# Heavy-tailed distributions in a stochastic gene autoregulation model

Pavol Bokes

June 3, 2021

**Abstract**

Synthesis of gene products in bursts of multiple molecular copies is an important source of gene expression variability. This paper studies large deviations in a Markovian drift–jump process that combines exponentially distributed bursts with deterministic degradation. Large deviations occur as a cumulative effect of many bursts (as in diffusion) or, if the model includes negative feedback in burst size, in a single big jump. The latter possibility requires a modification in the WKB solution in the tail region. The main result of the paper is the construction, via a modified WKB scheme, of matched asymptotic approximations to the stationary distribution of the drift–jump process. The stationary distribution possesses a heavier tail than predicted by a routine application of the scheme.

***Keywords:*** stochastic gene expression, bursting, WKB approximation, large deviations

***MSC 2020:*** 92C40; 60J76, 45D05, 41A60

## 1 Introduction

Bursty production of gene products (mRNA or protein molecules) makes an important contribution to the overall gene expression noise [1–4]. Bursts can be modelled as instantaneous jumps of a random process. Burst sizes have

been suggested to follow geometric (in a discrete process) or exponential (in a continuous process) distributions [5, 6]; we focus on the latter. Production of gene products is balanced by their degradation and/or dilution. Combining randomly timed and sized production bursts with deterministic decay leads to a Markovian drift–jump model of gene expression [7–10]. More fine-grained models of gene expression are based on a purely discrete [11–14] or a hybrid discrete–continuous state space [15–18]. The drift–jump model can be derived from the fine-grained processes using formal limit procedures [19–24].

In its basic formulation, the drift–jump model for gene expression admits a gamma stationary distribution [25]. The model possesses an explicit stationary distribution also in the presence of a Hill-type feedback in burst frequency [26]. Such regulation can result from common transcriptional control mechanisms [27]. In addition to feedback in burst frequency, there is evidence of feedback mechanisms that act on burst size or protein stability [28–30]. The explicit stationary solution to the drift–jump model has been extended to the case of feedback in protein stability [31]. However, in case of feedback in burst size, an explicit solution is unavailable, save for the special case of Michaelis–Menten-type response [32].

The near-deterministic regime of frequent and small bursts can be analysed using the Wentzel–Kramers–Brillouin (WKB) method; the WKB-approximate solutions closely agree with numerically obtained exact distributions even at moderate noise conditions [33]. Bursty production has been formulated and analysed with the WKB method also in the discrete state space [34–38]. Similar approaches have earlier been used in queueing systems [39, 40]. The standard WKB-type/diffusion-like results are guaranteed to apply for jump-size distributions with super-exponentially decaying tails [41]. Contrastingly, in the sub-exponential case, large deviations are driven by single big jumps [42]. The exponential case can combine both phenomena for random walks: the Cramer/WKB-type result applies in a region of sample space called the Cramer zone, while single big jumps contribute to deviations beyond the Cramer zone [43, 44].
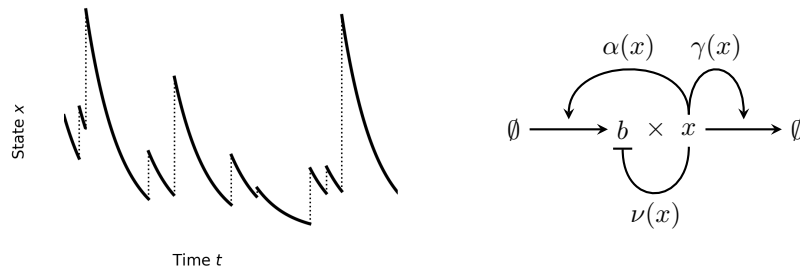
Figure 1: *Left:* Sketch of a typical sample path of the drift–jump gene-expression model. *Right:* Functions $\alpha(x)$, $\nu(x)$, and $\gamma(x)$ quantify feedback in burst frequency, burst size, and protein stability.

In this paper, the standard WKB-type approach will be shown to be suitable for the drift-jump gene expression model with positive feedback in burst size. If the feedback is negative, the WKB-approach will be shown to apply below a certain threshold (referred to, by analogy with random walks, as the Cramer zone), whereas beyond the threshold (referred to as the tail zone) single big jumps contribute to large deviations. Matched asymptotic approximations to the stationary distribution in the Cramer zone, in the tail zone, and on their boundary will be constructed using a formal singular perturbation approach [45–47].

The structure of the paper is as follows. Section 2 formulates the model. Section 3 presents the standard WKB approximation scheme. The core of the paper is Section 4, in which the modified WKB scheme is given. The boundary layer is treated in Section 5. The asymptotic results are cross-validated by simulations in Section 6. The paper is concluded in Section 7.

## 2 Model formulation

The drift–jump gene-expression model is a Markov process with piecewise continuous sample paths (Figure 1, left panel). The state $x$ of the process represents the concentration of a gene product (say a protein, for concreteness). The discontinuities in the sample path are the production bursts. Between bursts,

the protein concentration decays deterministically with rate constant $\gamma(x)$, i.e. as per $\dot{x} = -\gamma(x)$. Bursts occur with state-dependent frequency (propensity) $\varepsilon^{-1}\alpha(x)$. Burst sizes are drawn from an exponential distribution with rate parameter $\varepsilon^{-1}\nu(x)$, in which $x$ is the state of the process immediately before the burst; the reciprocal $\varepsilon/\nu(x)$ of the rate parameter gives the mean burst size. Decreasing the noise strength $\varepsilon$ makes bursts more frequent and smaller. The functions $\alpha(x)$, $\nu(x)$, and $\gamma(x)$ can implement feedback in burst frequency, burst size, and protein stability (Figure 1, right panel).

The probability density function $p(x,t)$ of being at state $x$ at time $t$ satisfies the integro–differential equation

$$\frac{\partial p}{\partial t} + \frac{\partial J}{\partial x} = 0, \tag{1}$$

$$J(x,t) = -\gamma(x)p(x,t) + \frac{1}{\varepsilon}\int_0^x p(y,t)\alpha(y)\exp\left(-\frac{\nu(y)(x-y)}{\varepsilon}\right)\mathrm{d}y. \tag{2}$$

In the conservation equation (1), $J = J(x,t)$ gives the flux of probability across a reference state $x$ at time $t$. By (2), it consists of a negative local flux due to deterministic decay and a positive non-local flux due to stochastic bursts. The non-local term integrates, over all states $y < x$, the probability $\varepsilon^{-1}p(y,t)\alpha(y)$ that a burst occurs multiplied by the exponential probability that the burst goes beyond the reference state $x$.

Estimating the integral in (2) by the Laplace method [48] as $\varepsilon \to 0$, we obtain $J \sim (\alpha(x)/\nu(x) - \gamma(x))p(x,t)$, which is the probability flux of a purely deterministic process

$$\frac{\mathrm{d}x}{\mathrm{d}t} = \frac{\alpha(x)}{\nu(x)} - \gamma(x). \tag{3}$$

Equation (3) is the deterministic limit of (1)–(2) (sometimes also referred to as the fluid limit or the law-of-large-numbers limit). Retaining a further term in the asymptotic expansion of the non-local term leads to an ad-hoc drift–diffusion approximation to the drift–jump process [49]. Such truncations exhibit different $\varepsilon \to 0$ asymptotics than the original problem [50].

4

Equating the flux in (2) to zero, we obtain a Volterra integral master equation

$$\gamma(x)p(x) = \frac{1}{\varepsilon} \int_0^x p(y)\alpha(y) \exp\left(-\frac{\nu(y)(x-y)}{\varepsilon}\right) \mathrm{d}y \qquad (4)$$

for the stationary distribution. Multiplying a solution $p(x)$ to (4) by a constant gives another solution. The multiplicative constant can be fixed by requiring that the total probability integrate to one. However, the dependence of the normalisation constant on $\varepsilon$ introduces unnecessary complications in the asymptotic expansions; we defer the normalisation until Section 6.

The principal aim of Sections 3–5 is to characterise the $\varepsilon \to 0$ asymptotics of solutions $p(x) = p(x; \varepsilon)$ to the Volterra master equation (4).

## 3   Standard WKB scheme

We seek an approximate solution to (4) in the WKB form

$$p(x; \varepsilon) = r(x; \varepsilon) \exp\left(-\frac{\Phi(x)}{\varepsilon}\right), \qquad (5)$$

where a regular dependence

$$r(x; \varepsilon) = r_0(x) + \varepsilon r_1(x) + O(\varepsilon^2) \qquad (6)$$

of the prefactor on $\varepsilon$ is postulated. The function $\Phi(x)$ in (5) is referred to as the quasipotential.

Inserting (5) into (4) gives

$$\gamma(x)r(x) \exp\left(-\frac{\Phi(x)}{\varepsilon}\right) = \frac{1}{\varepsilon} \int_0^x r(y)\alpha(y) \exp\left(-\frac{\Psi(x,y)}{\varepsilon}\right) \mathrm{d}y, \qquad (7)$$

where

$$\Psi(x,y) = \Phi(y) + \nu(y)(x-y). \qquad (8)$$

5

Differentiating (8) with respect to $y$ and setting $y = x$ gives relations

$$\Psi(x,x) = \Phi(x), \quad \partial_y \Psi(x,x) = \Phi'(x) - \nu(x) \quad \partial_y^2 \Psi(x,x) = \Phi''(x) - 2\nu'(x), \quad (9)$$

which tie up the local behaviour of $\Psi(x,y)$ near the boundary $y = x$ and that of the (yet unknown) quasipotential.

Provided that

$$\Psi(x,y) > \Psi(x,x) \quad \text{for} \quad y < x, \tag{10}$$

the dominant contribution to the integral on the right-hand side of (7) comes from an $O(\varepsilon)$-wide neighbourhood of the right boundary. Estimating the integral in (7) by the Laplace method and cancelling the common exponential term gives

$$\gamma(x)r(x) + \frac{\alpha(x)r(x)}{\partial_y \Psi(x,x)} = \varepsilon \frac{r(x)\alpha(x)\partial_y^2\Psi(x,x) - (r(x)\alpha(x))'\partial_y\Psi(x,x)}{(\partial_y\Psi(x,x))^3} + O(\varepsilon^2). \tag{11}$$

Inserting (6) and (9) into (11), and collecting $O(1)$ terms, yields the quasipotential

$$\Phi(x) = \int \nu(x) - \frac{\alpha(x)}{\gamma(x)} \mathrm{d}x, \tag{12}$$

while collecting $O(\varepsilon)$ terms determines the prefactor

$$r_0(x) = \frac{1}{\gamma(x)} \exp\left(\int \frac{\nu'(x)\gamma(x)}{\alpha(x)} \mathrm{d}x\right). \tag{13}$$

The constants of integration in the indefinite integrals in (12)–(13) add up to the normalisation constant in the probability distribution (5) and can be chosen arbitrarily.

The weak point of this section is the assumption (10). Combining (8) and (12), we see that

$$\partial_y\Psi(x,y) = -\frac{\alpha(y)}{\gamma(y)} + \nu'(y)(x - y). \tag{14}$$

If $\nu(x)$ is decreasing (positive feedback case), $\partial_y\Psi(x,y) < 0$ for $y \le x$, which

6

confirms (10) post hoc. If $\nu(x)$ is constant (no feedback in burst size), then $r_0(x) \exp(-\Phi(x)/\varepsilon)$ with (12)–(13) is the exact solution to (4) [31]. The case of negative feedback in burst size requires a subtler analysis, which is the subject of the rest of the paper.

# 4   Modified WKB scheme

From now on, we refer to the function $\Phi(x)$ defined by (12) as the local potential. The name reflects the fact that its derivation involved a local estimate of the integral in the Volterra master equation (7). We assume that the local potential satisfies

$$\Phi''(x) > 0, \quad \lim_{x \to 0} \Phi(x) = \infty, \quad \lim_{x \to \infty} \frac{\Phi(x)}{x} = \infty. \tag{15}$$

Assumptions (15) are satisfied e.g. by choosing

$$\alpha(x) = 1, \quad \gamma(x) = x, \quad \nu(x) = x^m, \quad m > 0. \tag{16}$$

Graphical examples in this section pertain to the parametric choice (16). The following subsection examines the behaviour of $\Psi(x, y)$ defined by (8) and constructs a modified potential.

## 4.1   Modified potential

For any fixed $y > 0$, equation $\Psi(x, y) = \Phi(x)$ in the unknown $x$ has two roots, the trivial one $x = y$, and a non-trivial one such that $x > y$ (Figure 2, left). Comparing the slopes of $\Phi(x)$ and $\Psi(x, y)$ at their non-trivial intersection, we obtain

$$\nu(x) - \frac{\alpha(x)}{\gamma(x)} > \nu(y) \quad \text{if } \Psi(x, y) \leq \Phi(x) \text{ and } y < x. \tag{17}$$

Let us look at the same equation but reverse the dependency between the two variables. For any fixed $x > 0$, equation $\Psi(x, y) = \Phi(x)$ in the unknown $y$ has a trivial root $y = x$, a non-trivial root $y = y_* < x_*$ if $x = x_*$, and two non-trivial
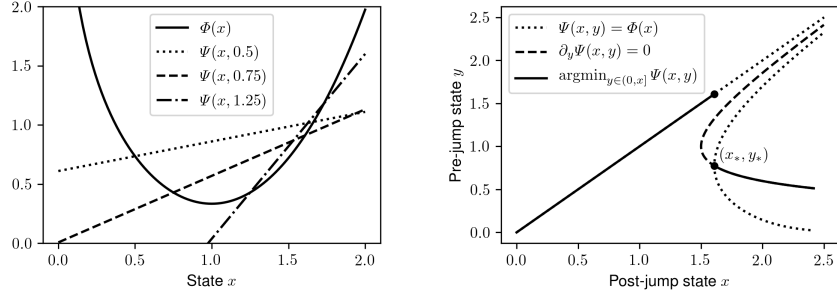
7

Figure 2: Properties of $\Psi(x, y)$ as defined by (8) and (12). The parametric choice (16) with $m = 2$ is used. *Left:* Solutions to $\Psi(x, y) = \Phi(x)$ in the unknown $x$. *Right:* Important curves in the domain of $\Psi(x, y)$.

roots if $x > x_*$ (Figure 2, right, dotted line); the critical pair $(x_*, y_*)$ satisfies

$$\Psi(x_*, y_*) = \Phi(x_*), \quad \partial_y \Psi(x_*, y_*) = 0. \tag{18}$$

Note that (17) implies that

$$\nu(x_*) - \frac{\alpha(x_*)}{\gamma(x_*)} > \nu(y_*). \tag{19}$$

The function $\Psi(x, y)$ is minimised by (cf. Figure 2, right, solid line)

$$\operatorname{argmin}_{y \in (0,x]} \Psi(x, y) = \begin{cases} x & \text{if } x \leq x_*, \\ y_\mathrm{m}(x) & \text{if } x \geq x_*, \end{cases} \tag{20}$$

where $y_\mathrm{m}(x)$ is the lower branch of the critical equation

$$\partial_y \Psi(x, y_\mathrm{m}(x)) = 0, \quad y_\mathrm{m}(x) \leq y_*. \tag{21}$$

We define the modified potential as

$$\tilde{\Phi}(x) = \min_{y \in (0,x]} \Psi(x, y) = \begin{cases} \Phi(x) & \text{if } x \leq x_*, \\ \Psi(x, y_m(x)) & \text{if } x \geq x_*. \end{cases} \tag{22}$$

8

The region $x < x_*$ will be referred to as the Cramer zone, and the complementary region $x > x_*$ as the tail zone. The derivative of the potential in the tail zone satisfies

$$\tilde{\Phi}'(x) = \partial_y \Psi(x, y_{\mathrm{m}}(x)) y_{\mathrm{m}}'(x) + \nu(y_{\mathrm{m}}(x)) = \nu(y_{\mathrm{m}}(x)) \text{ for } x > x_*. \tag{23}$$

Combining (12), (23), and (19), we find

$$\tilde{\Phi}'(x_*^-) = \nu(x_*) - \frac{\alpha(x_*)}{\gamma(x_*)} > \nu(y_*) = \tilde{\Phi}'(x_*^+), \tag{24}$$

meaning that the derivative of the modified potential is discontinuous at the boundary of the Cramer zone.

The purpose of the remainder of this section is to use the modified potential (22) as a basis for a WKB-type approximation to the solution $p(x, \varepsilon)$ to the integral equation (4). In the Cramer zone, condition (10) is satisfied and the standard procedure of Section 3 yields

$$p(x; \varepsilon) \sim r_0(x) \exp\left(-\frac{\Phi(x)}{\varepsilon}\right), \quad x < x_*, \tag{25}$$

where the prefactor is defined by (13). The next section argues that the modified potential (22) is appropriate outside the Cramer zone.

## 4.2 Dominant balance

If we look for a solution $p(x; \varepsilon)$ to (4) in a form that is logarithmically equivalent to $\exp(-\tilde{\Phi}(x)/\varepsilon)$, then the integrand on the right-hand side is logarithmically equivalent to $\exp(-\tilde{\Psi}(x, y)/\varepsilon)$, where

$$\tilde{\Psi}(x, y) = \tilde{\Phi}(y) + \nu(y)(x - y). \tag{26}$$

Let us investigate the behaviour of $\tilde{\Psi}(x, y)$ as function of $y \in (0, x]$ for a fixed $x > x_*$. The Cramer and the tail regions are thereby treated separately:

9

1. $y \le x_*$. Here we have

$$\tilde{\Psi}(x, y) = \Psi(x, y) \ge \Psi(x, y_{\mathrm{m}}(x)) = \tilde{\Phi}(x), \tag{27}$$

with equality in place if $y = y_{\mathrm{m}}(x)$.

2. $y \ge x_*$. Here

$$\tilde{\Psi}(x, y) = \Psi(y, y_{\mathrm{m}}(y)) + \nu(y)(x - y) \tag{28}$$

$$= \Phi(y_{\mathrm{m}}(y)) + \nu(y_{\mathrm{m}}(y))(y - y_{\mathrm{m}}(y)) + \nu(y)(x - y) \tag{29}$$

$$\ge \Phi(y_{\mathrm{m}}(y)) + \nu(y_{\mathrm{m}}(y))(x - y_{\mathrm{m}}(y)) = \Psi(x, y_{\mathrm{m}}(y)) \tag{30}$$

$$\ge \Psi(x, y_{\mathrm{m}}(x)) = \tilde{\Phi}(x), \tag{31}$$

where the estimate (30) holds for a non-decreasing $\nu(x)$ (negative feedback in burst size) and the estimate (31) follows from (20); both estimates become equalities if $y = x$.

The upshot of (27)–(31) is that

$$\tilde{\Psi}(x, y) \ge \tilde{\Phi}(x) \text{ for } y \in (0, x], \text{ with equality if } y \in \{y_{\mathrm{m}}(x), x\}. \tag{32}$$

The integral on the right-hand of (4) side will be logarithmically equivalent to $\exp(-\min_{y \in (0,x]} \tilde{\Psi}(x, y)/\varepsilon) = \exp(-\tilde{\Phi}(x)/\varepsilon)$, which is the asymptotics postulated for the solution. The use of the modified WKB potential (22) thus leads to a desired balance between the sides, at least to a logarithmic precision, of the master equation (4).

Important contributions to the integral term in (4) come from the neighbourhoods of the minimisers $y = y_{\mathrm{m}}(x)$ and $y = x$ of $\tilde{\Psi}(x, y)$ (as function of $y \in (0, x]$ for a fixed $x > x_*$). The function is locally parabolic near the internal minimiser $y = y_{\mathrm{m}}(x) < x_*$, but it is locally linear near the boundary minimiser $y = x > x_*$. By the Laplace method [48], an $O(\varepsilon^{1/2})$ neighbourhood of the parabolic minimiser, but only an $O(\varepsilon)$ neighbourhood of the linear minimiser,

10

contribute. In order to balance the contributions, we compensate at the level of prefactor, seeking the solution outside the Cramer zone in the form of

$$p(x, \varepsilon) \sim \varepsilon^{-1/2} \rho(x) \exp\left(-\frac{\tilde{\Phi}(x)}{\varepsilon}\right), \quad x > x_*. \tag{33}$$

The next subsection determines the prefactor $\rho(x)$ outside the Cramer zone.

## 4.3   The prefactor outside the Cramer zone

Inserting the WKB expansions (25) and (33) into the Volterra master equation (4), we find that for $\delta \gg \varepsilon^{1/2}$ we have

$$\gamma(x)\rho(x) \exp\left(-\frac{\tilde{\Phi}(x)}{\varepsilon}\right) = \varepsilon^{-1/2} \int_{y_{\mathrm{m}}(x)-\delta}^{y_{\mathrm{m}}(x)+\delta} \alpha(y) r_0(y) \exp\left(-\frac{\tilde{\Psi}(x, y)}{\varepsilon}\right) \mathrm{d}y$$

$$+ \varepsilon^{-1} \int_{x-\delta}^{x} \alpha(y)\rho(y) \exp\left(-\frac{\tilde{\Psi}(x, y)}{\varepsilon}\right) \mathrm{d}y + o\left(\exp\left(-\frac{\tilde{\Phi}(x)}{\varepsilon}\right)\right). \tag{34}$$

Estimating the integrals by the Laplace method, cancelling the common exponential term, and collecting at the leading order, we obtain

$$\gamma(x)\rho(x) = \left(\frac{2\pi}{\partial_y^2 \Psi(x, y_{\mathrm{m}}(x))}\right)^{1/2} \alpha(y_{\mathrm{m}}(x)) r_0(y_{\mathrm{m}}(x)) - \frac{\alpha(x)\rho(x)}{\partial_y \tilde{\Psi}(x, x)}. \tag{35}$$

Differentiating (26) with respect to $y$ and using (23) gives

$$\partial_y \tilde{\Psi}(x, y) = \nu(y_{\mathrm{m}}(y)) + \nu'(y)(x - y) - \nu(y), \quad y > x_*. \tag{36}$$

We set $y = x$ into (36) and insert the result into (35), arriving at

$$\rho(x) = \left(\frac{2\pi}{\partial_y^2 \Psi(x, y_{\mathrm{m}}(x))}\right)^{1/2} \frac{\alpha(y_{\mathrm{m}}(x)) r_0(y_{\mathrm{m}}(x))}{\gamma(x) - \frac{\alpha(x)}{\nu(x) - \nu(y_{\mathrm{m}}(x))}}. \tag{37}$$

Inequality (17) ensures that the denominator in (37) is positive (including at the boundary $x = x_*$).

In the next section, we tie up the loose ends in the approximation scheme by constructing an inner solution in a neighbourhood of the Cramer boundary

11

$x = x_*$ that matches (25) to the left and (33) to the right.

# 5   Boundary layer

The discontinuity in the potential derivative (24) and the mismatch of prefactor
magnitudes in (25) and (33) suggest the presence of a boundary layer near
$x = x_*$. We define the inner variable $\xi$ via the transformation

$$x = x_* + \kappa\varepsilon\ln\varepsilon + \varepsilon\xi, \tag{38}$$

where the constant $\kappa > 0$ will be specified later. Qualitatively, as $x$ increases
towards $x_*$, the integral in the Volterra equation begins to feel the "ghost" of the
internal minimum of $\Psi(x_*, y)$ (Figure 2, right panel): the local approximation
of Section 3 breaks down before $x_*$ is reached. The qualitative notion is made
quantitative in the rest of the section. Subsection 5.1 constructs the inner
solution that is valid in the boundary layer $\xi = O(1)$. Subsection 5.2 matches the
inner solution to the WKB approximations that are valid outside the boundary
layer.

## 5.1   Inner solution

The inner solution is sought to be proportional to a regular function of the inner
variable:

$$p(x_* + \kappa\varepsilon\ln\varepsilon + \varepsilon\xi; \varepsilon) \sim C(\varepsilon)f(\xi). \tag{39}$$

We divide the integration interval in (4) into $0 < y < x_\mathrm{o}$ and $x_\mathrm{o} < y < x$,
where $x_\mathrm{o}$ belongs to the overlap of the WKB approximation (25) and the inner
approximation (39).

In the first interval, the integral is estimated by means of the WKB approx-

12

imation (25) and the Laplace method as

$$
\begin{aligned}
\frac{1}{\varepsilon} & \int_0^{x_\circ} p(y)\alpha(y) \exp\left(-\frac{\nu(y)(x-y)}{\varepsilon}\right) \mathrm{d}y \\
&= \frac{1}{\varepsilon} \int_0^{x_\circ} p(y)\alpha(y) \exp\left(-\frac{\nu(y)(x_* - y)}{\varepsilon} - \nu(y)\xi\right) \varepsilon^{-\kappa\nu(y)} \mathrm{d}y \\
&\sim \frac{1}{\varepsilon} \int_0^{x_\circ} \alpha(y)r_0(y) \exp\left(-\frac{\Psi(x_*,y)}{\varepsilon} - \nu(y)\xi\right) \varepsilon^{-\kappa\nu(y)} \mathrm{d}y \\
&\sim \left(\frac{2\pi}{\partial_y^2 \Psi(x_*,y_*)}\right)^{1/2} \alpha(y_*)r_0(y_*) \exp\left(-\frac{\Phi(x_*)}{\varepsilon} - \nu(y_*)\xi\right) \varepsilon^{-\kappa\nu(y_*) - \frac{1}{2}}. \quad (40)
\end{aligned}
$$

In the second interval, the substitution $y = x_* + \kappa\varepsilon\ln\varepsilon + \varepsilon\eta$ and the inner approximation (39) give an asymptotic estimate

$$
\begin{aligned}
\frac{1}{\varepsilon} & \int_{x_\circ}^x p(y)\alpha(y) \exp\left(-\frac{\nu(y)(x-y)}{\varepsilon}\right) \mathrm{d}y \\
&\sim C(\varepsilon)\alpha(x_*) \int_{-\infty}^{\xi} f(\eta) \mathrm{e}^{-\nu(x_*)(\xi-\eta)} \mathrm{d}\eta. \quad (41)
\end{aligned}
$$

Requiring that (40) and (41) be of the same order implies

$$
C(\varepsilon) = \exp\left(-\frac{\Phi(x_*)}{\varepsilon}\right) \varepsilon^{-\kappa\nu(y_*) - \frac{1}{2}} \quad (42)
$$

for the proportionality constant in the inner solution (39).

Inserting (39), (40), and (41) into the Volterra master equation (4), and then dividing by $C(\varepsilon)$, yields

$$
\begin{aligned}
\gamma(x_*)f(\xi) = & \left(\frac{2\pi}{\partial_y^2 \Psi(x_*,y_*)}\right)^{1/2} \alpha(y_*)r_0(y_*)\mathrm{e}^{-\nu(y_*)\xi} \\
& + \alpha(x_*) \int_{-\infty}^{\xi} f(\eta)\mathrm{e}^{-\nu(x_*)(\xi-\eta)} \mathrm{d}\eta.
\end{aligned} \quad (43)
$$

Multiplying (43) by $\mathrm{e}^{\nu(x_*)\xi}$ and differentiating with respect to $\xi$ turns the inte-

gral equation (43) into a differential equation

$$
\begin{aligned}
\gamma(x_*)\frac{\mathrm{d}}{\mathrm{d}\xi}&\left(\mathrm{e}^{\nu(x_*)\xi}f(\xi)\right) = \alpha(x_*)\mathrm{e}^{\nu(x_*)\xi}f(\xi) \\
&+ \left(\frac{2\pi}{\partial_y^2\Psi(x_*,y_*)}\right)^{1/2}\alpha(y_*)r_0(y_*)(\nu(x_*)-\nu(y_*))\mathrm{e}^{(\nu(x_*)-\nu(y_*))\xi}.
\end{aligned}
\tag{44}
$$

Solving (44) yields

$$
f(\xi) = A\mathrm{e}^{-\left(\nu(x_*)-\frac{\alpha(x_*)}{\gamma(x_*)}\right)\xi} + B\mathrm{e}^{-\nu(y_*)\xi},
\tag{45}
$$

where

$$
B = \left(\frac{2\pi}{\partial_y^2\Psi(x_*,y_*)}\right)^{1/2}\frac{\alpha(y_*)r_0(y_*)}{\gamma(x_*)-\frac{\alpha(x_*)}{\nu(x_*)-\nu(y_*)}}
\tag{46}
$$

is found by the method of undetermined coefficients and $A$ is a constant of integration, which will be determined by asymptotic matching to the outer solution.

## 5.2   Matching

Two constants need to be determined to complete the inner solution, namely:

- the integration constant $A$ in (45);

- the constant $\kappa$ in the offset of the boundary layer (38).

These will be calculated in Section 5.2.2 by matching to the WKB solution (25) inside the Cramer zone. Before doing so, we demonstrate that the inner solution asymptotically matches the WKB solution (33) outside the Cramer zone.

### 5.2.1   Matching to the right

Owing to the inequality (17), the second term in the inner solution (45) dominates for $\xi \to \infty$; inserting it and (42) into (39) gives

$$
p(x;\varepsilon) \sim B\exp\left(-\frac{\Phi(x_*)}{\varepsilon}-\nu(y_*)\xi\right)\varepsilon^{-\kappa\nu(y_*)-\frac{1}{2}}
\tag{47}
$$

14

in the overlap of the inner solution and the outer solution to its right.

On the other hand, inserting the transformation (38) into the outer solution (33), re-expanding, and using (23) gives

$$
\begin{aligned}
p(x;\varepsilon) &\sim \varepsilon^{-1/2} \rho(x_*) \exp\left(-\frac{\Phi(x_*)}{\varepsilon} - \tilde{\Phi}'(x_*^+)(\kappa\ln\varepsilon + \xi)\right) \\
&= \rho(x_*) \exp\left(-\frac{\Phi(x_*)}{\varepsilon} - \nu(y_*)\xi\right) \varepsilon^{-\kappa\nu(y_*) - \frac{1}{2}}
\end{aligned}
\tag{48}
$$

in the overlap. Comparing (47) and (48), we find $B = \rho(x_*)$, which is consistent with (37) and (46).

### 5.2.2  Matching to the left

As $\xi \to -\infty$, the first term in (45) dominates, so that

$$
p(x;\varepsilon) \sim A \exp\left(-\frac{\Phi(x_*)}{\varepsilon} - \left(\nu(x_*) - \frac{\alpha(x_*)}{\gamma(x_*)}\right)\xi\right) \varepsilon^{-\kappa\nu(y_*) - \frac{1}{2}}
\tag{49}
$$

in the overlap of the inner solution and the outer solution to its left.

On the other hand, inserting (38) into the outer solution (25) gives

$$
\begin{aligned}
p(x;\varepsilon) &\sim r_0(x_*) \exp\left(-\frac{\Phi(x_*)}{\varepsilon} - \Phi(x_*^-)(\kappa\ln\varepsilon + \xi)\right) \\
&= r_0(x_*) \exp\left(-\frac{\Phi(x_*)}{\varepsilon} - \left(\nu(x_*) - \frac{\alpha(x_*)}{\gamma(x_*)}\right)\xi\right) \varepsilon^{-\kappa\left(\nu(x_*) - \frac{\alpha(x_*)}{\gamma(x_*)}\right)}.
\end{aligned}
\tag{50}
$$

Comparing (49) to (50) yields

$$
A = r_0(x_*), \quad \kappa = \frac{1}{2\left(\nu(x_*) - \frac{\alpha(x_*)}{\gamma(x_*)} - \nu(y_*)\right)};
\tag{51}
$$

inequality (19) thereby guarantees that $\kappa > 0$ as advertised at the beginning of the boundary-layer analysis. Equations (51) complete the inner solution and thus the asymptotic analysis of (4).

15

# 6    Numerical solution

Before being compared to a numerical solution, the asymptotic solutions are normalised by

$$\mathcal{N} = \int_0^{x_*} r_0(x) \exp\left(-\frac{\tilde{\Phi}(x)}{\varepsilon}\right) \mathrm{d}x + \varepsilon^{-1/2} \int_{x_*}^{\infty} \rho(x) \exp\left(-\frac{\tilde{\Phi}(x)}{\varepsilon}\right) \mathrm{d}x. \quad (52)$$

The integral of the WKB solution over the tail zone is exponentially smaller than the integral over the Cramer zone and can be neglected in (52). The Cramer-zone integral can in principle be estimated by the Laplace method by the local contribution from the minimiser of the potential $\Phi(x)$. However, practice shows that doing so introduces a relatively large numerical error. Instead, the normalisation constant can be calculated by numerical quadrature of (52).

For the numerical solution, sample paths $x_i(t)$, $i = 1, \ldots, N$, $0 \le t \le T$, subject to $x(0) = x_0$ are generated using the exact stochastic simulation algorithm (see the Appendix). The solution is constructed by the histogram method from the dataset of final-time values $\{x_i(T)\}_{i=1,\ldots,N}$. Specifically, we divide an interval $[0, x_{\max}]$ into $n$ equally sized bins, count the number of data in each bin, and divide the counts by $N x_{\max}/n$ so as to normalise into a probability density. The histogram estimate is close to the exact solution $p(x; \varepsilon)$ to the Volterra master equation (4) if the number of samples $N$ is large (so that the statistical error is small) and the simulation end time $T$ is large (so that the process equilibrates to steady state).

Figure 3 compares the three matched asymptotic approximations to the numerical solution for selected values of the noise strength $\varepsilon$. Decreasing $\varepsilon$ leads to a close agreement between the numerical solution and the asymptotic approximations in their respective regions of validity (Figure 3, top panels). As $\varepsilon$ decreases further (Figure 3, bottom panels), the Cramer-boundary and tail behaviour become exponentially improbable, and cannot be reliably estimated from a feasible number (say a billion) of samples. Nevertheless, the chosen examples demonstrate that the naive solution, which extends (25) outside the
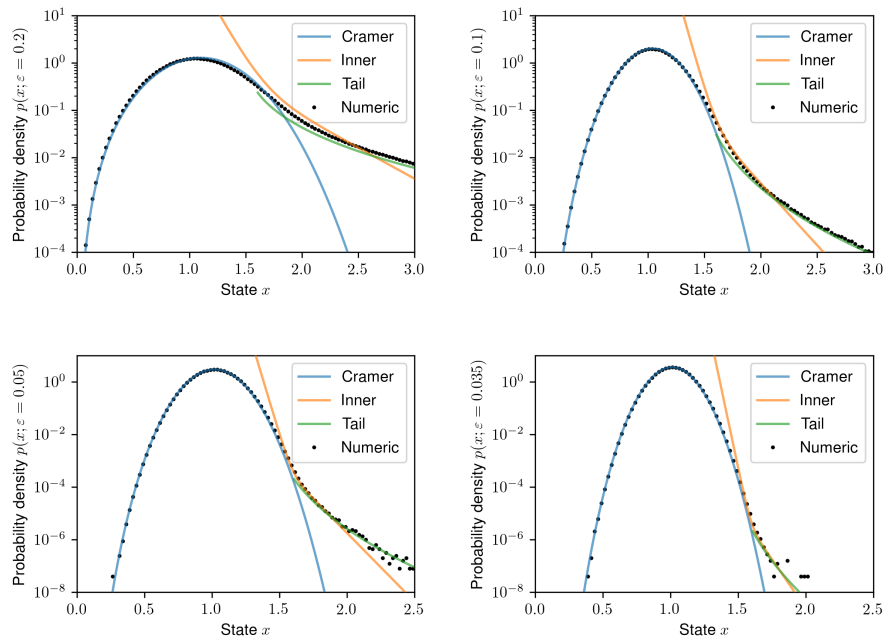
Figure 3: The simulation-based probability density (dots) is compared to the individual asymptotic approximations (solid lines), namely the WKB solution in the Cramer zone (25), the inner solution in the boundary layer (45), and the WKB solution in the tail zone (33). *Model parameters:* we use (16) with $m = 2$; values of $\varepsilon$ are specified in the label of the ordinate. *Numerical parameters:* $x_0 = 1$, $T = 30$, $N = 10^8$ (upper panels), $N = 10^9$ (lower panels), $n = 100$, $x_{\max} = 3$ (upper panels) and $x_{\max} = 2.5$ (lower panels).

Cramer zone, underestimates the tail of the stationary distribution, whereas the alternative approximations provide an adequate description.

# 7 Conclusion

This paper provides matched asymptotic approximations to the stationary distribution of a drift–jump model for stochastic gene expression. The analysis revolves around the estimation of the integral term in the Volterra master equation (4). The integral term represents the flux of probability due to production bursts through a reference state $x$. In the Cramer region ($x < x_*$), the flux consists solely from local contributions ($y \approx x$), whereas in the tail region ($x > x_*$), a contribution comes also from within the interval. The latter corresponds to the 'single big jumps' advertised in the abstract.

Negative feedback in burst size is a prerequisite for the singular behaviour in question. Conceptually, in the presence of negative feedback in burst size, it is 'cheaper' to hunker down and then take a giant leap, than to climb up with tiny steps. The result is thus in agreement with the broad principle that any large deviation occurs in the least unlikely of all the unlikely ways [51].

The analysis is formulated for general feedback responses satisfying certain constraints. A particular specimen, the power non-linearity $\nu(x) = x^m$, has been the main example throughout this text. The coefficient $m$ can be interpreted as the number of protein molecules that need to cooperate to repress the production burst. The solution to the Volterra equation (4) with a power non-linearity has previously been shown to satisfy $p(x) \sim c_1 x^{\frac{1}{\varepsilon} - 1}$ as $x \to 0$ and $p(x) \sim c_2 x^{-\frac{1}{\varepsilon m} - 1}$ as $x \to \infty$, where $c_1, c_2 > 0$ [52]. The same study provided a central-limit-theorem-type approximation that is valid as $\varepsilon \to 0$ for $|x - 1| = O(\varepsilon^{1/2})$. The current study thus contributes by approximations that apply as $\varepsilon \to 0$ throughout the state space $x > 0$. The popular Hill-type non-linearity $\nu(x) = 1/(1 + (x/K)^m)$ can be reduced to the power non-linearity by means of a simple transformation [52]. The conclusions arrived for the power non-linearity thus

easily extend to the Hill-type response.

Earlier studies argued that the subtleties that arise with feedback in burst size are an artefact of delay [32, 33]. Indeed, the memoryless property of the exponential distribution of burst sizes implies a lack of control at the infinitesimal timescale of burst growth. In light of this argument, the current results contribute to the understanding of the interplay between bursting and delay in biological systems [53–57].

# Appendix: Stochastic simulation algorithm

Here we provide an stochastic simulation algorithm that can be used to generate a sample path $x(t)$ of the process on a time interval $[0, T]$ subject to an initial condition $x(0) = x_0$. Similarly like the well-known Gibson–Bruck/Gillespie algorithm, the algorithm does not introduce truncation errors, but only statistical and round-off errors, and in this specific sense it is an exact simulation algorithm. For simplicity, we focus on the situation when the feedback acts only on burst size but not on burst frequency or protein stability; the general case is discussed in the end of the appendix.

Each sample path is generated iteratively as follows. Assume that the sample path $x(t)$ has already been generated on an interval $0 \leq t \leq t_{\mathrm{cur}}$ (initially $t_{\mathrm{cur}} = 0$ and $x(0) = x_0$ is an initial value). Assuming the absence of feedback in burst frequency ($\alpha(x) = 1$), the exponentially distributed waiting time until the coming burst is sampled by the inversion method as

$$\tau = -\varepsilon \ln \theta, \tag{53}$$

where $\theta$ is drawn from the uniform distribution in the unit interval. Assuming the absence of feedback in protein stability ($\gamma(x) = x$), the sample path decays exponentially until the coming burst:

$$x(t) = x(t_{\mathrm{cur}}) e^{-(t - t_{\mathrm{cur}})} \quad \text{for} \quad t_{\mathrm{cur}} < t < t_{\mathrm{cur}} + \tau. \tag{54}$$

19

At the time of the next burst the sample path is increased by the exponentially distributed burst size:

$$x(t) = x(t^-) - \frac{\varepsilon \ln \tilde{\theta}}{\nu(x(t^-))} \quad \text{for} \quad t = t_{\text{cur}} + \tau, \tag{55}$$

where $x(t^-) = x(t_{\text{cur}})\mathrm{e}^{-\tau}$ denotes the state of the sample path immediately before the burst; the variate $\tilde{\theta}$ is drawn from the uniform distribution in the unit interval independently of $\theta$. Thus one round of iteration via (53), (54), and (55) extends the sample path from the interval $[0, t_{\text{cur}}]$ to the interval $[0, t_{\text{cur}} + \tau]$. The algorithm is repeated until the state $x(T)$ at a required end time $T > 0$ is found.

The algorithm can be modified to account for feedback in burst frequency and protein stability. If feedback in burst frequency is present, the waiting time needs to be drawn from a distribution with a non-constant hazard function [8]. If feedback in protein stability is present, the sample path needs to be evolved as per $\dot{x} = -\gamma(x)$ between bursts.

# References

[1] R. D. Dar, B. S. Razooky, A. Singh, T. V. Trimeloni, J. M. McCollum, C. D. Cox, M. L. Simpson, and L. S. Weinberger, "Transcriptional burst frequency and burst size are equally modulated across the human genome," *P. Natl. Acad. Sci. USA*, vol. 109, pp. 17454–17459, 2012.

[2] Y. Wang, T. Ni, W. Wang, and F. Liu, "Gene transcription in bursting: a unified mode for realizing accuracy and stochasticity," *Biol. Rev.*, vol. 94, no. 1, pp. 248–258, 2019.

[3] J. Rodriguez and D. R. Larson, "Transcription in living cells: Molecular mechanisms of bursting," *Annu. Rev. Biochem.*, vol. 89, 2020.

[4] L. Schuh, M. Saint-Antoine, E. M. Sanford, B. L. Emert, A. Singh, C. Marr, A. Raj, and Y. Goyal, "Gene networks with transcriptional bursting reca-

20

pitulate rare transient coordinated high expression states in cancer," *Cell Syst.*, vol. 10, no. 4, pp. 363–378, 2020.

[5] J. Paulsson and M. Ehrenberg, "Random signal fluctuations can reduce random fluctuations in regulated components of chemical regulatory networks," *Phys. Rev. Lett.*, vol. 84, no. 23, pp. 5447–50, 2000.

[6] N. Friedman, L. Cai, and X. S. Xie, "Stochasticity in gene expression as observed by single-molecule experiments in live cells," *Israel J. Chem.*, vol. 49, no. 3-4, pp. 333–342, 2009.

[7] N. Friedman, L. Cai, and X. Xie, "Linking stochastic dynamics to population distribution: an analytical framework of gene expression," *Phys. Rev. Lett.*, vol. 97, p. 168302, 2006.

[8] P. Bokes, J. King, A. Wood, and M. Loose, "Transcriptional bursting diversifies the behaviour of a toggle switch: hybrid simulation of stochastic gene expression," *B. Math. Biol.*, vol. 75, pp. 351–371, 2013.

[9] M. Pájaro, I. Otero-Muras, C. Vázquez, and A. A. Alonso, "Transient hysteresis and inherent stochasticity in gene regulatory networks," *Nat. Commun.*, vol. 10, no. 1, pp. 1–7, 2019.

[10] J. Jedrak, M. Kwiatkowski, and A. Ochab-Marcinek, "Exactly solvable model of gene expression in a proliferating bacterial cell population with stochastic protein bursts and protein partitioning," *Phys. Rev. E*, vol. 99, no. 4, p. 042416, 2019.

[11] J. Holehouse, Z. Cao, and R. Grima, "Stochastic modeling of auto-regulatory genetic feedback loops: a review and comparative study," *Biophys. J.*, 2020.

[12] G. Giovanini, A. U. Sabino, L. R. Barros, A. F. Ramos, *et al.*, "A comparative analysis of noise properties of stochastic binary models for a self-

repressing and for an externally regulating gene," *Math. Biosci. Eng.*, vol. 17, no. 5, pp. 5477–5503, 2020.

[13] F. Veerman, N. Popović, and C. Marr, "Parameter inference with analytical propagators for stochastic models of autoregulated gene expression," *Int. J. Nonlinear Sci.*, 2021.

[14] M. K. Tonn, P. Thomas, M. Barahona, and D. A. Oyarzún, "Computation of single-cell metabolite distributions using mixture models," *Front. Cell Dev. Biol.*, vol. 8, p. 1596, 2020.

[15] J. Dattani and M. Barahona, "Stochastic models of gene transcription with upstream drives: exact solution and sample path characterization," *J. Roy. Soc. Interface*, vol. 14, p. 20160833, 2017.

[16] A. Kozdeba and A. Tomski, "Application of the goodwin model to autoregulatory feedback for stochastic gene expression," *Math. Biosci.*, vol. 327, p. 108413, 2020.

[17] P. Kurasov, D. Mugnolo, and V. Wolf, "Analytic solutions for stochastic hybrid models of gene regulatory networks," *Journal of Mathematical Biology*, vol. 82, no. 1, pp. 1–29, 2021.

[18] A. Crudu, A. Debussche, A. Muller, O. Radulescu, *et al.*, "Convergence of stochastic gene networks to hybrid piecewise deterministic processes," *Ann. Appl. Probab.*, vol. 22, pp. 1822–1859, 2012.

[19] P. Bokes, J. King, A. Wood, and M. Loose, "Multiscale stochastic modelling of gene expression," *J. Math. Biol.*, vol. 65, pp. 493–520, 2012.

[20] Y. T. Lin and C. R. Doering, "Gene expression dynamics with stochastic bursts: Construction and exact results for a coarse-grained model," *Phys. Rev. E*, vol. 93, p. 022409, 2016.

[21] Y. T. Lin and T. Galla, "Bursting noise in gene expression dynamics: linking microscopic and mesoscopic models," *J. Roy. Soc. Interface*, vol. 13, p. 20150772, 2016.

[22] C. Jia, M. Q. Zhang, and H. Qian, "Emergent lévy behavior in single-cell stochastic gene expression," *Phys. Rev. E*, vol. 96, no. 4, p. 040402, 2017.

[23] C. Jia, G. G. Yin, M. Q. Zhang, *et al.*, "Single-cell stochastic gene expression kinetics with coupled positive-plus-negative feedback," *Phys. Rev. E*, vol. 100, no. 5, p. 052406, 2019.

[24] X. Chen and C. Jia, "Limit theorems for generalized density-dependent markov chains and bursty stochastic gene regulatory networks," *J. Math. Biol.*, vol. 80, no. 4, pp. 959–994, 2020.

[25] L. Cai, N. Friedman, and X. Xie, "Stochastic protein expression in individual cells at the single molecule level," *Nature*, vol. 440, pp. 358–362, 2006.

[26] J. Jedrak and A. Ochab-Marcinek, "Influence of gene copy number on self-regulated gene expression," *J. Theor. Biol.*, vol. 408, pp. 222–236, 2016.

[27] L. Bintu, N. Buchler, H. Garcia, U. Gerland, T. Hwa, J. Kondev, and R. Phillips, "Transcriptional regulation by the numbers: models," *Curr. Opin. Genet. Dev.*, vol. 15, pp. 116–124, 2005.

[28] M. A. Hernandez, B. Patel, F. Hey, S. Giblett, H. Davis, and C. Pritchard, "Regulation of BRAF protein stability by a negative feedback loop involving the MEK–ERK pathway but not the FBXW7 tumour suppressor," *Cell. Signal.*, vol. 28, pp. 561–571, 2016.

[29] A. Sundqvist and J. Ericsson, "Transcription-dependent degradation controls the stability of the srebp family of transcription factors," *P. Natl. Acad. Sci. USA.*, vol. 100, pp. 13833–13838, 2003.

23

[30] M. A. Schikora-Tamarit, C. Toscano-Ochoa, J. D. Espinos, L. Espinar, and L. B. Carey, "A synthetic gene circuit for measuring autoregulatory feedback control," *Integr. Biol.*, vol. 8, pp. 546–555, 2016.

[31] P. Bokes and A. Singh, "Controlling noisy expression through auto regulation of burst frequency and protein stability," in *Češka M., Paoletti N. (eds) Hybrid Systems Biology. HSB 2019. Lecture Notes in Computer Science, vol 11705*, Springer, Cham, 2019.

[32] P. Bokes, "Maintaining gene expression levels by positive feedback in burst size in the presence of infinitesimal delay," *Discrete Cont. Dyn-B*, vol. 24, no. 10, p. 5539, 2019.

[33] P. Bokes, "Exact and WKB-approximate distributions in a gene expression model with feedback in burst frequency, burst size, and protein stability," *Discrete Cont. Dyn-B; doi: 10.3934/dcdsb.2021126*, 2021.

[34] S. Be'er and M. Assaf, "Rare events in stochastic populations under bursty reproduction," *J. Stat. Mech. Theory E.*, vol. 2016, p. 113501, 2016.

[35] M. Assaf and B. Meerson, "WKB theory of large deviations in stochastic populations," *J. Phys. A: Math. Theor.*, vol. 50, no. 26, p. 263001, 2017.

[36] J. Hertz, J. Tyrcha, and A. Correales, "Stochastic activation in a genetic switch model," *Phys. Rev. E*, vol. 98, no. 5, p. 052403, 2018.

[37] O. Vilk and M. Assaf, "Population extinction under bursty reproduction in a time-modulated environment," *Phys. Rev. E*, vol. 97, no. 6, p. 062114, 2018.

[38] P. Bokes, A. Borri, P. Palumbo, and A. Singh, "Mixture distributions in a stochastic gene expression model with delayed feedback: a WKB approximation approach," *J. Math. Biol.*, vol. 81, no. 1, pp. 343–367, 2020.

[39] C. Knessl, B. Matkowsky, Z. Schuss, and C. Tier, "Asymptotic analysis of

24

a state-dependent M/G/1 queueing system," *SIAM J. Appl. Math.*, vol. 46, no. 3, pp. 483–505, 1986.

[40] Z. Schuss, *Theory and applications of stochastic processes: an analytical approach.* Springer Science & Business Media, Berlin/Heidelberg, 2009.

[41] M. I. Freidlin and A. D. Wentzell, *Random perturbations of Dynamical Systems.* Springer, Heidelberg, 2012.

[42] A. Vezzani, E. Barkai, and R. Burioni, "Single-big-jump principle in physical modeling," *Phys. Rev. E*, vol. 100, no. 1, p. 012108, 2019.

[43] A. A. Borovkov and K. A. Borovkov, *Asymptotic analysis of random walks*, vol. 118. Cambridge University Press, 2008.

[44] A. A. Borovkov, *Probability Theory.* Springer, Heidelberg, 2013.

[45] R. Hinch and S. J. Chapman, "Exponentially slow transitions on a Markov chain: the frequency of calcium sparks," *Eur. J. Appl. Math.*, vol. 16, no. 04, pp. 427–446, 2005.

[46] J. Newby, "Bistable switching asymptotics for the self regulating gene," *J. Phys. A-math. Gen.*, vol. 48, p. 185001, 2015.

[47] P. C. Bressloff, *Stochastic processes in cell biology.* Springer, Heidelberg, 2014.

[48] A. H. Nayfeh, *Introduction to perturbation techniques.* John Wiley & Sons, New Jersey, 2011.

[49] N. van Kampen, *Stochastic Processes in Physics and Chemistry.* Elsevier, Amsterdam, 2006.

[50] J. Newby and S. J. Chapman, "Metastable behavior in Markov processes with internal states," *J. Math. Biol.*, vol. 69, no. 4, pp. 941–976, 2014.

[51] F. Den Hollander, *Large deviations*, vol. 14. American Mathematical Soc., 2008.

[52] P. Bokes, Y. Lin, and A. Singh, "High cooperativity in negative feedback can amplify noisy gene expression," *B. Math. Biol.*, vol. 80, pp. 1871–1899, 2018.

[53] J. M. Newby, "Spontaneous excitability in the Morris–Lecar model with ion channel noise," *SIAM J. Appl. Dyn. Syst.*, vol. 13, no. 4, pp. 1756–1791, 2014.

[54] E. Zavala and T. T. Marquez-Lago, "Delays induce novel stochastic effects in negative feedback gene circuits," *Biophys. J.*, vol. 106, no. 2, pp. 467–478, 2014.

[55] R. Martinez-Corral, E. Raimundez, Y. Lin, M. B. Elowitz, and J. Garcia-Ojalvo, "Self-amplifying pulsatile protein dynamics without positive feedback," *Cell Syst.*, vol. 7, no. 4, pp. 453–462, 2018.

[56] A. S. Sassi, M. Garcia-Alcala, M. J. Kim, P. Cluzel, and Y. Tu, "Filtering input fluctuations in intensity and in time underlies stochastic transcriptional pulses without feedback," *Proceedings of the National Academy of Sciences*, vol. 117, no. 43, pp. 26608–26615, 2020.

[57] J. Negrete, I. M. Lengyel, L. Rohde, R. A. Desai, A. C. Oates, and F. Jülicher, "Theory of time delayed genetic oscillations with external noisy regulation," *New J. Phys.*, vol. 23, no. 3, p. 033030, 2021.