



# Stochastic analysis of Chemical Reaction Networks using Linear Noise Approximation<sup>☆</sup>



Luca Cardelli<sup>a,b</sup>, Marta Kwiatkowska<sup>a</sup>, Luca Laurenti<sup>a,\*</sup>

<sup>a</sup> Department of Computer Science, University of Oxford, United Kingdom

<sup>b</sup> Microsoft Research, Cambridge, United Kingdom

## ARTICLE INFO

### Article history:

Received 27 November 2015

Received in revised form 8 July 2016

Accepted 1 September 2016

Available online 29 October 2016

### Keywords:

Chemical Reaction Networks

Linear Noise Approximation

Probabilistic logic

Model checking

## ABSTRACT

Stochastic evolution of Chemical Reactions Networks (CRNs) over time is usually analyzed through solving the Chemical Master Equation (CME) or performing extensive simulations. Analysing stochasticity is often needed, particularly when some molecules occur in low numbers. Unfortunately, both approaches become infeasible if the system is complex and/or it cannot be ensured that initial populations are small. We develop a probabilistic logic for CRNs that enables stochastic analysis of the evolution of populations of molecular species. We present an approximate model checking algorithm based on the Linear Noise Approximation (LNA) of the CME, whose computational complexity is independent of the population size of each species and polynomial in the number of different species. The algorithm requires the solution of first order polynomial differential equations. We prove that our approach is valid for any CRN close enough to the thermodynamical limit. However, we show on four case studies that it can still provide good approximation even for low molecule counts. Our approach enables rigorous analysis of CRNs that are not analyzable by solving the CME, but are far from the deterministic limit. Moreover, it can be used for a fast approximate stochastic characterization of a CRN.

© 2016 The Authors. Published by Elsevier Ireland Ltd. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

Chemical Reaction Networks (CRNs) and mass action kinetics are well studied formalisms for modelling biochemical systems (Chellaboina et al., 2009). In recent years, CRNs have also been successfully used as a formal programming language for biochemical systems (Soloveichik et al., 2010; Cardelli, 2013; Chen et al., 2013). There are two well established approaches for analysing chemical networks: deterministic and stochastic (Gillespie et al., 2013). The deterministic approach models the kinetics of a CRN as a system of ordinary differential equations (ODEs) and represents average behaviour, valid in the thermodynamic limit, when the molecular population is sufficiently high (Gillespie, 2009). The stochastic approach, on the other hand, is based on the Chemical Master Equation (CME) and models the CRN as a continuous-time Markov chain (CTMC) (Cardelli, 2008). The stochastic behaviour can

be analyzed by stochastic simulation (Gillespie et al., 2013) or by exhaustive probabilistic model checking of the CTMC, which can be performed, for example, by using PRISM (Kwiatkowska et al., 2011). Exhaustive analysis of the CTMC is able to find the best- and worst-case scenarios and is correct for any population size, but suffers from the state-space explosion problem (Kwiatkowska and Thachuk, 2014) and can only be used for relatively small systems. In contrast, deterministic methods are much more robust with respect to state-space explosion, but unable to represent stochastic fluctuations, which play a fundamental role when the system is not in thermodynamic equilibrium. As a consequence, approximate approaches to efficiently solve the CME are appealing. For instance, in Ammar et al. (2012) and Chinesta et al. (2015), the authors use proper generalized decomposition in order to efficiently derive a numerical solution of the CME. These approaches are based on the assumption that the probability of the system being in a particular state can be written as a finite sum of separable functions. Herein, we consider a different approach based on a continuous stochastic approximation of the CME (Van Kampen, 1992).

### 1.1. Contributions

In this paper we develop a novel approach for analysing the stochastic evolution of a CRN based on the Linear Noise

<sup>☆</sup> Authors order is alphabetical.

\* Corresponding author at: Department of Computer Science, University of Oxford, Wolfson Building, Parks Road, Oxford OX1 3BH, United Kingdom.

E-mail addresses: [luca@microsoft.com](mailto:luca@microsoft.com) (L. Cardelli),

[marta.kwiatkowska@cs.ox.ac.uk](mailto:marta.kwiatkowska@cs.ox.ac.uk) (M. Kwiatkowska), [luca.laurenti@cs.ox.ac.uk](mailto:luca.laurenti@cs.ox.ac.uk) (L. Laurenti).

Approximation (LNA) of the CME. We formulate SEL (Stochastic Evolution Logic), a probabilistic logic for CRNs that enables reasoning about probability, expectation and variance of linear combinations of populations of the species. Examples of properties that can be specified in our logic include “the maximum expected population size of species  $\lambda_1$  during the first 20 s is 75 molecules” and “the probability that the combined population of species  $\lambda_1$  and  $\lambda_2$  has degraded between 10 and 30 s is less than 0.1”. We propose an approximate model checking algorithm for the logic based on the LNA and implement it in Matlab and Java. We demonstrate that the complexity of model checking is polynomial in the initial number of species and independent of the initial molecule counts, thus ameliorating state-space explosion. Further, we show that model checking is exact when approaching the thermodynamic limit. Though the algorithm may not be accurate for systems far from the deterministic limit, this generally happens when the populations are small, in which case the analysis can be performed by transient analysis of the induced CTMC (Kwiatkowska et al., 2007). Our approach is essential for CRNs that cannot be analyzed by (partial) state space exploration, because of large or infinite state spaces. Moreover, it is useful for a fast (approximate) stochastic characterization of CRNs, since solving the LNA is much faster than solving the CME (Elf and Ehrenberg, 2003). We prove asymptotic correctness of LNA-based model checking and show on four examples that it is still possible to obtain very good approximations even for small population systems, compared to standard uniformisation (Kwiatkowska et al., 2007) and statistical model checking implemented in PRISM (Kwiatkowska et al., 2011).

## 1.2. Related work

The closest work to ours is by Bortolussi and Lanciani (2013), which uses the Central Limit Approximation (CLA) (essentially the same as the LNA) for checking restricted timed automata specifications, assuming a fixed population size. Wolf et al. (2010) develop a sliding window method to approximately verify infinite-state CTMCs, which applies to cases where most of the probability mass is concentrated in a confined region of the state space. Recently, Finite State Projection algorithms (FSP algorithms) for the solution or approximation of the CME have been introduced (Munsky and Khammash, 2006). Sliding window and FSP algorithms apply to the induced CTMC, but require at least partial exploration of the state space, and are thus not immune to state-space explosion. Moment closure techniques (Singh and Hespanha, 2008; Hespanha, 2008) improve scalability by estimating the first  $k \in \mathbb{N}$  moments of the distribution of the species over time. The LNA itself can be seen as a moment closure technique, as a Gaussian distribution is completely characterized by the first two moments. However, the LNA tells us more because it guarantees that, if certain conditions are satisfied, the distribution of the process is Gaussian.

Continuous Stochastic Logic (CSL), originally introduced in Aziz et al. (2000) and extended by Baier et al. (2003), is a logic widely used to perform model checking of continuous-time Markov chains. CSL combines temporal operators of the logic CTL with the probabilistic and steady-state operators, and is further extended with reward operators in Kwiatkowska et al. (2007). CSL model checking is based on solving the CME and proceeds through uniformisation of the CTMC, essentially a time discretisation, and thus involves traversal of the full state space. This can be partially ameliorated by fast adaptive uniformisation (Dannenberg et al., 2015) that does not consider states with negligible probability. An alternative is statistical model checking (SMC) which involves a key operator of CSL is probabilistic reachability, that is, computing the probability that a particular region of the state space is reached over a given time interval. Although SEL is endowed with a probabilistic operator, this operator gives the average value of the probability over time

and, if the time interval is not a singleton, this is not equivalent to probabilistic reachability. Nevertheless, as shown in Bortolussi et al. (2016), SEL and our approximate model checking algorithm can be extended to express reachability, but currently lacks reward operators. The CSL steady-state operator of CSL cannot be added to CSL because LNA is accurate only for finite time. PRISM (Kwiatkowska et al., 2011) implements CSL model checking using uniformisation, fast adaptive uniformisation and statistical model checking.

Hybrid Automata Stochastic Logic HASL (Ballarini et al., 2011a) is an expressive specification formalism for stochastic Petri nets based on linear hybrid automata that is employed by the tool Cosmos (Ballarini et al., 2011b). CRNs have a natural interpretation in terms of stochastic Petri nets, see e.g. Barbot and Kwiatkowska (2015). The HASL formalism is more expressive than SEL and CSL, and can express CSL probabilistic reachability and expected reward properties. HASL model checking proceeds through statistical model checking of the product of a HASL specification automaton and the Petri net, and is implemented in Cosmos. In contrast, SEL model checking follows through approximating the solution of the CME with the Gaussian process induced by the LNA.

## 1.3. Structure of the paper

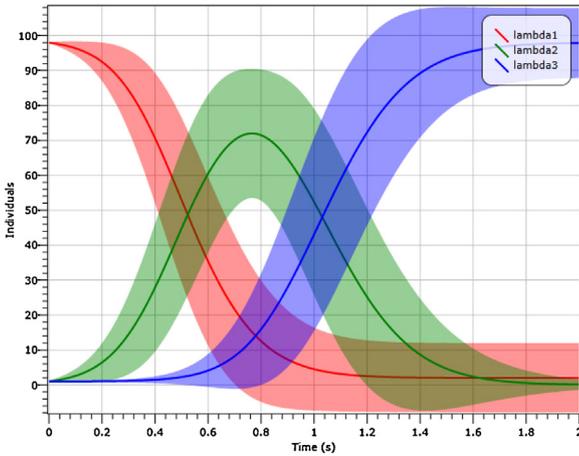
In Section 2 we summarize the deterministic and stochastic modelling approaches for CRNs, and in Section 3 we describe the Linear Noise Approximation method. Section 4 introduces the logic SEL and the corresponding model checking algorithm based on the LNA. In Section 5 we demonstrate our approach on four networks taken from the literature. Section 6 concludes the paper.

## 2. Chemical Reaction Networks

A *Chemical Reaction Network (CRN)*  $C = (\Lambda, R)$  is a pair of finite sets, where  $\Lambda$  is the set of *chemical species* and  $R$  the set of reactions.  $|\Lambda|$  denotes the size of the set of species. A *reaction*  $\tau \in R$  is a triple  $\tau = (r_\tau, p_\tau, k_\tau)$ , where  $r_\tau, p_\tau \in \mathbb{N}^{|\Lambda|}$  and  $k_\tau \in \mathbb{R}_{>0}$ .  $r_\tau$  and  $p_\tau$  represent the stoichiometry of reactants and products and  $k_\tau$  is the coefficient associated to the rate of the reaction; its dimension is  $s^{-1}$ . We often write reactions as  $\lambda_1 + \lambda_3 \rightarrow^{k_1} 2\lambda_2$  instead of  $\tau_1 = ([1, 0, 1]^T, [0, 2, 0]^T, k_1)$ , where  $\cdot^T$  indicates the transpose of a vector. We define the *net change* associated to a reaction  $\tau$  by  $\nu_\tau = p_\tau - r_\tau$ . For example, for  $\tau_1$  as above, we have  $\nu_{\tau_1} = [-1, 2, -1]^T$ .

We make the assumption that the system is well stirred, that is, the probability of the next reaction occurring between two molecules is independent of the location of those molecules. We consider fixed volume  $V$  and temperature; under these assumptions a *configuration* or *state*  $x \in \mathbb{N}^{|\Lambda|}$  of the system is given by the number of molecules of each species. We define  $[x] = x/N$ , the vector of the species *concentration* in  $x$  for a given  $N$ , where  $N = V \cdot N_A$  is the volumetric factor,  $V$  is the volume of the solution and  $N_A$  Avogadro's number. The physical dimension of  $N$  is  $\text{mol}^{-1} L$ , where  $\text{mol}$  indicates mole and  $L$  is litre. Given  $\lambda_i \in \Lambda$  then  $\#\lambda_i \cdot x \in \mathbb{N}$  represents the number of molecules of  $\lambda_i$  in  $x$  and  $[\lambda_i] \cdot x \in \mathbb{R}$  the concentration of  $\lambda_i$  in the same configuration. In some cases we elide  $x$ , and we simply write  $\#\lambda_i$  and  $[\lambda_i]$  instead of  $\#\lambda_i \cdot x$  and  $[\lambda_i] \cdot x$ . They are related by  $[\lambda_i] = \#\lambda_i / N$ . The dimension of  $[\lambda_i]$  is  $\text{mol} L^{-1}$ .

The propensity  $\alpha_{n,\tau}$  of a reaction  $\tau$  in terms of the number of molecules (here subscript  $n$  stands for the number of molecules) is a function of the current configuration of the system  $x$  such that  $\alpha_{n,\tau}(x)dt$  is the probability that a reaction event occurs in the next infinitesimal interval  $dt$ . In this paper we assume as valid the stochastic form of the law of mass action, so the propensity rates are proportional to the number of molecules that participate in the reaction (Cardelli, 2008). Stochastic models consider the system in terms of numbers of molecules, while deterministic ones,



**Fig. 1.** Expected number and standard deviation of species of the CRN of Example 2.1 for the given initial conditions, calculated by simulating the CME.

generally, in terms of concentrations, denoted  $\alpha_{c,\tau}(x)$  where subscript  $c$  stands for concentrations, and the relationship is as follows. For a reaction  $\tau = (r_\tau, p_\tau, k_\tau)$ , given the configuration  $x$  and  $r_{\tau,i}$ , the  $i$ th component of  $r_\tau$ , then  $\alpha_{c,\tau}(x) = k_\tau \prod_{i=1}^{|\Lambda|} ([\lambda_i]_x)^{r_{\tau,i}}$  is the propensity function expressed in terms of concentrations as given by the deterministic law of mass action. It is possible to show that, for any order of reaction,  $\alpha_{n,\tau}(x) \approx N\alpha_{c,\tau}(x)$  if  $N$  is sufficiently large (Anderson and Kurtz, 2011). Note that  $\alpha_{c,\tau}$  is independent of  $N$ . In this paper we are interested only in finite time horizon, because of the problematic character of studying solutions of ODEs for infinite time horizon (Bortolussi et al., 2013).

**Example 2.1.** Consider the CRN  $C = (\{\lambda_1, \lambda_2, \lambda_3\}, R)$ , where  $R = \{(\lambda_1 + \lambda_2 \rightarrow 10\lambda_2 + \lambda_2), (\lambda_2 + \lambda_3 \rightarrow 10\lambda_3 + \lambda_3)\}$ , with initial conditions  $\#\lambda_1 = 98, \#\lambda_2 = 1, \#\lambda_3 = 1$ , for a system with  $N = 1000$ . Fig. 1 plots the expectation and standard deviation of population sizes. We may wish to check if the maximum expected value of  $\#\lambda_2$  remains smaller than 75 molecules during the first 2 s. However, the system is stochastic, so we also need to analyze whether the variance is limited enough when  $\#\lambda_2$  reaches the maximum. Sometimes, analysis of first and second moments does not suffice, so it could be of interest to check the probability of some events, for instance, is the probability that  $\#\lambda_2 - (\#\lambda_1 + \#\lambda_3) > 0$ , between  $t_1 = 0.5$  s and  $t_2 = 1.0$  s, greater than 0.6?

### 2.1. Deterministic semantics

Let  $C = (\Lambda, R)$  be a CRN. The deterministic model approximates the concentration of the species of the system over time as a set of autonomous polynomial first order differential equations:

$$\frac{d\Phi(t)}{dt} = F(\Phi(t)) \quad (1)$$

where  $F(\Phi(t)) = \sum_{\tau=(r_\tau, p_\tau, k_\tau) \in R} \nu_\tau \alpha_{c,\tau}(\Phi(t))$  and  $\alpha_{c,\tau}(\Phi(t)) = k_\tau \prod_{i=1}^{|\Lambda|} \Phi_i(t)^{r_{\tau,i}}$ . Function  $\Phi: \mathbb{R}_{\geq 0} \rightarrow \mathbb{R}^{|\Lambda|}$  describes the behaviour of the system as a set of deterministic equations assuming a continuous state-space semantics, and therefore  $\Phi(t) \in \mathbb{R}^{|\Lambda|}$  is the vector of the species concentrations at time  $t$ . Assuming  $t_0 = 0$ , the initial condition is  $\Phi(0) = [x_0]$ , expressed as a concentration. Note that  $F(\Phi(t))$  is Lipschitz continuous, so  $\Phi$  exists and is unique (Ethier and Kurtz, 2009).  $\Phi$  represents the evolution of the system deterministically, neglecting stochastic fluctuations. However, it is often the case that stochasticity cannot be neglected. This is true especially when there are species with small populations. As a consequence, in such cases a stochastic model is needed.

### 2.2. Stochastic semantics

CRNs are well represented by CTMCs, whose transient analysis can be performed via the Chemical Master Equation (CME) (Van Kampen, 1992).

**Definition 1.** Given a CRN  $C = (\Lambda, R)$  and the volumetric factor  $N$ , we define a time-homogeneous CTMC (Cinlar, 2013; Pinsky and Karlin, 2010)  $(X^N(t), t \in \mathbb{R}_{\geq 0})$  with state space  $S \subseteq \mathbb{N}^{|\Lambda|}$ . Given  $x_0 \in S$ , the initial configuration of the system, then  $P(X^N(0) = x_0) = 1$ . The transition rate from state  $x_i$  to state  $x_j$  is defined as  $r(x_i, x_j) = \sum_{\{\tau \in R | x_j = x_i + \nu_\tau\}} N\alpha_{c,\tau}(x_i)$ .

$X^N(t)$  describes the stochastic evolution of molecular populations of each species at time  $t$ . For  $x \in S$ , we define  $P^{(t)}(x) = P(X^N(t) = x | X^N(0) = x_0)$ , where  $x_0$  is the initial configuration. The CME describes the time evolution of  $X^N$  as:

$$\frac{d}{dt} (P^{(t)}(x)) = \sum_{\tau \in R} [N\alpha_{c,\tau}(x - \nu_\tau) P^{(t)}(x - \nu_\tau) - N\alpha_{c,\tau}(x) P^{(t)}(x)]. \quad (2)$$

The CME can be equivalently defined in terms of the infinitesimal generator matrix (Wolf et al., 2010), which admits computing an approximation of the CME using, for example, fast adaptive uniformisation (Didier et al., 2009; Dannenberg et al., 2015) or the sliding window method (Wolf et al., 2010).

We also define the CTMC  $(X^N(t)/N, t \in \mathbb{R}_{\geq 0})$  with state space  $S \subseteq \mathbb{Q}^{|\Lambda|}$ . If  $[x_0] \in S$  is the initial configuration, then  $P((X^N(0)/N) = [x_0]) = 1$ . The transition rate from state  $[x_i]$  to  $[x_j]$  is defined as  $r([x_i], [x_j]) = \sum_{\{\tau \in R | [x_j] = [x_i] + (\nu_\tau/N)\}} N\alpha_{c,\tau}(x_i)$ .  $(X^N(t)/N)$  is the random vector describing the system at time  $t$  in terms of concentrations. In Anderson and Kurtz (2011) and Ethier and Kurtz (2009) it is proved that  $\limsup_{N \rightarrow \infty} \sup_{t' \leq t} \|((X^N(t')/N) - \Phi(t'))\| = 0$  almost

surely for any time  $t$ . This explains the relationship between the two different semantics, where the deterministic solution can be viewed as a limit of the stochastic solution, valid when close enough to the thermodynamic limit (i.e., for large molecular counts).

### 3. Linear Noise Approximation

The solution of the CME can be computationally expensive, or even infeasible, because the set of reachable states can be huge or infinite. The Linear Noise Approximation (LNA) has been introduced by Van Kampen as a second order approximation of the system size expansion of the CME (Van Kampen, 1992). It permits a stochastic characterization of the evolution of a CRN, still maintaining scalability comparable to that of the deterministic models.

In what follows, we introduce the LNA following the derivation in Wallace et al. (2012). This approach clearly shows that, assuming mass action kinetics, the LNA is always accurate for any CRN, if  $N$  is large enough, at least for a limited time. In fact, the original derivation of Van Kampen (1992) does not shield much light on the validity of such an approximation because truncating the expansion of the CME works fine only if terms of higher order are well behaved, and this is not always the case (Wallace et al., 2012).

In order to derive the LNA, we first consider the following conditions, namely the *leap conditions*. Given a CRN  $C = (\Lambda, R)$ ,  $X^N$  satisfies the leap conditions at time  $t$ , if for any  $\tau \in R$ , it holds that there exists a finite time interval  $dt$  such that:

$$\alpha_{n,\tau}(X^N(t)) \text{ is constant in } [t, t + dt] \text{ and}$$

$$\alpha_{n,\tau}(X^N(t)) \cdot dt \gg 1.$$

In Gillespie (2000), Gillespie shows that if these conditions are satisfied then the CME can be approximated by the following *Chemical*

Langevin Equation (CLE):

$$X^N(t + dt) = X^N(t) + \sum_{\tau \in R} \nu_{\tau} \alpha_{n,\tau}(X^N(t)) dt + \sum_{\tau \in R} \nu_{\tau} \sqrt{\alpha_{n,\tau}(X^N(t))} N_{\tau}(0, 1) \sqrt{dt} \quad (3)$$

where  $N_{\tau}(0, 1)$  are a set of independent normally distributed random variables with expected value 0 and variance 1. It is possible to show that, assuming mass action kinetics, for  $N$  large enough, the leap conditions can always be satisfied, and so Eq. (3) can be considered as a valid approximation of the real Markov process, at least for finite time: the Gaussian nature of the CLE makes it impossible to handle rare events. Since stochastic fluctuations depend on the volumetric factor,  $N$ , of the system, and, specifically, for average concentrations are of the order of  $N^{(1/2)}$  (Elf and Ehrenberg, 2003), we can assume that Eq. (3) has a solution of the form:

$$X^N(t) \approx N\Phi(t) + N^{(1/2)}Z(t) \quad (4)$$

where  $Z(t) = (Z_1(t), Z_2(t), \dots, Z_{|\Lambda|}(t))$  is a random vector, independent of  $N$ , representing the stochastic fluctuations at time  $t$  and  $\Phi(t)$  is given by the solution of Eq. (1). This assumption can also be justified by considering the work of Ethier and Kurtz (2009). Assuming such a structure for the solution of Eq. (3), then the probability distribution of  $Z(t)$  can be approximated by the following linear Fokker–Plank equations (Risken, 1984):

$$\frac{\partial P(Z, t)}{\partial t} = - \sum_{i=1}^{|\Lambda|} \sum_{j=1}^{|\Lambda|} \frac{\partial F_j(\Phi(t))}{\partial \Phi_i} \frac{\partial (Z_j P(Z, t))}{\partial Z_j} + \frac{1}{2} \sum_{i=1}^{|\Lambda|} \sum_{j=1}^{|\Lambda|} G_{i,j}(\Phi(t)) \frac{\partial^2 P(Z, t)}{\partial Z_i \partial Z_j} \quad (5)$$

where  $G(\Phi(t)) = \sum_{\tau \in R} \nu_{\tau} \alpha_{c,\tau}(\Phi(t))$  and  $F_j(\Phi(t))$  is the  $j$ th component of  $F(\Phi(t))$ . Solving a general Fokker–Planck equation generally cannot be done in closed form (Ammar et al., 2006). However, it is well known that the solution of Eq. (5) yields a Gaussian process (Van Kampen, 1992). For every time  $t$ ,  $Z(t)$  has a multivariate normal distribution whose expected value,  $E[Z(t)]$ , and covariance matrix,  $C[Z(t)]$ , are the solution of the following equations (Elf and Ehrenberg, 2003):

$$\frac{dE[Z(t)]}{dt} = J_F(\Phi(t))E[Z(t)] \quad (6)$$

$$\frac{dC[Z(t)]}{dt} = J_F(\Phi(t))C[Z(t)] + C[Z(t)]J_F^T(\Phi(t)) + G(\Phi(t)) \quad (7)$$

where  $J_F(\Phi(t))$  is the Jacobian of  $F(\Phi(t))$ . We consider as initial conditions  $E[Z(0)] = 0$  and  $C[Z(0)] = 0$ . This means that  $E[Z(t)] = 0$  for every  $t$ .

The LNA is therefore obtained as an approximation of the CLE, and so it yields all its conditions of validity. However, we need to check that Hypothesis (4) is effectively satisfied, which implies the need to check that  $C[Z(t)]$  remains bounded. This ensures that the LNA can be an accurate approximation, at least for a limited time, for any CRN. The following theorem ensures that, for  $N \rightarrow \infty$ , the LNA is always an accurate approximation.

**Theorem 1.** (Ethier and Kurtz, 2009) Let  $C = (\Lambda, R)$  be a CRN in a system of size  $N$  and  $X^N$  the CTMC induced by  $C$ . Let  $\Phi(t)$  be the solution of Eq. (1) with initial condition  $\Phi(0) = \frac{x_0}{N}$  and  $Z$  be the Gaussian process with expected value and variance given by Eqs. (6) and (7). Call  $\bar{X}^N = N^{1/2} \left( \frac{X^N(t)}{N} - \Phi(t) \right)$ . Then, for any  $x \in \mathbb{R}^{|\Lambda|}$ ,

$$\lim_{N \rightarrow \infty} \mathbb{P}_{\bar{X}^N(t)}(x) = \mathbb{P}_{Z(t)}(x), \quad (8)$$

where  $\mathbb{P}_{\bar{X}^N(t)}$  and  $\mathbb{P}_{Z(t)}$  are the cumulative distribution functions of random variables  $\bar{X}^N(t)$  and  $Z(t)$ , respectively.

Even if both LNA and CLE give rise to a Gaussian process, solving the LNA is much simpler than solving the CLE as explained in Wallace et al. (2012). As a consequence, it can be used for a fast stochastic characterization of the stochastic semantics of any CRN. LNA is exact in the limit of high populations, but can also be used for quite small populations. In fact, if the species of interest present a unimodal distribution, and the molecular count is such that a continuous approximation can be reasonable, then the LNA is generally surprisingly accurate.

To compute the LNA it is necessary to solve  $O(|\Lambda|^2)$  first order differential equations, but the complexity is independent of the initial number of molecules of each species. Therefore, one can avoid the exploration of the state space that methods based on uniformization rely upon.

### 3.1. Probabilistic analysis of CRNs

The LNA thus permits approximation of the probability distribution of  $X^N(t)$  with the probability distribution of  $Y^N(t) = N\Phi(t) + N^{1/2}Z(t)$ , where  $Z$ , and hence  $Y^N$ , are Gaussian processes. As a consequence,  $Y^N(t)$  has a multivariate Gaussian distribution, so it is completely characterized by its expected value and covariance matrix, whose values are respectively  $E[Y^N(t)] = N\Phi(t)$  and  $C[Y^N(t)] = N^{(1/2)}C[Z(t)]N^{(1/2)} = N C[Z(t)]$ .

Since  $Y^N$  has a multivariate normal distribution, then every linear combination of its components is normally distributed. Therefore, given  $B = [b_1, b_2, \dots, b_{|\Lambda|}]$  where  $b_1, b_2, \dots, b_{|\Lambda|} \in \mathbb{Z}$ , we can consider the random variable  $BY^N(t)$ , which defines a linear combination of the species at time  $t$ . For every  $t$ ,  $BY^N(t)$  is a normal random variable, whose expected value and variance are:

$$E[BY^N(t)] = BE[Y^N(t)] \quad (9)$$

$$C[BY^N(t)] = BC[Y^N(t)]B^T. \quad (10)$$

For a specific time  $t_k$ , it is possible to calculate the probability that  $BY^N(t_k)$  is within a set  $I$  of closed, disjoint real intervals  $[l_i, u_i]$ , where  $l_i, u_i \in \mathbb{R} \cup \{+\infty, -\infty\}$ . This probability  $\Omega_{Y^N, B, I}(t_k)$  is given by:

$$\Omega_{Y^N, B, I}(t_k) = \sum_{[l_i, u_i] \in I} \int_{l_i}^{u_i} g(x|E[BY^N(t_k)], C[BY^N(t_k)]) dx \quad (11)$$

where  $g(x|EV, \sigma^2)$  is the Gaussian distribution with expected value  $EV$  and covariance  $\sigma^2$ . We recall that it is possible to find numerical solution of Eq. (11) in constant time using the Z table (Patel and Read, 1996).

**Example 3.1.** Consider the CRN of Example 2.1, then we can obtain the probability that  $\#\lambda_1 - 2\#\lambda_3$  is at least 10 at time 20 by defining  $B' = [1, 0, -2]$ ,  $I' = \{[10, +\infty]\}$  and calculating  $\Omega_{Y^N, B', I'}(20)$ .

The following theorems are consequences of results in Wallace et al. (2012) and Ethier and Kurtz (2009), which can be generalized for reactions with a finite number of reagents and products. They show asymptotic pointwise convergence of expected value, variance and probability.

**Theorem 2.** Let  $C = (\Lambda, R)$  be a CRN. Suppose the solution of Eq. (7) is bounded, then, for any finite instant of time  $t_i$

$$\lim_{N \rightarrow \infty} \|\Omega_{Y^N, B, I}(t_i) - \tilde{\Omega}_{X^N, B, I}(t_i)\| = 0, \quad (12)$$

where  $\tilde{\Omega}_{X^N, B, I}(t_i)$  is the probability that  $B(X^N)$  is within  $I$  at time  $t_i$ .

**Theorem 3.** Suppose the solution of Eq. (7) is bounded, then, approaching the thermodynamic limit, for any finite instant of time  $t_k$ :

$$\lim_{N \rightarrow \infty} \|C[BY^N(t_k)] - C[BX^N(t_k)]\| = 0 \quad (13)$$

$$\lim_{N \rightarrow \infty} \|E[BY^N(t_k)] - E[BX^N(t_k)]\| = 0. \quad (14)$$

To solve the differential equations (6) and (7), it is necessary to use a numerical method such as the adaptive Runge–Kutta algorithm (Butcher, 1987). This yields the solution for a finite set of sampling times  $\Sigma = [t_1, \dots, t_{|\Sigma|}] \in \mathbb{R}^{|\Sigma|}$ , where  $t_1 \leq \dots \leq t_k \leq \dots \leq t_{|\Sigma|}$  and  $|\Sigma|$  is the sample size. Assuming  $Y^N$  is separable, that is, it is possible to completely define the behaviour of  $Y^N$  by only considering a countable number of points, we can calculate  $\Omega_{Y^N, B, I}$  for any point in  $\Sigma$  and, if points are dense enough, then this set exhaustively describes the probability that  $BX^N$  is within  $I$  over time. This restriction is not a limitation since for any stochastic process there exists a separable modification of it (Itô, 2006).

#### 4. Stochastic Evolution Logic (SEL)

Let  $C = (\Lambda, R)$  be a CRN with initial state  $x_0$ , in a system of size  $N$ . We now define the logic SEL (Stochastic Evolution Logic) which enables evaluation of the probability, variance and expectation of linear combinations of populations of the species of  $C$ .

The syntax of SEL is given by:

$$\eta := P_{\sim p}[B, I]_{[t_1, t_2]} \mid Q_{\sim v}[B]_{[t_1, t_2]} \mid \eta_1 \wedge \eta_2 \mid \eta_1 \vee \eta_2$$

where  $Q = \{supV, infV, supE, infE\}$ ,  $\sim \in \{<, >\}$ ,  $p \in [0, 1]$ ,  $v \in \mathbb{R}$ ,  $B \in \mathbb{Z}^{|\Lambda|}$ ,  $I = \{[l_i, u_i] \mid l_i, u_i \in \mathbb{R} \cup \{+\infty, -\infty\} \wedge [l_i, u_i] \cap [l_j, u_j] = \emptyset, i \neq j\}$  and  $[t_1, t_2]$  is a closed interval, with the constraint that  $t_1 \leq t_2$  and  $t_1, t_2 \in \mathbb{R}$ . If  $t_1 = t_2$  the interval reduces to a singleton.

Formulae  $\eta$  describe global properties of the stochastic evolution of the system.  $(B, I)$  specifies a linear combination of the species of  $C$  and a set of intervals, where  $B \in \mathbb{Z}^{|\Lambda|}$  is the vector defining the linear combination and  $I$  represents a set of disjoint closed real intervals.  $P_{\sim p}[B, I]_{[t_1, t_2]}$  is the probabilistic operator, which specifies the probability that the linear combination defined by  $B$  falls within the range  $I$  over the time interval  $[t_1, t_2]$ .  $supE, infE, supV, infV$  respectively yield the supremum and infimum of expected value and variance of the random variables associated to  $B$  within the specified time interval.

**Example 4.1.** Consider the CRN of Example 2.1. Checking if the variance of  $\#\lambda_1$  remains smaller than  $K_1$  within  $[t_j, t_k]$  can be expressed as  $supV[[1, 0, 0]]_{[t_j, t_k]}$ . Another example is

checking if, in the same interval,  $(\#\lambda_1 - \#\lambda_2)$  is at least  $K_2$  or within  $[K_3, K_4]$ , with  $K_3 < K_4 < K_2$ , with probability greater than 0.95:  $P_{>0.95}[[1, -1, 0], ([K_3, K_4], [K_2, \infty))]_{[t_j, t_k]}$ . Equivalently, instead of writing  $B$ , we write directly the linear combination it defines. For example, in the latter case we have  $P_{>0.95}[(\#\lambda_1 - \#\lambda_2), ([K_3, K_4], [K_2, \infty))]_{[t_j, t_k]}$ .

We now comment about expressiveness of SEL in relation to CSL, the logic typically employed to specify properties of the CTMCs induced by CRNs. Though SEL includes the probabilistic operator  $P_{\sim p}[B, I]_{[t_1, t_2]}$ , this is different from the probabilistic reachability operator of CSL. As shown in Bortolussi et al. (2016), under some restrictions SEL can be endowed with the probabilistic reachability operator, but reward operators for SEL have not been studied. The steady-state operator of CSL cannot be handled by LNA because it is accurate only for finite time.

#### 4.1. Semantics

Given a CRN  $C = (\Lambda, R)$  with initial configuration  $x_0$  in a system of fixed volumetric factor  $N$ , its stochastic behaviour is described by the CTMC  $X^N$  of Definition 1. We define a path of CTMC  $X^N$  as a sequence  $\omega = x_0 t_1 x_1 t_1 x_2 \dots$  where  $x_i$  is a state and  $t_i \in \mathbb{R}_{>0}$  is the time spent in the state  $x_i$ . A path is finite if there is a state  $x_k$  that is absorbing.  $\omega \otimes t$  is the state of the path at time  $t$ .  $Path(X^N, x_0)$  is the set of all (finite and infinite) paths of the CTMC starting in  $x_0$ . We work with the standard probability measure  $Prob$  over paths  $Path(X^N, x_0)$  defined using cylinder sets (Kwiatkowska et al., 2007).

We first define when a path  $\omega$  satisfies  $(B, I)$  at time  $t$ :

$$\omega, t \models (B, I) \leftrightarrow \exists [l_i, u_i] \in I \cdot l_i \leq B(\omega \otimes t) \leq u_i.$$

Note that  $B(\omega \otimes t)$  is well defined because  $\omega \otimes t \in \mathbb{N}^{|\Lambda|}$ .

We now define  $Pr_{B, I}^{X^N}(t) = Prob\{\omega \in Path(X^N, x_0) \mid \omega, t \models (B, I)\}$ , then if the time interval is a singleton the satisfaction relation for the probabilistic operator is:

$$X^N, x_0 \models P_{\sim p}[B, I]_{[t_1, t_1]} \leftrightarrow Pr_{B, I}^{X^N}(t_1) \sim p$$

Instead, for  $t_1 < t_2$  we have:

$$X^N, x_0 \models P_{\sim p}[B, I]_{[t_1, t_2]} \leftrightarrow \frac{1}{t_2 - t_1} \int_{t_1}^{t_2} Pr_{B, I}^{X^N}(t) dt \sim p$$

$Pr_{B, I}^{X^N}(t)$  is the probability of the set of paths of  $X^N$  such that the linear combination of the species defined by  $B$  falls within  $I$ . It is well defined since we have previously defined the probability measure  $Prob$  on  $Path(X^N, x_0)$ . To define the satisfaction relation of the probabilistic operator we simply take the average value of  $Pr_{B, I}^{X^N}(t)$  during the interval  $[t_1, t_2]$ . For the remaining operators the satisfaction relation is defined as:

$$X^N, x_0 \models supV[B]_{[t_1, t_2]} \leftrightarrow \sup_{\sim v} (C[B(X^N)], [t_1, t_2]) \sim v$$

$$X^N, x_0 \models infV[B]_{[t_1, t_2]} \leftrightarrow \inf_{\sim v} (C[B(X^N)], [t_1, t_2]) \sim v$$

$$X^N, x_0 \models supE[B]_{[t_1, t_2]} \leftrightarrow \sup_{\sim v} (E[B(X^N)], [t_1, t_2]) \sim v$$

$$X^N, x_0 \models infE[B]_{[t_1, t_2]} \leftrightarrow \inf_{\sim v} (E[B(X^N)], [t_1, t_2]) \sim v$$

$$X^N, x_0 \models \eta_1 \wedge \eta_2 \leftrightarrow X^N, x_0 \models \eta_1 \wedge X^N, x_0 \models \eta_2$$

$$X^N, x_0 \models \eta_1 \vee \eta_2 \leftrightarrow X^N, x_0 \models \eta_1 \vee X^N, x_0 \models \eta_2$$

where  $inf(\cdot, [t_1, t_2])$  and  $sup(\cdot, [t_1, t_2])$  respectively denote the infimum and supremum within  $[t_1, t_2]$ .

#### 4.2. LNA-based Approximate Model Checking for CRNs

Stochastic model checking of CRNs is usually achieved by transient analysis of the CTMC  $X^N$  (Kwiatkowska et al., 2007), which involves solving the CME and thus suffers from the state-space explosion problem. We propose an approximate model checking algorithm based on LNA. The inputs are a SEL formula  $\eta$ , the stochastic process  $X^N$  induced by the CRN and initial state  $x_0$ . The output is *true* in case the formula is verified, and otherwise *false*.

The algorithm proceeds by induction on the structure of formula  $\eta$ , successively computing whether each subformula is satisfied or not. We assume that Eqs. (6) and (7) are solved numerically where  $\Sigma$  is the finite set of sample points on which their solution is defined

and that  $t_0$ , initial time, and  $t_{max}$ , final time, are always sampling points.

#### 4.2.1. Probabilistic operator

To evaluate  $P_{\sim p}[(B, I)]_{[t_1, t_2]}$  we construct the function  $Prob_{(B, I)}(t) = \Omega_{Y^N, B, I}(t_i)$  for  $t \in [t_i, t_{i+1})$ ,  $t_i, t_{i+1} \in \Sigma$  (alternatively, it can be constructed as the interpolation of the values of  $\Omega_{Y^N, B, I}$  over  $\Sigma$  points).

**Lemma 1.**  $Prob_{(B, I)}$  is integrable on  $\mathbb{R}_{\geq 0}$ .

**Proof.**  $Prob_{(B, I)}$  is a bounded function with at most  $|M_I|$  discontinuities, where  $|M_I| \in \mathbb{N}_{>0}$ . Therefore, the set of discontinuities is a countable set, and countable sets have measure 0. Hence, the conditions of the Lebesgue criterion for integrability holds. This concludes the proof.  $\square$

**Theorem 2** guarantees the pointwise correctness of  $Prob_{(B, I)}$  and its integrability allows us to compute the following approximation, then compare to threshold  $p$  to decide the truth value. If  $t_2 \neq t_1$  then  $1/(t_2 - t_1) \int_{t_1}^{t_2} Pr_{B, I}^{X^N}(t) dt \approx 1/(t_2 - t_1) \int_{t_1}^{t_2} Prob_{B, I}(t) dt$  else if  $t_1 = t_2$  then  $Pr_{B, I}^{X^N}(t_1) \approx Prob_{B, I}(t_1)$ .

#### 4.2.2. Expectation and variance operators

To evaluate  $sup(C[B(X^N)], [t_1, t_2])$ ,  $inf(C[B(X^N)], [t_1, t_2])$ ,  $sup(E[B(X^N)], [t_1, t_2])$  and  $inf(E[B(X^N)], [t_1, t_2])$  we use the LNA, namely, compute the expected value and variance of Eqs. (9) and (10). **Theorem 3** guarantees the quality of the approximation. We can now compute the following approximations, then compare to the threshold  $v$ :

$$sup(C[B(X^N)], [t_1, t_2])$$

$$\approx \max\{C[BY^N(t_k)] \mid (t_k \in \Sigma \wedge t_1 \leq t_k \leq t_2) \vee (t_k \in L_{[t_1, t_2]})\}$$

$$inf(C[B(X^N)], [t_1, t_2])$$

$$\approx \min\{C[BY^N(t_k)] \mid (t_k \in \Sigma \wedge t_1 \leq t_k \leq t_2) \vee (t_k \in L_{[t_1, t_2]})\}$$

and similarly for the expected value.  $L_{[t_1, t_2]} = \{t_i \in \Sigma \mid \#t_j \in \Sigma \cdot |t_1 - t_j| < |t_1 - t_i|\}$  ensures that for any time interval there is at least one sampling point, even if the interval is a singleton. Note that, for each sub-formula, the algorithm involves the calculation of some quantity, so one can define a quantitative semantics for SEL as in [Donzé and Maler \(2010\)](#). In our implementation, the syntax for obtaining this numerical value is by replacing the bounds in the  $Q$  and  $P$  operators with “=” as shown in the case studies in Section 5.

LNA-based model checking can also be used for systems far from the thermodynamic limit, at a cost of some loss of precision. LNA assumes continuous state space, and it is not possible to justify this assumption for very small populations. However, if the distributions of interest are not multi-modal and the noise term is finite and approximated by a Gaussian distribution, then LNA gives very good approximation even for quite small systems. It is clear that model checking accuracy increases as  $N$  grows. We emphasize that the model checking algorithm we have presented is also able to handle CRNs whose stochastic semantics is an infinite CTMC, which occur frequently in biological models.

#### 4.2.3. Complexity of LNA-based approximate model checking

The time complexity for model checking formula  $\eta$  against a CRN  $C = (\Lambda, R)$  is linear in  $|\eta|$ . In the worst case, analysis of a single operator requires the solution of  $O(|\Lambda|^2)$  polynomial differential equations for a bounded time. However, an efficient implementation can solve the  $O(|\Lambda|^2)$  ODEs only once for the interval  $[0, t_{max}]$ , and then reuse this result for every operator, where  $t_{max}$  is the greatest (finite) time of interest. Note that ODEs are solved in terms of

concentrations (a value between 0 and 1 by convention), ensuring independence from the number of molecules of each species, although stiffness can slow down the solution of the LNA.

## 5. Experimental results

We implemented the methods in a framework based on Matlab and Java. The experiments were run on an Intel Dual Core i7 machine with 8 GB of RAM. To solve the differential equations, we use *Matlab ode45*, a variable step Runge–Kutta algorithm. We employ LNA-based model checking for the analysis of four biological reaction networks: a Phosphorelay Network ([Csikász-Nagy et al., 2011](#)), a Gene Expression Model ([Thattai and Van Oudenaarden, 2001](#); [Mateescu, 2011](#)), the FGF pathway ([Heath et al., 2008](#)) and the GW network ([Cardelli, 2014](#)). For every network, the CRN and parameters have been taken from the referenced papers. We coded the same CRNs in PRISM ([Kwiatkowska et al., 2011](#)) in order to compare accuracy and time of execution with standard uniformisation of the CME ([Kwiatkowska et al., 2007](#)) and statistical model checking (SMC) techniques (confidence interval method) as implemented in PRISM. For the FGF and GW case studies, global analysis and SMC cannot be used, because the state space is too large for direct analysis, and SMC requires many time-consuming simulations to obtain good accuracy.

### 5.1. Phosphorelay network

We consider a three-layer phosphorelay network whose structure is derived from [Csikász-Nagy et al. \(2011\)](#). Phosphorelay networks are extended two-component signalling systems found in diverse bacteria, lower eukaryotes and plants. Each layer of the network,  $(L1, L2, L3)$ , can be found in phosphorylate form  $(L1p, L2p, L3p)$ . We consider the initial condition  $\#L1p = \#L2p = \#L3p = 0$ ,  $\#L1 = \#L2p = \#L3p = Init$ , where  $Init \in \mathbb{N}$ . Then we analyze the ligand  $B$ , whose initial condition is  $\#B = 3 * Init$ . We are interested in checking the following SEL property:

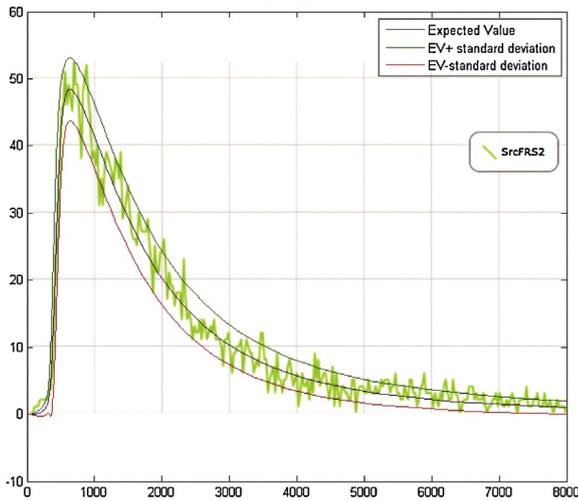
$$P_{>0.7}[(\#L1p - \#L3p), [0, +\infty)]_{[0, 100]}$$

$$\wedge P_{>0.98}[(\#L3p - \#L1p), [0, +\infty)]_{[300, 600]}$$

which is verified if, in the first interval, the probability that  $\#L1p$  is greater than  $\#L3p$  is  $>0.7$  and if, between 300 and 600, with probability  $>0.98$ ,  $\#L3p$  is greater than  $\#L1p$ . We evaluate this formula in three different initial conditions, firstly  $Init = 32$  and  $N = 5000$ , then  $Init = 64$  and  $N = 10,000$ , and finally  $Init = 100$  and  $N = 15,625$ , so the same concentration but different numbers of molecules. In all cases, the LNA-based model checking evaluates the formula as true. To understand the quality of the approximation, we check the following quantitative formula  $P_{\sim \gamma}[(\#L3p - \#L1p), [0, +\infty)]_{[T, T]}$  for  $T \in [0, 600]$  (recall that in our implementation “=” gives the quantity calculated by model checking the operator). We compare the results with the evaluation of the corresponding CSL formula using standard uniformisation (*Unif*) with error  $10^{-7}$  ([Kwiatkowska et al., 2007](#)). The following table shows the results. *MaxErr* is the maximum error computed by LNA-based approach compared to standard uniformisation and *AvgErr* is the average error; *Time(.)* stands for execution time.

Init	Time (LNA)	Time (Unif)	MaxErr	AvgErr
20	0.22 s	2 min	0.0675	0.0519
32	0.23 s	5 min	0.059	0.02
64	0.26 s	>2 h	0.0448	0.0027
100	0.3 s	>2 h	0.03	0.0011

Note that as *Init* increases the error of our method decreases, while the execution time is practically independent of the molecular count. LNA-based algorithms are faster in all cases. Thus our



**Fig. 2.** Expected number and standard deviation of species of  $\#Src:FRS2$  in the FGF pathway during the first 8000s estimated by our method is compared with a stochastic simulation of the same species.

approach can be used even for quite small population systems, giving fast approximate stochastic characterization.

## 5.2. Gene expression

We consider a simple CRN that models the transcription of a gene into a mRNA molecule, and the translation of the latter into a protein. The CRN, rates and initial conditions are the same as in [Mateescu \(2011\)](#). The stochastic semantics of the reaction network is an infinite CTMC, and we use this model to show that our method can handle infinite state-space processes. We consider the quantitative property  $supE[\#mRNA]_{[T,T]}$ , which gives the number of molecules of  $mRNA$  in the system at time  $T$ . We compare our method with SMC estimation of the same property by using 50,000 simulations, for  $T = \{300, 600, 900, 1200\}$ , and in the following tables we compare the results in terms of execution time ( $Time(\cdot)$ ) and estimated expected value of  $\#mRNA$  ( $ExpVal(\cdot)$ ). LNA-based model checking is several orders of magnitude faster without loss of accuracy.

$T$	Time (LNA)	Time (Simul)	ExpVal (LNA)	ExpVal (Simul)
300	0.52 s	75 s	100.17	$100.14 \pm 0.1$
600	0.54 s	198 s	142.15	$142.11 \pm 0.1$
900	0.54 s	337 s	159.73	$159.74 \pm 0.1$
1200	0.56 s	483 s	167.1	$167.1 \pm 0.1$

## 5.3. Fibroblast growth factor (FGF) pathway

Fibroblast growth factors (FGF) are a family of proteins which play a key role in the process of cell signalling in a variety of contexts, like wound healing and skeletal development. We consider the model of FGF signalling pathway developed in [Heath et al. \(2008\)](#), which is composed of more than 50 reactions and species. We consider the system with initially 105 molecules for species with non-zero initial concentration. Analysis of the model reveals that the phosphorylated form of  $FRS2$  can bind the protein  $Src$ , and then this new complex,  $Src:FRS2$ , can relocate out. We want to check if the expected value of  $\#Src:FRS2$  during the first 3000s reaches a maximum value greater than 40. We do that by checking the property  $supE[\#Src : FRS2]_{[0,3000]}$ . The formula evaluates to true, and in [Fig. 2](#) we analyze the expected value and standard deviation of  $\#Src:FRS2$ . We obtain these values directly from the logic

considering the quantitative interpretation of  $supE[\#Src : FRS2]_{[T,T]}$  and  $supV[\#Src : FRS2]_{[T,T]}$  for  $T \in [0, 3000]$ . It is possible to see that, after an initial peak, relocation causes exponential decay.

In the same figure we show a single stochastic simulation of the system for the same initial conditions, confirming our evaluation. Moreover, the approximation can be justified theoretically.  $\#Src:FRS2$  converges to zero necessarily and this demonstrates the unimodality of the distribution of the species; we note that the variance is finite, so Eq. (4) holds.

## 5.4. DNA strand displacement of the GW network

GW is a network related to the G2-M cell cycle switch ([Novak and Tyson, 1993](#)). Under particular initial conditions, it has been shown that GW can emulate the Approximate Majority algorithm ([Cardelli, 2014](#)). Here, we consider the two-domain DNA strand-displacement implementation of GW ([Cardelli, 2013](#)). The corresponding CRN is composed of 340 species and 240 reactions. For our analysis the species of interest are  $R$  and  $P$ , whose initial conditions are  $\#R = 90$  and  $\#P = 10$ . These species model the switch for the activating phosphatase  $Cdc25$ ; initial conditions of other species are taken from the referenced papers. We check the property  $P_{>0.6}[\#R - \#P, [65, +\infty]]_{[0,T]}$  for  $T = \{1000, 2000, 3500, 5000\}$ , in a system of size  $N = 45,000$ . The results are reported in the following table.

$T$	Time execution	Quantitative value	Qualitative value
1000	420 s	0.4297	False
2000	780 s	0.5313	False
3500	1380 s	0.6535	True
5000	2120 s	0.7349	True

## 6. Concluding remarks

We presented a novel probabilistic logic (SEL) for analysing stochastic behaviour of CRNs and proposed an approximate model checking algorithm of the CME based on the LNA. We have implemented the algorithm and demonstrated on four non-trivial examples that LNA-based model checking enables analysis of CRNs with hundreds of species, and even infinite CTMCs, at a cost of some loss of accuracy. It would be interesting to find bounds on the approximation error when the system is far from the thermodynamic limit. However, the error is not only dependent on the value of  $N$ , but also on the structure of the CRN, the rates, and the property. As recently shown in [Cardelli et al. \(2016\)](#), it is possible to formulate a stochastic hybrid approach, in which the LNA is used only for a subset of species, namely, those for which the leap conditions are satisfied. Other species are treated as a discrete-state Markov process. This improves precision of the stochastic analysis of CRNs when multimodality is present. One of the most attractive features of the LNA is that it enables a stochastic analysis of a CRN by solving a set of ODEs quadratic in the number of species. We aim to exploit this feature in order to enable synthesis and symbolic analysis of CRNs. It would also be interesting to extend SEL with reward operators as in, e.g., [Baier et al. \(2000\)](#) and [Kwiatkowska et al. \(2007\)](#), with which one can express properties such as the expected number of molecules and variance of a species when certain events happen.

## Acknowledgements

The authors would like to thank Luca Bortolussi for helpful discussions. This research is supported by a Royal Society Research Professorship and ERC AdG VERIWARE.

## References

- Ammar, A., Mokdad, B., Chinesta, F., Keunings, R., 2006. A new family of solvers for some classes of multidimensional partial differential equations encountered in kinetic theory modeling of complex fluids. *J. Non-Newton. Fluid Mech.* 139 (3), 153–176.
- Ammar, A., Cueto, E., Chinesta, F., 2012. Reduction of the chemical master equation for gene regulatory networks using proper generalized decompositions. *Int. J. Numer. Methods Biomed. Eng.* 28 (9), 960–973.
- Anderson, D.F., Kurtz, T.G., 2011. Continuous time Markov chain models for chemical reaction networks. In: *Design and Analysis of Biomolecular Circuits*. Springer, pp. 3–42.
- Aziz, A., Sanwal, K., Singhal, V., Brayton, R., 2000. Model-checking continuous-time Markov chains. *ACM Trans. Comput. Logic* 1 (1), 162–170.
- Baier, C., Haverkort, B., Hermanns, H., Katoen, J.-P., 2000. On the logical characterisation of performability properties. In: *International Colloquium on Automata, Languages, and Programming*. Springer, pp. 780–792.
- Baier, C., Haverkort, B., Hermanns, H., Katoen, J.-P., 2003. Model-checking algorithms for continuous-time Markov chains. *IEEE Trans. Soft. Eng.* 29 (6), 524–541.
- Ballarini, P., Djafri, H., Dufloy, M., Haddad, S., Pekergin, N., 2011a. COSMOS: a statistical model checker for the hybrid automata stochastic logic. In: *Proc. 8th Int. Conf. Quantitative Evaluation of Systems (QEST'11)*. IEEE Computer Society Press, pp. 143–144.
- Ballarini, P., Djafri, H., Dufloy, M., Haddad, S., Pekergin, N., 2011b. Cosmos: a statistical model checker for the hybrid automata stochastic logic. In: *2011 Eighth International Conference on Evaluation of Systems (QEST)*. IEEE, pp. 143–144.
- Barbot, B., Kwiatkowska, M., 2015. On quantitative modelling and verification of DNA walker circuits using stochastic petri nets. In: *Devillers, R., Valmari, A. (Eds.), Application and Theory of Petri Nets and Concurrency, Vol. 9115 of Lecture Notes in Computer Science*. Springer International Publishing, pp. 1–32.
- Bortolussi, L., Lanciani, R., 2013. Model checking Markov population models by central limit approximation. In: *International Conference on Quantitative Evaluation of Systems*. Springer, pp. 123–138.
- Bortolussi, L., Hillston, J., Latella, D., Massink, M., 2013. Continuous approximation of collective system behaviour: a tutorial. *Perform. Eval.* 70 (5), 317–349.
- Bortolussi, L., Cardelli, L., Kwiatkowska, M., Laurenti, L., 2016. Approximation of probabilistic reachability for chemical reaction networks using the linear noise approximation. In: *Agha, G., Van Houdt, B. (Eds.), Proceedings Quantitative Evaluation of Systems: 13th International Conference (QEST 2016)*. Springer, QC, Canada, pp. 72–88, [http://dx.doi.org/10.1007/978-3-319-43425-4\\_5](http://dx.doi.org/10.1007/978-3-319-43425-4_5). ISBN: 978-3-319-43425-4.
- Butcher, J.C., 1987. *The Numerical Analysis of Ordinary Differential Equations: Runge–Kutta and General Linear Methods*. Wiley-Interscience.
- Cardelli, L., 2008. On process rate semantics. *Theor. Comput. Sci.* 391 (3), 190–215.
- Cardelli, L., 2013. Two-domain DNA strand displacement. *Math. Struct. Comput. Sci.* 23 (02), 247–271.
- Cardelli, L., 2014. Morphisms of reaction networks that couple structure to function. *BMC Syst. Biol.* 8 (1), 84.
- Cardelli, L., Kwiatkowska, M., Laurenti, L., 2016. A stochastic hybrid approximation for chemical kinetics based on the linear noise approximation. In: *Proceedings International Conference on Computational Methods in Systems Biology*. Springer, pp. 147–167.
- Chellaboina, V., Bhat, S., Haddad, M., Bernstein, D.S., 2009. Modeling and analysis of mass-action kinetics. *IEEE Control Syst.* 29 (4), 60–78.
- Chen, Y.-J., Dalchau, N., Srinivas, N., Phillips, A., Cardelli, L., Soloveichik, D., Seelig, G., 2013. Programmable chemical controllers made from DNA. *Nat. Nanotechnol.* 8 (10), 755–762.
- Chinesta, F., Magnin, M., Roux, O., Ammar, A., Cueto, E., 2015. Kinetic theory modeling and efficient numerical simulation of gene regulatory networks based on qualitative descriptions. *Entropy* 17 (4), 1896–1915.
- Cinlar, E., 2013. *Introduction to Stochastic Processes*. Courier Corporation.
- Csikász-Nagy, A., Cardelli, L., Soyer, O.S., 2011. Response dynamics of phosphorelays suggest their potential utility in cell signalling. *J. R. Soc. Interface* 8 (57), 480–488.
- Dannenber, F., Hahn, E.M., Kwiatkowska, M., 2015. Computing cumulative rewards using fast adaptive uniformization. *ACM Trans. Model. Comput. Simul. (TOMACS)* 25 (2), 9.
- Didier, F., Henzinger, T.A., Mateescu, M., Wolf, V., 2009. Fast adaptive uniformization of the chemical master equation. In: *International Workshop on High Performance Computational Systems Biology, 2009, HIBI'09*. IEEE, pp. 118–127.
- Donzé, A., Maler, O., 2010. *Robust Satisfaction of Temporal Logic Over Real-Valued Signals*. Springer.
- Elf, J., Ehrenberg, M., 2003. Fast evaluation of fluctuations in biochemical networks with the linear noise approximation. *Genome Res.* 13 (11), 2475–2484.
- Ethier, S.N., Kurtz, T.G., 2009. *Markov Processes: Characterization and Convergence*, vol. 282. John Wiley & Sons.
- Gillespie, D.T., 2000. The chemical Langevin equation. *J. Chem. Phys.* 113 (1), 297–306.
- Gillespie, D.T., 2009. Deterministic limit of stochastic chemical kinetics. *J. Phys. Chem. B* 113 (6), 1640–1644.
- Gillespie, D.T., Hellander, A., Petzold, L.R., 2013. Perspective: Stochastic algorithms for chemical kinetics. *J. Chem. Phys.* 138 (17), 170901.
- Heath, J., Kwiatkowska, M., Norman, G., Parker, D., Tymchyshyn, O., 2008. Probabilistic model checking of complex biological pathways. *Theor. Comput. Sci.* 391 (3), 239–257.
- Hespanha, J., 2008. Moment closure for biochemical networks. In: *3rd International Symposium on Communications, Control and Signal Processing, 2008, ISCCSP 2008*. IEEE, pp. 142–147.
- Itô, K., 2006. *Essentials of Stochastic Processes*, vol. 231. American Mathematical Soc.
- Kwiatkowska, M., Thachuk, C., 2014. Probabilistic model checking for biology. *Softw. Syst. Saf.* 36, 165.
- Kwiatkowska, M., Norman, G., Parker, D., 2007. Stochastic model checking. In: *International School on Formal Methods for the Design of Computer, Communication and Software Systems*. Springer, pp. 220–270.
- Kwiatkowska, M., Norman, G., Parker, D., 2011. Prism 4.0: Verification of probabilistic real-time systems. In: *International Conference on Computer Aided Verification*. Springer, pp. 585–591.
- Mateescu, M.-E.-C., 2011. *Propagation models for biochemical reaction networks*. Ph.D. thesis, EPFL.
- Munsky, B., Khammash, M., 2006. The finite state projection algorithm for the solution of the chemical master equation. *J. Chem. Phys.* 124 (4), 044104.
- Novak, B., Tyson, J.J., 1993. Numerical analysis of a comprehensive model of M-phase control in *Xenopus oocyte* extracts and intact embryos. *J. Cell Sci.* 106 (4), 1153–1168.
- Patel, J.K., Read, C.B., 1996. *Handbook of the normal distribution*, vol. 150. CRC Press.
- Pinsky, M., Karlin, S., 2010. *An Introduction to Stochastic Modeling*. Academic Press.
- Risken, H., 1984. *Fokker–Planck Equation*. Springer.
- Singh, A., Hespanha, J.P., 2006. Lognormal moment closures for biochemical reactions. In: *Proceedings of the 45th IEEE Conference on Decision and Control*. IEEE, pp. 2063–2068.
- Soloveichik, D., Seelig, G., Winfree, E., 2010. DNA as a universal substrate for chemical kinetics. *Proc. Natl. Acad. Sci. U. S. A.* 107 (12), 5393–5398.
- Thattai, M., Van Oudenaarden, A., 2001. Intrinsic noise in gene regulatory networks. *Proc. Natl. Acad. Sci. U. S. A.* 98 (15), 8614–8619.
- Van Kampen, N.G., 1992. *Stochastic Processes in Physics and Chemistry*, vol. 1. Elsevier.
- Wallace, E., Gillespie, D., Sanft, K., Petzold, L., 2012. Linear noise approximation is valid over limited times for any chemical system that is sufficiently large. *IET Syst. Biol.* 6 (4), 102–115.
- Wolf, V., Goel, R., Mateescu, M., Henzinger, T.A., 2010. Solving the chemical master equation using sliding windows. *BMC Syst. Biol.* 4 (1), 42.