarrow \tilde{\mathbf{n}}}(t) - p_{\leftarrow \tilde{\mathbf{n}}}(t) \right). \quad (11)$$

Here we use the effective transition rates

$$\tilde{q}_{\tilde{\mathbf{m}} \leftarrow \tilde{\mathbf{n}}}(t) = \lim_{\delta t \rightarrow 0} \frac{1}{\delta t} P(\tilde{q}(t + \delta t) = \tilde{\mathbf{m}} | \tilde{q}(t) = \tilde{\mathbf{n}}, q_0) \quad (\tilde{\mathbf{m}} \neq \tilde{\mathbf{n}}). \quad (12)$$

If  $\tilde{q}$  is not Markovian, the transition probability from a state  $\tilde{\mathbf{n}}$  at time  $t$  will be affected by the history of the process, and hence on the initial distribution  $q_0$ . If  $\tilde{q}$  is Markovian, Eq. (12) reduces to the classical transition rate which is independent of  $q_0$ .

Looking at Eq. (9) we make two important observations. As intimated in the discussion of the telegraph example of the previous section, if any trajectory under  $\tilde{q}$  has jumps which are not allowed under  $p$ , i.e. if  $\tilde{q}_{\tilde{\mathbf{m}} \leftarrow \tilde{\mathbf{n}}} \neq 0$  while  $p_{\tilde{\mathbf{m}} \leftarrow \tilde{\mathbf{n}}} = 0$ , then the cross-entropy (9) will be  $\infty$ : the reduced model  $p$

is not flexible enough to model  $\tilde{q}$ . On the other hand, if  $p$  contains transitions which are impossible under  $\tilde{q}$ , i.e. if  $p_{\tilde{m} \leftarrow \tilde{n}} \neq 0$  while  $\tilde{q}_{\tilde{m} \leftarrow \tilde{n}} = 0$ , then these can only increase the cross-entropy (9). This means that the optimal reduced model  $p$  does not contain transitions beyond those in  $\tilde{q}$ , and in the context of the CME we can therefore assume that  $p$  and  $\tilde{q}$  have the same reactions, with different propensities (Markovian for  $p$ , not necessarily for  $\tilde{q}$ ).

If the  $i$ -th reaction in  $p$  has propensity function  $\rho_i(\tilde{\mathbf{n}}; t)$  and net stoichiometry  $s_i$  we obtain the following decomposition of the cross-entropy rate at time  $t$ :

$$H(\tilde{q}; p)_t = - \sum_i \sum_{\tilde{\mathbf{n}}} \tilde{q}_t(\tilde{\mathbf{n}}) [\tilde{q}_{\tilde{n}+s_i \leftarrow \tilde{n}} \log \rho_i(\tilde{\mathbf{n}}; t) - \rho_i(\tilde{\mathbf{n}}; t)], \quad (13)$$

where the first sum is over all reactions in  $p$ , or equivalently all visible reactions in  $q$ . The total cross-entropy is obtained by integrating Eq. (13) over  $[0, T]$ , and we can find the optimal  $p$  by optimising the cross-entropy for each reaction separately.

We can minimize Eq. (13) analytically if the full model  $q$  is Markovian. Assume there is precisely one reaction in  $q$  with net stoichiometry  $\tilde{s}_i$  in the projection, and let  $\sigma_i$  be its propensity function. Differentiating the above equation with respect to  $\rho_i(\tilde{\mathbf{n}}; t)$  and setting the derivative to zero we obtain

$$\rho_i^*(\tilde{\mathbf{n}}; t) = \sum_z q_t(z | \tilde{\mathbf{n}}) \sigma_i(\tilde{\mathbf{n}}, z; t) = \mathbb{E}_z[\sigma_i(\tilde{\mathbf{n}}, z; t) | \tilde{\mathbf{n}}; t]. \quad (14)$$

The optimal propensity for a reaction under  $p$  is the expected propensity under the original model conditioned on the observed state  $\tilde{\mathbf{n}}$ . In particular, if the propensity of the original reaction does not depend on unobserved species it can be taken over directly. We explore the implications of these results on the moments of the reduced system in Section 2.5.4.

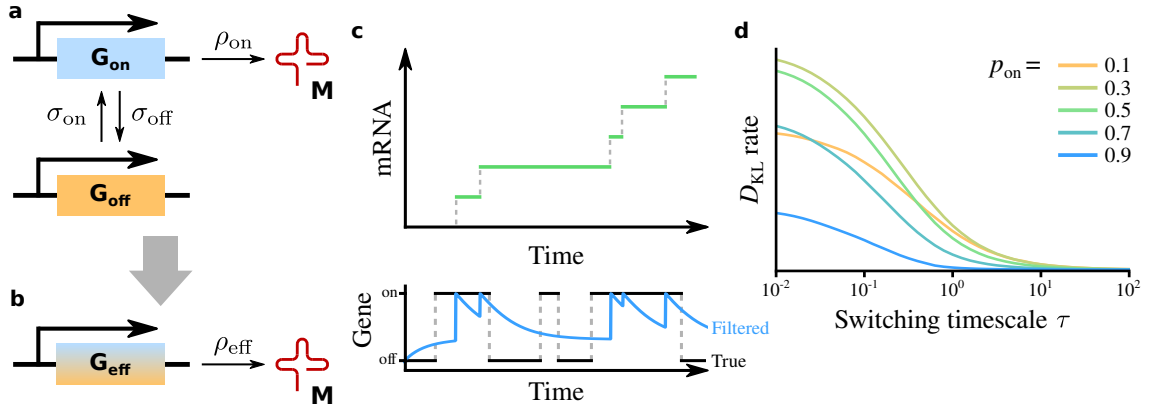
In practice we often place constraints on the reduced propensities  $\rho_i(\tilde{\mathbf{n}}; t)$  such as time-homogeneity and mass-action kinetics, which result in constrained optima. For example, if  $\rho_i$  is taken to be independent of time we have to integrate Eq. (13) over  $[0, T]$ , and as  $T \rightarrow \infty$  it can be verified that the total cross-entropy is minimised for

$$\rho_i^*(\tilde{\mathbf{n}}) = \sum_z q_\infty(z | \tilde{\mathbf{n}}) \sigma_i(\tilde{\mathbf{n}}, z) = \mathbb{E}_z[\sigma_i(\tilde{\mathbf{n}}, z) | \tilde{\mathbf{n}}; t = \infty], \quad (15)$$

which is the steady-state version of Eq. (14).

For other forms of propensity functions the optimum can generally be obtained by a similar averaging procedure. In all cases the main difficulties lie in estimating the relevant conditional expectations, which can be done numerically using the Monte-Carlo approach presented in the previous section. Eq. (14) is perhaps the central result of this paper, and we will see that it coincides with the propensities obtained using the QSSA, the QEA and the LMA in when these are applied. These methods can therefore all be seen as special cases of our variational approach based on minimising the KL divergence. As an aside, our results provide a new derivation of the marginal CME proposed in [30, 34] and the nontrivial result that Eq. (14) and Alg. 1 target equivalent objectives.

In particular Alg. 1 can be bypassed by computing certain conditional moments of a system, which provides an alternative way to perform model reduction. For many systems these conditional moments can be efficiently computed using moment equations, possibly involving moment closure approximations, which have attracted a large amount of literature [4]. The LMA, for example, involves a bootstrapping approach to estimate the relevant moments. Alternatively the conditional expectations could be computed by numerically integrating the CME, although we do not expect this to scale well to large systems. In consequence Alg. 1, while a convenient way to implement our approach in full generality and empirically verify its performance, can often be bypassed.



**Figure 3:** Model reduction illustrated using the telegraph model without degradation. **(a)** Schematic of the model. **(b)** The reduced model is equivalent to a Poisson process with rate  $\rho_{\text{eff}}$ . **(c)** Example trajectory from the telegraph model, showing mRNA numbers (top) and the gene state (bottom) in time. If we only observe mRNA numbers we can infer the current gene state by filtering, obtaining the probability for the gene being on at a given time (blue line). When mRNA numbers increase the on probability jumps to 1 since transcription only happens in that state. **(d)** KL divergence rate between the telegraph model and its Poisson reduction, for various choices of the switching timescale  $\tau = \sigma_{\text{on}} + \sigma_{\text{off}}$  and on probability  $p_{\text{on}} = \sigma_{\text{on}} / (\sigma_{\text{on}} + \sigma_{\text{off}})$ , assuming a fixed transcription rate  $\rho_{\text{on}} = 1$ . The Poisson approximation becomes more accurate as the switching timescale increases compared to the transcription timescale. KL divergences were estimated numerically using Alg. 2.

### Example: Telegraph Model, cont.

Consider the telegraph model from the previous section. If we make the reduced model  $p$  more flexible by allowing a nonlinear and time-dependent propensity  $\rho(m; t) = \rho_{\text{eff}}(m; t)$ , then the cross-entropy between  $\tilde{q}$  and  $p$  can be computed as

$$H(\tilde{q}; p) = \int_0^T \sum_m \tilde{q}_t(m) (\tilde{q}_{m+1 \leftarrow m}(t) \log \rho_{\text{eff}}(m; t) - \rho_{\text{eff}}(m; t)) dt. \quad (16)$$

Note that both  $\tilde{q}$  and  $p$  have the same initial conditions. The effective propensities of  $\tilde{q}$  can be computed as

$$\tilde{q}_{m+1 \leftarrow m}(t) = \rho_{\text{on}} q_t(g = \text{on} | m), \quad (17)$$

which is  $\rho_{\text{on}}$  weighed by the probability that the gene is on at time  $t$  given that there are  $m$  mRNA molecules present. This is the optimal propensity for the reduced model  $p$  in accordance with Eq. (14), featuring a nonlinear dependence on  $m$ , and this can be shown to be the optimal Markovian approximation to  $\tilde{q}$ . We will see in Section 2.5.4 that this approximation exactly matches the mRNA distribution of the full model at all times.

If we require mass-action kinetics for the reduced model, i.e. if we let  $\rho(m; t) = \rho_{\text{eff}}(t)$  for an effective rate constant  $\rho_{\text{eff}}(t)$ , minimising the cross-entropy (16) yields

$$\rho_{\text{eff}}(t) = \rho_{\text{on}} q_t(g = \text{on}). \quad (18)$$

The effective transcription rate learned by  $p$  is the expected transcription rate of  $q$  at time  $t$ , this time averaged over all  $m$ .

Finally if we require  $\rho(m; t) = \rho_{\text{eff}}$  to be time-independent we obtain by a similar procedure

$$\rho_{\text{eff}} = \rho_{\text{on}} \cdot \frac{1}{T} \int_0^T q_t(g = \text{on}) dt, \quad (19)$$

which is just the previous result averaged over  $t$ . This is the special case considered in the previous section, where  $p$  was restricted to be a Poisson process with constant intensity. As  $T \rightarrow \infty$ , the optimum (19) converges to the expected transcription rate of  $q$  at steady state.

It can be checked that the propensity functions (17) and (18) exactly preserve the mean mRNA numbers of the full system. Similarly, the time-independent propensity function (19) preserves the mean mRNA numbers at steady state. For mass-action propensities, however, the variances will be underestimated. This is because the reduced models do not model gene switching, which increases the noise in the system. This is a common occurrence with model reduction and is well-known for the telegraph model, since mRNA numbers are always distributed following a Poisson mixture distribution, which will always have a higher variance than a single Poisson distribution with the same mean [39]. We will further explore how well reduced models capture marginal distributions of the original in Section 2.5.4.

## 2.4 Computing KL Divergences

To this point we have been occupied with minimising the KL divergence (1), or rather the cross-entropy (3), with respect to  $p$ . Perhaps counter-intuitively, minimising the KL divergence is easier than computing it, since the entropy

$$H(\tilde{q})_{[0,T]} = -\mathbb{E}_{\tilde{\mathbf{n}}_{[0,T]}}[\log q(\tilde{\mathbf{n}}_{[0,T]})], \quad (20)$$

can be ignored during optimisation. In order to assess the performance of the reduced model obtained this way we want to compute the full KL divergence, and hence the entropy, explicitly. As for the cross-entropy the integral in Eq. (20) is generally intractable, and we will approximate the entropy by simulating  $N$  samples  $\tilde{\mathbf{n}}_{[0,T]}^{(i)}$  from  $\tilde{q}$  and computing

$$\widehat{H}(\tilde{q}) = -\frac{1}{N} \sum_{i=1}^N \log \tilde{q}(\tilde{\mathbf{n}}_{[0,T]}^{(i)}). \quad (21)$$

Here we face another difficulty, for if  $\tilde{q}$  is not Markovian we cannot use Eq. (4) to compute log-likelihood of a trajectory  $\tilde{\mathbf{n}}_{[0,T]}$ . Since by assumption  $q$  is a Markov process, the projection  $\tilde{q}$  is just a partially observed Markov process and we can use the so-called forward algorithm for Hidden Markov Models [40] to compute the log-likelihood of a trajectory.

The forward algorithm computes the joint probability  $q(\tilde{\mathbf{n}}_{[0,t]}, \mathbf{z}_t)$  for the observed trajectory up to time  $t$  and the current state of the hidden species,  $\mathbf{z}_t$ . As shown in Appendix C, the forward algorithm in this case yields the following set of jump ODEs:

$$q(\tilde{\mathbf{n}}_0, \mathbf{z}_0) = q_0(\tilde{\mathbf{n}}_0, \mathbf{z}_0), \quad (22)$$

$$\begin{aligned} \frac{d}{dt} q(\tilde{\mathbf{n}}_{[0,t]}, \mathbf{z}_t) &= \sum_{\mathbf{z}' \neq \mathbf{z}_t} (q_{(\tilde{\mathbf{n}}_t, \mathbf{z}_t) \leftarrow (\tilde{\mathbf{n}}_t, \mathbf{z}')} q(\tilde{\mathbf{n}}_{[0,t]}, \mathbf{z}') - q_{(\tilde{\mathbf{n}}_t, \mathbf{z}') \leftarrow (\tilde{\mathbf{n}}_t, \mathbf{z}_t)} q(\tilde{\mathbf{n}}_{[0,t]}, \mathbf{z}_t)) \\ &\quad - \sum_i \sigma_i(\tilde{\mathbf{n}}_t, \mathbf{z}_t) q(\tilde{\mathbf{n}}_{[0,t]}, \mathbf{z}_t), \end{aligned} \quad (23)$$

$$\lim_{t \searrow t_k} q(\tilde{\mathbf{n}}_{[0,t]}, \mathbf{z}_t = \mathbf{z}) = \lim_{t \nearrow t_k} (\sigma_{j_k}(\tilde{\mathbf{n}}_k, \tilde{\mathbf{z}}_t = \mathbf{z}) q(\tilde{\mathbf{n}}_{[0,t]}, \mathbf{z}')), \quad (24)$$

with jumps at the jump times of the observed trajectory  $\tilde{\mathbf{n}}_{[0,T]}$ . The first sum in Eq. (23) represents hidden reactions (those that do not affect  $\tilde{\mathbf{n}}$ ), and the second sum is over all visible reactions in  $q$ , with propensity functions given by the  $\sigma_i$ . The marginal likelihood of the observed trajectory can then be computed by summing over the unobserved variables  $\mathbf{z}$

$$\tilde{q}(\tilde{\mathbf{n}}_{[0,T]}) = \sum_{\mathbf{z}} q(\tilde{\mathbf{n}}_{[0,T]}, \mathbf{z}). \quad (25)$$

The above equations take the form of a modified CME where the observed variables are fixed. They can be solved using an adaptation of the Finite State Projection algorithm [5], and in cases where the unobserved state space is finite, exact solutions can be computed numerically. We summarise the procedure to estimate the KL divergence defined in Eq. (1) in Algorithm 2.

---

**Algorithm 2** Computing the projected KL divergence between a full and a reduced model.

---

**Inputs:** number of simulations  $N$ , simulation length  $T$ , full model  $q$ , reduced model  $p$

**Output:** KL divergence estimate  $\widehat{\text{KL}}(\tilde{q} \| p)$

---

```

 $\hat{L} := 0$ 
for all  $i = 1, \dots, N$  do
  sample  $n_{[0,T]}$  from  $q$ 
  project  $n_{[0,T]}$  to  $\tilde{\mathbf{n}}_{[0,T]}$ 
  compute  $p(\tilde{\mathbf{n}}_{[0,T]})$  using (4)
  compute  $\tilde{q}(\tilde{\mathbf{n}}_{[0,T]})$  using (22)–(24), (25).
   $\hat{L} += \tilde{q}(\tilde{\mathbf{n}}_{[0,T]}) - p(\tilde{\mathbf{n}}_{[0,T]})$ 
return  $\hat{L}/N$ 

```

---

As further shown in Appendix C, a corollary of the above equations is the following jump ODE for the marginal distribution itself:

$$\log \tilde{q}(\tilde{\mathbf{n}}_0) = \log \tilde{q}_0(\tilde{\mathbf{n}}(0)), \quad (26)$$

$$\frac{d}{dt} \log \tilde{q}(\tilde{\mathbf{n}}_{[0,t]}) = - \sum_i \mathbb{E}_{z_t} [\sigma_i(\tilde{\mathbf{n}}(t), z_t) | \tilde{\mathbf{n}}_{[0,t]}, q_0; t], \quad (27)$$

$$\lim_{t \searrow t_k} \log \tilde{q}(\tilde{\mathbf{n}}_{[0,t]}) = \lim_{t \nearrow t_k} \left( \log \tilde{q}(\tilde{\mathbf{n}}_{[0,t]}) + \log \mathbb{E}_{z_t} [\sigma_{j_k}(\tilde{\mathbf{n}}_k, \tilde{z}_t) | \tilde{\mathbf{n}}_{[0,t]}, q_0; t] \right). \quad (28)$$

These equations generalise Eqs. (4) to the presence of unobserved species  $\mathbf{z}$ . These equations can be found in various places in the literature, e.g. [33, 35, 41–44] although we are not aware of a consistent nomenclature for this modification of the CME.

The equations for the log-likelihood are linear and can be split into a sum of contributions for each reaction in  $\tilde{q}$ . Integrating over all reduced trajectories  $\tilde{\mathbf{n}}_{[0,T]}$  shows that the entropy  $H(\tilde{q})$  can thus be expressed as a sum of contributions from each reaction in  $\tilde{q}$ . Eq. (13) shows that the same decomposition holds for the cross-entropy  $H(\tilde{q}; p)$ , so entire KL divergence  $\text{KL}(\tilde{q} \| p)$  can be decomposed into contributions coming from the individual reactions, which allows us to differentially analyse the accuracy of the reduced model  $p$  for each reaction separately. We note that all reactions together determine the distribution  $q$  on trajectories over which we integrate, so this decomposition is not strict in practice.

### Example: Telegraph Model (cont. 2)

We can solve the filtering problem for the telegraph model exactly as the unobserved state space is 2-dimensional. Let  $F(i; t) = q(g(t) = i, m_{[0,t]})$ . The telegraph model consists of three reactions, the

hidden switching reactions with parameters  $\sigma_{\text{on}}$  and  $\sigma_{\text{off}}$ , and the visible transcription reaction  $\rho_{\text{on}}$ . Specialising Eq. (23) to this system we obtain

$$\frac{d}{dt} F(\text{on}; t) = -(\rho_{\text{on}} + \sigma_{\text{off}}) F(\text{on}; t) + \sigma_{\text{on}} F(\text{off}; t), \quad (29)$$

$$\frac{d}{dt} F(\text{off}; t) = \sigma_{\text{off}} F(\text{on}; t) - \sigma_{\text{on}} F(\text{off}; t), \quad (30)$$

where  $m_{[0,t]}$  denotes the observed mRNA trajectory. At each mRNA production event we update the probabilities according to Eq. (24), which in this case reads:

$$\lim_{t \searrow t_k} F(\text{on}; t) = \rho_{\text{on}} \cdot \lim_{t \nearrow t_k} F(\text{on}; t), \quad (31)$$

$$\lim_{t \searrow t_k} F(\text{off}; t) = 0. \quad (32)$$

The marginal likelihood of the observed mRNA trajectory is then given by

$$q(m_{[0,t]}) = F(\text{on}; t) + F(\text{off}; t). \quad (33)$$

These derivations can also be found in [33].

As a byproduct the above equations yield the conditional distribution over the gene state at each time  $t$  given the trajectory prior to that time point. These are the filtered distributions and illustrated in Fig. 3c, but we only need the marginal likelihood for our purposes.

Using the above equations to compute the log-likelihood for a large number of trajectories sampled from the telegraph model  $q$  yields a numerical approximation of the entropy  $H(\tilde{q})$ , and together with Eq. (6) we obtain an estimate of the KL divergence  $\text{KL}(\tilde{q} \| p)$ . For time-homogeneous propensities we empirically found the KL divergence to be asymptotically proportional to  $T$ , and we call the proportionality factor the (steady-state) KL divergence rate. Note that the cross-entropy rate can be defined directly using Eq. (11).

In Fig. 3d we show how this KL divergence rate between the telegraph model and its Poisson approximation changes for various choices of  $\sigma_{\text{on}}$  and  $\sigma_{\text{off}}$ . We observe that the KL divergence rate tends to 0 as the switching rates increase; indeed the Poisson approximation can be obtained from the QEA for this example. We further note that the KL divergence is maximal when the gene spends substantial amounts of time in each state; for  $p_{\text{on}}$  close to 1 the gene is almost always active and the system approaches constitutive expression with rate  $\rho$ , while for  $p_{\text{on}}$  close to 0 the gene only activates sporadically and resembles a constitutive gene with a very low expression rate  $\rho \cdot p_{\text{on}}$ . Visually inspecting mRNA trajectories generated in both regimes show that the results become harder to distinguish from the Poisson process.

Our results on the telegraph model without degradation apply almost verbatim when degradation is included: adding the reaction  $M \rightarrow \emptyset$  to the full and reduced models does not affect the KL divergences considered. *A priori* adding degradation changes the probability distribution over trajectories, which can now exhibit decreasing mRNA numbers, but the log-likelihood contributed by the degradation reaction is the same between the telegraph model and its reduction, and therefore the difference cancels in Eqs. (26)–(28) and Eq. (13). Since the log-likelihood contributed by the gene switching and transcription reactions does not depend on mRNA numbers, which do change in the presence of degradation, the total KL divergence is unaffected. In the presence of feedback, degradation would indirectly affect the KL divergence via its effect on mRNA numbers; we refer the interested reader to the recent paper [35] for an analysis of this type of information flow in stochastic biochemical reaction networks. We can readily verify that our observations on the mean and variance of mRNA numbers under the telegraph model and its reduction remain valid in the presence of degradation.



## 2.5 Relationship with Known Approaches

### 2.5.1 The Quasi-Steady State Approximation

In deterministic chemical kinetics the Quasi-Steady State Approximation (QSSA) is a model reduction technique that can be applied when a system can be partitioned into so-called slow species  $\mathbf{n}^s$  and fast species  $\mathbf{n}^f$  such that the fast species  $\mathbf{n}^f$  evolve on a much faster time scale than the slow species. On the timescale on which the slow species evolve, the fast species can therefore be assumed to reach their steady state almost instantaneously. Thus one assumes  $\frac{d}{dt}\mathbf{n}^f = 0$ , which allows for simplification of the remaining equations for the slow species. The QSSA has famously been applied to Michaelis-Menten enzyme kinetics with the enzyme-substrate complex  $ES$  as the fast species, resulting in the classical Michaelis-Menten propensity for product formation (see 3.2 for more information).

The QSSA was extended to the stochastic case in [8]. Here one assumes that conditioned on  $\mathbf{n}^s$ , the fast species reach steady state nearly instantaneously compared to the timescale of interest. In our formulation the reduced state space consists of the slow species,  $\tilde{\mathbf{n}} = \mathbf{n}^s$ , and the unobserved species are the fast species,  $\mathbf{z} = \mathbf{n}^f$ . As shown in [8] the QSSA yields the following reduced CME for the approximation  $p$  (Eqs. (10), (11) in [8]):

$$\frac{d}{dt}p_t(\mathbf{n}^s) = \sum_{i=1}^r [\tilde{\rho}_i(\mathbf{n}^s - \mathbf{S}_i^s; t) p_t(\mathbf{n}^s - \mathbf{S}_i^s) - \tilde{\rho}_i(\mathbf{n}^s; t) p_t(\mathbf{n}^s)], \quad (34)$$

where the reduced propensities  $\tilde{\rho}_i$  are defined as

$$\tilde{\rho}_i(\mathbf{n}^s; t) = \mathbb{E}_{\mathbf{n}^f}[\rho_i(\mathbf{n}^s, \mathbf{n}^f) | \mathbf{n}^s; t]. \quad (35)$$

This agrees precisely with Eq. (14), which illustrates how the QSSA can be seen as minimising the KL divergence (B4). We remark that Eq. (34) is a special case of the phenomenon discussed in Section 2.5.4. We can compare (35) with the true propensities of the projection  $\tilde{q}$ , which depend on the entire history of the observed trajectory and the initial distribution  $q_0$ :

$$\rho_i^{\text{exact}}(\mathbf{n}; t | q_0, \mathbf{n}_{[0,t]}^s) = \mathbb{E}_{\mathbf{n}^f}[\rho_i(\mathbf{n}^s, \mathbf{n}^f) | \mathbf{n}_{[0,t]}^s, q_0; t]. \quad (36)$$

This equation, which describes the instantaneous reaction rate of the non-Markovian process  $\tilde{q}$ , has previously been considered e.g. in [8, 33].

Some intuition for the relationship between timescale separation and our approach can be gained from equations Eqs. (26)–(28), which describe the probability of a reduced trajectory  $\tilde{\mathbf{n}}_{[0,T]}$  under  $\tilde{q}$ . Reaction propensities under  $q$  in general depend on the unobserved species, whose distribution is correlated with the history of the current trajectory. If we assume that the timescale of the unobserved species  $\mathbf{n}^f$  is very fast, then this correlation will decay very quickly and the time-dependence of the conditional distributions in Eqs. (26)–(28) will be negligible, so we can replace these equations by

$$\log \tilde{q}(\tilde{\mathbf{n}}_0) = \log \tilde{q}_0(\tilde{\mathbf{n}}(0)), \quad (37)$$

$$\frac{d}{dt} \log \tilde{q}(\tilde{\mathbf{n}}_{[0,t]}) = - \sum_i \mathbb{E}_{\mathbf{z}_t}[\sigma_i(\tilde{\mathbf{n}}(t), \mathbf{z}_t) | \tilde{\mathbf{n}}_t], \quad (38)$$

$$\lim_{t \searrow t_k} \log \tilde{q}(\tilde{\mathbf{n}}_{[0,t]}) = \lim_{t \nearrow t_k} \left( \log \tilde{q}(\tilde{\mathbf{n}}_{[0,t]}) + \log \mathbb{E}_{\mathbf{z}_t}[\sigma_{j_k}(\tilde{\mathbf{n}}_k, \mathbf{z}_t) | \tilde{\mathbf{n}}_t] \right). \quad (39)$$

which is another way of showing that  $\tilde{q}$  has propensities given by (35).

### 2.5.2 The Quasiequilibrium Approximation

While the QSSA relies on a reaction network being divisible into slow and fast species that evolve on two different timescales, in practice it is more frequently the case that some *reactions* in a system will be fast and that some will be slow. The Quasiequilibrium Approximation (QEA) developed in [17, 45, 46] modifies the QSSA to model this scenario by reformulating the system using extents  $\mathbf{a} = (a_1, \dots, a_r)$ , where the extent  $a_i$  is defined the number of times reaction  $i$  has occurred.

The extents themselves form a Markov chain, and the state of a system can be obtained from its extents as  $\mathbf{n}(t) = \mathbf{n}_0 + \mathbf{S}\mathbf{a}(t)$ , where  $\mathbf{S}$  is the stoichiometry matrix. As derived e.g. in [17], the CME of this system is

$$\frac{d}{dt}q_i(\mathbf{a}) = \sum_i [\rho_i(\mathbf{S}\mathbf{a} - \mathbf{S}_i)q_t(\mathbf{a} - \mathbf{e}_i) + \rho_i(\mathbf{S}\mathbf{a})q_t(\mathbf{a})]. \quad (40)$$

Here the sum is over reactions and  $\mathbf{e}_i$  is the vector with a 1 in the  $i$ -th position and 0 elsewhere. The QEA assumes that the extent vector can be divided into slow and fast components  $\mathbf{a}^s$  and  $\mathbf{a}^f$ , respectively, corresponding to slow and fast reactions, and we define our reduced system to consist only of the slow reactions. From here we can proceed analogously to the QSSA discussed above and obtain that the reduced propensities are given by

$$\tilde{\rho}_i(\mathbf{a}^s; t) = \mathbb{E}_{\mathbf{a}^f}[\rho_i(\mathbf{S}\mathbf{a}) | \mathbf{a}^s; t], \quad (41)$$

if the fast reactions can be assumed to equilibrate instantaneously.

A subtlety of this argument is the fact that the fast extents  $\mathbf{a}^f$  can only increase in time and therefore will not admit a steady-state distribution in general. The conditional means, however, will converge in many cases, e.g. if  $\mathbf{a}^f$  consists of both directions of a reversible reaction. If this is the case the QEA yields a well-defined reduction, which moreover agrees with (14) and therefore minimises the KL divergence to the full model on the space of extents.

### 2.5.3 The Linear Mapping Approximation

The Linear Mapping Approximation (LMA) [25] replaces a bimolecular reaction of the form  $G + X \xrightarrow{\sigma} G^*$ , where  $G$  is a binary species representing a gene state, by a reaction  $G \xrightarrow{\bar{\sigma}} G^*$  with effective propensity  $\bar{\sigma}$ . Assuming mass action kinetics, the propensity function of the bimolecular reaction is  $\rho(\mathbf{n}) = gx\sigma$ , where  $g$  is the state of  $G$  and  $x$  the abundance of the species  $X$ , and the propensity of the linearised version is  $\tilde{\rho}(\mathbf{n}) = \bar{\sigma}g$  for an appropriate choice of  $\bar{\sigma}$ . Taking the derivative with respect to  $\bar{\sigma}$  of the cross-entropy rate (13) at time  $t$  yields

$$\frac{\partial H(\tilde{q}; p)_t}{\partial \bar{\sigma}} = \sum_{\mathbf{n}} q_t(\mathbf{n}) \left( g - \frac{\sigma gx}{\bar{\sigma}} \right), \quad (42)$$

which vanishes if and only if

$$\bar{\sigma} = \sigma \cdot \mathbb{E}[x | g = 1; t]. \quad (43)$$

This provides a mathematical justification for the mean-field assumption underlying [25]. The LMA approximates this conditional expectation by imposing a self-consistency condition on the linearised system, thereby deriving an effective approximation to (43).

As an aside we can compute the KL divergence rate between  $q$  and the optimal reduction  $p$  given by (43) analytically to obtain

$$\text{KL}(q \| p)_t = \sigma \cdot P(g = 1) \cdot (\mathbb{E}[x \log x | g = 1] - \mathbb{E}[x | g = 1] \mathbb{E}[\log x | g = 1]). \quad (44)$$

Here the dependence of the expectations on the time  $t$  are suppressed. This expression for the discrepancy incurred by the LMA, which resembles the definition of the variance, quantifies the intuition that the LMA should be accurate if fluctuations of  $X$  in the unbound gene state are small.

The above derivation is not entirely rigorous as the linearisation has a different net stoichiometry than the original reaction and KL divergence between  $q$  and  $p$  is infinite. To remedy this we can either neglect fluctuations in  $X$  due to binding in the original model or consider the KL divergence on a reaction-by-reaction basis as discussed in 2.4, since the two reactions correspond despite their different net stoichiometries.

If fluctuations in  $X$  are absent or can be neglected, the optimal reduction using Eq. (43) preserves first-order moments for all species, assuming that all other reactions are linear. Indeed, replacing the reaction  $G \xrightarrow{\sigma[X]} G^*$  by its linearisation only changes the moment equation for  $\mathbb{E}[g]$ . For the full model, the latter is given by

$$\frac{d}{dt}\mathbb{E}[g; t] = \sigma \mathbb{E}[gx; t] + \dots, \quad (45)$$

while the reduction has

$$\frac{d}{dt}\mathbb{E}[g; t] = \sigma \mathbb{E}[x | g = 1; t] \mathbb{E}[g; t] + \dots, \quad (46)$$

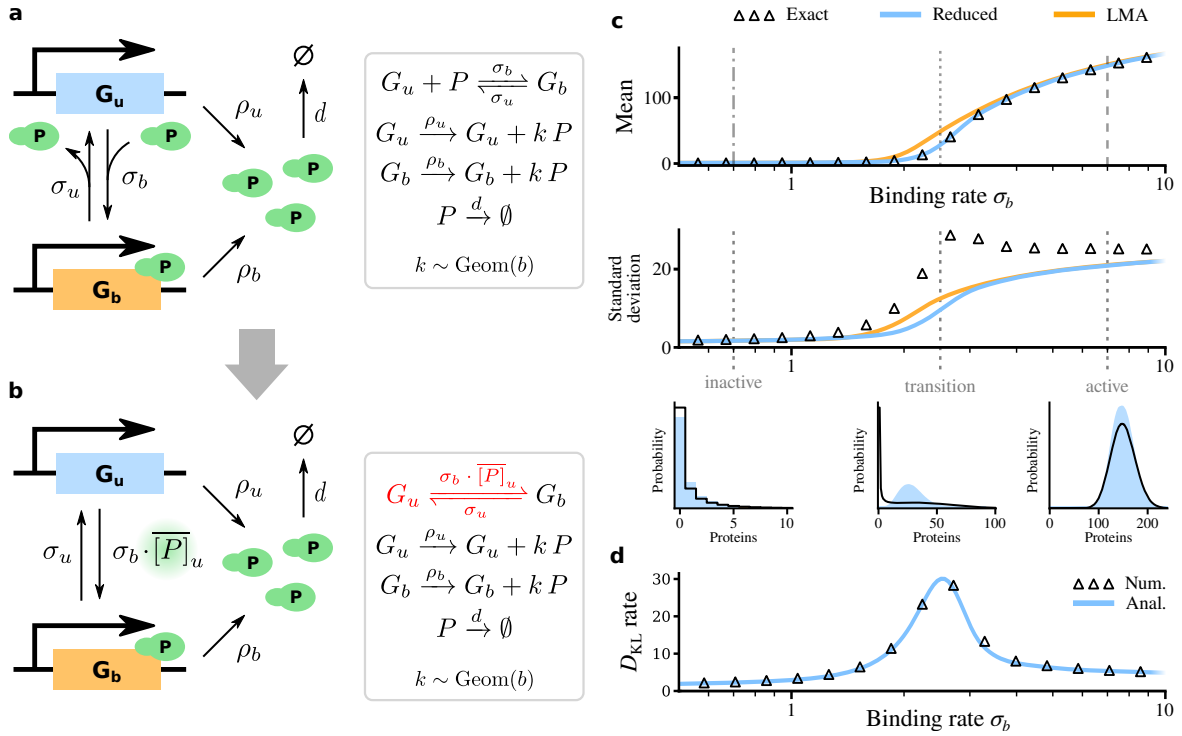
which is equivalent for a binary species  $g$  (the remaining terms only depend on species means by our linearity assumption). Second and higher-order moments will in generally differ between the two as we will see. This sheds light on a more subtle aspect of the LMA: the self-consistent approach in [25] involves computing the second-order moment  $\mathbb{E}[gx; t]$ , which can differ in the reduction. For this reason the original approach, while generally quite accurate, can lead to slightly suboptimal propensities as we shall see in Section 3.1.

#### 2.5.4 Moment-Matching

In Appendix D we show that the optimal reduced model  $p$ , defined using the propensity functions (14), exactly preserves the marginal distributions of the full model. In particular, if we place no constraints on the reduced propensity functions the resulting system  $p$  will have the same means, variances and higher-order moments for the observed species as the full system  $q$ . Thus our approach is also related to moment-matching, in particular to the proposed moment-based methods in [11, 47].

The fact that the reduced model  $p$  can exactly capture the probability distribution of the projected model  $\tilde{q}$  at any time  $t$  does not mean that the two are equivalent. The difference can be seen on the level of individual trajectories. The reduced model  $p$  is always memoryless, whereas this is not usually the case for  $\tilde{q}$ , resulting in different dynamical dynamical properties such as autocorrelations. This ability to discriminate models on a trajectory-level illustrates an advantage of our path-based approach over purely moment-based model reduction.

We emphasise that exact preservation of moments only holds for the globally optimal propensity functions defined by Eq. (14), which can depend on  $\tilde{n}$  and  $t$  in a complicated way. In practice one often use parametric propensity functions, such as mass-action propensities or Hill functions, and as a consequence this feature is lost. As with the LMA or the telegraph model it may still be possible to preserve some moments (such as means) exactly with the right choice of propensity function. Our results imply that discrepancies in the moments between the full and reduced models result from the fact that Eq. (14) is only computed approximately for methods like the QSSA or the QEA.



**Figure 4:** Reduction of an autoregulatory feedback loop using the Linear Mapping Approximation. **(a)** Schematic of the full reaction network. The number  $k \geq 0$  of proteins produced at each reaction is geometrically distributed with mean  $b$ . **(b)** Reduced form of the same reaction system, where the protein-gene binding reaction is replaced by a mean-field approximation. **(c)** For the chosen parameter values the system exhibits critical behaviour around  $\sigma_b = 2.5$ , transitioning from a mostly unbound to a mostly bound state. Near the transition protein fluctuations increase and the mean-field assumption in the LMA breaks down. A comparison of steady-state moments and distributions for the full model, the optimal reduction and the LMA shows that the effective reaction rate computed by the latter is generally close to optimal. **(d)** KL divergence rate at steady state between the full model and the reduced version, computed analytically using Eq. (44) and using Monte Carlo simulations. The peak around the transition matches the observed discrepancy between the full and the reduced model. The remaining model parameters are  $\sigma_u = 400$ ,  $\rho_u = 0.3$ ,  $\rho_b = 105$ ,  $b = 2$ ,  $d = 1$ .

### 3 Numerical Experiments

#### 3.1 Autoregulatory Feedback Loop

In this section we analyse a simple model of stochastic gene expression featuring positive autoregulation (see Fig. 4a). The system in question consists of a single gene found in two states, bound ( $G_b$ ) and unbound ( $G_u$ ), as well as the coded protein  $P$  which can bind to the gene to increase its own transcription rate. For simplicity we do not consider mRNA dynamics explicitly, instead modelling protein production as occurring in geometrically distributed bursts (see [48] for a derivation). The linearised version of the system, using the optimal propensity derived in Section 2.5.3, is shown in Fig. 4b. In Fig. 4c we compare the steady-state distributions of the full model with its linearisation, where the effective binding rate is computed numerically using Alg. 1, and the LMA, where the binding rate is approximated using the self-consistent approach in [25].

This system exhibits critical behaviour for some parameter values (see Fig. 4c), and it was shown

in [49] that computing the moments of this system is difficult near points of criticality as most moment closure techniques as well as the LMA yield inaccurate results. Our results show that for all parameters the self-consistent equations of the LMA provides results close to the numerically computed optimum. As discussed in Section 2.5.3, the exact reduction very closely reproduces mean protein numbers, with the LMA incurring a small bias near the critical point. In contrast, neither reduction is able to capture the increased fluctuations near the critical point, significantly underestimating the variance in protein numbers. Comparing the protein distributions for the full and the optimally reduced model shows a large discrepancy near the critical point, compared to parameters far from it.

We compute the KL divergence rate at steady state between the full model and its linearisation using (44) and via Monte Carlo estimation (see Fig. 4d). The steady-state KL divergence rate exhibits a notable peak near the critical point, coinciding with the parameter regime where the linearisation fails to capture the full model. As we move away in either direction from the critical point the KL divergence decreases in accordance with the better approximation of the system by its linearised version. This shows how the KL divergence can be used to assess how well model reduction works for different parameter regimes.

### 3.2 Michaelis-Menten Kinetics

We next consider one of the most-studied reaction networks in biology, the Michaelis-Menten model of enzyme kinetics, consisting of an enzyme  $E$  and a substrate  $S$  that reversibly bind to form an enzyme-substrate complex  $ES$  which converts the substrate into the product  $P$  (see Fig. 5a). This system is often treated as an example application of both the QSSA and the QEA, which remain the standard model reduction techniques applied to this example.

The QSSA classically models the enzyme  $E$  and the enzyme-substrate complex  $ES$  as the fast species, and replaces the full system by a one-step reaction where the substrate is converted to the product with a Hill-type propensity function (see Fig. 5b). This approximation is known to be valid when

$$E_T \ll [S] + K_m, \quad (47)$$

where  $E_T$  is the total number of enzymes and  $K_m$  is the Michaelis-Menten constant  $(k_{-1} + k_2)/k_1$  [50, 51]. Note that the reduced model is invariant under the scaling

$$E_T \mapsto \varepsilon \cdot E_T \quad k_i \mapsto \varepsilon^{-1} \cdot k_i, \quad i = -1, 1, 2, \quad (48)$$

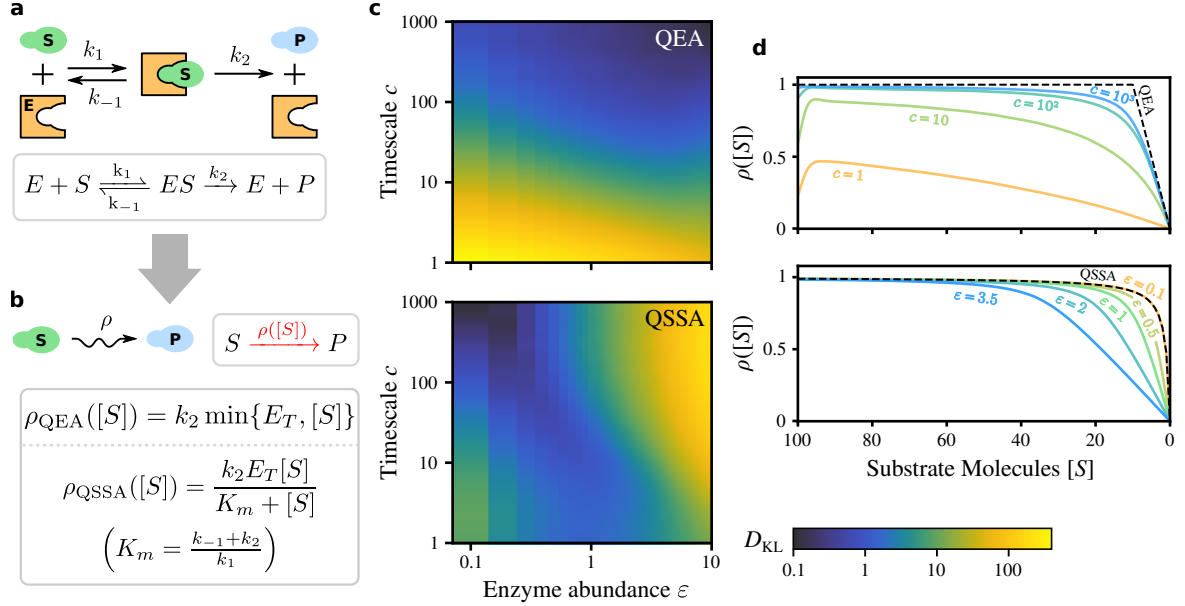
where the parameter  $\varepsilon$  can be seen as controlling the accuracy of the reduction: the rescaled QSSA condition (47) reads

$$\varepsilon E_T \ll [S] + K_m, \quad (49)$$

which suggests that small values of  $\varepsilon$  should result in the full model more closely resembling the QSSA reduction.

The stochastic QEA for this system is analysed in [17] and also results in a one-step reaction where the propensity function is now piecewise linear (see Fig. 5b). As opposed to the QSSA, the QEA is accurate when substrate binding and unbinding is fast, i.e.  $k_1, k_{-1}$  are sufficiently large. Note that the QEA reduction is invariant under the scaling

$$k_1 \mapsto c \cdot k_1 \quad k_{-1} \mapsto c \cdot k_{-1}, \quad (50)$$



**Figure 5:** Comparison of the QEA and the QSSA for Michaelis-Menten kinetics. **(a)** Schematic of the Michaelis-Menten system. We assume mass action kinetics. **(b)** Reduction of the Michaelis-Menten system under the QEA and the QSSA yields a one-step process with effective (non-mass action) propensities that depend on the method used. **(c)** Total KL divergence between the full model and the QEA and QSSA for different values of the timescale parameter  $c$  in (50) and the enzyme abundance parameter  $\varepsilon$  in (48). The QEA generally becomes more accurate as  $c$  increases, and the QSSA becomes more accurate as  $\varepsilon$  decreases. Note the small copy number effects that become apparent in the QSSA for low values of  $\varepsilon$  and  $c$ . We used the analytically obtained propensities for the QEA and the QSSA, and numerically estimated the KL divergences using Alg. 2. **(d)** Comparison of the effective reaction propensities computed according to Eq. (14) with those predicted by the QEA and the QSSA. Here substrates include enzyme-bound substrates. The top figure shows the effective propensities for  $\varepsilon = 1$ : as  $c \rightarrow \infty$  the effective propensities converge to the function predicted by the QEA. The bottom figure shows the effective reaction propensities for  $c = 1000$ , and as  $\varepsilon$  decreases the propensities approach the function predicted by the QSSA. The remaining parameters were fixed to  $k_1 = k_{-1} = 0.001$ ,  $k_2 = 0.1$ ,  $E_T = 10$  and  $S_T = 100$ .



where  $c$  controls the timescale at which quasiequilibrium is reached, and therefore the accuracy of the QEA. We stress that (48) and (50) are two independent scalings and can be performed simultaneously, as will be the case in Fig. 5c.

In general it may be difficult to predict which of the two approaches is more accurate unless one is clearly in one of the limiting regimes. Based on the KL divergence between the full model and either of the two reactions we can investigate this question for a range of parameters. The reduced model consists of two species,  $S$  and  $P$ , and since  $[S] + [P]$  is conserved we can describe it in terms of either of the two. The correct projection for this example identifies the species  $S$  and  $ES$  in the full model, ie. we define  $[S]_{\text{red}} = [S] + [ES]$  in order to eliminate the binding and unbinding reactions from the projection. This lumping of two rapidly equilibrating species is standard when using the QEA (see e.g. [24]), but it applies equally well to the QSSA in this case.

In Fig. 5c we use Alg. 2 to estimate the total KL divergence over  $[0, \infty]$  between the full model and both reductions for a fixed number of substrate molecules. The total KL divergence is finite since all trajectories enter the same absorbing state defined by  $[S] = 0$  in a finite amount of time. As expected for the QEA its accuracy increases with  $c$ , but whereas the QSSA tends to become less accurate for large  $E_T$ , in the low enzyme regime we observe a similar decrease in accuracy that is not explained by the deterministic theory.

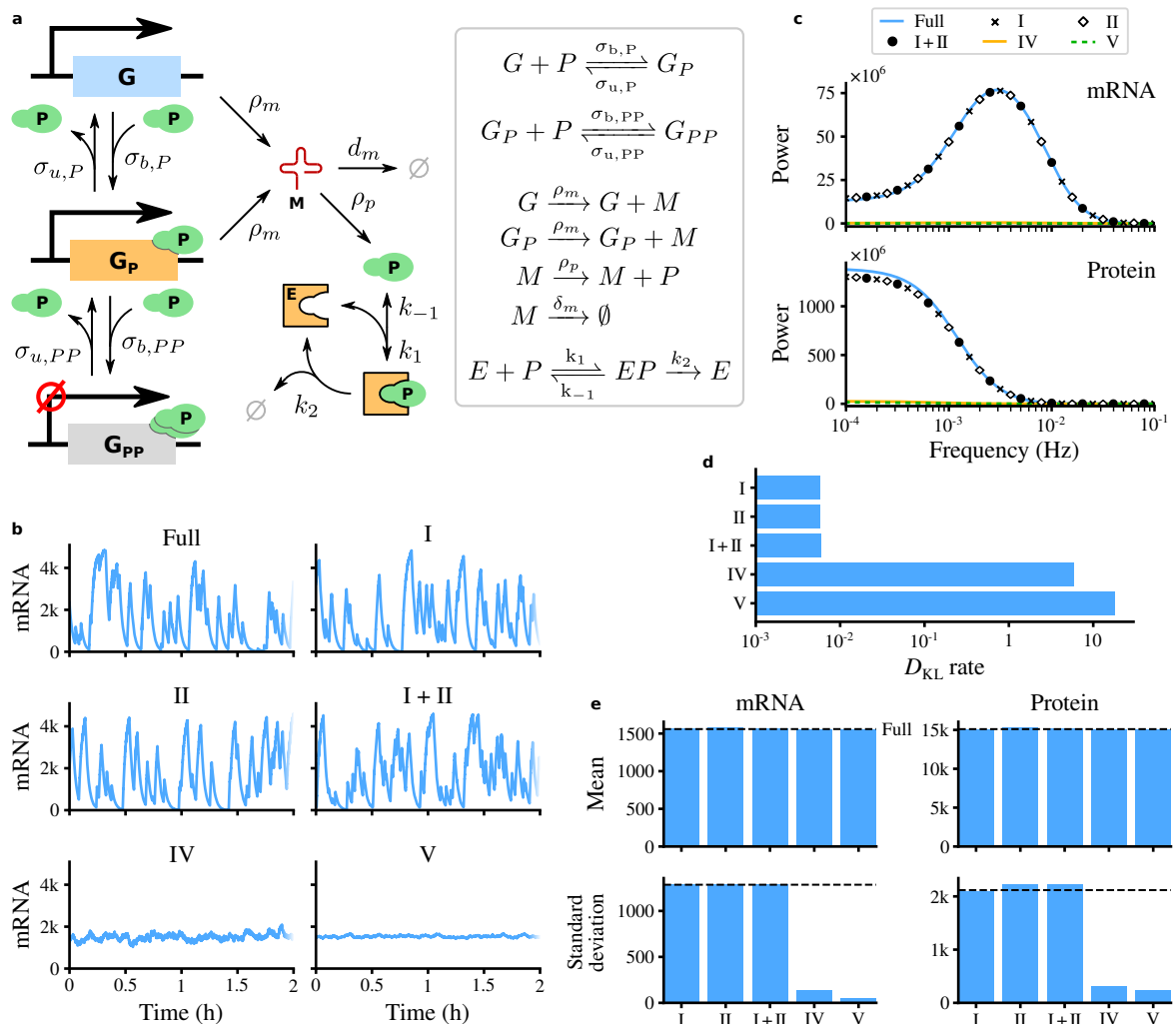
The observed decrease is a small copy number effect caused by the fact that the waiting time distribution between two productions of  $P$  is not exponential: for very small  $E_T$  the unbinding of an enzyme immediately after such an event implies that a significant fraction of such productions occur as a two-step process (where the free enzyme binds another substrate and then converts it) [52]. For large enough  $c$  the binding step is very fast, and assuming that  $k_2 \gg k_{-1}$  so that unbinding is unlikely to occur again, the conversion of substrates to products can be viewed as an effective one-step process. The effect of nonexponential waiting time distributions has previously been analysed in [53]; we see that the KL divergence on trajectories can be sensitive to subtle dynamical effects such as waiting time distribution that are not visible when considering e.g. the moments of a system, which are well predicted by the QSSA for small values of  $c$  and  $\varepsilon$ .

Overall we can see that for the chosen parameter values the QSSA generally performs better than the QEA in the regime of small  $\varepsilon$ , corresponding to low enzyme numbers, and the QEA is most accurate for  $c$ , keeping in mind the scalings in (48) and (50). Neither approximation is satisfactory if the number of total enzymes is larger than the amount of substrates, but the binding and unbinding rates are small. In this scenario other approximations, such as the total QSSA [54], which we do not investigate here, will generally be more accurate.

Figure 5d compares the effective propensities for the full system, as a function of the unconverted substrate abundance, with the predictions made by the QEA resp. the QSSA. In the case of the QEA we see that the effective propensities slowly converge to their limit as the timescale parameter  $c$  increases. In contrast the QSSA provides a good approximation to the effective propensity as long as the number of substrate molecules is larger than the number of enzymes. While the effective propensities are the optimal choice for the reduced model, the actual quality of the approximation is affected both by the size of the fluctuations of the actual propensities around their mean as well as the degree to which the waiting times of the original system follow an exponential distribution.

### 3.3 Genetic Oscillator

Our final example is the gene expression system shown in Fig. 6a. This system consists of a gene, mRNA and protein, as well as a Michaelis-Menten type protein degradation mechanism. Up to two protein molecules can bind the gene, and in the twice bound state the gene pauses transcription. This model has been considered in [11] as an example system where naive reduction of the CME using



**Figure 6:** Comparison of different reductions for the oscillatory gene network. **(a)** Schematic of the reaction system. We assume mass action kinetics for all reactions in the full model. **(b)** Example trajectories for the full model and its reductions. Note the oscillatory behaviour in the original model that is not preserved in reductions IV and V. We used Alg. 1 to compute the effective parameters for each reduction (cf. Table 2) from a single long trajectory probing steady-state behaviour. **(c)** Power spectra of mRNA and protein concentrations (defined as copy numbers normalised by the system size  $\Omega$ ). Models I, II and I+II closely reproduce the oscillatory behaviour of the full model while IV and V do not show sustained oscillations. **(d)** KL divergence rates between the full model and all reductions, estimated using Alg. 2. **(e)** Steady-state means and standard deviations for mRNA and protein numbers, estimated numerically using the Gillespie algorithm. While all reductions closely approximate the mean, Models IV and V do not match the variance of the full model. The parameters for the full model, taken from [11], are  $\Omega = 1000$ ,  $E_T = 10$ ,  $\rho_m = 50$ ,  $\rho_p = 0.0045$ ,  $d_m = 0.01$ ,  $k_1 = 0.1$ ,  $k_{-1} = 10$ ,  $k_2 = 10$ ,  $\sigma_{b,P} = 0.001$ ,  $\sigma_{u,P} = 100$ ,  $\sigma_{b,PP} = 1000$ ,  $\sigma_{u,PP} = 1$ .

Hill-type effective propensities is unable to accurately capture the noise in the original. The true system can exhibit oscillations in mRNA numbers that are caused by the two-step negative feedback and are not present in the reduced version. In this section we want to analyse what reductions can be performed on the full model while still keeping oscillatory behaviour.

We consider five different reductions for this model which are listed in Table 2. Each reduction removes or combines several species such as gene states, thereby reducing the dimensionality of the system. The rightmost column describes the correct projection for each model, derived using an analogous argument as in the Michaelis-Menten example. We fit the unknown (effective) parameters for each model numerically by minimising the KL divergence from the full model using Alg. 1. Typical mRNA trajectories for all models can be seen in Fig. 6b. The oscillations in the full model are clearly visible, and they are inherited by reductions I, II and I+II. In contrast, simplifying the first binding step as in IV and V results in a non-oscillatory system (see Fig. 6c).

In Fig. 6d we compare the KL divergence rates between the full model and reductions I-V (note that the different reductions are defined on different state spaces). While models I, II and I+II have comparatively low KL divergences from the original there is a sharp increase with models IV and V. From this alone one could expect that the latter perform significantly worse, which is indeed the case as shown in Fig. 6b. Fig. 6e presents a comparison of the mean and standard deviations for mRNA and protein abundances; while all models approximate the means very closely, the predicted standard deviations for model IV and V are very far from their true values consistent with the lack oscillatory behaviour. This resembles the results in Section 3.1 which also show excellent agreement of the full and reduced model on the mean level, independent of the total approximation quality.

This comparison of various reductions shows how variational model reduction can be employed on a more sophisticated scale, where multiple reductions are possible. Given a list of possible simplifications we can automatically find optimal parameters for each reduction and compute the KL divergence from the true model. Reductions which do not significantly affect the output will generally lead to very low KL divergences compared to those which do, as in this example where a large gap between KL divergence rates separates reductions I, II and I+II from IV and V. The latter do not retain much of the moments or oscillatory dynamics of the full model compared to the other reductions.

## 4 Discussion

In this work we presented an information-theoretic approach to model reduction for the CME based on minimising the KL divergence between the full model and the proposed reduction. Based on this variational principle we determine the optimal reaction propensities for the reduced model in Eq. (14) and show that it underlies some of the most common approaches to model reduction for the CME: the QSSA [8–10], the QEA [16, 17] and the LMA [25]. As a consequence, we establish that these methods can be seen as special cases of our approach. We furthermore obtain a general justification for the mean-field arguments proposed in the literature and connect them with information theory and likelihood-based approaches. We provide a numerical algorithm for automated fitting of a reduced model based on minimising the KL divergence via stochastic gradient descent, and a numerical algorithm for estimating this KL divergence in order to assess model fit. While the KL divergence between Markov chains has been considered before in e.g. [19–21, 31, 32, 34, 35, 55], to our knowledge it has not been studied in connection with standard model reduction methods as done in this paper. Our numerical results show how the KL divergence can be computed in practice and used in the context of model reduction.

Model	Full reactions	Reduced reactions	Projection map
<b>I</b>	$E + P \rightleftharpoons EP \longrightarrow E$	$P \xrightarrow{d_{\text{eff}}} \emptyset$	$[P]_{\text{red}} := [P] + [EP]$
<b>II</b>	$G + P \rightleftharpoons G_P$ $G_P + P \rightleftharpoons G_{PP}$ $G_P \longrightarrow G_P + P$	$G + P \xrightleftharpoons[\sigma_u / (K_u + [P])]{\sigma_{b,P}} G_{P/PP}$ $G_{P/PP} \xrightarrow{\rho_{\text{eff}}} G_{P/PP} + P$	$[P]_{\text{red}} := [P] + [G_{PP}]$ $[G_{P/PP}]_{\text{red}} := [G_P] + [G_{PP}]$
<b>I + II</b>	(I & II above)	(I & II above)	$[P]_{\text{red}} := [P] + [EP] + [G_{PP}]$ $[G_{P/PP}]_{\text{red}} := [G_P] + [G_{PP}]$
<b>IV</b>	$G + P \rightleftharpoons G_P$ $G_P + P \rightleftharpoons G_{PP}$ $G \longrightarrow G + P$ $G_P \longrightarrow G_P + P$	$G_{-}/P + P \xrightleftharpoons[\sigma_{u,PP}]{\sigma_b[P]/(K_b + [P])} G_{PP}$ $G_{-}/P \xrightarrow{\rho_m} G_{-}/P + P$	$[P]_{\text{red}} := [P] + [G_P] + [G_{PP}]$ $[G_{-}/P]_{\text{red}} := [G] + [G_P]$
<b>V</b>	$G + P \rightleftharpoons G_P$ $G_P + P \rightleftharpoons G_{PP}$ $G \longrightarrow G + P$ $G_P \longrightarrow G_P + P$	$\emptyset \xrightarrow{\frac{\alpha + \beta[P]}{1 + \gamma[P] + \delta[P]^2}} P$	$[P]_{\text{red}} := [P] + [G_P] + 2[G_{PP}]$

**Table 2:** Five possible reductions for the oscillator model. Model I+II is the combination of the reductions in Models I and II. All propensities are shown without their mass action terms. Species that are lumped or removed in the reduced model are shown in gray, and species that do not exist in the full model in blue. Parameters that are introduced in the reduced model are shown in red and found by minimisation of the KL divergence between the models; all remaining parameters are taken over from the full model. The functional forms of the propensity functions for Models II, IV and V were derived using the QEA.

Using three biologically relevant examples we illustrated how these model reduction techniques can be analysed from this variational perspective. For the autoregulatory feedback loop we showed that the KL divergence provides a useful metric of approximation quality and can detect parameter regimes where the mean-field approximation behind the LMA fails. Using the Michaelis-Menten system we demonstrated how our approach can be used to decide between possible reductions of a model, particularly in the non-asymptotic regime where neither the QSSA nor the QEA are strictly valid. Finally we used the genetic oscillator in [11] as an example to show how different reductions of a given model can be fit automatically, separately and in combination, and assesses which steps in a putative model reduction procedure would impact approximation quality more than others.

In this paper we only focused on discrete stochastic models, mainly those described using the Chemical Master Equation. This approach disregards continuum approximations such as the Chemical Langevin Equation and the System Size Expansion [56], which are frequently used in practice, and it is worth asking whether information-theoretic approaches can be extended to this context. A priori this appears difficult since the KL divergence and related quantities tend to diverge in the continuum limit, such as the divergence between Brownian motions with different diffusion coefficients. This similarly applies to the wide class of hybrid approaches [30, 42, 57, 58]. In particular, abundance separation, as opposed to time-scale separation with the QSSA and the QEA, does not fall under the umbrella approach we considered. In this scenario other perspectives on model reduction are likely needed.

A major drawback of the general approach we presented is that the KL divergences we use cannot usually be optimised analytically. As we have shown, the minimum typically corresponds to computing conditional expectations as in (14), which is rarely possible exactly, e.g. in the case of the LMA, where moment equations are used to approximate conditional means of protein numbers, or the QSSA, where deterministic rate equations are commonly used to derive the approximate propensities. The Monte-Carlo approach we proposed in Alg. 1 has the advantage of being completely generic, but unlike most existing techniques it requires numerical optimisation and will inevitably be slower than those. The optimisation problem in question typically being low-dimensional and convex, most of the computational time is required to generate sample trajectories and compute gradients, with variable computational complexity that heavily depends on the system.

One important aspect of model reduction we have not addressed is that of choosing an appropriate architecture for a reduced model. As we investigated in the case of the genetic oscillator, while we can always optimise the parameters given an architecture, the quality of the approximation can greatly depend on which reductions are performed. Although methods for automatically choosing from a predefined set of approximations of a system exist in special cases [15], a fully general algorithm that computes a maximally reduced version of a given model within a given threshold is still missing from the literature. In combination with other machine learning approaches to model reduction such as [59], however, we hope that our approach may be one of several steps towards such a procedure.

The KL divergence is a well-studied quantity in statistics but it lacks some desirable properties in the context of model reduction. While we were able to establish that optimising the KL divergence in theory leads to a reduced system with the same marginal distributions, establishing a relationship between the KL divergence and dynamical information such as the autocorrelation and power spectra is more difficult. Our investigations suggest that the KL divergence on the space of trajectories is a useful quantity that relates to both the marginal distributions and dynamical properties of the system, and can be a useful metric of model fit. We are optimistic that further investigations in this area will demarcate more clearly which properties of a model are well captured by the KL divergence and which are not.

## Code availability

Code implementing the methods introduced in this paper can be found at <https://github.com/kaandocal/modred>.

## Acknowledgments

This work was supported by the EPSRC Centre for Doctoral Training in Data Science (EPSRC grant EP/L016427/1) and the University of Edinburgh for K. Ö. and a Leverhulme Trust grant (Grant No. RPG-2020-327) for R. G. The authors are grateful to the anonymous reviewers for their useful feedback and suggestions.

## References

- [1] M. Kærn, T. C. Elston, W. J. Blake, and J. J. Collins, “Stochasticity in gene expression: from theories to phenotypes,” *Nat. Rev. Genet.*, **6**(6): 451–464, 2005.
- [2] R. Satija and A. K. Shalek, “Heterogeneity in immune responses: from populations to single cells,” *Trends Immunol.*, **35**(5): 219–229, 2014.
- [3] T. Lipniacki, B. Hat, J. R. Faeder, and W. S. Hlavacek, “Stochastic effects and bistability in T cell receptor signaling,” *J. Theor. Biol.*, **254**(1): 110–122, 2008.
- [4] D. Schnoerr, G. Sanguinetti, and R. Grima, “Approximation and inference methods for stochastic biochemical kinetics - a tutorial review,” *J. Phys. A*, **50**(9): 093001, 2017.
- [5] B. Munsky and M. Khammash, “The Finite State Projection algorithm for the solution of the Chemical Master Equation,” *J. Chem. Phys.*, **124**(4): 044104, 2006.
- [6] N. Ali Eshthewy and L. Scholz, “Model reduction for kinetic models of biological systems,” *Symmetry*, **12**(5): 863, 2020.
- [7] T. J. Snowden, P. H. van der Graaf, and M. J. Tindall, “Methods of model reduction for large-scale biological systems: a survey of current methods and trends,” *Bull. Math. Biol.*, **79**(7): 1449–1486, 2017.
- [8] C. V. Rao and A. P. Arkin, “Stochastic chemical kinetics and the quasi-steady-state assumption: application to the Gillespie algorithm,” *J. Chem. Phys.*, **118**(11): 4999–5010, 2003.
- [9] J. K. Kim, K. Josić, and M. R. Bennett, “The relationship between stochastic and deterministic quasi-steady state approximations,” *BMC Syst. Biol.*, **9**(1): 87, 2015.
- [10] J. K. Kim, K. Josić, and M. R. Bennett, “The validity of quasi-steady-state approximations in discrete stochastic simulations,” *Biophys. J.*, **107**(3): 783–793, 2014.
- [11] P. Thomas, A. V. Straube, and R. Grima, “The slow-scale Linear Noise Approximation: an accurate, reduced stochastic description of biochemical networks under timescale separation conditions,” *BMC Syst. Biol.*, **6**(1): 39, 2012.
- [12] P. Thomas, A. V. Straube, and R. Grima, “Limitations of the stochastic quasi-steady-state approximation in open biochemical reaction networks,” *J. Chem. Phys.*, **135**(18): 181103, 2011.
- [13] H.-W. Kang, W. R. KhudaBukhsh, H. Koepl, and G. A. Rempala, “Quasi-steady-state approximations derived from the stochastic model of enzyme kinetics,” *Bull. Math. Biol.*, **81**(5): 1303–1336, 2019.



- [14] J. Eilertsen, K. Srivastava, and S. Schnell, “Stochastic enzyme kinetics and the quasi-steady-state reductions: Application of the slow-scale Linear Noise Approximation à la Fenichel,” *J. Math. Biol.*, **85**(1): 3, 2022.
- [15] Y. M. Song, H. Hong, and J. K. Kim, “Universally valid reduction of multiscale stochastic biochemical systems using simple non-elementary propensities,” *PLoS Comput Biol*, **17**(10): e1008952, 2021.
- [16] E. L. Haseltine and J. B. Rawlings, “Approximate simulation of coupled fast and slow reactions for stochastic chemical kinetics,” *J. Chem. Phys.*, **117**(15): 6959–6969, 2002.
- [17] J. Goutsias, “Quasiequilibrium approximation of fast reaction kinetics in stochastic biochemical systems,” *The Journal of Chemical Physics*, **122**(18): 184102, 2005.
- [18] L. Cardelli, I. Pérez Verona, M. Tribastone, *et al.*, “Exact maximal reduction of stochastic reaction networks by species lumping,” *Bioinform.*, **37**: 2021.
- [19] R. A. Amjad, C. Blöchl, and B. C. Geiger, “A generalized framework for Kullback-Leibler Markov aggregation,” *IEEE Trans. Automat. Contr.*, **65**(7): 3068–3075, 2020.
- [20] K. Deng, P. G. Mehta, and S. P. Meyn, “Optimal Kullback-Leibler aggregation via spectral theory of Markov chains,” *IEEE Trans. Automat. Contr.*, **56**(12): 2793–2808, 2011.
- [21] B. C. Geiger, T. Petrov, G. Kubin, and H. Koepl, “Optimal Kullback-Leibler aggregation via information bottleneck,” *IEEE Trans. Automat. Contr.*, **60**(4): 1010–1022, 2015.
- [22] S. Bo and A. Celani, “Multiple-scale stochastic processes: decimation, averaging and beyond,” *Phys. Rep.*, **670**: 1–59, 2017.
- [23] J. Holehouse, A. Sukys, and R. Grima, “Stochastic time-dependent enzyme kinetics: closed-form solution and transient bimodality,” *J. Chem. Phys.*, **153**(16): 164113, 2020.
- [24] C. Jia and R. Grima, “Dynamical phase diagram of an auto-regulating gene in fast switching conditions,” *J. Chem. Phys.*, **152**(17): 174110, 2020.
- [25] Z. Cao and R. Grima, “Linear Mapping Approximation of gene regulatory networks with stochastic dynamics,” *Nat. Commun.*, **9**(1): 3305, 2018.
- [26] T. Jahnke and W. Huisinga, “Solving the Chemical Master Equation for monomolecular reaction systems analytically,” *J. Math. Biol.*, **54**(1): 1–26, 2007.
- [27] C. Gadgil, C. H. Lee, and H. G. Othmer, “A stochastic analysis of first-order reaction networks,” *Bull. Mathem. Biol.*, **67**(5): 901–946, 2005.
- [28] T. Zhou and J. Zhang, “Analytical results for a multistate gene model,” *SIAM J. Appl. Math.*, **72**(3): 789–818, 2012.
- [29] Y. Li, D.-Q. Jiang, and C. Jia, “Steady-state joint distribution for first-order stochastic reaction kinetics,” *Phys. Rev. E*, **104**(2): 024408, 2021.
- [30] T. Jahnke, “On Reduced Models for the Chemical Master Equation,” *Multiscale Model. Simul.*, **9**(4): 1646–1676, 2011.
- [31] M. Opper and G. Sanguinetti, “Variational inference for Markov jump processes,” in *Advances in Neural Information Processing Systems*, vol. 20, 2007.
- [32] C. Wildner and H. Koepl, “Moment-based variational inference for Markov jump processes,” in *Proceedings of the 36th International Conference on Machine Learning*, 2019, 6766–6775.
- [33] C. Zechner and H. Koepl, “Uncoupled analysis of stochastic reaction networks in fluctuating environments,” *PLOS Computational Biology*, **10**(12): e1003942, 2014.

- [34] L. Bronstein and H. Koepl, “Marginal process framework: a model reduction tool for Markov jump processes,” *Phys. Rev. E*, **97**(6): 062147, 2018.
- [35] A.-L. Moor and C. Zechner, *Dynamic information transfer in stochastic biochemical networks*, 2022. arXiv: [2208.04162](https://arxiv.org/abs/2208.04162) [q-bio.MN].
- [36] W. R. KhudaBukhsh, A. Auddy, Y. Disser, and H. Koepl, “Approximate lumpability for Markovian agent-based models using local symmetries,” *J. Appl. Probab.*, **56**(3): 647–671, 2019.
- [37] K. P. Murphy, *Probabilistic Machine Learning: Advanced Topics*. The MIT Press, 2023.
- [38] J. Peccoud and B. Ycart, “Markovian modeling of gene-product synthesis,” *Theor. Popul. Biol.*, **48**(2): 222–234, 1995.
- [39] J. K. Kim and J. C. Marioni, “Inferring the kinetics of stochastic gene expression from single-cell RNA-sequencing data,” *Genome Biol.*, **14**(1): R7, 2013.
- [40] L. Rabiner, “A tutorial on Hidden Markov Models and selected applications in speech recognition,” *Proc. IEEE*, **77**(2): 257–286, 1989.
- [41] F. Confortola and M. Fuhrman, “Filtering of continuous-time Markov chains with noise-free observation and applications,” *Stochastics*, **85**: 216–251, 2010.
- [42] J. Hasenauer, V. Wolf, A. Kazeroonian, and F. J. Theis, “Method of conditional moments (MCM) for the Chemical Master Equation: a unified framework for the method of moments and hybrid stochastic-deterministic models,” *J Math Biol*, **69**(3): 687–735, 2014.
- [43] L. Duso and C. Zechner, “Selected-node stochastic simulation algorithm,” *The Journal of Chemical Physics*, **148**(16): 164108, 2018.
- [44] M. Rathinam and M. Yu, “State and parameter estimation from exact partial state observation in stochastic reaction networks,” *The Journal of Chemical Physics*, **154**(3): 034103, 2021.
- [45] Y. Cao, D. T. Gillespie, and L. R. Petzold, “The slow-scale Stochastic Simulation Algorithm,” *J. Chem. Phys.*, **122**(1): 014116, 2005.
- [46] E. L. Haseltine and J. B. Rawlings, “On the origins of approximations for stochastic chemical kinetics,” *J. Chem. Phys.*, **123**(16): 164115, 2005.
- [47] D. W. Kim, H. Hong, and J. K. Kim, “Systematic inference identifies a major source of heterogeneity in cell signaling dynamics: The rate-limiting step number,” *Science Advances*, **8**(11): eab14598, 2022.
- [48] C. Jia and R. Grima, “Small protein number effects in stochastic models of autoregulated bursty gene expression,” *J. Chem. Phys.*, **152**(8): 084115, 2020.
- [49] K. Öcal, R. Grima, and G. Sanguinetti, “Parameter estimation for biochemical reaction networks using Wasserstein distances,” *J. Phys. A*, **53**(3): 034002, 2019.
- [50] L. A. Segel, “On the validity of the steady state assumption of enzyme kinetics,” *Bull. Math. Biol.*, **50**(6): 579–593, 1988.
- [51] D. Gillespie, K. Sanft, and L. Petzold, “Legitimacy of the stochastic Michaelis–Menten approximation,” *IET Syst. Biol.*, **5**(1): 58–69, 2011.
- [52] R. Grima and A. Leier, “Exact product formation rates for stochastic enzyme kinetics,” *J. Phys. Chem. B*, **121**(1): 13–23, 2017.
- [53] D. T. Gillespie, Y. Cao, K. R. Sanft, and L. R. Petzold, “The subtle business of model reduction for stochastic chemical kinetics,” *J. Chem. Phys.*, **130**(6): 064103, 2009.

- [54] S. MacNamara, A. M. Bersani, K. Burrage, and R. B. Sidje, “Stochastic chemical kinetics and the total quasi-steady-state assumption: application to the Stochastic Simulation Algorithm and Chemical Master Equation,” *J. Chem. Phys.*, **129**(9): 095105, 2008.
- [55] Z. Rached, F. Alajaji, and L. Campbell, “The Kullback-Leibler divergence rate between Markov sources,” *IEEE Trans. Automat. Contr.*, **50**(5): 917–921, 2004.
- [56] N. van Kampen, *Stochastic Processes in Physics and Chemistry*, 3rd. Elsevier, 2007.
- [57] A. Hellander and P. Lötstedt, “Hybrid method for the Chemical Master Equation,” *Journal of Computational Physics*, **227**(1): 100–122, 2007.
- [58] S. Smith, C. Cianci, and R. Grima, “Model reduction for stochastic chemical systems with abundant species,” *J. Chem. Phys.*, **143**(21): 214105, 2015.
- [59] G. Caravagna, L. Bortolussi, and G. Sanguinetti, “Matching models across abstraction levels with Gaussian processes,” in *Computational Methods in Systems Biology*, E. Bartocci, P. Lio, and N. Paoletti, Eds., Cham: Springer International Publishing, 2016, 49–66.
- [60] D. T. Gillespie, “A general method for numerically simulating the stochastic time evolution of coupled chemical reactions,” *J. Comput. Phys.*, **22**(4): 403–434, 1976.

## A Likelihoods of Fully Observed Trajectories

Let  $p$  be a continuous-time Markov chain defined on the space  $\mathcal{X}$  with initial distribution  $p_0$ . If  $p$  has time-independent transition rates the Stochastic Simulation Algorithm [60], reproduced in Algorithm 3, returns exact samples from  $p$ . We can inspect the SSA to obtain the probability of drawing any given trajectory as follows:

- The first state  $\mathbf{n}(0)$  is drawn from the initial distribution  $p_0$  and has probability  $p_0(\mathbf{n}(0))$
- The time  $\tau_1$  until the next jump is drawn from an exponential distribution with rate  $p_{\leftarrow \mathbf{n}_0}$  and has probability distribution function  $p_{\leftarrow \mathbf{n}_0} \exp(-\tau_1 p_{\leftarrow \mathbf{n}_0})$ .
- The next state  $\mathbf{n}_1$  is drawn from the transition distribution  $p_{\mathbf{n}_1 | \mathbf{n}_0}$  and has probability  $p_{\mathbf{n}_1 \leftarrow \mathbf{n}_0} / p_{\leftarrow \mathbf{n}_0}$
- The time  $\tau_2$  until the next jump is drawn from an exponential distribution with rate  $p_{\leftarrow \mathbf{n}_1}$  and has probability distribution function  $p_{\leftarrow \mathbf{n}_1} \exp(-\tau_2 p_{\leftarrow \mathbf{n}_1})$ , etc.

We thus arrive at the probability

$$p(\mathbf{n}_{[0,T]}) = p_0(\mathbf{n}(0)) \cdot p_{\leftarrow \mathbf{n}_0} \cdot \exp(-\tau_1 p_{\leftarrow \mathbf{n}_0}) \cdot \frac{p_{\mathbf{n}_1 \leftarrow \mathbf{n}_0}}{p_{\leftarrow \mathbf{n}_0}} \cdot \exp(-\tau_2 p_{\leftarrow \mathbf{n}_1}) \cdot \dots \cdot \exp(-\tau_k p_{\leftarrow \mathbf{n}_{k-1}}) \quad (\text{A1})$$

where the last term corresponds to the probability that  $t_{k+1} > T$ , ie. that the next jump occurs after time  $T$  (note the absence of the prefactor  $p_{\leftarrow \mathbf{n}_k}$ ). Each term  $p_{\leftarrow \mathbf{n}_i}$  in the above expression cancels with the denominator in the next term, and upon taking logarithms we arrive at Eq. (4).

---

**Algorithm 3** The Stochastic Simulation Algorithm to simulate samples from a Markov chain  $p$ .

---

**Input:** simulation length  $T$ , Markov chain  $p$ , initial distribution  $p_0$

**Output:**  $\mathbf{n}_{[0,T]}$  - sampled trajectory

---

**sample**  $\mathbf{n}_0 \sim p_0$

$t \leftarrow 0, i \leftarrow 0$

**while**  $t < T$  **do**

**sample**  $\tau_{i+1} \sim \text{Exp}(1/p_{\leftarrow \mathbf{n}_i})$

**sample**  $\mathbf{n}_{i+1} = \mathbf{m}$  **with probability**  $p_{\mathbf{m} \leftarrow \mathbf{n}_i} / p_{\leftarrow \mathbf{n}_i}$

$t += \tau_{i+1}$

**return** states  $(\mathbf{n}_0, \mathbf{n}_1, \dots)$ , jump times  $(t_1, t_2, \dots)$

---

## B Kullback-Leibler Divergence between Markov Chains

In [31] the authors derive an expression for the Kullback-Leibler divergence between two continuous-time Markov chains  $q$  and  $p$  defined on a common state space  $\mathcal{X}$ . For any integer  $N > 1$  let  $q^{(N)}$  resp.  $p^{(N)}$  be the time discretisations of  $q$  and  $p$  with time step  $\delta t = T/N$ . These define probability distributions on  $\mathcal{X}^{N+1}$  with KL divergence

$$\text{KL}(q^{(N)} \parallel p^{(N)}) = \sum_{\mathbf{n}} q_0(\mathbf{n}) \log \frac{q_0(\mathbf{n})}{p_0(\mathbf{n})} + \sum_{i=1}^N \sum_{\mathbf{n}, \mathbf{m}} q_{i-1}^{(N)}(\mathbf{n}) q_i^{(N)}(\mathbf{m} | \mathbf{n}) \log \frac{q_i^{(N)}(\mathbf{m} | \mathbf{n})}{p_i^{(N)}(\mathbf{m} | \mathbf{n})} \quad (\text{B1})$$

The transition rates are related to the discrete-time transition probabilities as follows:

$$q_i^{(N)}(\mathbf{n} | \mathbf{n}) = 1 - (\delta t) q_{\leftarrow \mathbf{n}}(t^{(i)}) + o(\delta t) \quad (\text{B2})$$

$$q_i^{(N)}(\mathbf{m} | \mathbf{n}) = (\delta t) q_{\mathbf{m} \leftarrow \mathbf{n}}(t^{(i)}) + o(\delta t) \quad (\mathbf{m} \neq \mathbf{n}) \quad (\text{B3})$$

where  $t^{(i)} = iT/N$  is the  $i$ -th discretisation time point. Using these identities for  $q$  and  $p$  it can be verified that the KL divergence (B1) converges to the following as  $N \rightarrow \infty$ :

$$\begin{aligned} \text{KL}(q \| p) &= \text{KL}(q_0 \| p_0) - \int_0^T dt \sum_{\mathbf{n}} q_t(\mathbf{n}) (q_{\leftarrow \mathbf{n}}(t) - p_{\leftarrow \mathbf{n}}(t)) \\ &\quad + \int_0^T dt \sum_{\mathbf{m} \neq \mathbf{n}} q_t(\mathbf{n}) q_{\mathbf{m} \leftarrow \mathbf{n}}(t) (\log q_{\mathbf{m} \leftarrow \mathbf{n}}(t) - \log p_{\mathbf{m} \leftarrow \mathbf{n}}(t)) \end{aligned} \quad (\text{B4})$$

In general the KL divergence can be written as the difference of the cross-entropy  $H(q; p)$  and the entropy  $H(q; q)$ , and we define the cross-entropy of two continuous-time Markov chains as

$$H(q; p) = H(q_0; p_0) + \int_0^T dt \sum_{\mathbf{n}} q_t(\mathbf{n}) p_{\leftarrow \mathbf{n}}(t) - \int_0^T dt \sum_{\mathbf{m} \neq \mathbf{n}} q_t(\mathbf{n}) q_{\mathbf{m} \leftarrow \mathbf{n}}(t) \log p_{\mathbf{m} \leftarrow \mathbf{n}}(t). \quad (\text{B5})$$

The above derivation is valid only if  $q_{\mathbf{m} \leftarrow \mathbf{n}} \neq 0$  implies  $p_{\mathbf{m} \leftarrow \mathbf{n}} \neq 0$ ; otherwise the cross-entropy and KL divergence are both infinite.

This definition of the cross-entropy generalises to the case where  $q$  is not Markovian, e.g. when it is the projection of a Markov process. In the discrete-time case we have

$$H(\tilde{q}^{(N)}; p^{(N)}) = H(\tilde{q}_0; p_0) - \sum_{i=1}^N \sum_{\tilde{\mathbf{n}}, \tilde{\mathbf{m}}} \tilde{q}_{i-1}^{(N)}(\tilde{\mathbf{n}} | q_0) \tilde{q}_i^{(N)}(\tilde{\mathbf{m}} | \tilde{\mathbf{n}}, q_0) \log p_i^{(N)}(\tilde{\mathbf{m}} | \tilde{\mathbf{n}}), \quad (\text{B6})$$

where the marginal and conditional distributions of  $\tilde{q}$  depend on the initial distribution of  $q$  since  $\tilde{q}$  is no longer assumed to be memoryless. Expressing the transition probabilities for  $p$  in terms of transition rates using Eqs. (B2) and (B3), as well as the identities

$$\tilde{q}_i^{(N)}(\tilde{\mathbf{n}} | \tilde{\mathbf{n}}, q_0) = 1 - (\delta t) \tilde{q}_{\leftarrow \tilde{\mathbf{n}}}(t^{(i)}) + o(\delta t), \quad (\text{B7})$$

$$\tilde{q}_i^{(N)}(\tilde{\mathbf{m}} | \tilde{\mathbf{n}}, q_0) = (\delta t) \tilde{q}_{\tilde{\mathbf{m}} \leftarrow \tilde{\mathbf{n}}}(t^{(i)}) + o(\delta t), \quad (\tilde{\mathbf{m}} \neq \tilde{\mathbf{n}}) \quad (\text{B8})$$

which follow from the definition of the marginal transition rates in Eq. (12), we obtain

$$\begin{aligned} H(\tilde{q}^{(N)}; p^{(N)}) &= H(\tilde{q}_0; p_0) + \sum_{i=1}^N \sum_{\tilde{\mathbf{n}}} \tilde{q}_{i-1}^{(N)}(\tilde{\mathbf{n}} | q_0) \left( 1 - (\delta t) \sum_{\tilde{\mathbf{m}} \neq \tilde{\mathbf{n}}} \tilde{q}_{\tilde{\mathbf{m}} \leftarrow \tilde{\mathbf{n}}}(t^{(i)}) \right) \log \left( 1 - (\delta t) \sum_{\tilde{\mathbf{m}} \neq \tilde{\mathbf{n}}} p_{\tilde{\mathbf{m}} \leftarrow \tilde{\mathbf{n}}} \right) \\ &\quad - (\delta t) \sum_{i=1}^N \sum_{\tilde{\mathbf{n}}} \tilde{q}_{i-1}^{(N)}(\tilde{\mathbf{n}} | q_0) \sum_{\tilde{\mathbf{m}} \neq \tilde{\mathbf{n}}} \tilde{q}_{\tilde{\mathbf{m}} \leftarrow \tilde{\mathbf{n}}}(t^{(i)}) \log p_{\tilde{\mathbf{m}} | \tilde{\mathbf{n}}}(t^{(i)}) - N \log(\delta t) + o(\delta t). \end{aligned} \quad (\text{B9})$$

The term  $N \log(\delta t)$  is independent of  $p$  and we will renormalise it to 0; this is standard when deriving continuum-limit versions of many information-theoretic quantities and was implicitly used in our definition of the cross-entropy in Eq. (B5). Proceeding with the above equation to get

$$\begin{aligned} H(\tilde{q}^{(N)}; p^{(N)}) &\approx H(\tilde{q}_0; p_0) + (\delta t) \sum_{i=1}^N \sum_{\tilde{\mathbf{n}}} \tilde{q}_{i-1}^{(N)}(\tilde{\mathbf{n}} | q_0) (\tilde{q}_{\leftarrow \tilde{\mathbf{n}}}(t^{(i)}) - p_{\leftarrow \tilde{\mathbf{n}}}(t^{(i)})) \\ &\quad - (\delta t) \sum_{i=1}^N \sum_{\tilde{\mathbf{n}}} \tilde{q}_{i-1}^{(N)}(\tilde{\mathbf{n}} | q_0) \sum_{\tilde{\mathbf{m}} \neq \tilde{\mathbf{n}}} \tilde{q}_{\tilde{\mathbf{m}} \leftarrow \tilde{\mathbf{n}}}(t^{(i)}) \log p_{\tilde{\mathbf{m}} | \tilde{\mathbf{n}}}(t^{(i)}) - N \log(\delta t) + o(\delta t). \end{aligned}$$

which, after dropping the  $N \log(\delta t)$  term and rearranging, yields Eq. (9) in the limit.

## C Likelihoods of Reduced Trajectories

Let  $q$  be a continuous-time Markov chain defined on the discrete state space  $\mathcal{X}$  with initial distribution  $q_0$ , and let  $\tilde{q}$  be its projection onto  $\tilde{\mathcal{X}}$ . Computing the log probability of a trajectory  $\tilde{\mathbf{n}}_{[0,T]}$  under  $\tilde{q}$  requires integrating over all possible full trajectories  $\mathbf{n}_{[0,T]} = (\tilde{\mathbf{n}}_{[0,T]}, \mathbf{z}_{[0,T]})$  that are compatible with  $\tilde{\mathbf{n}}_{[0,T]}$ , that is,

$$\tilde{q}(\tilde{\mathbf{n}}_{[0,T]}) = \int_{\mathbf{z}_{[0,T]}} q(\tilde{\mathbf{n}}_{[0,T]}, \mathbf{z}_{[0,T]}) d\mathbf{z}_{[0,T]}. \quad (\text{C1})$$

In order to compute the integral on the right-hand side we consider the time-discretisations of  $q$  and  $\tilde{q}$ . The marginal likelihood in this case can be obtained using the well-known forward algorithm for HMMs, which sequentially computes the joint probabilities  $q^{(N)}(\tilde{\mathbf{n}}_{0:i}, \mathbf{z}_i)$ :

$$q^{(N)}(\tilde{\mathbf{n}}_0, \mathbf{z}_0) = q_0(\tilde{\mathbf{n}}_0, \mathbf{z}_0), \quad (\text{C2})$$

$$q^{(N)}(\tilde{\mathbf{n}}_{0:i+1}, \mathbf{z}_{i+1}) = \sum_{\mathbf{z}} q^{(N)}(\tilde{\mathbf{n}}_{0:i}, \mathbf{z}_i = \mathbf{z}) q^{(N)}(\tilde{\mathbf{n}}_{i+1}, \mathbf{z}_{i+1} | \tilde{\mathbf{n}}_i, \mathbf{z}_i = \mathbf{z}). \quad (\text{C3})$$

Using (B2), (B3) we can write the second equation as

$$\begin{aligned} q^{(N)}(\tilde{\mathbf{n}}_{0:i+1}, \mathbf{z}_{i+1}) &= q^{(N)}(\tilde{\mathbf{n}}_{0:i}, \mathbf{z}_i = \mathbf{z}_{i+1}) \left(1 - (\delta t) q_{\leftarrow(\tilde{\mathbf{n}}_i, \mathbf{z}_{i+1})}(\mathbf{t}^{(i)})\right) \\ &\quad + (\delta t) \sum_{\mathbf{z} \neq \mathbf{z}_{i+1}} q^{(N)}(\tilde{\mathbf{n}}_{0:i}, \mathbf{z}_i = \mathbf{z}) q_{(\tilde{\mathbf{n}}_i, \mathbf{z}_{i+1}) \leftarrow (\tilde{\mathbf{n}}_i, \mathbf{z})}(\mathbf{t}^{(i)}) + o(\delta) \quad (\mathbf{n}_{i+1} = \mathbf{n}_i), \end{aligned} \quad (\text{C4})$$

$$q^{(N)}(\tilde{\mathbf{n}}_{0:i+1}, \mathbf{z}_{i+1}) = (\delta t) \sum_{\mathbf{z}} q^{(N)}(\tilde{\mathbf{n}}_{0:i}, \mathbf{z}_i = \mathbf{z}) q_{(\tilde{\mathbf{n}}_{i+1}, \mathbf{z}_{i+1}) \leftarrow (\tilde{\mathbf{n}}_i, \mathbf{z})}(\mathbf{t}^{(i)}) + o(\delta) \quad (\mathbf{n}_{i+1} \neq \mathbf{n}_i). \quad (\text{C5})$$

Here the two cases correspond to no visible jump and a visible jump occurring at the  $i$ -th step, respectively. Eq. (C4) can be reorganised as

$$\begin{aligned} q^{(N)}(\tilde{\mathbf{n}}_{0:i+1}, \mathbf{z}_{i+1}) &= q^{(N)}(\tilde{\mathbf{n}}_{0:i}, \mathbf{z}_i = \mathbf{z}_{i+1}) + (\delta t) \sum_{\mathbf{z} \neq \mathbf{z}_{i+1}} q^{(N)}(\tilde{\mathbf{n}}_{0:i}, \mathbf{z}_i = \mathbf{z}) q_{(\tilde{\mathbf{n}}_i, \mathbf{z}_{i+1}) \leftarrow (\tilde{\mathbf{n}}_i, \mathbf{z})}(\mathbf{t}^{(i)}) \\ &\quad - (\delta t) \sum_{\mathbf{z} \neq \mathbf{z}_{i+1}} q^{(N)}(\tilde{\mathbf{n}}_{0:i}, \mathbf{z}_i = \mathbf{z}_{i+1}) q_{(\tilde{\mathbf{n}}_i, \mathbf{z}) \leftarrow (\tilde{\mathbf{n}}_i, \mathbf{z}_{i+1})}(\mathbf{t}^{(i)}) \\ &\quad - (\delta t) \sum_{\mathbf{n} \neq \mathbf{n}_{i+1}} \sum_{\mathbf{z}} q^{(N)}(\tilde{\mathbf{n}}_{0:i}, \mathbf{z}_i = \mathbf{z}_{i+1}) q_{(\tilde{\mathbf{n}}, \mathbf{z}) \leftarrow (\tilde{\mathbf{n}}_i, \mathbf{z}_{i+1})}(\mathbf{t}^{(i)}) + o(\delta), \end{aligned} \quad (\text{C6})$$

Here the first and second sums correspond to hidden reactions that do not change  $\tilde{\mathbf{n}}$  and the last double represents visible reactions that affect  $\tilde{\mathbf{n}}$ . Expressing the transition rates for the visible reactions in terms of the reaction propensities  $\sigma_j$  and taking the continuum limit  $\delta t \rightarrow 0$  now yields Eqs. (22)-(24).

To arrive at Eqs. (26)-(28) we sum Eqs. (22)-(24) over all  $\mathbf{z}$ . Marginalising Eq. (22) immediately yields Eq. (26), and for the other two we obtain

$$\begin{aligned} \frac{d}{dt} \sum_{\mathbf{z}_t} q(\tilde{\mathbf{n}}_{[0,t]}, \mathbf{z}_t) &= \sum_{\mathbf{z}_t} \sum_{\mathbf{z}' \neq \mathbf{z}_t} (q_{(\tilde{\mathbf{n}}_t, \mathbf{z}_t) \leftarrow (\tilde{\mathbf{n}}_t, \mathbf{z}')} q(\tilde{\mathbf{n}}_{[0,t]}, \mathbf{z}') - q_{(\tilde{\mathbf{n}}_t, \mathbf{z}') \leftarrow (\tilde{\mathbf{n}}_t, \mathbf{z}_t)} q(\tilde{\mathbf{n}}_{[0,t]}, \mathbf{z}_t)) \\ &\quad - \sum_{i \text{ vis.}} \sum_{\mathbf{z}_t} \sigma_i(\tilde{\mathbf{n}}_t, \mathbf{z}_t) q(\tilde{\mathbf{n}}_{[0,t]}, \mathbf{z}_t), \end{aligned} \quad (\text{C7})$$

$$\lim_{t \searrow t_k} \sum_{\mathbf{z}_t} q(\tilde{\mathbf{n}}_{[0,t]}, \mathbf{z}_t) = \lim_{t \nearrow t_k} \sum_{\mathbf{z}_t} \sigma_{j_k}(\tilde{\mathbf{n}}_k, \tilde{\mathbf{z}}_t) q(\tilde{\mathbf{n}}_{[0,t]}, \mathbf{z}_t). \quad (\text{C8})$$

The first double sum in Eq.(C7), corresponding to all hidden reactions, vanishes. Write the joint distributions in terms of conditional distributions as

$$q_t(\tilde{\mathbf{n}}_{[0,t]}, \mathbf{z}) = q_t(\tilde{\mathbf{n}}_{[0,t]}) q_t(\mathbf{z} | \tilde{\mathbf{n}}_{[0,t]}, q_0) \quad (\text{C9})$$

where we explicitly include the dependence on the initial distribution. We arrive at

$$\frac{d}{dt} q(\tilde{\mathbf{n}}_{[0,t]}) = - \sum_{i \text{ vis.}} \sum_{\mathbf{z}_t} \sigma_i(\tilde{\mathbf{n}}_t, \mathbf{z}_t) q(\mathbf{z}_t | \tilde{\mathbf{n}}_{[0,t]}, q_0) q(\tilde{\mathbf{n}}_{[0,t]}), \quad (\text{C10})$$

$$\lim_{t \searrow t_k} q(\tilde{\mathbf{n}}_{[0,t]}) = \lim_{t \nearrow t_k} \sum_{\mathbf{z}_t} \sigma_{j_k}(\tilde{\mathbf{n}}_k, \tilde{\mathbf{z}}_t) q(\mathbf{z}' | \tilde{\mathbf{n}}_{[0,t]}, q_0) q(\tilde{\mathbf{n}}_{[0,t]}). \quad (\text{C11})$$

We recognise the sum over  $\mathbf{z}_t$  as a conditional expectation:

$$\frac{d}{dt} q(\tilde{\mathbf{n}}_{[0,t]}) = - \sum_{i \text{ vis.}} \mathbb{E}_{\mathbf{z}_t} [\sigma_i(\tilde{\mathbf{n}}(t), \mathbf{z}_t) | \tilde{\mathbf{n}}_{[0,t]}, q_0; t] q(\tilde{\mathbf{n}}_{[0,t]}), \quad (\text{C12})$$

$$\lim_{t \searrow t_k} q(\tilde{\mathbf{n}}_{[0,t]}) = \lim_{t \nearrow t_k} \mathbb{E}_{\mathbf{z}_t} [\sigma_i(\tilde{\mathbf{n}}(t), \mathbf{z}_t) | \tilde{\mathbf{n}}_{[0,t]}, q_0; t] q(\tilde{\mathbf{n}}_{[0,t]}). \quad (\text{C13})$$

The two equations rearrange to yield (27) and (28).

## D Marginal distributions after reduction

In this section we compute the marginal distributions of a projected reaction network  $\tilde{q}$ :

$$\frac{d}{dt} \tilde{q}_t(\tilde{\mathbf{n}}) = \frac{d}{dt} \sum_{\mathbf{z}} q_t(\tilde{\mathbf{n}}, \mathbf{z}) = \sum_i \sum_{\mathbf{z}} (\sigma_i(\tilde{\mathbf{n}} - \tilde{\mathbf{S}}_i, \mathbf{z} - \mathbf{S}_i^z) q_t(\tilde{\mathbf{n}} - \tilde{\mathbf{S}}_i, \mathbf{z} - \mathbf{S}_i^z) - \sigma_i(\tilde{\mathbf{n}}) q_t(\tilde{\mathbf{n}}, \mathbf{z})) \quad (\text{D1})$$

$$= \sum_i \sum_{\mathbf{z}} (\sigma_i(\tilde{\mathbf{n}} - \tilde{\mathbf{S}}_i, \mathbf{z} - \mathbf{S}_i^z) q_t(\mathbf{z} - \mathbf{S}_i^z | \tilde{\mathbf{n}} - \tilde{\mathbf{S}}_i) \tilde{q}_t(\tilde{\mathbf{n}} - \tilde{\mathbf{S}}_i) - \sigma_i(\tilde{\mathbf{n}}) q_t(\mathbf{z} | \tilde{\mathbf{n}}) \tilde{q}_t(\tilde{\mathbf{n}})) \quad (\text{D2})$$

$$= \sum_i (\mathbb{E} [\sigma_i(\tilde{\mathbf{n}} - \tilde{\mathbf{S}}_i, \mathbf{z} - \mathbf{S}_i^z) | \tilde{\mathbf{n}} - \tilde{\mathbf{S}}_i; t] \tilde{q}_t(\tilde{\mathbf{n}} - \tilde{\mathbf{S}}_i) - \mathbb{E} [\sigma_i(\tilde{\mathbf{n}}) | \tilde{\mathbf{n}}; t] \tilde{q}_t(\tilde{\mathbf{n}})). \quad (\text{D3})$$

This corresponds exactly to the CME for the reduced system with propensities given by Eq. (14). As a result, the marginal distributions of the optimal reduction  $p$  and  $\tilde{q}$  agree at all times.