



Inferring the spatiotemporal DNA replication program from noisy data

A. Baker and J. Bechhoefer

Department of Physics, Simon Fraser University, Burnaby, British Columbia, Canada V5A 1S6

(Received 13 December 2013; published 6 March 2014)

We generalize a stochastic model of DNA replication to the case where replication-origin-initiation rates vary locally along the genome and with time. Using this generalized model, we address the inverse problem of inferring initiation rates from experimental data concerning replication in cell populations. Previous work based on curve fitting depended on arbitrarily chosen functional forms for the initiation rate, with free parameters that were constrained by the data. We introduce a nonparametric method of inference that is based on Gaussian process regression. The method replaces specific assumptions about the functional form of the initiation rate with more general prior expectations about the smoothness of variation of this rate, along the genome and in time. Using this inference method, we recover, with high precision, simulated replication schemes from noisy data that are typical of current experiments.

DOI: [10.1103/PhysRevE.89.032703](https://doi.org/10.1103/PhysRevE.89.032703)

PACS number(s): 87.10.-e, 87.14.gk, 87.18.Vf, 82.60.Nh

I. INTRODUCTION

Cells must accurately duplicate their DNA content at every cell cycle. Depending on the organism, DNA replication can initiate at one or at multiple sites called *origins of replication*. The DNA is copied by a pair of oppositely moving *replication forks* that propagate away from each origin. These forks actively copy the genome away from the origin until they encounter another replication fork. DNA replication can thus be modeled as a process of initiation, growth, and coalescences occurring in an asynchronous, parallel way until the whole genome is copied. In this process, initiation has been observed to be a stochastic process [1–6], while fork propagation, at the large scales (10–100 kilobase) between origins, is largely deterministic, and often constant [7]. Fork stalls from DNA damage and other causes can alter the replication program [8], but we do not consider such effects here.

The elements of stochastic initiation, deterministic growth, and coalescence are formally equivalent to the processes of nucleation, growth, and coalescence in crystallization kinetics, and this equivalence has inspired efforts to model DNA replication kinetics using the formalism developed in the 1930s by Kolmogorov, Johnson, Mehl, and Avrami (KJMA) for crystallization kinetics [9]. Of course, DNA replication takes place in a space that is topologically one dimensional, a fact that allows one to take advantage of exact solutions to the KJMA equations in one dimension [10].

The rate of initiation of origins is typically highly variable, both in space, along the genome, and in time, throughout *S phase*, the part of the cell cycle in which the genome is duplicated. In many cases, we can describe the initiation process by a rate $I(x, t)$, where $I(x, t) dx dt$ gives the probability of initiation to occur in $(x, x + dx)$ at $(t, t + dt)$ given that x is unreplicated up until time t . Loosely, we will say that $I(x, t)$ is the probability for an origin to initiate, or “fire”, at (x, t) .

In addition to its intrinsic theoretical interest, describing replication stochastically can help biologists understand better the biological dynamics underlying replication. As we discuss below, experiments have recently begun to deliver large amounts of data concerning cell populations undergoing replication. For example, it is now possible to measure the fraction of cells $f(x, t)$ that have replicated the locus x along

the genome by a time t after the beginning of *S phase* [11]. In contrast to the case of crystallization kinetics, there is little fundamental understanding of the structure of the initiation function $I(x, t)$. Since direct observation of initiations *in vivo* has not been possible, the task is to estimate, or infer, $I(x, t)$ from data such as the replication fraction $f(x, t)$ or—more conveniently, it will turn out—the unreplicated fraction $s(x, t) = 1 - f(x, t)$, which is also the probability that the locus x is unreplicated at time t .

In this paper, we have two goals: The first, presented in Secs. II and III, is to collect and generalize previous results on the application of the KJMA formalism to DNA replication. Previous work has focused on special cases: models of replication in *Xenopus laevis* (frog) embryos were based on experiments that averaged data from the whole genome [12] and thus could neglect spatial variations. Conversely, in recent experiments on a small section of a mouse genome, spatial variations dominated and temporal variations could be neglected. In budding yeast, origins are restricted to specific sites along the genome [13], which also leads to a restricted form of the initiation function. In general, however, both spatial and temporal variations are important, and we extend here the full KJMA formalism to handle such cases. Section IV gives a brief example that illustrates the kinds of results and insights that this approach to modeling replication can provide.

The second goal is to present a different way to infer initiation rates $I(x, t)$ from replication data such as $s(x, t)$. Replication timing data are increasingly available for a variety of organisms and cell types [11, 14–17], and advances in experimental techniques now allow the determination of the probability distribution of genome-wide replication timing at fine spatial and temporal scales. For instance, in yeast, the unreplicated fraction profiles have been determined at 1 kb resolution in space and 5 min resolution in time [18]. The increasing availability of data makes the ability to infer initiation rates important.

Our main result, presented in Sec. V, is to adapt the technique of *Gaussian process regression* to “invert” experimental replication data and estimate the initiation function $I(x, t)$ and fork velocity v . Previous approaches have mainly used curve fitting, a technique that postulates a suitable functional form

for $I(x,t)$, with free parameters that are then constrained by fitting to the data. This technique was used to infer initiation functions in frog embryos [12], budding yeast [18–21], and limited regions of human somatic cells [22].

Although the above efforts were successful, curve-fitting methods are time consuming, requiring considerable effort to generate initial guesses that are close enough to the final inference. The situation is even more difficult if one wants to describe replication over the whole genome of higher eukaryotes. In these organisms, initiations are not limited to well-positioned replication origins but also occur in large extended initiation zones whose functional form is not known *a priori*. Furthermore, the mapping of well-positioned replication origins and extended initiation zones along the genome is difficult [23], and not much is known about the firing-time distributions. These added uncertainties make curve-fitting approaches to local genomic data in higher eukaryotes problematic.

Given the difficulty of extending and automating curve-fit approaches, we explore here an alternative that does not depend on knowing *a priori* the functional form of the initiation function. The technique, Gaussian process regression, is based on the Bayesian approach to data analysis and gives a systematic way to infer the initiation rate without making detailed assumptions about its functional form in the way required of curve-fit methods. Although Gaussian process regression is more powerful than curve-fitting methods, it can be simpler to apply. Because no detailed tuning of initiation conditions is required, the method can in principle be automated. In contrast, curve-fitting methods require a good technical understanding to use successfully.

II. GENERAL REPLICATION PROGRAM

We begin by establishing relationships that must be obeyed by any spatiotemporal replication program with a constant fork velocity. (Scenarios with variable or stochastic fork velocities are straightforward but beyond the scope of this paper.) We then show that many quantities of interest, such as the densities of right- and left-moving forks or the initiation and termination densities, are related to derivatives of the unreplicated fraction profiles. Then we describe briefly how to use these relationships to characterize the replication program.

A. DNA replication kinetics quantities

If the replication fork velocity v is constant, the replication program in one cell cycle is completely specified by the genomic positions and firing times of the replication origins. From each origin, two divergent forks propagate at constant velocity until they meet and coalesce with a fork of the opposite direction at a replication terminus [Fig. 1(a)]. The spatial and temporal coordinates of replication termini, as well as the propagation lines of the replication forks and the replication timing (the time at which a locus is replicated), can all be derived from the genomic positions and firing times of the replication origins. Note that the inherent stochasticity of the replication program implies that the number of activated origins, along with their positions and firing times, change from one cell cycle to another, as depicted in Fig. 1(b).

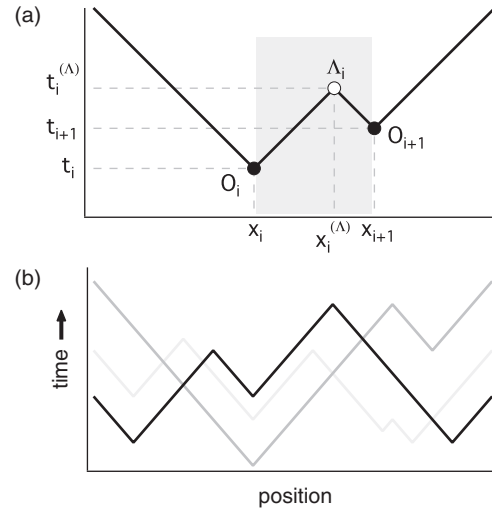


FIG. 1. Spatiotemporal representation of the replication program. (a) Replication program in one cell cycle. From each replication origin O_i (filled disk), two replication forks propagate at constant velocity v until they meet a fork of the opposite direction at a replication terminus Λ_i (hollow disk). The *replication timing curve*—the time at which a locus is replicated—is given by the intersecting set of propagation lines of the replication forks (dark zigzag line). The shaded area shows the domain of terminus Λ_i . (b) Replication program in several cell cycles. The number of activated origins, their genomic positions, and firing times change from one cell cycle to another.

Consequently, the number of terminations and initiations, the number of forks, and the replication timing curve all change from one cell cycle to another.

Let us define several quantities describing a stochastic DNA replication program. The *initiation and termination densities* $\rho_{\text{init}}(x,t)$ and $\rho_{\text{ter}}(x,t)$ give the (ensemble average) number of initiation and termination events observed in any given spatiotemporal region. The corresponding spatial densities are given by $\rho_{\text{init}}(x) = \int_0^{t_{\infty}} dt \rho_{\text{init}}(x,t)$ and $\rho_{\text{ter}}(x) = \int_0^{t_{\infty}} dt \rho_{\text{ter}}(x,t)$. Note that although the integration formally is to $t = \infty$, the end of replication for a finite genome of length L will at a finite (but stochastic) time t_{end} [24]. Often, $\rho_{\text{init}}(x)$ is called the *efficiency* of the locus x , as it equals the fraction of cells where locus x has initiated.

In this paper, we use the compact notation (\pm) to distinguish right-moving forks (velocity $+v$) from left-moving forks (velocity $-v$). The *fork densities* $\rho_{\pm}(x,t)$ give the spatial densities of (\pm) forks at a given time t . In other words, the (ensemble average) number of (\pm) forks in a genomic region $[x_1, x_2]$ at time t is given by $\int_{x_1}^{x_2} dx \rho_{\pm}(x,t)$. Also, as forks propagate at velocity $\pm v$, the number of (\pm) forks crossing the locus x during $[t_1, t_2]$ is given by $\int_{t_1}^{t_2} v dt \rho_{\pm}(x,t)$. Consequently, the proportions of cell cycles where the locus x is replicated by a (\pm) fork is given by $p_{\pm}(x) = \int_0^{t_{\infty}} v dt \rho_{\pm}(x,t)$. The *replication fork polarity* $p(x) = p_+(x) - p_-(x)$ measures the average directionality of the fork replicating the locus x .

Replication timing—the time when a locus is replicated—changes from one cell cycle to another. The variations can be intrinsic, due to stochastic initiation in an individual cell, and extrinsic, due to a population of cells. These variations lead

to a probability distribution $P(x,t)$ for the replication timing at locus x . The closely related *unreplicated fraction* $s(x,t)$ is defined to be the fraction of cells where x is unreplicated at time t . Since $s(x,t)$ equals the probability that replication at x occurs after t , we see that $P(x,t) = -\partial_t s(x,t)$. The ensemble average of the replication timing, or *mean replication timing*, is then¹

$$T(x) = \int_0^{t_\infty} dt P(x,t)t = \int_0^{t_\infty} dt s(x,t). \quad (1)$$

B. Derivatives of the unreplicated fraction profiles

We can establish a number of relations among the quantities defined in Sec. II A. In particular,

$$v[\rho_+(x,t) + \rho_-(x,t)] = -\partial_t s(x,t), \quad (2a)$$

$$\rho_+(x,t) - \rho_-(x,t) = \partial_x s(x,t), \quad (2b)$$

$$\rho_\pm(x,t) = -\frac{1}{2} \left(\frac{1}{v} \partial_t \mp \partial_x \right) s(x,t), \quad (2c)$$

$$\rho_{\text{init}}(x,t) - \rho_{\text{ter}}(x,t) = -\frac{1}{2} v \square s(x,t), \quad (2d)$$

$$p(x) = vT'(x), \quad (2e)$$

$$\rho_{\text{init}}(x) - \rho_{\text{ter}}(x) = \frac{1}{2} vT''(x) = \frac{1}{2} p'(x), \quad (2f)$$

where $\square = \frac{1}{v^2} \partial_t^2 - \partial_x^2$ is the d'Alembertian operator. See the Appendix for a proof of these relations.

From Eq. (2c), the densities of right- and left-moving forks are directly given by derivatives of the unreplicated fraction. The sum of the fork densities in Eq. (2a) is related to $P(x,t) = -\partial_t s(x,t)$, the probability distribution of replication timing at locus x . Equations (2e) and (2f), previously derived in [25,26], and, in special cases, in [27,28], show that the shape of the mean replication timing curve $T(x)$ gives direct information about the fork polarity and the relative densities of initiation and termination in a region. For instance, the replication fork polarity profile $p(x)$ was estimated in the human genome using Eq. (2e) and shown to be the key determinant of the compositional and mutational strand asymmetries generated by the replication process [25,29,30].

Contrary to intuition [11], the above equations show that there need not be a direct correspondence between well-positioned replication origins and timing-curve minima [27]. Around a fixed, isolated origin i located at position x_i , the initiation density profile reduces to a Dirac delta function, $\rho_{\text{init}} = E_i \delta(x - x_i)$, where the height E_i is the observed efficiency of origin i (the fraction of cells where origin i has initiated). Equation (2f) shows that the isolated origin i produces a jump discontinuity of height $2E_i$ in the fork polarity profile. Equation (2e) shows that at a minimum in $T(x)$, the fork polarity $p(x)$ must change sign. Mathematically, the efficiency E_i of the origin may or may not be large enough to produce a sign shift in $p(x)$ corresponding to a minimum of

the $T(x)$ curve. More intuitively, a weak origin (one that rarely fires in a cell cycle) in a region that is almost always replicated by a nearby strong origin may not affect the timing curve enough to produce a local minimum. As a result, even fixed, isolated origins do not necessarily imply minima in the mean replicating time curve [27]. Indeed, in budding yeast, about one origin in three is not associated with a local minimum of the timing curve [21].

III. INDEPENDENT ORIGIN FIRING

The results of Sec. II B are valid for any initiation rule. If, also, origins fire independently, then the whole spatiotemporal replication program is analytically solvable. ‘‘Independence’’ here means that an initiation event neither impedes nor favors origin initiation at another locus and implies that we can define a local initiation rate of unreplicated DNA, $I(x,t)$. The local initiation rate then completely specifies the stochastic replication program. Most models of the replication program proposed so far [19–21,27,31,32] assume the independent firing of replication origins and are thus special cases of the general formalism presented here. (An exception is [33].) The replication program is then formally analogous to a one-dimensional nucleation-and-growth process with time- and space-dependent nucleation (initiation) rate. In the 1930s, the kinetics of nucleation-and-growth processes were analytically derived for crystallization by Kolmogorov, Johnson, Mehl, and Avrami in the *KJMA theory* of phase transition kinetics [9]. Here, we will prove that the quantities describing DNA replication—the unreplicated fraction profiles and the probability distribution of the replication timing curve, the density of initiation and termination and of forks—can all be analytically derived from the local initiation rate.

The KJMA formulation of the replication program is an exactly solvable model, as all higher-moment correlation functions can also be analytically derived, for example, the joint probability distribution of replication timing at different loci, or the joint densities of initiations at different loci. We will show that, even when origins fire independently, the propagation of forks creates correlations in nearby replication times and in nearby initiation events.

Many of these relationships were previously derived for the special case of well-positioned replication origins [21,28]. The present formalism is more general, as it can include extended initiation zones, and offers a more compact and elegant derivation of these relationships.

A. Unreplicated fraction

We first note that the locus x is unreplicated at time t if and only if (iff) no initiations occur in the past ‘‘cone’’ $V_{(x,t)}[v]$ of (x,t) [gray area in Fig. 2(a)] defined by

$$V_{(x,t)}[v] = \{(x',t') : |x - x'| \leq v(t - t')\}. \quad (3)$$

When the context is unambiguous, we will use the more compact notation $X = (x,t)$ and $V_X \equiv V_{(x,t)}[v]$. The unreplicated fraction then equals the probability that no initiations occur in V_X (Kolmogorov’s argument [9]). As initiations occur independently with an initiation rate $I(x,t)$, this probability is given by a Poisson distribution with time- and space-dependent

¹Our definition differs from that of [27,28] in that we neglect the very small probability that no initiations occur on a chromosome. Replication is then not well defined.

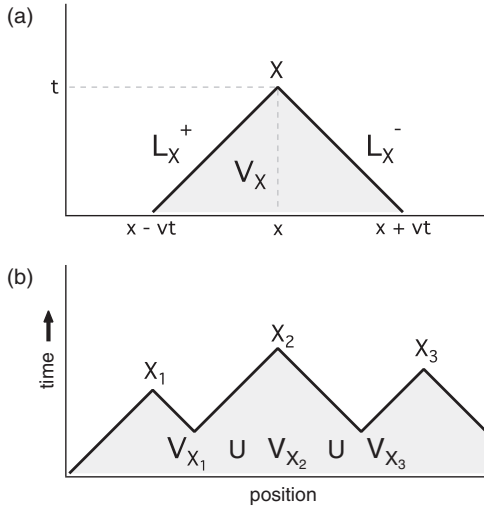


FIG. 2. Kolmogorov's argument. (a) A locus x is unreplicated at time t iff no initiation occurs in the past cone V_X of $X = (x, t)$, the gray region demarcated by the lines L_X^\pm . (b) The loci x_1, x_2, x_3 are all unreplicated at times t_1, t_2, t_3 iff no initiation occurs in $V_{X_1} \cup V_{X_2} \cup V_{X_3}$ (gray region).

rate [34]. Thus, the unreplicated fraction is given by [35]

$$s(x, t) = e^{-\int_{V_X} dt' dx' I(x', t')}. \quad (4)$$

B. Replication timing and fork densities

We can extend Kolmogorov's argument to find the fork densities. From Eqs. (2c) and (4), we find

$$\rho_\pm(x, t) = \left[\int_{L_X^\pm} I \right] s(x, t), \quad (5)$$

where the integrals of I over the lines L_X^+ and L_X^- in Fig. 2(a) are defined as

$$\int_{L_X^\pm} I \equiv \int_0^t dt' I[x \mp v(t - t'), t']. \quad (6)$$

The interpretation of Eq. (5) is straightforward: a (\pm) fork passes by x at time t iff no initiation occurs in V_X and one initiation occurs along L_X^\pm .

Similarly, from Eq. (2a),

$$\begin{aligned} P(x, t) &= v[\rho_+(x, t) + \rho_-(x, t)] = -\partial_t s(x, t), \\ &= v \left[\int_{L_X^+} I + \int_{L_X^-} I \right] s(x, t). \end{aligned} \quad (7)$$

In words: to have replication at $X = (x, t)$, no initiation occurs in V_X and an initiation along either the line L_X^+ or the line L_X^- causes a fork of velocity v to sweep by.

C. Initiation and termination densities

The initiation rate $I(x, t)$ gives the number of initiations at an unreplicated site. The initiation density $\rho_{\text{init}}(x, t)$ is then determined by the rate of initiation at (x, t) times the probability that no initiations occurred previously in the triangular area V_X

defined in Fig. 2(a):

$$\rho_{\text{init}}(x, t) = I(x, t)s(x, t). \quad (8)$$

From Eqs. (2d), (4), and (8), the density of terminations is

$$\rho_{\text{ter}}(x, t) = 2v \left[\int_{L_X^+} I \right] \left[\int_{L_X^-} I \right] s(x, t). \quad (9)$$

A termination at $X = (x, t)$ implies that no initiation occurs in V_X , one initiation occurs along L_X^+ , and one along L_X^- .

D. Rate equations for fork densities

From the above formalism, we can easily recover the rate-equation formalism proposed in [36] for fork densities. First, using Eq. (2c), the relation (2d) can be rewritten as a rate equation for the density of right- or left-moving forks,

$$(\partial_t \pm v\partial_x)\rho_\pm(x, t) = \rho_{\text{init}}(x, t) - \rho_{\text{ter}}(x, t). \quad (10)$$

Then, from Eqs. (5), (8), and (9) we find [36]

$$(\partial_t \pm v\partial_x)\rho_\pm(x, t) = Is - 2v \frac{\rho_+\rho_-}{s}. \quad (11)$$

Intuitively, fork densities change either because forks enter or leave a region (transport) or because there is initiation (birth) or termination (death).

E. Correlations in replication timing

As discussed in [21], the observation that neighboring loci tend to have similar replication times can be fully consistent with the independent-firing assumption. To more precisely quantify the correlation between replication times at different loci, we introduce the N -point unreplicated fraction $s(X_1, \dots, X_N)$, where X_i denotes the spacetime point (x_i, t_i) . We define s to be the fraction of cells where each of the N loci x_i is unreplicated at time t_i . The joint probability distribution of replication timing at loci x_1, \dots, x_N is then given by

$$P(X_1, \dots, X_N) = (-1)^N \partial_{t_1} \cdots \partial_{t_N} s(X_1, \dots, X_N). \quad (12)$$

In Fig. 2(b) we note that each locus x_i is unreplicated at time t_i iff no initiations occur in $V_{X_1} \cup \dots \cup V_{X_N}$, the union of past cones. Therefore,

$$s(X_1, \dots, X_N) = e^{-\int_{V_{X_1} \cup \dots \cup V_{X_N}} dX' I(X')}. \quad (13)$$

In [37], Sekimoto derived an equivalent expression in the more general setting of a time-dependent growth law.

To see why replication-fork propagation creates correlations between the replication times at different loci, consider the $N = 2$ case. Since $V_{X_1} \cup V_{X_2} = V_{X_1} + V_{X_2} - V_{X_1} \cap V_{X_2}$, the two-point unreplicated fraction is equal to

$$s(X_1, X_2) = s(X_1)s(X_2)e^{+\int_{V_{X_1} \cap V_{X_2}} dX' I(X')}. \quad (14)$$

If the replication times at loci x_1 and x_2 were uncorrelated, both their probability distributions and their cumulative distributions would factor: $P(X_1, X_2) = P(X_1)P(X_2)$ and $s(X_1, X_2) = s(X_1)s(X_2)$. It is clear from Eq. (14) that replication times at loci x_1 and x_2 are correlated because initiation events may occur in their common past cone $V_{X_1} \cap V_{X_2}$. Indeed, if $I(X)$ is not identically zero in $V_{X_1} \cap V_{X_2}$, then $s(X_1, X_2) \neq s(X_1)s(X_2)$. However, if the loci x_1 and x_2 are

sufficiently far apart—that is, if $|x_1 - x_2| \geq 2vt_{\text{end}}$, where t_{end} is the duration of the S phase—then their past cones do not intercept, and the replication times at x_1 and x_2 are indeed uncorrelated.

F. The joint density of initiation

In Sec. III E, we saw that the propagation of replication forks creates correlations in the timing of replication: a location near an origin will tend to replicate soon after that origin fires. A less obvious kind of correlation also exists in the initiation densities, where, again, we argue that apparent correlations can sometimes be deceptive. Indeed, experimental observations of apparent origin synchrony [38] or of sequential firing, as observed in temporal transition regions [39], suggest that initiations may be temporally and spatially correlated, contradicting the independent-firing assumption. Here, we will see that inferring independence from such observations can be subtle.

In order to quantify the correlations observed in the distribution of initiations, we introduce the N -point joint density of initiations $\rho_{\text{init}}(X_1, \dots, X_N)$, defined as the probability of observing, during the same cell cycle, an initiation at each X_i . Let us first assume that no X_i belongs to the past cone of another X_j , as depicted in Fig. 2(b). Then, an initiation at each X_i implies also that no initiation has occurred in $V_{X_1} \cup \dots \cup V_{X_N}$. Since the origins fire independently, the joint density of initiation is

$$\rho_{\text{init}}(X_1, \dots, X_N) = I(X_1) \cdots I(X_N) s(X_1, \dots, X_N). \quad (15)$$

To illustrate why replication-fork propagation necessarily creates correlations in the joint density of initiation, we rewrite these expressions for $N = 2$:

$$\rho_{\text{init}}(X_1, X_2) = \rho_{\text{init}}(X_1) \rho_{\text{init}}(X_2) e^{\int_{V_{X_1} \cap V_{X_2}} dX' I(X')}. \quad (16)$$

As in Eq. (14), initiation densities at X_1 and X_2 are correlated because of possible origin firing in their common past cone $V_{X_1} \cap V_{X_2}$. To prove that neighboring initiations influence each other then takes more than the observation of initiation clusters or of sequential firing of nearby origins. Only a clear departure from Eq. (16) would provide definitive evidence.

Finally, if one of the X_i belongs to the past cone of another X_j , $\rho_{\text{init}}(X_1, \dots, X_N)$ is necessary null. As re-replication is not allowed, we cannot observe an initiation in the future cone of another origin firing. The joint density of initiation must satisfy this trivial correlation.

G. Well-positioned replication origins

In organisms such as the budding yeast *Saccharomyces cerevisiae*, origins initiate at predefined sites called *potential origins*. The local initiation rate then has the form [21]

$$I(x, t) = \sum_i \delta(x - x_i) I_i(t), \quad (17)$$

where x_i is the position of potential origin i and $I_i(t)$ its initiation rate. All the analytical formulas derived in [21, 28] are recovered as particular cases of the more general and compact

expressions Eqs. (4)–(16), with the local initiation rate given by Eq. (17).²

Let us specify the expressions for $s(x, t)$ and $\rho_{\text{init}}(x)$ in the case of well-positioned origins. From Eqs. (4) and (17), the unrepliated fraction can be written

$$s(x, t) = \prod_i s_i \left(t - \frac{|x - x_i|}{v} \right), \quad (18)$$

$$\text{where } s_i(t) \equiv e^{-\int_0^t dt' I_i(t')} \quad (19)$$

is the probability that the potential origin i has yet not initiated at time t . In words, the locus x is unrepliated a time t iff each origin i has not initiated before time $t - |x - x_i|/v$. From Eqs. (8) and (17), the initiation density profile will have sharp peaks at potential-origin sites:

$$s(x, t) = \prod_i s_i \left(t - \frac{|x - x_i|}{v} \right), \quad (20)$$

$$\text{where } s_i(t) \equiv e^{-\int_0^t dt' I_i(t')} \quad (21)$$

where E_i , the *observed efficiency* of origin i , is defined as the fraction of cells where the origin i has activated before the end of the S phase. The observed efficiency of the origin i depends on its initiation properties but is also affected by the initiation properties of neighboring origins [21, 28]. Indeed, when the locus x_i is replicated by a fork coming from a neighboring origin, the potential origin i will not be activated during this cell cycle, and the potential origin is *passively replicated*. It is then interesting to consider the *potential efficiency* of a replication origin—the probability that the origin would activate during the S phase if passive replication by neighboring origins is prevented. The potential efficiency q_i of origin i , denoted *origin competence* in [19, 28], is equal to

$$q_i = 1 - e^{-\int_0^{\infty} dt' I_i(t')} = 1 - s_i(t_{\infty}), \quad (22)$$

as $s_i(t)$ is the probability that the origin i has not yet initiated at time t . Contrary to a claim in [28], the KJMA formalism does not assume 100% competent origins; in general, $q_i < 1$.³ In budding yeast, passive replication has a strong impact on

²To make the connection with Refs. [21, 28] more explicit, let us specify some of the quantities introduced in those references in terms of the local initiation rate. The initiation probability density $\phi_i(t)$ in [21], or origin activation time probability density $p_i(t)$ in [28], is given by $\phi_i(t) = p_i(t) = -\partial_t s_i(t) = I_i(t) s_i(t)$. Note that this is not a normalized probability distribution, as $\int_0^{\infty} dt \phi_i(t) = q_i < 1$ is the potential efficiency, or competence, of origin i . In [28], we also have $M_i(x, t) = s_i(t - |x - x_i|/v)$ and $p_i(x, t) = I_i(t - |x - x_i|/v) s_i(t - |x - x_i|/v)$. In terms of the local initiation rate, the combinatorial expressions in [28] simplify greatly; for instance, $p_i(x, t) \prod_{j \neq i} M_j(x, t) = I_i(t - |x - x_i|/v) s(x, t)$.

³For example, we can have $q_i < 1$ if origins fail to be licensed prior to the start of the S phase [19, 28]. Let the *licensing probability* for origin i be L_i . Then $s_i(t) = (1 - L_i) + L_i e^{-I' t}$, where I' is the initiation rate *if licensed*, assumed, for simplicity to be constant for all origins and all time. From Eq. (21), $I_i(t) = -\frac{d}{dt} \ln s_i(t) \sim I' \left(\frac{L_i}{1 - L_i} \right) e^{-I' t}$ for times $t \gg 1/I'$. A finite licensing probability thus cuts off the effective initiation rate at long times, and the failure to

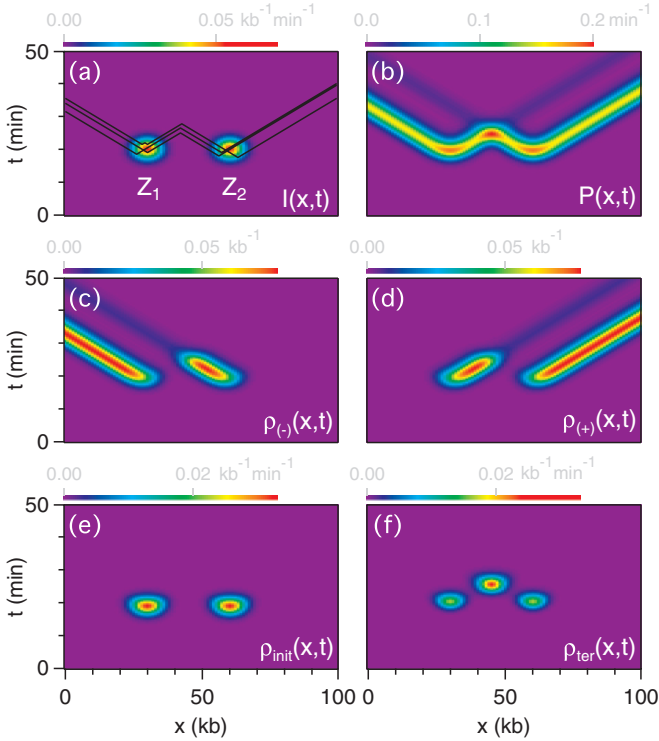


FIG. 3. (Color online) Replication program with two extended initiation zones Z_1 and Z_2 . (a) Heat map of the local initiation rate $I(x,t)$. The black lines correspond to single cell cycle realization of the replication program, obtained by Monte Carlo simulation. (b) Replication distribution Eq. (7). (c),(d) Densities of left- and right-moving forks, Eq. (5). (e),(f) Densities of initiation [Eq. (8)] and termination [Eq. (9)].

the efficiencies of replication origins: the observed efficiency is usually much smaller than the potential efficiency [21].

IV. EXAMPLE REPLICATION PROGRAM

Let us now illustrate the formalism developed in the two preceding sections on an artificial replication program that consists of two extended initiation zones Z_1 and Z_2 . In Fig. 3(a), the spatiotemporal representation of the local initiation rate is color coded by a heat map. To give an idea of the resulting stochasticity, we sample by Monte Carlo simulation five realizations of the replication program, represented by the black lines on Fig. 3(a). Several aspects of the replication program, analytically derived from the local initiation rate using the results of Secs. II and III, are represented on Figs. 3(b)–3(f).

Notice how Fig. 3 reveals many fine details about the replication process. For example, the density of termination events in Fig. 3(f) shows three zones. At the center is the strongest one, representing the case where forks from the two origin regions collide after propagating roughly to the midpoint between the initiation zones Z_1 and Z_2 . The two

weaker termination zones overlap with the initiation zones and represent cases where two or more initiation events within the same zone lead to a fork collision soon after the initiation event. Solving the analytical model allows us to detect and quantify the probability for these different scenarios to occur.

V. INFERRING THE LOCAL INITIATION RATE

In Secs. II–IV, we showed how to solve the *forward* problem of replication: given an initiation rate $I(x,t)$, calculate various quantities of interest for the replication process, for example the unreplicated fraction $s(x,t)$. Now we consider the *inverse* problem: given a noisy measurement of $s(x,t)$, can we infer $I(x,t)$? In particular, we advance a nonparametric method that avoids having to define a model structure for $I(x,t)$.

To test the method under well-controlled circumstances, we will focus on inverting simulated data based on the spatiotemporal replication program presented in Sec. IV. The data will have a spacetime resolution comparable to that of present experiments and will include noise levels that are also typical.

We begin by first reviewing past attempts to solve this inverse problem, including fitting strategies and analytic approaches based on expressing the initiation rate $I(x,t)$ as a function of the nonreplicated fraction $s(x,t)$. After discussing the limitations of previous attempts, we then propose a Bayesian, nonparametric approach to infer $I(x,t)$ from replication timing data. We will test this inference scheme on the artificial data set described above and show that near-perfect reconstruction of the replication program (with negligible posterior uncertainty) is attained for many quantities of interest, such as the unreplicated fraction, the densities of replication forks, and the densities of initiation and termination. The local initiation rate is also inferred with low posterior uncertainty in most regions except at the end of the S phase, where the unreplicated fraction, already close to zero, is insensitive to large variations in the initiation rate.

A. Curve-fitting strategies

As discussed in the Introduction (Sec. I), the replication fork velocity v and initiation function $I(x,t)$ can be estimated by curve fitting [12,18–22]. The main issue is that one must make strong assumptions about the prior functional form for $I(x,t)$, for example, whether origins are localized along the genome, the type of time dependence, etc. Besides requiring *a priori* knowledge about the biology that is not always available, the underlying forms may not really be what is assumed. Also, the number of parameters needed is not clear in advance. For example, the number of detectable origins in budding yeast is an output of the inference process. In addition, in the most commonly used implementations of curve fitting, one needs to provide initial values for all parameters. It is often not easy to find initial values that are close enough to the best estimate to converge. Methods such as genetic algorithms that can optimize the fit globally improve upon this aspect, albeit at the cost of greater computation [18,19].

For all these reasons, a successful curve fit requires both *a priori* knowledge and a good level of technical expertise. Below, we will explore a strategy that requires only vague *a priori* expectations and that can, in principle, be automated.

license origins can be absorbed into the effective initiation rate. Note that if $L_i = 1$, we recover $I_i = I'$.

B. Exact inverse

Recently, we showed how to invert explicitly the KJMA formula Eq. (4), thereby determining analytically $I(x,t)$ from $s(x,t)$ [40]:⁴

$$I(x,t) = -\frac{1}{2}v\Box \ln s(x,t). \quad (23)$$

Because Eq. (23) gives an exact expression for $I(x,t)$, it would seem to provide an alternative to curve-fit approaches: rather than guess the form of $I(x,t)$, we can simply calculate it from the data $s(x,t)$. Unfortunately, the analytical inverse is numerically unstable: taking two derivatives amplifies noise tremendously. Thus, Eq. (23) can be naively applied only if essentially noise-free data for $s(x,t)$ are available. For example, in [40], we used Eq. (23) to invert simulations that had negligible numerical noise. When applied directly to low-resolution experimental data with realistic amounts of noise, Eq. (23) gives unphysical results such as negative initiation rates [12]. Simple *ad hoc* fixes, such as smoothing $s(x,t)$ over fixed space and time scales [40], lead to unacceptable distortion in the estimate of $I(x,t)$ and also do not give uncertainties in estimated initiation rates. All of these shortcomings motivate a more fundamental approach.

C. Bayesian inference

Here, we will adopt a Bayesian, nonparametric approach to more properly infer $I(x,t)$ from replication timing data. Bayesian methods offer a consistent and conceptually well-founded framework for inference, where all assumptions are explicitly stated [42].

1. Introduction

The Bayesian formulation is well adapted to parameter-estimation problems [42]. In our case, the goal is to infer the parameter I (the local initiation rate) from the data d (a noisy measurement of the unreplicated fraction). We recall that the posterior probability of I , given data d , is determined by Bayes' theorem, which is derived from the product and sum rules of probability theory [42]:

$$\underbrace{P(I|d,\beta)}_{\text{posterior}} = \frac{1}{P(d|\beta)} \underbrace{P(d|I,\beta)}_{\text{likelihood}} \underbrace{P(I|\beta)}_{\text{prior}}, \quad (24)$$

where the normalizing factor, the *evidence*, is given by

$$\underbrace{P(d|\beta)}_{\text{evidence}} = \int dI P(d|I,\beta) P(I|\beta). \quad (25)$$

In Eq. (24), the likelihood follows the noise model for the data, while the prior encodes any available information—even vague—about the parameter to infer; in replication, for instance, we know that initiation rates $I(x,t)$ must be positive. We also expect that temporal and spatial variations of $I(x,t)$

are smooth, although we may not know the smoothness scales. Below, we will describe in more detail the probabilistic model used for inference given such vague priors.

Often, the specification of a probabilistic model for the likelihood and the prior requires an additional set of parameters, called *hyperparameters*, symbolized by β in Eqs. (24) and (25). In our case, the hyperparameters comprise the fork velocity v , which affects the relationship Eq. (4) between the unreplicated fraction data and the initiation rate, the noise level affecting the data, and additional parameters encoding prior information about the initiation rate, for example the temporal scale of smoothness. These hyperparameters can themselves be inferred by another application of Bayes' theorem [43]:

$$P(\beta|d) = \frac{1}{P(d)} P(d|\beta) P(\beta). \quad (26)$$

The posterior probability of the hyperparameters is thus proportional to the evidence and the prior probability of the hyperparameters. Given the posterior $P(\beta|d)$, we can eliminate the hyperparameters by *marginalization*, or “integrating out.” For example,

$$P(I|d) = \int P(I|d,\beta) P(\beta|d) d\beta. \quad (27)$$

The Bayesian formulation is also well adapted to model selection. Given data and candidate theories, Bayes' theorem allows one to estimate the most probable model [43]. For instance, we could compare the probabilistic model presented here and the fitting procedure (which can easily be reformulated in a Bayesian framework) employed in yeast. We could even compare to a theoretical model that extends the KJMA formalism to take into account correlations in the origin firing. Such model comparisons are beyond the scope of the present paper.

The inference task here is complicated by the nonlinear relationship Eq. (4) between the data (the unreplicated fraction profiles) and the initiation rate we seek to infer and by the positivity constraint on the initiation rate. Indeed, if the relationship were linear and no positivity constraint needed to be enforced, then we would be able to derive the posterior Eq. (24) analytically. Below, we will approximate the posterior probability distribution by its mode, the *maximum a posteriori* (MAP) approximation, which requires a high-dimensional nonlinear optimization algorithm. To estimate the width of the posterior, we will sample directly the posterior by Markov chain Monte Carlo (MCMC) techniques. Finally, to estimate the evidence, we will use the *Laplace approximation*, which is the analog of the saddle-point approximation in statistical physics.

2. Likelihood

We model the data as a noisy version of the unreplicated fraction s , sampled in time and space:

$$d_k = s(x_k, t_k) + \xi_k, \\ \text{with } s(x,t) = e^{-\int_{v_X} dx' dt' I(x',t')}, \quad (28)$$

with noise described by independent, identically distributed (i.i.d.) Gaussian random variables of standard deviation σ_d .

⁴Generalizing Eq. (23) to a space- and time-dependent velocity field $v(x,t)$ is straightforward, albeit cumbersome [41]: $I(x,t) = \frac{1}{2}[(\frac{1}{v}\partial_t v)\frac{1}{v}\partial_t + (\partial_x v)\partial_x - v\Box] \ln s(x,t)$. Because it is not at present clear whether systematic (as opposed to random) variation of fork velocities is important, we focus on the constant- v case.

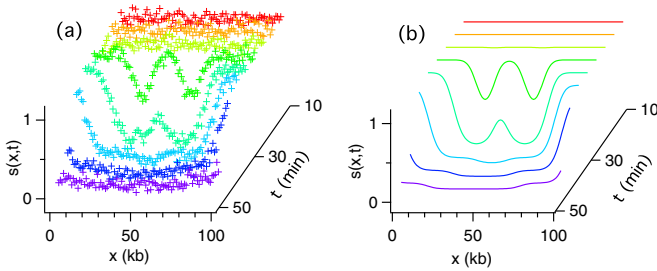


FIG. 4. (Color online) Simulation of the replication program with extended initiation zones (Fig. 3). (a) Artificial data set generated by adding Gaussian noise of standard deviation $\sigma_d = 0.05$ to the true unreplicated fractions in (b). In (a) and (b), the unreplicated replicated fraction is given at every 1 kb from $x = 1$ to 100 kb and every 5 min from $t = 10$ to 50 min.

Thus, $\xi_k \sim \mathcal{N}(0, \sigma_d^2)$, and the likelihood is

$$\begin{aligned} P(d|I, v, \sigma_d) &= P_{\text{noise}}(d - s) \\ &= \prod_k \frac{1}{\sqrt{2\pi\sigma_d^2}} e^{-(1/2\sigma_d^2)[d_k - s(x_k, t_k)]^2}, \end{aligned} \quad (29)$$

where the product is over all data points k .

In the artificial data set shown in Fig. 4, the noisy unreplicated fractions are sampled every 1 kb in a fragment of 100 kb and every 5 min from $t = 10$ min to $t = 50$ min. These resolutions match that of the recent budding-yeast experiments described above. We chose $\sigma_d = 0.05$, again typical of current experiments [17,21]. Note that Fig. 4 can also be interpreted as a plot of replicated fraction $f = 1 - s$ from times of 10 to 50 min.

Comment on the noise model. Although we model the noise by i.i.d. Gaussian random variables of standard deviation σ_d , it is straightforward to substitute any noise model in Eq. (29), including correlations, time- or space-dependent variance, or non-Gaussian distributions. As a real-world example, the analysis of data on budding yeast showed a variance that increased throughout the S phase and a noise distribution, that while Gaussian for small fluctuations, was exponential for larger ones [21]. In general, small deviations from the Gaussian form will not affect the analysis much.

3. Prior

A key advantage of the Bayesian formulation is that we can specify the prior, the set of possible initiation rate functions, without having to impose a particular functional form. Nevertheless, we do have some vague prior knowledge about $I(x, t)$ that should be used to constrain the set of possible initiation functions: it must be positive and its temporal variations are smooth. In some cases, spatial variations are also smooth.

To ensure the positivity of the initiation rate, we change variables, defining

$$I(x, t) \equiv I_0 \times 10^{m(x, t)}. \quad (30)$$

In other words, rather than trying to infer the initiation rate I directly, we will infer its logarithm m .

To enforce smooth variations in the initiation rate, we will use a Gaussian process prior [44] on $m = \log_{10}(I/I_0)$:

$$m \sim \mathcal{GP}(0, \Sigma), \quad (31)$$

with a homogeneous, squared-exponential covariance function that depends on the spatial separation Δx and the temporal separation Δt :

$$\Sigma(x, t; x + \Delta x, t + \Delta t) = \sigma_0^2 e^{-(\Delta x/\ell_0)^2} e^{-(\Delta t/\tau_0)^2}. \quad (32)$$

A Gaussian process m can be viewed as the infinite-dimensional analog of the multivariate normal distribution; it defines a probability distribution over functions. The precise definition is that the values of m at an arbitrary set of points (X_1, \dots, X_N) are distributed according to the multivariate normal distribution $[m(X_1), \dots, m(X_N)] \sim \mathcal{N}(0, \Sigma)$, with covariance matrix $\Sigma_{ij} = \Sigma(X_i, X_j)$. In our case, we would like to infer the initiation rate at a spatial resolution of $\delta x = 1$ kb and a temporal resolution of $\delta t = 0.5$ min (we set δt in order to have $\delta x = v \delta t$, with a fork velocity equal to $v = 2$ kb min⁻¹). This defines the grid of points $X \equiv (x, t)$ where m should be evaluated. The prior distribution on $m = \{m(x, t)\}$ is therefore the multivariate normal

$$P(m|\sigma_0, \tau_0, \ell_0) = \frac{1}{\sqrt{\det(2\pi\Sigma)}} e^{-(1/2)m\Sigma^{-1}m}, \quad (33)$$

with the covariance function Σ evaluated at the grid of points (x, t) using Eq. (32). In the covariance function Eq. (32), σ_0 quantifies the prior expectations about the range of values taken by m . The square-exponential decay as a function of the time interval separating two points, on a characteristic time scale τ_0 , enforces the smoothness of the function m on the same time scale, and similarly for the spatial scale ℓ_0 . The limit $\ell_0 \rightarrow 0$ means that m values at different genomic positions are uncorrelated. It is obtained by replacing the squared exponential in Eq. (32) by a Dirac delta function $\delta(\Delta x)$. The squared-exponential form of the covariance matrix the defined in Eq. (32) is the standard choice in the Gaussian process literature [44]. It yields smooth functions that are differentiable to all orders. Other choices are possible, when appropriate. For example, correlations that decay exponentially lead to functions that are continuous but not differentiable.

In Gaussian process regression, the task is to go from a Gaussian-prior representation of $m(x, t)$ [Eq. (31)] to a posterior representation that incorporates the noisy observations d_k . Note that many authors define a Gaussian process regression to be one where the posterior distribution for m is also a Gaussian process (that is, they assume that the data are related to m by a linear transformation). Here, the data and m are nonlinearly related, and the resulting distribution for m is non-Gaussian. For simplicity, we also refer to this case as Gaussian process regression, but we will need to use special techniques to deal with the non-Gaussian nature of the posterior distribution.

4. Hyperparameters

As discussed earlier, we can estimate the hyperparameters from the data set itself. Here, instead of carrying out this procedure for all of them, we will do so only for the most interesting ones, the fork velocity v and the spatial scale ℓ_0 for $I(x, t)$ variations. The latter is especially delicate, in

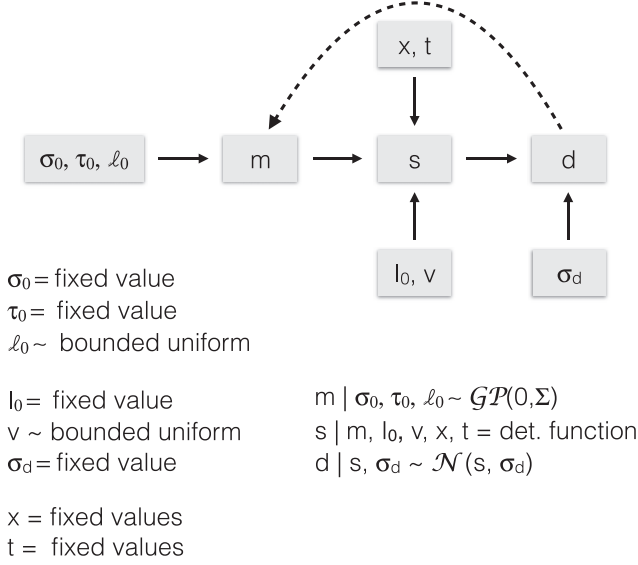


FIG. 5. Diagram summarizing the forward replication model $m \rightarrow d$ evaluated at grid points (x, t) and its hyperparameters $(\sigma_0, \tau_0, \ell_0, l_0, v, \sigma_d)$. The symbol “ \sim ” means “distributed as,” and the dashed arrow denotes the inference $d \rightarrow m$.

that some organisms, such as budding yeast, have near- δ -function initiation sites, while others, such as frog embryos, permit initiation anywhere and have slowly varying densities. Accordingly, we will carry out the self-consistent selection for these parameters below.

We first fix the hyperparameters of lesser interest. For example, σ_0 and l_0 set the range of values allowed for the initiation rate. Their precise value should not matter much, as long as the allowed range of values is larger than the actual range of values taken by the initiation rate. Here, we choose $l_0 = 10^{-4} \text{ kb}^{-1} \text{ min}^{-1}$ and $\sigma_0 = 3$, to allow for a very wide range of values for the initiation rate. This choice allows a 1σ range of initiation rates of between 10^{-1} and $10^{-7} \text{ kb}^{-1} \text{ min}^{-1}$.

The temporal scale τ_0 defines how quickly $I(x, t)$ can vary. Although in principle as interesting as the spatial scale ℓ_0 , the evidence to date suggests that the experimental range of values is much narrower. For example, previous analysis of the replication kinetics in yeast [21] is consistent with $\tau_0 \approx 10 \text{ min.}$, about 1/4 the duration of the S phase, and we used this value in the inference procedure.

The complete probabilistic model is summarized in Fig. 5. Below, we first use the model to infer the logarithmic initiation rate $m = \log(I/I_0)$ from the data d , assuming the hyperparameters to be known. In the last section, we will solve for m_{MAP} over a grid of values for v and ℓ_0 and find that the posterior is almost entirely concentrated at the correct (simulation) values.

5. Posterior

The posterior $P(m|d, \beta)$ for the logarithmic initiation rate $m = \log_{10}(I/I_0)$ is given by Bayes’ theorem Eq. (24), with the likelihood given by Eq. (29), the prior given by Eq. (33), and the hyperparameters $\beta = \{v, \sigma_d, l_0, \sigma_0, \tau_0, \ell_0\}$. Note that the parameter m to infer is evaluated at a resolution of 1 kb in space

and 0.5 min in time and thus forms an $(N_x N_t = 100 \times 100)$ -dimensional vector. Thus, the posterior for m is a probability distribution defined on a very high-dimensional (10^4) space. Below, we will consider both replacing the distribution by its mode (maximum *a posteriori* approximation) and sampling the posterior by MCMC techniques.

6. Maximum a posteriori approximation

The mode of the posterior distribution, which gives the *maximum a posteriori* estimate, can be found by minimizing the “energy” functional [45,46]

$$\begin{aligned}
 E(m) &= -\ln P(m, d|\beta) \\
 &= -\ln P(d|m, \beta) - \ln P(m|\beta) \\
 &= \frac{1}{2\sigma_d^2} \sum_k [d_k - s(x_k, t_k)]^2 + \frac{1}{2} N_d \ln 2\pi \sigma_d^2 \\
 &\quad + \frac{1}{2} m \Sigma^{-1} m + \frac{1}{2} \ln \det(2\pi \Sigma)
 \end{aligned}
 \tag{34}$$

$$\text{with } s(x, t) = e^{-\int_{v_x} dx' dt' I_0 \times 10^{m(x', t')}}.
 \tag{34}$$

The quantity $E(m)$ is the negative logarithm of the joint posterior, with N_d the number of data points. The MAP estimate $m_{\text{MAP}} = \text{argmin } E(m)$ can be interpreted as a compromise between minimizing the least-square fit $\frac{1}{2\sigma_d^2} (d - s)^2$ (the “energy”) and minimizing $\frac{1}{2} m \Sigma^{-1} m$ (the “entropy”), where smoother states have lower entropy because they are compatible with fewer data sets. Alternatively, we can view the minimization as a regularized “Tikhonov” inverse [47], where the compromise is between finding the m that best reproduces the data d and minimizing the Tikhonov penalty term, which favors smooth m on the spatial scale ℓ_0 and temporal scale τ_0 .

We minimized E in Eq. (34) numerically via the Newton conjugate gradient algorithm [48]. Although we minimize in a 10^4 -dimensional space, the program converges in less than a minute on a regular laptop.

The MAP approximation is to replace the posterior distribution by a Dirac δ function at its mode,

$$P(m|d, \beta) \simeq \delta(m - m_{\text{MAP}}).
 \tag{35}$$

That is, we simply substitute m_{MAP} into the analytical expression of the initiation rate and into all other quantities of interest. As shown in Fig. 6(a), the estimated local initiation rate $I_{\text{MAP}} = I_0 \times 10^{m_{\text{MAP}}}$ is very close to the true initiation rate, Fig. 3(a). Similarly, the estimated unreplicated fraction, Fig. 6(b), density of right- and left-moving forks, Figs. 6(c) and 6(d), as well as the density of initiation, Fig. 6(e), and termination, Fig. 6(f), obtained by simply substituting I_{MAP} in the analytical expressions of Sec. III are indistinguishable from their true values in Fig. 3. Finally, note that all those quantities are reconstructed at the desired temporal resolution of 0.5 min, while the original data d in Fig. 4 have only a temporal resolution of 5 min. This interpolation is possible because the temporal smoothness scale $\tau_0 = 10 \text{ min.}$

7. MCMC sampling of the posterior

The MAP approximation Eq. (35) would seem to be a rather crude one, as it neglects the posterior uncertainty

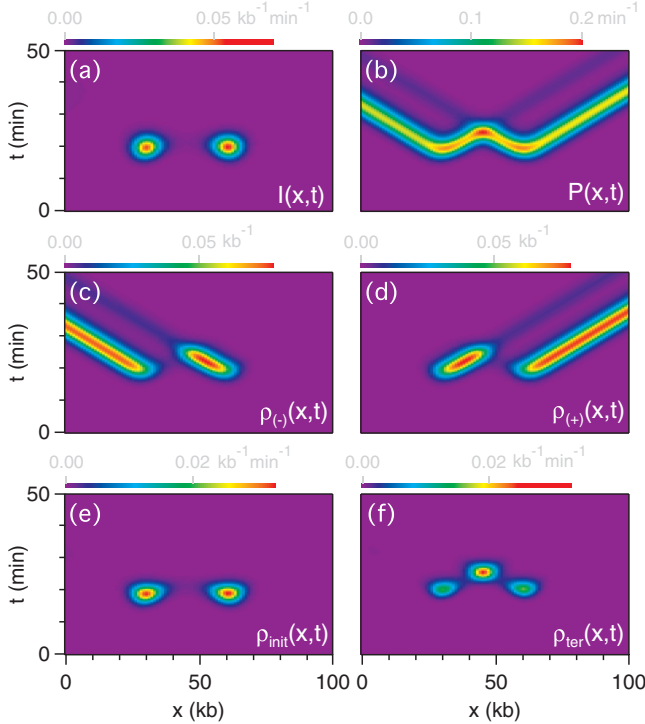


FIG. 6. (Color online) Near-perfect reconstruction of the replication program in Fig. 3. All characteristics of the replication program are reconstructed using the MAP estimate m_{MAP} of $m = \log_{10}(I/I_0)$. (a) Local initiation rate $I_{\text{MAP}} = I_0 \times 10^{m_{\text{MAP}}}$. (b) Replication distribution Eq. (7). (c),(d) Densities of left- and right-moving forks, Eq. (5). (e),(f) Densities of initiation [Eq. (8)] and termination [Eq. (9)].

for m . Moreover, the MAP estimate m_{MAP} is usually not a representative sample from the posterior, and its value is not invariant under reparametrization [43]. However, in our particular case, the MAP estimate m_{MAP} does yield a very accurate reconstruction of the replication program: Since, as we will see below, the posterior uncertainty for most quantities turns out to be negligible, samples from the posterior distribution are almost always close to the MAP value.

To estimate the width of the posterior distribution Eq. (24), we used Markov chain Monte Carlo sampling. We first implemented the classic Metropolis-Hastings algorithm, but it was very slow. We then tried instead the Hamiltonian Monte Carlo algorithm [43], which was about 100 times faster. We initialized the Markov chain at the MAP estimate in order to skip the burn-in phase and used the Hessian of the energy $E(m)$ as a preconditioning matrix for the momentum. We generated an effectively independent sample (i.e., an evaluation over the entire spacetime grid) every 5 s on a regular laptop. Ten samples from the posterior distribution are given in Fig. 7, as well as the 90% credible interval. We see that the posterior uncertainty for the unreplicated fraction, Fig. 7(c), the replication fork polarity, Fig. 7(d), the density of initiation, Fig. 7(e), and of termination, Fig. 7(f), are negligible, with a small posterior uncertainty at the boundaries. The local initiation rate has low posterior uncertainty, Fig. 7(a), except at the end of the S phase, Fig. 7(b). The large uncertainty on the initiation rate at the end of the S phase is easily understandable: At the end of the S phase, the unreplicated fractions are close

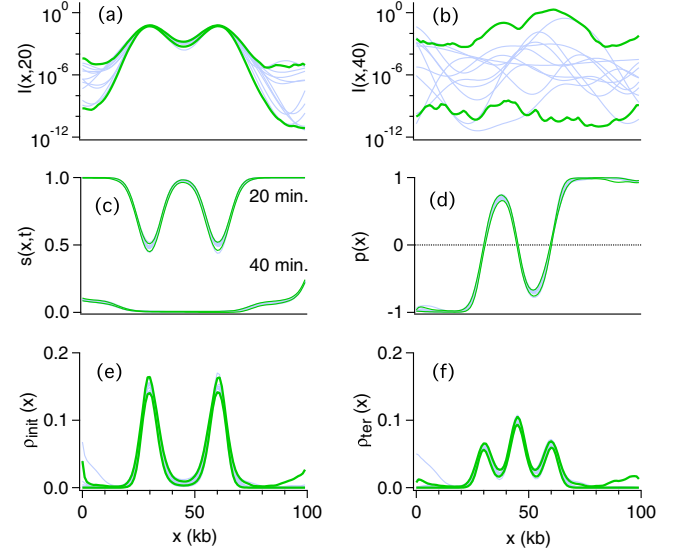


FIG. 7. (Color online) Negligible posterior uncertainty, except for the initiation rate at the end of the S phase. Ten MCMC samples (light blue lines) from the posterior probability distribution Eq. (24), and the 90% credible interval (between heavy green lines). Local initiation rate at (a) $t = 20$ min and (b) 40 min. (c) Unreplicated fractions at $t = 20$ and 40 min. (d) Replication fork polarity. Spatial density of (e) initiation and (f) termination.

to zero; thus, even large variations of the local initiation rate result in minor variations in the unreplicated fractions that will be much smaller than the noise level. The initiation rate thus cannot be accurately inferred in these regions. However, as we have seen in Figs. 7(c)–7(f), the large uncertainty in the local initiation rate at the end of the S phase results in negligible uncertainty for other quantities of interest.

8. Inferring v and ℓ_0

We inferred the most important hyperparameters, the fork velocity v and the spatial smoothness scale ℓ_0 , directly from the data. By Bayes' theorem applied to the hyperparameters in Eq. (26), the posterior distribution for v and ℓ_0 is given by

$$P(v, \ell_0 | d, \beta') = \frac{1}{P(d | \beta')} P(d | v, \ell_0, \beta') P(v, \ell_0), \quad (36)$$

where $\beta' = \{\sigma_d, I_0, \sigma_0, \tau_0\}$ contains the remaining hyperparameters. If we assume a flat prior on v and ℓ_0 , the posterior $P(v, \ell_0 | d, \beta')$ is simply proportional to the evidence $P(d | v, \ell_0, \beta') = P(d | \beta)$. From Eq. (25), the evidence $P(d | \beta)$ is evaluated by integrating the joint posterior $P(m, d | \beta)$ over m , a 10^4 -dimensional vector. Such a high-dimensional integration cannot be performed numerically. In the Laplace approximation [43], the joint posterior is approximated by a Gaussian around its maximum (the MAP estimate m_{MAP}):

$$P(m, d | \beta) \simeq e^{-E_{\text{MAP}} - (1/2)(m - m_{\text{MAP}}) \cdot E''_{\text{MAP}}(m - m_{\text{MAP}})}, \quad (37)$$

where E_{MAP} is the energy Eq. (34) at the MAP, and E''_{MAP} is the Hessian of the energy evaluated at the MAP. As the distribution is a Gaussian, the integration over m can be done

analytically. The logarithmic evidence is then

$$\ln P(d|\beta) \simeq \frac{1}{2} \ln \det(2\pi E''_{\text{MAP}}) - E_{\text{MAP}}. \quad (38)$$

This formula corresponds to the saddle-point approximation often encountered in statistical physics.

We then evaluated the Laplace approximation of the logarithmic evidence on a grid of values for v and ℓ_0 , spanning $v = 1$ kb/min to $v = 3$ kb/min every 0.1 kb/min for the fork velocity, and $\ell_0 = 0$ kb to $\ell_0 = 20$ kb every 5 kb for the spatial smoothing scale. We found that the value of the evidence at $v = 2$ kb/min and $\ell_0 = 15$ kb (the true values of the artificial data set) was several orders of magnitude larger than the evidence at other values. In other words, the posterior probability for v and ℓ_0 is, at the resolution considered, almost equal to 1 at the true values of (v, ℓ_0) and zero elsewhere. Therefore, for the data set considered here, we can infer accurately (at a resolution of 0.1 kb/min and 5 kb) the fork velocity and the spatial scale with near certainty.

9. Scaling up to a genome-wide analysis

The inference procedure described here was performed on a small fragment of length 100 kb. The Gaussian-process method involves, in effect, inverting a matrix whose size scales with the number N_d of $s(x, t)$ data points. Since matrix inversion varies as N_d^3 and since mammalian chromosomes are 10–100 times longer, we would need to increase the computation times by a factor of 10^3 – 10^6 , which would be demanding. However, the N_d^3 scaling is a worst-case scenario that holds when all data points are correlated. But only data on length scales of ℓ_0 and time scales t_0 are actually correlated, meaning that there are roughly independent blocks of size $\ell_0 \times t_0$, within which correlations must be accounted for. Consequently, the matrix is *sparse*, and computations should scale as N_d , not N_d^3 . The Gaussian-processes literature includes many specific examples where such sparse approximations have been successful [49].

VI. CONCLUSIONS

In this article, we have generalized the forward analysis of the DNA replication problem to the case of arbitrary initiation rates $I(x, t)$. We then introduced an inference procedure based on a Gaussian-process prior that avoids the need of earlier curve-fitting methods to specify the form of $I(x, t)$ in advance. We then showed that a small test case (100 kb genome) with typical replication parameters and typical experimental noise and resolution could be successfully inverted, with very small errors for all replication quantities of interest, except in cases where the experimental data were not very informative. (These cases were typically the end of the S phase and the edges of the sample.) The method may in principle be generalized to handle realistic genome sizes.

Assuming that the method does scale up and can successfully reproduce earlier analyses, we will then have a powerful method for learning about DNA replication in multiple organisms. Further, while we have focused on microarray and sequencing experiments, our methods should be compatible with the numerous other experimental methods, including fluorescence-activated cell sorting (FACS) [50], molecular combing [50], and Okazaki-fragment mapping [51]. Moreover, while the analysis is conceptually more complicated than curve

fitting, it can be automated and thus has the potential to be more widely used in the biological community.

From a more theoretical point of view, Gaussian-process regression [44] can be regarded as the equivalent of a free-field theory, in that the objects of interest are fields (defined over space and time) and are supposed to always show Gaussian fluctuations. In our case, the nonlinear relation between the replication data and the initiation rate of interest meant that our result was far from Gaussian. Although we used MCMC methods to sample the resulting non-Gaussian distributions, it would be interesting to explore other approaches to data analysis. In one approach, the parameter space of the probabilistic model defines a Riemannian manifold, allowing one to formulate a search algorithm for the MAP estimate [52] or MCMC exploration [53] in geometric terms. Taking a geometric approach can speed up the numerical algorithms discussed here. Alternatively, one can use the equivalent of interacting field theories and not assume Gaussian distributions. In this regard, the work of EnBlin and collaborators on *information field theory* [54] is an especially promising approach.

ACKNOWLEDGMENT

We thank Scott Yang and Nick Rhind for their suggestions. This work was funded by NSERC (Canada).

APPENDIX

We prove Eqs. (2a)–(2f) by first considering the replication program in one cell cycle. Then we show that the results derived for a single cell cycle generalize straightforwardly to the ensemble average for a stochastic or variable replication program.

1. In one cell cycle

Consider N origins O_1, \dots, O_N located at genomic positions $x_1 < \dots < x_N$ and initiated at times t_1, \dots, t_N , with fork velocities $\pm v$. From simple geometry [Fig. 1(a)], we see that each pair of origins (O_i, O_{i+1}) leads to a single termination event Λ_i at location $x_i^{(\Lambda)}$ and time $t_i^{(\Lambda)}$, where

$$\begin{aligned} x_i^{(\Lambda)} &= \frac{1}{2}(x_{i+1} + x_i) + \frac{1}{2}v(t_{i+1} - t_i), \\ t_i^{(\Lambda)} &= \frac{1}{2v}(x_{i+1} - x_i) + \frac{1}{2}(t_{i+1} + t_i). \end{aligned} \quad (A1)$$

The spatiotemporal densities of initiation and termination are therefore given by

$$\begin{aligned} \rho_{\text{init}}(x, t) &= \sum_i \delta(x - x_i) \delta(t - t_i), \\ \rho_{\text{ter}}(x, t) &= \sum_i \delta(x - x_i^{(\Lambda)}) \delta(t - t_i^{(\Lambda)}), \end{aligned} \quad (A2)$$

where $\delta(\cdot)$ is the Dirac delta function. Integrating over time gives the corresponding spatial densities:

$$\rho_{\text{init}}(x) = \sum_i \delta(x - x_i), \quad \rho_{\text{ter}}(x) = \sum_i \delta(x - x_i^{(\Lambda)}). \quad (A3)$$

The replication timing curve $T(x)$ is defined as the time at which the locus x is replicated and is represented as the solid line in Fig. 1(a). Let us define the *domain of origin* O_i to be $x \in [x_{i-1}^{(\Lambda)}, x_i^{(\Lambda)}]$. Within the domain, the replication timing curve is given by

$$T(x) = t_i + \frac{|x - x_i|}{v}. \quad (\text{A4})$$

The straight lines about each origin are one-dimensional analogs of the “light cones” of relativity that radiate from a source. In the similarly defined *domain of terminus* Λ_i , defined as $x \in [x_{i-1}, x_i]$ and illustrated in Fig. 1(a), the replication timing curve is given by the “past cones” from Λ_i :

$$T(x) = t_i^{(\Lambda)} - \frac{|x - x_i^{(\Lambda)}|}{v}. \quad (\text{A5})$$

The *unreplicated fraction* $s(x, t)$ is given by

$$s(x, t) = H[T(x) - t], \quad (\text{A6})$$

where H is the Heaviside step function.

In Fig. 1(a) in the domain of origin O_i , right-moving and left-moving replication forks have densities that are given by

$$\rho_{\pm}(x, t) = H[\pm(x - x_i)] \frac{1}{v} \delta[t - T(x)]. \quad (\text{A7})$$

Equivalently, in the domain of terminus Λ_i , the fork densities are given by

$$\rho_{\pm}(x, t) = H[\mp(x - x_i^{(\Lambda)})] \frac{1}{v} \delta[t - T(x)]. \quad (\text{A8})$$

Note that $p_{\pm}(x) = \int_0^{t_{\infty}} v dt \rho_{\pm}(x, t)$ equals 1 if the locus x is replicated by a \pm fork. Thus, the replication fork polarity $p(x) = p_+(x) - p_-(x) = \pm 1$ gives the directionality (\pm) of

the fork replicating the locus x . In the domain of origin O_i , the replication fork polarity is equal to

$$p(x) = \text{sgn}(x - x_i). \quad (\text{A9})$$

It is then straightforward, using the theory of distribution [55] and the above definitions Eqs. (A2)–(A9), to differentiate $s(x, t)$ and check the relations Eqs. (2a)–(2f).

2. Ensemble average

Because of the stochasticity of the replication program [1,4–6], the number of activated origins, their positions, and their firing times, all change from one cell cycle to another [Fig. 1(b)]. This variability may also reflect heterogeneity in the population of cells considered. For instance, mixtures of different cell types or cells of the same cell type but with different epigenetic states can give different stochastic replication programs. The ensemble average then corresponds to a superimposition of the different replication programs. A clear-cut example of the latter is the replication program in the human female X chromosome, where the ensemble average of replication seems to be “biphasic,” superposing the replication programs from the active and inactive X chromosomes [16,56].

The unreplicated fraction $s(x, t)$, the densities of initiation $\rho_{\text{mit}}(x, t)$ and termination $\rho_{\text{ter}}(x, t)$, the fork densities $\rho_{\pm}(x, t)$, the fork polarity $p(x)$, and the mean replication timing $T(x)$ defined in Sec. II A all correspond to the ensemble averages of their one-cell-cycle counterparts given in Sec. A 1. We proved in Sec. A 1 that the relations Eqs. (2a)–(2f) were true in each cell cycle. As derivatives and averages commute, we can straightforwardly extend Eqs. (2a)–(2f) to the ensemble average.

-
- [1] K. L. Friedman, B. J. Brewer, and W. L. Fangman, *Genes Cells* **2**, 667 (1997).
- [2] J. Herrick, P. Stanislawski, O. Hyrien, and A. Bensimon, *J. Mol. Biol.* **300**, 1133 (2000).
- [3] I. Lucas, M. Chevrier-Miller, J. M. Sogo, and O. Hyrien, *J. Mol. Biol.* **296**, 769 (2000).
- [4] P. K. Patel, B. Arcangioli, S. P. Baker, A. Bensimon, and N. Rhind, *Mol. Biol. Cell* **17**, 308 (2006).
- [5] N. Rhind, *Nat. Cell Biol.* **8**, 1313 (2006).
- [6] D. M. Czajkowsky, J. Liu, J. L. Hamlin, and Z. Shao, *J. Mol. Biol.* **375**, 12 (2008).
- [7] M. D. Sekedat, D. Fenyő, R. S. Rogers, A. J. Tackett, J. D. Aitchison, and B. T. Chait, *Mol. Syst. Biol.* **6**, 353 (2010).
- [8] M. G. Gauthier, J. Herrick, and J. Bechhoefer, *Phys. Rev. Lett.* **104**, 218104 (2010).
- [9] A. N. Kolmogorov, *Bull. Acad. Sci. URSS* **3**, 335 (1937); W. A. Johnson and P. Mehl, *Trans. AIME* **135**, 416 (1939); M. Avrami, *J. Chem. Phys.* **7**, 1103 (1939); **8**, 212 (1940); **9**, 177 (1941).
- [10] K. Sekimoto, *Int. J. Mod. Phys. B* **05**, 1843 (1991); E. Ben-Naim and P. L. Krapivsky, *Phys. Rev. E* **54**, 3562 (1996).
- [11] M. K. Raghuraman, E. A. Winzeler, D. Collingwood, S. Hunt, L. Wodicka, A. Conway, D. J. Lockhart, R. W. Davis, B. J. Brewer, and W. L. Fangman, *Science* **294**, 115 (2001).
- [12] J. Herrick, S. Jun, J. Bechhoefer, and A. Bensimon, *J. Mol. Biol.* **320**, 741 (2002).
- [13] C. A. Nieduszynski, S. I. Hiraga, P. Ak, C. J. Benham, and A. D. Donaldson, *Nucleic Acids Res.* **35**, D40 (2007).
- [14] D. Schübeler, D. Scalzo, C. Kooperberg, B. van Steensel, J. Delrow, and M. Groudine, *Nat. Genet.* **32**, 438 (2002).
- [15] I. Hiratani, T. Ryba, M. Itoh, T. Yokochi, M. Schwaiger, C.-W. Chang, Y. Lyou, T. M. Townes, D. Schubeler, and D. M. Gilbert, *PLoS Biol.* **6**, e245 (2008).
- [16] R. S. Hansen, S. Thomas, R. Sandstrom, T. K. Canfield, R. E. Thurman, M. Weaver, M. O. Dorschner, S. M. Gartler, and J. A. Stamatoyannopoulos, *Proc. Natl. Acad. Sci. USA* **107**, 139 (2010).
- [17] C. A. Müller, M. Hawkins, R. Retkute, S. Malla, R. Wilson, M. J. Blythe, R. Nakato, M. Komata, K. Shirahige, A. P. S. de Moura, and C. A. Nieduszynski, *Nucleic Acids Res.* **42**, e3 (2014).
- [18] M. Hawkins, R. Retkute, C. A. Müller, N. Saner, T. U. Tanaka, A. P. de Moura, and C. A. Nieduszynski, *Cell Rep.* **5**, 1132 (2013).
- [19] A. P. S. de Moura, R. Retkute, M. Hawkins, and C. A. Nieduszynski, *Nucleic Acids Res.* **38**, 5623 (2010).

- [20] H. Luo, J. Li, M. Eshaghi, J. Liu, and R. K. M. Karuturi, *BMC Bioinformatics* **11**, 247 (2010).
- [21] S. C.-H. Yang, N. Rhind, and J. Bechhoefer, *Mol. Syst. Biol.* **6**, 404 (2010).
- [22] A. Demczuk, M. G. Gauthier, I. Veras, S. Kosiyatrakul, C. L. Schildkraut, M. Busslinger, J. Bechhoefer, and P. Norio, *PLoS Biol.* **10**, e1001360 (2012).
- [23] D. M. Gilbert, *Nat. Rev. Genet.* **11**, 673 (2010).
- [24] J. Bechhoefer and B. Marshall, *Phys. Rev. Lett.* **98**, 098105 (2007); S. C.-H. Yang and J. Bechhoefer, *Phys. Rev. E* **78**, 041917 (2008).
- [25] A. Baker, B. Audit, C. Chen, B. Moindrot, A. Leleu, G. Guilbaud, A. Rappailles, C. Vaillant, A. Goldar, F. Mongelard, Y. d'Aubenton-Carafa, O. Hyrien, C. Thermes, and A. Arneodo, *PLoS Comput. Biol.* **8**, e1002443 (2012).
- [26] B. Audit, A. Baker, C.-L. Chen, A. Rappailles, G. Guilbaud, H. Julienne, A. Goldar, Y. d'Aubenton-Carafa, O. Hyrien, C. Thermes, and A. Arneodo, *Nat. Protoc.* **8**, 98 (2013).
- [27] R. Retkute, C. A. Nieduszynski, and A. P. S. de Moura, *Phys. Rev. Lett.* **107**, 068103 (2011).
- [28] R. Retkute, C. A. Nieduszynski, and A. de Moura, *Phys. Rev. E* **86**, 031916 (2012).
- [29] A. Baker, C. L. Chen, H. Julienne, B. Audit, Y. d'Aubenton-Carafa, C. Thermes, and A. Arneodo, *Eur. Phys. J. E* **35**, 123 (2012).
- [30] A. Baker, H. Julienne, C. L. Chen, B. Audit, Y. d'Aubenton Carafa, C. Thermes, and A. Arneodo, *Eur. Phys. J. E* **35**, 92 (2012).
- [31] J. Lygeros, K. Koutroumpas, S. Dimopoulos, I. Legouras, P. Kouretas, C. Heichinger, P. Nurse, and Z. Lygerou, *Proc. Natl. Acad. Sci. USA* **105**, 12295 (2008).
- [32] J. J. Blow and X. Q. Ge, *EMBO Rep.* **10**, 406 (2009).
- [33] S. Jun, J. Herrick, A. Bensimon, and J. Bechhoefer, *Cell Cycle* **3**, 211 (2004).
- [34] N. Van Kampen, *Stochastic Processes in Physics and Chemistry*, 3rd ed. (North-Holland, Amsterdam, 2007).
- [35] S. Jun, H. Zhang, and J. Bechhoefer, *Phys. Rev. E* **71**, 011908 (2005).
- [36] M. G. Gauthier, P. Norio, and J. Bechhoefer, *PLoS One* **7**, e32053 (2012).
- [37] K. Sekimoto, *Physica A* **135**, 328 (1986).
- [38] R. Berezney, D. D. Dubey, and J. A. Huberman, *Chromosoma* **108**, 471 (2000).
- [39] G. Guilbaud, A. Rappailles, A. Baker, C.-L. Chen, A. Arneodo, A. Goldar, Y. d'Aubenton-Carafa, C. Thermes, B. Audit, and O. Hyrien, *PLoS Comput. Biol.* **7**, e1002322 (2011).
- [40] A. Baker, B. Audit, S. C.-H. Yang, J. Bechhoefer, and A. Arneodo, *Phys. Rev. Lett.* **108**, 268101 (2012).
- [41] A. Baker, Ph.D. thesis, Ecole Normale Supérieure de Lyon, 2011.
- [42] E. T. Jaynes, *Probability Theory: The Logic of Science* (Cambridge University Press, Cambridge, 2003).
- [43] D. MacKay, *Information Theory, Inference and Learning Algorithms* (Cambridge University Press, Cambridge, 2003).
- [44] C. E. Rasmussen and C. K. Williams, *Gaussian Processes for Machine Learning* (MIT Press, Cambridge, MA, 2006).
- [45] J. C. Lemm, *Bayesian Field Theory* (The Johns Hopkins University Press, Baltimore, 2003).
- [46] W. Bialek, *Biophysics: Searching for Principles* (Princeton University Press, Princeton, NJ, 2012).
- [47] E. T. Jaynes, in *Inverse Problems*, edited by D. W. McLaughlin, SIAM-AMS Proceedings No. 14 (AMS, Providence, RI, 1984), pp. 151–166.
- [48] J. Nocedal and S. Wright, *Numerical Optimization* (Springer, New York, 2006).
- [49] E. Snelson and Z. Ghahramani, in *Advances in Neural Information Processing Systems*, edited by Y. Weiss, B. Schölkopf, and J. Platt (MIT Press, Cambridge, MA, 2006), Vol. 18, pp. 1259–1266.
- [50] E. Ma, O. Hyrien, and A. Goldar, *Nucleic Acids Res.* **40**, 2010 (2012).
- [51] S. R. McGuffee, D. J. Smith, and I. Whitehouse, *Mol. Cell* **50**, 123 (2013).
- [52] M. K. Transtrum, B. B. Machta, and J. P. Sethna, *Phys. Rev. Lett.* **104**, 060201 (2010); *Phys. Rev. E* **83**, 036701 (2011).
- [53] M. Girolami and B. Calderhead, *J. R. Stat. Soc. B* **73**, 123 (2011).
- [54] N. Oppermann, M. Selig, M. R. Bell, and T. A. Enßlin, *Phys. Rev. E* **87**, 032136 (2013).
- [55] W. Rudin, *Functional Analysis*, 2nd ed. (McGraw-Hill, New York, 1991).
- [56] A. Koren and S. A. McCarroll, *Genome Res.* **24**, 64 (2013).