

Statistical distributions of sequencing by synthesis
with probabilistic nucleotide incorporation

Yong Kong

Department of Molecular Biophysics and Biochemistry

W.M. Keck Foundation Biotechnology Resource Laboratory

Yale University

333 Cedar Street, New Haven, CT 06510

email: yong.kong@yale.edu

Abstract

Sequencing by synthesis is used in many next-generation DNA sequencing technologies. Some of the technologies, especially those exploring the principle of single-molecule sequencing, allow incomplete nucleotide incorporation in each cycle. We derive statistical distributions for sequencing by synthesis by taking into account the possibility that nucleotide incorporation may not be complete in each flow cycle. The statistical distributions are expressed in terms of nucleotide probabilities of the target sequences and the nucleotide incorporation probabilities for each nucleotide. We give exact distributions both for fixed number of flow cycles and for fixed sequence length. Explicit formulas are derived for the mean and variance of these distributions. The results are generalizations of our previous work for pyrosequencing. Incomplete nucleotide incorporation leads to significant change in the mean and variance of the distributions, but still they can be approximated by normal distributions with the same mean and variance. The results are also generalized to handle sequence context dependent incorporation. The statistical distributions will be useful for instrument and software development for sequencing by synthesis platforms.

1 Introduction

The next-generation DNA sequencing technology is changing biological research in many aspects and opening up new directions of experimental designs. The underlining technology of many of the current available platforms and those that are still under development can be categorized as sequencing by synthesis (SBS). In SBS, nucleotides are added to the reactions repeatedly in a pre-determined cyclic manner. Those nucleotides that are complementary to the template sequence will be potentially incorporated at each step, usually driven by enzymes. The presence or absence of any signal at each step reveals the nature of the template being sequenced.

In some of these SBS technologies, such as the Illumina Genome Analyzer, a mixture of four nucleotides is added at each cycle and the identity of the nucleotide is distinguished by the four different fluorescent labels. The nucleotides are modified so that only one base can be incorporated in each cycle. In this case, the length of sequencing reads has a simple relation with the number of sequencing cycles: they are equal to each other. In some other SBS technologies, such as pyrosequencing, however, only one kind of nucleotides is added to the reaction in each cycle. The length of sequencing reads will not have a fixed relation with the number of sequencing cycles. Rather, the read length is dependent on the sequence context and the nucleotide flow order. For the traditional pyrosequencing, the nucleotide incorporation is optimized to be as complete as possible in each cycle to avoid errors caused by dephasing (discussed below). At each extension cycle, ideally the nucleotide should be incorporated 100%

if it is complementary to the target sequence (including the homopolymer regions). Analytical statistical distributions have been derived for this situation in a previous paper (Kong, 2008). In the present paper we deal with statistical distributions for the case where the nucleotide incorporation is not achieved 100% at each cycle. For the emerging single-molecule DNA sequencing technology (Gupta, 2008), sometimes it is desirable to *not* achieve 100% nucleotide incorporation at each cycle in order to maximize the resolution of homopolymer regions and increase the accuracy of nucleotide incorporation (Harris *et al.*, 2008).

In this paper we derive statistical distributions for SBS with probabilistic, or incomplete, nucleotide incorporation. Mathematically these results are generalizations of the results obtained previously for 100% nucleotide incorporation at each cycle (Kong, 2008). These distributions are useful in different stages in the use and development of the next-generation DNA sequencing technology, such as instrument development and testing, algorithm and software development, and the everyday machine performance monitoring and trouble-shooting.

The paper is organized as follows. First in the remaining of this *Introduction* section we give a brief description of the single-molecule DNA sequencing technology (SMS), for which the theoretical results of present work will find most useful. We also define the necessary notation here. The derivation of the main results, which are exact under the assumptions of the sequence model, will be presented in the *Bivariate Generating Functions* section. After that, we present the explicit formulas for the mean and variance of the distributions, for both fixed number of cycles and fixed sequences length.

After that, we generalize the results to the situation where nucleotide incorporation depends on sequence context. The main results are summarized in Eqs. (6), (10), (16), (17)), and (18).

1.1 Single-molecule DNA sequencing

Single-molecule DNA sequencing (SMS) technology becomes commercially available only in 2008 (Harris *et al.*, 2008), but major efforts are underway to develop different kinds of SMS technologies and it is expected that more SMS platforms will be available in the next few years (Gupta, 2008). Some of these SMS technologies employ SBS as the underlining principle.

One of the advantages of SMS compared with bulk sequencing technologies is that, unlike other bulk sequencing platforms, SMS does not have a clonal amplification step for the target sequences. Instead, a single target sequence molecule is used as the template. This avoids the problems of bias and errors introduced by PCR in the amplification step. Another advantage of SMS is that the problem of dephasing associated with bulk sequencing can be circumvented.

In the bulk sequencing, a significant copy numbers of the target sequence have to be present to obtain a detectable signal. During the sequencing process, the synchronization between identical individual templates will be lost gradually, leading to signal decay and sequencing errors. To avoid dephasing, the reaction reagents and synthesis chemistry are usually tuned to drive the enzymatic incorporation to completion at

each cycle. As a consequence, the misincorporation rate will increase. For SMS, since each molecule is monitored individually, the dephasing problem does not exist. The reaction kinetics can be controlled to adjust the rate of incorporation to the benefit of sequencing accuracy. For example, slow reaction kinetics can be used to limit incorporation to two or three bases (Harris *et al.*, 2008). For a homopolymer region GGG, the bulk sequencing will try to incorporate three C's in one cycle; for SMS, however, zero, one, two, or three C's can be incorporated in a single synthesis cycle. This flexibility in incorporation rate can be utilized to increase the resolution of homopolymer region (Harris *et al.*, 2008).

As a result, for SMS the consecutive zeros in the signal (where there is no incorporation) is not limited to 3 as in the pyrosequencing case. Instead, if the nucleotide complementary to the template is not incorporated in the current cycle, it still has chance to be incorporated in the next cycle, or the next-next cycle, without detrimental effect on sequencing quality. The relation between the read length and flow cycle in SMS will also be different from that of traditional pyrosequencing. In the following we will derive an analytical expression for the distribution, which will depend on nucleotide composition of the target sequences as well as the incorporation probability of each nucleotide.

1.2 Notation and definitions

To avoid the unnecessary specification of the detailed names of the four kinds of nucleotides, in the following we will use a , b , c , and d to represent any permutations of the usual nucleotides A , C , G , and T , as we did previously (Kong, 2008). Throughout the paper we assume that the nucleotides in the target sequence are independent of each other. Note that this is different from the sequence-context dependence of *nucleotide incorporation*, which will be addressed in Section 4. The probabilities for the four nucleotides in the target sequence are denoted as p_a , p_b , p_c , and p_d .

The nucleotide incorporation probabilities are denoted as α_{ij} , where $i = a, b, c$, and d stands for the kind of nucleotides, $j = 0, 1, 2, \dots$ denotes the cycle number, with the current cycle as 0. For example, when it is complementary to the template, if the chance for nucleotide b to be incorporated in the current cycle, the next cycle, and the next-next cycle is $1/3$, $1/2$, and $1/6$, respectively, then $\alpha_{b0} = 1/3$, $\alpha_{b1} = 1/2$, $\alpha_{b2} = 1/6$, and $\alpha_{bj} = 0$ for $j > 2$. Obviously, the pyrosequencing discussed previously is a special case in this notation with $\alpha_{i0} = 1$ and $\alpha_{ij} = 0$ for $j > 0$. By definition $\sum_{j=0}^{\infty} \alpha_{ij} = 1$. The generating functions (GFs) of the nucleotide incorporation probabilities $g_i(x)$ ($i = a, b, c$, and d) and their associated elementary symmetric functions (ESFs) $t_i(x)$ ($i = 1, 2, 3$, and 4) will be defined in the next section (Eqs. (4) and (8)). The α_{ij} will be generalized to be sequence-context dependent in Section 4.

We also use the same definition of nucleotide *flow cycle number* as before (Kong, 2008, Table 1), which is defined as the “quad cycle” of successive four nucleotides

Table 1: The relation between nucleotide flow and cycle number.

cycle number	\dots	f	$f + 1$	$f + 2$	$f + 3$	\dots
nucleotide flow	\dots	a b c d	a b c d	a b c d	a b c d	\dots

$\{abcd\}$. The cycle number is denoted as f in the following. We will use n for the length of a sequence. The relation between nucleotide flow and cycle number is illustrated in Table 1.

To extract coefficients from the expansion of GFs, we use the notation $[x^n]f(x)$ to denote the coefficient of x^n in the series of $f(x)$ in powers of x . Similarly, we use $[x^n y^m]f(x, y)$ to denote the coefficient of $x^n y^m$ in the bivariate $f(x, y)$.

2 Bivariate Generating Functions

In this section we first set up a detailed example of the recurrence relations between probabilities of sequence length and flow cycle number, and then establish a set of equations for the general case. The set of equations of probabilities cannot be solved in closed forms. However, if we transform the equations of probabilities into their corresponding generating functions (GFs), then these GFs can be easily obtained.

2.1 Recurrences

Let $L_i(f, n)$, $i = a, b, c$, and d denote the probability (up to a normalization factor, see below) of sequences with a length of n that is synthesized in f flow cycles with the last nucleotide being i . First let's look at one example. In this example, we assume

that for each nucleotide only the first three incorporation probabilities are non-zero, i.e., $\alpha_{ij} = 0$ for $j > 2$. The following four recurrence relations can be established. It might be helpful to refer back to Table 1 for the understanding of these recurrences:

$$\begin{aligned}
L_a(f+3, n+1) &= p_a [L_a(f+3, n)\alpha_{a0} + L_a(f+2, n)\alpha_{a1} + L_a(f+1, n)\alpha_{a2} \\
&\quad + L_b(f+2, n)\alpha_{a0} + L_b(f+1, n)\alpha_{a1} + L_b(f, n)\alpha_{a2} \\
&\quad + L_c(f+2, n)\alpha_{a0} + L_c(f+1, n)\alpha_{a1} + L_c(f, n)\alpha_{a2} \\
&\quad + L_d(f+2, n)\alpha_{a0} + L_d(f+1, n)\alpha_{a1} + L_d(f, n)\alpha_{a2}], \tag{1a}
\end{aligned}$$

$$\begin{aligned}
L_b(f+3, n+1) &= p_b [L_a(f+3, n)\alpha_{b0} + L_a(f+2, n)\alpha_{b1} + L_a(f+1, n)\alpha_{b2} \\
&\quad + L_b(f+3, n)\alpha_{b0} + L_b(f+2, n)\alpha_{b1} + L_b(f+1, n)\alpha_{b2} \\
&\quad + L_c(f+2, n)\alpha_{b0} + L_c(f+1, n)\alpha_{b1} + L_c(f, n)\alpha_{b2} \\
&\quad + L_d(f+2, n)\alpha_{b0} + L_d(f+1, n)\alpha_{b1} + L_d(f, n)\alpha_{b2}], \tag{1b}
\end{aligned}$$

$$\begin{aligned}
L_c(f+3, n+1) &= p_c [L_a(f+3, n)\alpha_{c0} + L_a(f+2, n)\alpha_{c1} + L_a(f+1, n)\alpha_{c2} \\
&\quad + L_b(f+3, n)\alpha_{c0} + L_b(f+2, n)\alpha_{c1} + L_b(f+1, n)\alpha_{c2} \\
&\quad + L_c(f+3, n)\alpha_{c0} + L_c(f+2, n)\alpha_{c1} + L_c(f+1, n)\alpha_{c2} \\
&\quad + L_d(f+2, n)\alpha_{c0} + L_d(f+1, n)\alpha_{c1} + L_d(f, n)\alpha_{c2}], \tag{1c}
\end{aligned}$$

$$\begin{aligned}
L_d(f+2, n+1) &= p_d [L_a(f+2, n)\alpha_{d0} + L_a(f+1, n)\alpha_{d1} + L_a(f, n)\alpha_{d2} \\
&\quad + L_b(f+2, n)\alpha_{d0} + L_b(f+1, n)\alpha_{d1} + L_b(f, n)\alpha_{d2} \\
&\quad + L_c(f+2, n)\alpha_{d0} + L_c(f+1, n)\alpha_{d1} + L_c(f, n)\alpha_{d2} \\
&\quad + L_d(f+2, n)\alpha_{d0} + L_d(f+1, n)\alpha_{d1} + L_d(f, n)\alpha_{d2}]. \tag{1d}
\end{aligned}$$

In general, for arbitrary nucleotide incorporation probabilities α_{ij} , we can write down the following four recurrences:

$$L_a(f, n) = p_a \left\{ \sum_{j=0}^f L_a(f-j, n-1) \alpha_{aj} + \sum_{j=0}^{f-1} L_b(f-j-1, n-1) \alpha_{aj} \right. \\ \left. + \sum_{j=0}^{f-1} L_c(f-j-1, n-1) \alpha_{aj} + \sum_{j=0}^{f-1} L_d(f-j-1, n-1) \alpha_{aj} \right\}, \quad (2a)$$

$$L_b(f, n) = p_b \left\{ \sum_{j=0}^f L_a(f-j, n-1) \alpha_{bj} + \sum_{j=0}^f L_b(f-j, n-1) \alpha_{bj} \right. \\ \left. + \sum_{j=0}^{f-1} L_c(f-j-1, n-1) \alpha_{bj} + \sum_{j=0}^{f-1} L_d(f-j-1, n-1) \alpha_{bj} \right\}, \quad (2b)$$

$$L_c(f, n) = p_c \left\{ \sum_{j=0}^f L_a(f-j, n-1) \alpha_{cj} + \sum_{j=0}^f L_b(f-j, n-1) \alpha_{cj} \right. \\ \left. + \sum_{j=0}^f L_c(f-j, n-1) \alpha_{cj} + \sum_{j=0}^{f-1} L_d(f-j-1, n-1) \alpha_{cj} \right\}, \quad (2c)$$

$$L_d(f, n) = p_d \left\{ \sum_{j=0}^f L_a(f-j, n-1) \alpha_{dj} + \sum_{j=0}^f L_b(f-j, n-1) \alpha_{dj} \right. \\ \left. + \sum_{j=0}^f L_c(f-j, n-1) \alpha_{dj} + \sum_{j=0}^f L_d(f-j, n-1) \alpha_{dj} \right\}. \quad (2d)$$

The recurrences cannot be solved in closed forms. However, if we transform these recurrences into their corresponding GFs, then these GFs can be solved in compact forms. The bivariate GFs of $L_i(f, n)$ are defined as

$$G_i(x, y) = \sum_{n=1}^{\infty} \sum_{f=1}^{\infty} L_i(f, n) x^f y^n, \quad i = a, b, c, d. \quad (3)$$

We also define the GFs for nucleotide incorporation probabilities as

$$g_i(x) = p_i \sum_{j=0}^{\infty} \alpha_{ij} x^j \quad i = a, b, c, d. \quad (4)$$

Since $\sum_{j=0}^{\infty} \alpha_{ij} = 1$, we have

$$g_i(1) = p_i \quad i = a, b, c, d.$$

To transform the system of recurrence equations in Eq. (2) to a system of equations of GFs, we need to use the following identities:

$$\sum_{f=1}^{\infty} \sum_{n=1}^{\infty} \sum_{j=0}^{f-1} L_i(f-j, n-1) \alpha_{ij} x^f y^n = y [G_i(x, y) + G_{i0}(x)] g_i(x), \quad i = a, b, c, d$$

and

$$\sum_{f=1}^{\infty} \sum_{n=1}^{\infty} \sum_{j=0}^{f-2} L_i(f-j-1, n-1) \alpha_{ij} x^f y^n = xy [G_i(x, y) + G_{i0}(x)] g_i(x), \quad i = a, b, c, d$$

where

$$G_{i0}(x) = \sum_{f=1}^{\infty} L_i(f, 0) x^f, \quad i = a, b, c, d.$$

By definition, we have

$$G_{a0}(x) = x,$$

$$G_{b0}(x) = 0,$$

$$G_{c0}(x) = 0,$$

$$G_{d0}(x) = 0.$$

The system of equations of GFs after the transform is:

$$G_a = [G_a + x + xG_b + xG_c + xG_d]yg_a, \quad (5a)$$

$$G_b = [G_a + x + G_b + xG_c + xG_d]yg_b, \quad (5b)$$

$$G_c = [G_a + x + G_b + G_c + xG_d]yg_c, \quad (5c)$$

$$G_d = [G_a + x + G_b + G_c + G_d]yg_d \quad (5d)$$

which can be solved as

$$G_a(x, y) = \frac{g_a xy}{H} F, \quad (6a)$$

$$G_b(x, y) = \frac{g_b xy}{H} [1 - (g_c + g_d)(1 - x)y + g_c g_d (1 - x)^2 y^2], \quad (6b)$$

$$G_c(x, y) = \frac{g_c xy}{H} [1 - g_d(1 - x)y], \quad (6c)$$

$$G_d(x, y) = \frac{g_d xy}{H}, \quad (6d)$$

where

$$H = 1 - t_1 y + t_2(1 - x)y^2 - t_3(1 - x)^2 y^3 + t_4(1 - x)^3 y^4, \quad (7)$$

and

$$F = [1 - (g_b + g_c + g_d)(1 - x)y + (g_b g_c + g_b g_d + g_c g_d)(1 - x)^2 y^2 - g_b g_c g_d (1 - x)^3 y^3].$$

Here we use *elementary symmetric functions* (ESFs) $t_i(x)$ of the nucleotide incorporation probabilities GFs $g_i(x)$ to put the solutions of $G_i(x, y)$ into a more compact form.

These ESFs $t_i(x)$ of four variables are defined as:

$$\begin{aligned}
t_1(x) &= g_a + g_b + g_c + g_d, \\
t_2(x) &= g_a g_b + g_a g_c + g_a g_d + g_b g_c + g_b g_d + g_c g_d, \\
t_3(x) &= g_a g_b g_c + g_a g_b g_d + g_a g_c g_d + g_b g_c g_d, \\
t_4(x) &= g_a g_b g_c g_d.
\end{aligned} \tag{8}$$

If we put $x = 1$ into $t_i(x)$, we get back the ESFs s_i of nucleotide probabilities p_i in the target sequences, which we used in the previous work (Kong, 2008):

$$\begin{aligned}
s_1 = t_1(1) &= p_a + p_b + p_c + p_d = 1, \\
s_2 = t_2(1) &= p_a p_b + p_a p_c + p_a p_d + p_b p_c + p_b p_d + p_c p_d, \\
s_3 = t_3(1) &= p_a p_b p_c + p_a p_b p_d + p_a p_c p_d + p_b p_c p_d, \\
s_4 = t_4(1) &= p_a p_b p_c p_d.
\end{aligned} \tag{9}$$

For complete nucleotide incorporation, $g_i(x) = p_i$. In this case we have $t_i(x) = s_i$ as a constant, instead of a function of x .

We see from the expressions of $G_i(x, y)$ in Eq. (6) that they are identical in *forms* to the solutions in the traditional pyrosequencing where 100% incorporation is assumed (Kong, 2008), with ESFs $t_i(x)$ of the nucleotide incorporation probabilities GFs $g_j(x)$ replacing the ESFs s_i of nucleotide probabilities p_j , $i = 1, 2, 3, 4$ and $j = a, b, c, d$.

From the expressions of the GFs in Eq. (6) we can see that they are not symmetric with respect to the nucleotide incorporation probability GFs $g_i(x)$. If we only consider the nucleotide flows that end up in the same “quad cycle” (see Table 1), then we can

add the four GFs $G_i(x, y)$ together to obtain a symmetric GF

$$\begin{aligned} G(x, y) &= \sum_{n=1}^{\infty} \sum_{f=1}^{\infty} L(f, n) x^f y^n = G_a + G_b + G_c + G_d \\ &= \frac{xy}{H} [t_1 - t_2(1-x)y + t_3(1-x)^2 y^2 - t_4(1-x)^3 y^3]. \end{aligned} \quad (10)$$

Here $L(f, n)$ is the (unnormalized) probability of a sequence with a length of n which can be synthesized in f flow cycles.

2.2 Normalization factors

To treat $L_i(f, n)$ as real probabilities, they have to be normalized. As we did in the previous work, these normalization factors can be obtained by setting $x = 1$ and $y = 1$ in the GFs in Eq. (6) or Eq. (10). If we set $x = 1$, then we obtained the normalization factors for $L_a(f, n)$ when the sequence length is fixed at n . From Eq. (6) we get

$$G_i(1, y) = \frac{p_i y}{1 - y},$$

from which we obtain the normalization factors u_i for $L_i(f, n)$ as

$$u_i = \sum_{f=1}^{\infty} L_i(f, n) = [y^n] G_i(1, y) = p_i \quad (11)$$

The normalization factor for $L(f, n) = \sum_i L_i(f, n)$ is simply the sum of u_i :

$$u = u_a + u_b + u_c + u_d = 1.$$

For fixed cycle, by setting $y = 1$, the denominator of $G_i(x, 1)$ becomes

$$1 - t_1 + t_2(1-x) - t_3(1-x)^2 + t_4(1-x)^3.$$

Since $t_1(1) = 1$, we see that $x = 1$ is a root of the denominator of $G_i(x, 1)$, and it is the dominant part in the expansion of $G_i(x, 1)$:

$$G_i(x, 1) = \frac{\beta_{i1}}{1-x} + \dots$$

The coefficients of the expansion can be evaluated as

$$\beta_{i1} = (1-x)G_i(x, 1)|_{x=1} = \frac{p_i}{s_2 + t'_1(1)},$$

from which come the normalization factors for $L_i(f, n)$ when the number of flow cycles is fixed at f :

$$v_i = \sum_{n=1}^{\infty} L_a(f, n) = [x^f]G_i(x, 1) \approx \frac{p_i}{s_2 + t'_1(1)}. \quad (12)$$

Here $t'_1(x) = \partial t_1(x)/\partial x$ and

$$t'_1(1) = \left. \frac{\partial t_1(x)}{\partial x} \right|_{x=1}$$

stands for the value of the first derivative of $t_1(x)$ evaluated at $x = 1$.

The normalization factor for $L(f, n) = \sum_i L_i(f, n)$ when the number of flow cycles is fixed at f is

$$v = v_a + v_b + v_c + v_d \approx \frac{1}{s_2 + t'_1(1)}. \quad (13)$$

As we will see below, v is an important part in the expression of mean and variance for $L_i(f, n)$. When the normalization factor v is compared with that of the complete incorporation case (Kong, 2008), the only difference is the extra term of $t'_1(1)$.

2.3 Mean and variance

The availability of GFs makes it easy to derive the mean and variance for the distributions of $L_i(f, n)$. When the sequence length is fixed, the mean and variance are given by

$$\bar{f}(n) = \frac{1}{u} [y^n] \frac{\partial G(x, y)}{\partial x} \Big|_{x=1}, \quad (14a)$$

$$\sigma_{\bar{f}}^2(n) = \frac{1}{u} [y^n] \frac{\partial^2 G(x, y)}{\partial x^2} \Big|_{x=1} + \bar{f}(n) - \bar{f}^2(n). \quad (14b)$$

Similar formulas apply to the individual nucleotide flow GFs G_a , G_b , etc, with their corresponding normalization factors u_i as shown in Eq. (11).

When the number of flow cycles is fixed, the mean and variance are given by

$$\bar{n}(f) = \frac{1}{v} [x^f] \frac{\partial G(x, y)}{\partial y} \Big|_{y=1}, \quad (15a)$$

$$\sigma_{\bar{n}}^2(f) = \frac{1}{v} [x^f] \frac{\partial^2 G(x, y)}{\partial y^2} \Big|_{y=1} + \bar{n}(f) - \bar{n}^2(f). \quad (15b)$$

Similar formulas apply to the individual nucleotide flow GFs, with their corresponding normalization factors v_i shown in Eq. (12).

3 Distributions at fixed sequence length and fixed number of flow cycles

With the help of the expression of GF $G(x, y)$ in Eq. (10), we can calculate *exact* distribution of $L(f, n)$ at a fixed sequence length n or a fixed number of cycles f , by expanding $G(x, y)$ and extracting x^f or y^n respectively. By using formulas Eqs. (14)

and (15) discussed above, we can calculate the mean and variance of these distributions. We will discuss the two cases separately in the following.

3.1 Fixed sequence length: distribution of flow cycles

When the length the target sequences is fixed at n , the mean $\bar{f}(n)$ and variance $\sigma_f^2(n)$ of the number of flow cycles f that is needed to determine the sequences are calculated from Eq. (14) as

$$\bar{f}(n) = [s_2 + t'_1]n - s_2 + 1, \quad (16a)$$

$$\begin{aligned} \sigma_f^2(n) = [s_2 - 3s_2^2 + 2s_3 + t'_1 + 2t'_2 + t''_1 - 4t'_1s_2 - t_1^2]n \\ + [5s_2^2 - s_2 - 4s_3 - 2t'_2 + 4t'_1s_2]. \end{aligned} \quad (16b)$$

Here and in the following we use the abbreviations such as t'_1 to denote $t'_1(1)$, t''_1 as $t''_1(1) = \partial^2 t_1(x)/\partial x^2|_{x=1}$, etc. They are the derivatives of $t_i(x)$ with respect to x evaluated at the value of $x = 1$.

From Eqs. (16a) and (16b) we can see that both the mean $\bar{f}(n)$ and the variance $\sigma_f^2(n)$ of flow cycles increase linearly with the sequence length n . Compared with the complete incorporation situation (Kong, 2008), the difference is the extra terms of the derivatives of $t_i(x)$ evaluated at the value of $x = 1$. When the nucleotide incorporation is complete at each cycle, $t_i(x) = p_i$ so all these derivatives vanish, leading back to the original results of pyrosequencing (Kong, 2008).

In Figure 1 the exact distributions of flow cycles are shown for a fixed sequence length of $n = 100$ base pairs with an artificial example of unequal nucleotide compo-

sition probabilities p_i and nucleotide incorporation probabilities α_{ij} . The nucleotide probabilities used here are $p_a = 3/10 = 0.3$, $p_b = 1/5 = 0.2$, $p_c = 1/5 = 0.2$, and $p_d = 3/10 = 0.3$. The non-zero nucleotide incorporation probabilities α_{ij} are (from $j = 0$ and up) $\alpha_{aj} = [1/10, 1/5, 2/5, 3/10]$, $\alpha_{bj} = [3/10, 1/5, 1/10, 1/10, 1/10, 1/10, 1/10]$, $\alpha_{cj} = [3/10, 3/10, 3/10, 1/10]$, and $\alpha_{dj} = [2/5, 1/5, 1/5, 1/10, 1/10]$. To keep numerical precision, we used exact calculation throughout by using either integers or exact fractions for all the coefficients in the expansion of the GFs, as we did in the previous work (Kong, 2008). The expansion was done using PARI/GP, a computer algebra system (The PARI Group, 2008).

Also shown in Figure 1 in continuous curve is the normal distribution $N(\bar{f}(n), \sigma_f^2(n))$, the mean $\bar{f}(n)$ and variance $\sigma_f^2(n)$ of which are calculated from Eqs. (16a) and (16b). It is evident that, as in the complete nucleotide incorporation situation, the exact distributions can be approximated quite accurately by normal distributions with the same mean and variance. For our example here, the normal distribution is $N(201.63, 211.5197)$. Compared with the complete nucleotide incorporation case, the incomplete incorporation significantly increases the number of flow cycles needed to cover the sequences of the same length. At the same time, the variance also increases significantly. For complete nucleotide incorporation with the same nucleotide composition probabilities p_i in the target sequences and at the same fixed sequence length $n = 100$, the mean and variance of flow cycle are 37.63 and 8.0045, respectively (Kong, 2008).

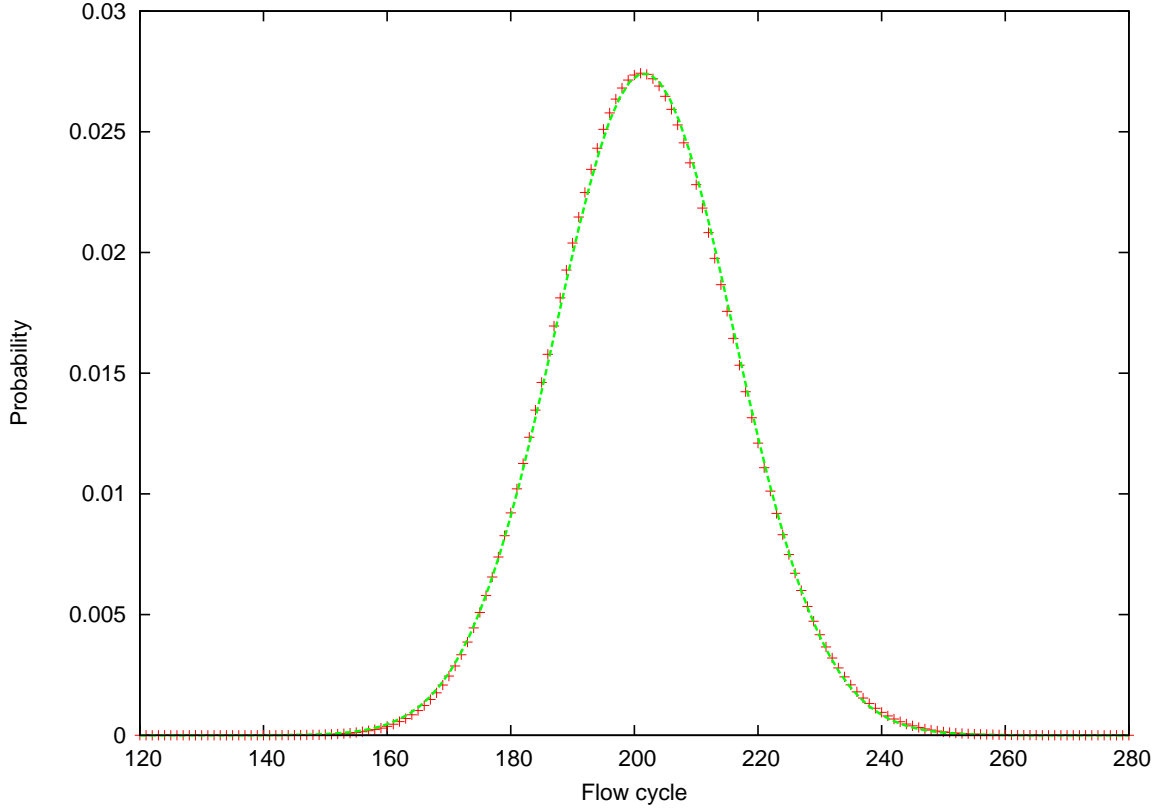


Figure 1: The distribution of flow cycles for a fixed sequence length of $n = 100$ base pairs. The nucleotide composition probabilities used here are $p_a = 3/10 = 0.3$, $p_b = 1/5 = 0.2$, $p_c = 1/5 = 0.2$, and $p_d = 3/10 = 0.3$. The non-zero nucleotide incorporation probabilities are $\alpha_{aj} = [1/10, 1/5, 2/5, 3/10]$, $\alpha_{bj} = [3/10, 1/5, 1/10, 1/10, 1/10, 1/10, 1/10]$, $\alpha_{cj} = [3/10, 3/10, 3/10, 1/10]$, and $\alpha_{dj} = [2/5, 1/5, 1/5, 1/10, 1/10]$. The exact distribution is plotted as '+' and is calculated from Eq. (10). The continuous curve is the normal distribution $N(\bar{f}(n), \sigma_f^2(n))$ of the same mean and variance as those of the exact distribution, where $\bar{f}(n)$ and $\sigma_f^2(n)$ are calculated from Eqs. (16a) and (16b). The normal distribution shown here is $N(201.63, 211.5197)$.

3.2 Fixed flow cycle: distribution of sequence length

When the number of flow cycles f is fixed, the mean $\bar{n}(f)$ and variance $\sigma_n^2(f)$ of the length of the sequences that can be determined by these flow cycles are calculated by Eqs (15a) and (15b) as:

$$\bar{n}(f) \approx vf - v^2[2s_2^2 - 2s_3 + 3t'_1s_2 + t_1'' - 2t_2'], \quad (17a)$$

$$\begin{aligned} \sigma_n^2(f) \approx v^3wf - v^4[2s_2(s_2s_3 + 3s_4) - 2s_3t'_1(3t'_1 + 2s_2) + 6t'_1s_4 - 2(2t'_2 + t_1'' + 2s_3)^2 \\ + (3t_2'' + t_1''' + 6t'_3 + t_1''(5s_2 + t'_1) - 2t_2'(t'_1 - 3s_2))/v + t_1's_2/v^2], \quad (17b) \end{aligned}$$

where

$$w = s_2 - 3s_2^2 + 2s_3 + t'_1 + 2t'_2 + t_1'' - 4t'_1s_2 - t_1'^2.$$

Eqs. (17a) and (17b) shows that both the average sequence length $\bar{n}(f)$ and the variance $\sigma_n^2(f)$ increase linearly with the number of flow cycle f . As discussed in section 2.2, the small extra terms are ignored in the above expressions.

In Figure 2 the exact distribution of sequence length in base pairs is shown for a fixed number of flow cycles $f = 50$ (200 nucleotide flows). The nucleotide composition probabilities p_i and nucleotide incorporation probabilities α_{ij} used here are the same as in the previous section. These exact distribution is calculated from Eq. (10) in section 2.

Also shown in Figure 2 in continuous curve is the normal distribution $N(\bar{n}(f), \sigma_n^2(f))$, the mean $\bar{n}(f)$ and variance $\sigma_n^2(f)$ of which are calculated from Eqs. (17a) and (17b).

Just like the distributions of the number of flow cycles at fixed sequence length as

discussed in the previous section, the exact distributions of sequence length at a fixed number of flow cycles can also be approximated well with normal distributions with the same mean and variance as those of the exact distributions. For the example here, the normal distributions is $N(25.0856, 13.1454)$. The introduction of incomplete nucleotide incorporation significantly reduces the mean and variance of the read length of sequences that can be determined at a given number of flow cycles. For complete nucleotide incorporation, with all other parameters being the same as above, the mean and variance of sequence length for this example would be 134.0117 and 78.5114, respectively (Kong, 2008).

Compared with the fit between the exact and normal distributions in Figure 1, the curve of normal distribution in Figure 2 shows some disagreements with the exact distribution, with the exact distribution having a slightly longer tails on the right and a slightly shorter tail on the left when compared to the normal distribution. The discrepancy is similar to that found in the complete nucleotide incorporation situation (Kong, 2008, Figure 2).

4 Generalization to sequence context dependent incorporation

In the previous discussions we assume that the nucleotide incorporation does not depend on sequence context. This may be only a first order approximation to the real

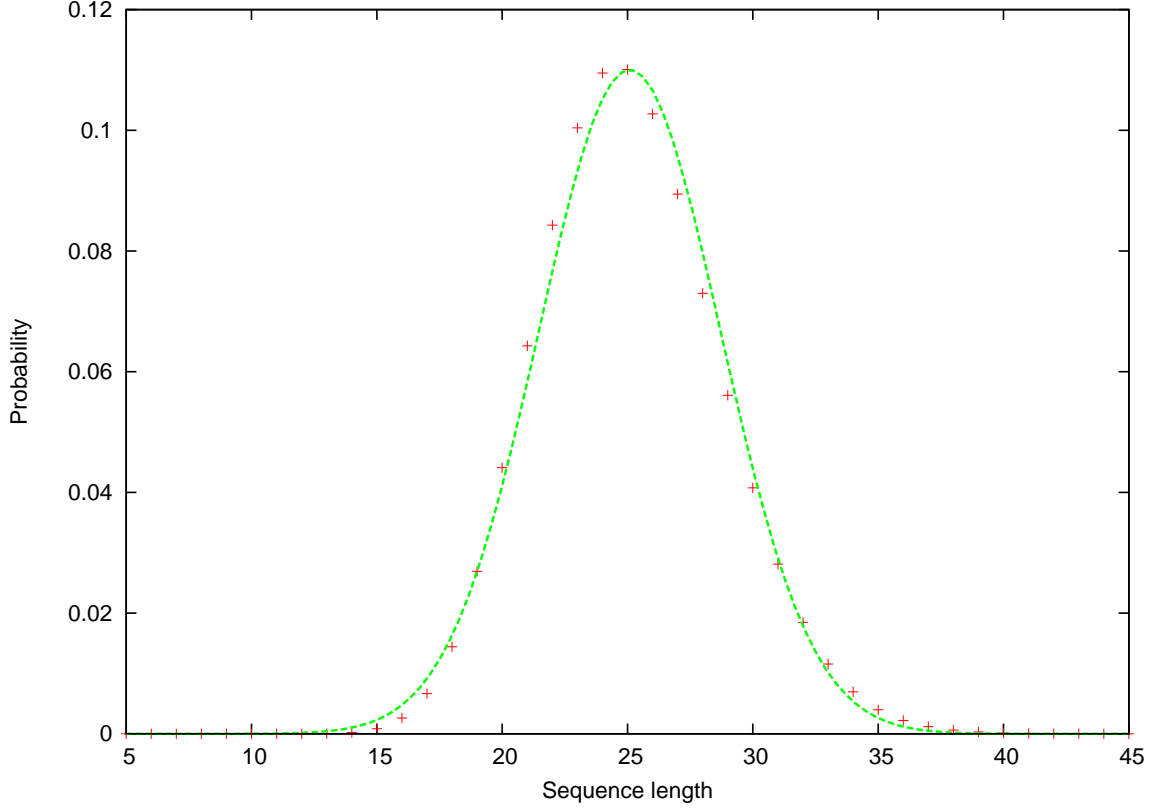


Figure 2: The distributions of sequence length in base pairs for a fixed number of flow cycles $f = 50$. The nucleotide composition probabilities p_i and nucleotide incorporation probabilities α_{ij} used here are the same as in Figure 1. The exact distribution is plotted as '+' and is calculated from Eq. (10). The continuous curve is the normal distribution $N(\bar{n}(f), \sigma_n^2(f))$ with the same mean and variance as those of the exact distribution, where $\bar{n}(f)$ and $\sigma_n^2(f)$ are calculated from Eqs. (17a) and (17b). The normal distribution shown here is $N(25.0856, 13.1454)$.

situation (Harris *et al.*, 2008). The results in previous sections, however, can be generalized to take into account sequence context dependent nucleotide incorporation. Instead of using α_{ij} , which describes the probability of nucleotide i being incorporated in the j th cycle, we can introduce α_{kij} , which is the probability of nucleotide i being incorporated in the j th cycle if the previous incorporated nucleotide is k . Correspondingly, the nucleotide incorporation GFs will become

$$g_{ki}(x) = p_i \sum_{j=0}^{\infty} \alpha_{kij} x^j, \quad i, k = a, b, c, d.$$

Instead of four $g_i(x)$, we now have 16 $g_{ki}(x)$. The system of equations of GFs in Eq. (5) will become

$$G_a = \left[(G_a + x)g_{aa} + xg_{ba}G_b + xg_{ca}G_c + xg_{da}G_d \right] y, \quad (18a)$$

$$G_b = \left[(G_a + x)g_{ab} + g_{bb}G_b + xg_{cb}G_c + xg_{db}G_d \right] y, \quad (18b)$$

$$G_c = \left[(G_a + x)g_{ac} + g_{bc}G_b + g_{cc}G_c + xg_{dc}G_d \right] y, \quad (18c)$$

$$G_d = \left[(G_a + x)g_{ad} + g_{bd}G_b + g_{cd}G_c + g_{dd}G_d \right] y. \quad (18d)$$

The GFs $G_i(x, y)$ can be solved in terms of $g_{ki}(x)$. The solution of $G = G_a + G_b + G_c + G_d$ is in the same form as Eq. (10):

$$G(x, y) = \frac{xy [\mathcal{S}_1 - \mathcal{S}_2 y + \mathcal{S}_3 y^2 - \mathcal{S}_4 y^3]}{1 - \mathcal{T}_1 y + \mathcal{T}_2 y^2 - \mathcal{T}_3 y^3 + \mathcal{T}_4 y^4}$$

where $\mathcal{S}_i(x)$ and $\mathcal{T}_i(x)$ are functions of $g_{ki}(x)$, and $\mathcal{S}_4 = \mathcal{T}_4$.

For a particular set of nucleotide incorporation probabilities α_{kij} (and hence the 16 nucleotide incorporation GFs g_{ki}) Eq. (18) can be used to solve for $G_i(x, y)$ for the

exact distributions, and Eqs. (14) and (15) can be used to calculate the mean and variance for the approximate normal distributions. The explicit expressions of mean and variance like those in Eqs. (16) and (17), however, seem difficult to obtain in compact forms for the sequence context dependent incorporation.

5 Discussion

In this paper we derived the statistical distributions for SBS with probabilistic nucleotide incorporation. The solutions are generalizations of the results obtained previously for pyrosequencing, where nucleotide incorporation is assumed to be complete for each flow cycle. Exact distributions can be obtained from the GFs (Eq. (6)), and these exact distributions can be approximated by normal distributions, the mean and variance of which are calculated from the explicit formulas derived from GFs (Eqs. (16) and (17)). The exact distributions can also be obtained when the nucleotide incorporation is sequence context dependent (Eq. (18)).

Probabilistic, or incomplete nucleotide incorporation, although a thing to avoid in traditional bulk sequencing, bring benefits for SMS, such as increased accuracy for incorporation and higher resolution for homopolymer regions. In SMS each template molecule is monitored individually, which makes it possible to bypass the phasing problem faced by the bulk sequencing technologies. The potentials of higher throughput and lower cost make it possible that SMS technologies will become a major biological and biomedicine research tool in the near future. The statistical distribution derived

here will be useful for instrument and software development of the next-generation sequencing platforms, including the SMS technologies.

Acknowledgment

This work was supported by Yale School of Medicine.

References

- Gupta, P. K., 2008. Single-molecule DNA sequencing technologies for future genomics research. *Trends in Biotechnology* In Press, Corrected Proof. URL <http://dx.doi.org/10.1016/j.tibtech.2008.07.003>.
- Harris, T. D., Buzby, P. R., Babcock, H., Beer, E., Bowers, J., Braslavsky, I., Causey, M., Colonell, J., Dimeo, J., Efcavitch, W. J., Giladi, E., Gill, J., Healy, J., Jarosz, M., Lapen, D., Moulton, K., Quake, S. R., Steinmann, K., Thayer, E., Tyurina, A., Ward, R., Weiss, H., and Xie, Z., 2008. Single-molecule dna sequencing of a viral genome. *Science* 320, 106–109.
- Kong, Y., 2008. Statistical distributions of pyrosequencing. *Journal of Computational Biology* In Press.
- The PARI Group, 2008. *PARI/GP, version 2.3.4*. Bordeaux. Available from <http://pari.math.u-bordeaux.fr/>.