

On Pattern Frequency Occurrences in a Markovian Sequence¹

M. Régnier² and W. Szpankowski³

Abstract. Consider a given pattern H and a random text T generated by a Markovian source. We study the frequency of pattern occurrences in a random text when overlapping copies of the pattern are counted separately. We present exact and asymptotic formulae for moments (including the variance), and probability of r pattern occurrences for three different regions of r , namely: (i) $r = O(1)$, (ii) central limit regime, and (iii) large deviations regime. In order to derive these results, we first construct certain language expressions that characterize pattern occurrences which are later translated into generating functions. We then use analytical methods to extract asymptotic behaviors of the pattern frequency from the generating functions. These findings are of particular interest to molecular biology problems (e.g., finding patterns with unexpectedly high or low frequencies, and gene recognition), information theory (e.g., second-order properties of the relative frequency), and pattern matching algorithms (e.g., q -gram algorithms).

Key Words. Frequency of pattern occurrences, Markov source, Autocorrelation polynomials, Languages, Generating functions, Asymptotic analysis, Large deviations.

1. Introduction. Repeated patterns and related phenomena in words (also called sequences or strings) are known to play a central role in many facets of computer science, telecommunications, and molecular biology. One of the most fundamental questions arising in such studies is the frequency of pattern occurrences in another string known as the *text*. Applications of these results include wireless communications (see [1]), approximate pattern matching (see [23] and [37]), molecular biology (see [32]), code synchronization (see [18]–[20]), and source coding (see [8]). In fact, this work and the one by Fudos et al. [14] were motivated by problems arising in approximate pattern matching by q -grams (see [23] and [37]), developing performance models for database systems in wireless communications (see [1]), and gene recognition in a DNA sequence (see [32]), respectively. Actually, one of the earliest applications appears to be in code synchronization (see [18]).

We study the problem in a probabilistic framework in which the text is generated randomly either by a memoryless source (the so-called *Bernoulli model*) or by a Markovian source (the so-called *Markovian model*). In the former, every symbol of a finite

¹ This paper was presented in part at the 1997 International Symposium on Information Theory, Ulm, Germany. This research was supported by NATO Collaborative Grant CRG.950060. Part of this work was done during authors visits at Purdue University and at INRIA, Rocquencourt. The first author was additionally supported by ESPRIT LTR Project No. 20244 (ALCOM-IT) and GREG “Motifs dans les Sequences.” The second author was additionally supported by NSF Grants CCR-9201078, NCR-9206315, and NCR-9415491.

² INRIA, Rocquencourt, 78153 Le Chesnay Cedex, France. Mireille.Regnier@inria.fr.

³ Department of Computer Science, Purdue University, West Lafayette, IN 47907, USA. spa@cs.purdue.edu.

alphabet is created independently of the other symbols, and the probabilities of symbol generation are not the same (if all probabilities of symbol generation are the same, the model is called the *symmetric* Bernoulli model). In the Markovian model, the next symbol depends on a finite number of previous symbols.

The problem of pattern occurrences in a random string is a classical one. Feller in 1968 already suggested a solution in his book [10]. Several other authors also contributed to this problem: e.g., see [3], [5], [22], [27], [31] and references therein. However, the most important recent contributions belong to Guibas and Odlyzko, who in a series of papers (see [18]–[20]) laid the foundations of the analysis for the symmetric Bernoulli model. In particular, the authors of [20] computed the moment generating function for the number of strings of length n that do *not* contain any one of a given set of patterns. Certainly, this suffices to estimate the probability of at least one pattern occurrence in a random string generated by the symmetric Bernoulli model. Furthermore, Guibas and Odlyzko [20] in a passing remark also presented some basic results for several pattern occurrences in a random text for the symmetric Bernoulli model, and for the probability of no occurrence of a given pattern in the asymmetric Bernoulli model. Recently, Fudos et al. [14] computed the probability of exactly r occurrences of a pattern in a random text in the *asymmetric* Bernoulli model, just directly extending the results of Guibas and Odlyzko. The Markovian model was tackled by Li [27], and Chrysaphinou and Papastavridis [5] who extended the Guibas and Odlyzko results of no pattern occurrence to Markovian texts. Prum et al. [33] (see also [36]) obtained the limiting distribution for the number of pattern occurrences in the Markovian model but without an explicit computation of the variance. Recently, Flajolet et al. [12] considered pattern occurrences in a random tree. Some other contributions are [3], [7], [16], [24], [25], [30], [32], and [39].

In this paper we provide a complete characterization of the frequency of pattern occurrences in a random text generated according either to the Bernoulli model or the Markovian model. Our method of analysis treats both models uniformly, and therefore we concentrate on discussing the Markovian model. Let O_n denote the number of occurrences of a given pattern H in a random text when *overlapping* copies of the pattern are counted separately. In Theorem 2.1 we present the generating function of O_n which can be used to compute exactly the probability of r pattern occurrences in the text. Furthermore, this allows for an easy computation of all moments, using for instance a symbolic computation system. In this paper we present explicit formulae for the mean and the variance of O_n . We observe that the evaluation of the variance was quite challenging as pointed out in [32] and [33]. It turns out that the variance depends on the internal structure of the pattern through the so-called *autocorrelation polynomial*. We should point out that Prum et al. [33] proposed two statistical methods to estimate the variance which should be compared with our computations (see Theorem 2.2 and Section 3).

We also estimate asymptotically the probability of exactly r occurrences of the pattern for three different ranges of r (see Theorem 2.2); namely, (i) $r = O(1)$, (ii) $r = EO_n + x\sqrt{n}$ for $x = O(1)$ (i.e., central limit regime), and (iii) $r = (1 + \delta)EO_n$ (i.e., large deviations regime). For our results to hold we assume that $nP(H) \rightarrow \infty$ (see [16] for other regimes of $nP(H)$). However, for a *given* pattern H it is natural to assume that the length of the pattern is constant with respect to n (and we adopt this assumption throughout).

Our results should be of particular interest to molecular biology, pattern matching algorithms, and information theory (e.g., relative frequency, code synchronization, coding, etc.). Two problems of molecular biology can benefit from these results, namely: finding patterns with unexpectedly (high or low) frequencies (the so-called contrast words) [15], and recognizing genes by using statistical properties [11]. Statistical methods have been successfully used from the early '80s to extract information from sequences of DNA. In particular, identifying deviant short motifs, the frequency of which is either too high or too low, might point out unknown biological information (see [11] and others for the analysis of functions of contrast words in DNA texts). From this perspective, our results give estimates for the statistical significance of deviations of word occurrences from the expected values and allow a biologist to build a dictionary of contrast words in genetic texts. Recently, Coward [6] used our results in the search of exceptional patterns in the yeast genome.

Another biological problem for which our results might be useful is gene recognition. Most gene recognition techniques rely on the observation that the statistics of patterns (motifs/codon) occurrences in coding and noncoding regions are different. Our findings allow the estimation of the statistical significance of such differences, and the construction of the confidence interval for pattern occurrences.

These results can also be used to recognize statistical properties of various other information sources such as images, text, etc. In information theory, *relative frequency* defined as $\Delta_n = O_n/(n - m + 1)$, where m is the length of the pattern, is often used to assess statistics of information sources. It is well known [8], [29] that Δ_n converges almost surely to the probability $P(H)$ of the pattern H , but much less is known about second-order properties of Δ_n such as the limiting distribution, large deviations, and rate of convergence. The rate of convergence to the source entropy—which is related to the rate of convergence of the relative frequency [29]—has recently appeared in the formulation of some results on data compression (see [28], [38], and [41]). Marton and Shields [29] proved that Δ_n converges exponentially fast to $P(H)$ for sources satisfying the so-called *blow-up property* (e.g., Markov sources, hidden Markov, etc.). Our results characterize precisely such a convergence in the central limit regime and the large deviations regime for Markovian sources.

In the accompanying paper [34], we extended our results to *approximate* pattern occurrences or a set of pattern occurrences. Such extension is vital to some approximate pattern matching algorithms. Recently, Sutinen and Szpankowski [37] used these results for performance evaluation of q -gram filtration algorithms.

This paper is organized as follows. In the next section we present our main results and their consequences. The proofs are delayed until the last section. Our derivation in Section 3.1 use a language approach, thus is also valid for Markovian models since no probabilistic assumption is made. In Section 3.2 we translate language relationships into associated generating functions, and finally we use analytical tools in Section 3.3 to derive asymptotic results.

2. Main Results. We consider two strings, a pattern string $H = h_1 h_2 \cdots h_m$ and a text string $T = t_1 t_2 \cdots t_n$ of respective lengths equal to m and n over an alphabet \mathcal{S} of size V . We write $\mathcal{S} = \{1, 2, \dots, V\}$ to simplify the presentation. Throughout, we assume that

the pattern string is *fixed* and given, while the text string is random. More precisely, we consider the following two probabilistic models of text generation:

(B) **BERNOULLI MODEL.** The text is a realization of a sequence of independently, identically distributed (i.i.d.) random variables, such that a symbol $s \in \mathcal{S}$ occurs with probability $P(s)$.

(M) **MARKOVIAN MODEL.** The text is a realization of a *stationary* Markov sequence of order K , that is, the probability of the next symbol occurrence depends on the K previous symbols. In most derivations we deal only with the first-order Markov chain ($K = 1$), and then we define the transition matrix $\mathbf{P} = \{p_{i,j}\}_{i,j \in \mathcal{S}}$ where $p_{i,j} = \Pr\{t_{k+1} = j | t_k = i\}$. By $\boldsymbol{\pi} = (\pi_1, \dots, \pi_V)$ we denote the stationary distribution satisfying $\boldsymbol{\pi}\mathbf{P} = \boldsymbol{\pi}$, and we write Π for the stationary matrix that consists of V identical rows equal to $\boldsymbol{\pi}$. Finally, by \mathbf{Z} we denote the **fundamental matrix** $\mathbf{Z} = (\mathbf{I} - (\mathbf{P} - \Pi))^{-1}$ where \mathbf{I} is the identity matrix.

Throughout the paper we use systematically the following notation: lowercase bold letters represents row vectors (e.g., $\boldsymbol{\pi}$), while uppercase bold letters denote matrices (e.g., Π). To extract a particular element, say with index (i, j) , from a matrix, say \mathbf{P} , we write $[\mathbf{P}]_{i,j} = p_{i,j}$. Finally, we recall that $(\mathbf{I} - \mathbf{P})^{-1} = \sum_{k \geq 0} \mathbf{P}^k$ provided the inverse matrix exists (i.e., $\det(\mathbf{I} - \mathbf{P}) \neq 0$ or $\|\mathbf{P}\| < 1$ for any matrix norm $\|\cdot\|$). Below, we write $P(H_i^j) = \Pr\{T_{i+k}^{j+k} = H_i^j\}$ for the probability of an occurrence of substring $H_i^j = h_i \cdots h_j$ in the random text T_{i+k}^{j+k} between symbols $i+k$ and $j+k$ for any k (in particular, $P(H)$ denotes the probability of H appearing in the text).

Our goal is to estimate the frequency of multiple pattern occurrences in the text assuming either the Bernoulli or the Markovian model. We find it convenient and useful to express our findings in terms of languages. A language \mathcal{L} is a collection of words satisfying some properties. We associate with a language \mathcal{L} a generating function defined below.

DEFINITION 1. For any language \mathcal{L} we define its generating function $L(z)$ as

$$(1) \quad L(z) = \sum_{w \in \mathcal{L}} P(w)z^{|w|},$$

where $P(w)$ is the stationary probability of word w occurrence, $|w|$ is the length of w , and we adopt a usual convention that $P(\varepsilon) = 1$, where ε is the empty word. In addition, we define the H -conditional generating function of \mathcal{L} as

$$(2) \quad L_H(z) = \sum_{w \in \mathcal{L}} P(w|w_{-m} = h_1 \cdots w_{-1} = h_m)z^{|w|} = \sum_{w \in \mathcal{L}} P\{w|w_{-m}^{-1} = H\}z^{|w|},$$

where w_{-i} stands for a symbol preceding the first character of w at distance i .

It turns out that several properties of pattern occurrences depend on the so-called *autocorrelation polynomial* that we define next:

DEFINITION 2. Given a string H , we define the *autocorrelation set* \mathcal{A} as

$$(3) \quad \mathcal{A} = \{H_{k+1}^m : H_1^k = H_{m-k+1}^m\},$$

and by HH we denote the set of positions k satisfying $H_1^k = H_{m-k+1}^m$. The generating function of language \mathcal{A} is denoted as $A(z)$ and we call it the *autocorrelation polynomial*. Its H -conditional generating function is denoted $A_H(z)$. In particular,

$$(4) \quad A_H(z) = \sum_{k \in HH} P(H_{k+1}^m | H_k^k) z^{m-k}$$

for a Markov chain of order $K = 1$.

Before we proceed, we present a simple example illustrating the definitions we have introduced so far.

EXAMPLE (Autocorrelation Functions). Assume that $H = 101$ over a binary alphabet $\mathcal{S} = \{0, 1\}$. Observe that $HH = \{1, 3\}$ and $\mathcal{A} = \{\varepsilon, 01\}$, where ε is the empty word. Thus, for the symmetric Bernoulli model (both symbols occur with the same probability equal to 0.5) we have $A(z) = 1 + z^2/4$, while for the Markovian model of order 1, we obtain $A_{101}(z) = 1 + p_{10}p_{01}z^2$.

We can now formulate our main results. In what follows, we denote by $O_n(H)$ (or simply by O_n) a random variable representing the number of occurrences of H in a random text T of size n . We introduce the generating function of the language \mathcal{T}_r of words that contain exactly r occurrences of H , namely, $T^{(r)}(z) = \sum_{n \geq 0} \Pr\{O_n(H) = r\} z^n$. We also define a bivariate generating function as follows:

$$(5) \quad T(z, u) = \sum_{r=1}^{\infty} T^{(r)}(z) u^r = \sum_{r=1}^{\infty} \sum_{n=0}^{\infty} \Pr\{O_n(H) = r\} z^n u^r$$

for $|z| \leq 1$ and $|u| \leq 1$.

Our main results are summarized in the following two theorems. The first theorem presents exact formulae for the generating functions $T^{(r)}(z)$ and $T(z, u)$, and can be used to compute exactly all parameters related to the pattern occurrence $O_n(H)$. In the second theorem, we provide asymptotic formulae for $\Pr\{O_n(H) = r\}$ for three regimes of r , namely: (i) $r = O(1)$, (ii) $r = EO_n + x\sqrt{\text{Var } O_n}$ when $x = O(1)$ (i.e., local central limit), (iii) $r = (1 + \delta)EO_n$ for some δ (i.e., large deviations). All proofs are presented in the next section: Section 3.2 contains the proof of Theorem 2.1 while the proof of Theorem 2.2 can be found in Section 3.3.

THEOREM 2.1. Let H be a given pattern of size m , and let T be a random text of length n generated according to a stationary Markov chain of order 1 over a V -ary alphabet \mathcal{S} . The generating functions $T^{(r)}(z)$ and $T(z, u)$ become

$$(6) \quad T^{(r)}(z) = R(z)M_H^{r-1}(z)U_H(z), \quad r \geq 1,$$

$$(7) \quad T(z, u) = R(z) \frac{u}{1 - uM_H(z)} U_H(z),$$

where

$$(8) \quad M_H(z) = 1 + \frac{z - 1}{D_H(z)},$$

$$(9) \quad U_H(z) = \frac{1}{D_H(z)},$$

$$(10) \quad R(z) = z^m P(H) \frac{1}{D_H(z)},$$

with

$$(11) \quad D_H(z) = (1 - z)A_H(z) + z^m P(H)(1 + (1 - z)F(z)).$$

The function $F(z)$ is defined for $|z| \leq R$ where $R = 1/\|\mathbf{P} - \Pi\|$ as follows:

$$(12) \quad F(z) = \frac{1}{\pi_{h_1}} [(\mathbf{P} - \Pi)(\mathbf{I} - (\mathbf{P} - \Pi)z)^{-1}]_{h_m, h_1},$$

where h_1 and h_m are the first and last symbols of H , respectively. In the Bernoulli model, $F(z) = 0$, and hence

$$D_H(z) = (1 - z)A_H(z) + z^m P(H),$$

with the other formulae as above.

Theorem 2.1 is the starting point of our next finding that deals with asymptotics for $n \rightarrow \infty$ when $nP(H) \rightarrow \infty$. The results below are derived in Section 3.3 using analytical tools.

THEOREM 2.2. *Let the hypotheses of Theorem 2.1 be fulfilled and $nP(H) \rightarrow \infty$.*

(i) **Moments.** *The expectation $EO_n(H)$ satisfies, for $n \geq m$,*

$$(13) \quad EO_n(H) = P(H)(n - m + 1),$$

while the variance becomes, for some $r > 1$,

$$(14) \quad \text{Var } O_n(H) = nc_1 + c_2 + O(r^{-n}),$$

where

$$(15) \quad c_1 = P(H)(2A_H(1) - 1 - (2m - 1)P(H) + 2P(H)E_1),$$

$$(16) \quad c_2 = P(H)((m - 1)(3m - 1)P(H) - (m - 1)(2A_H(1) - 1) - 2A'_H(1)) - 2(2m - 1)P(H)^2E_1 + 2E_2P(H)^2,$$

and the constants E_1, E_2 are

$$(17) \quad E_1 = \frac{1}{\pi_{h_1}} [(\mathbf{P} - \Pi)\mathbf{Z}]_{h_m, h_1},$$

$$(18) \quad E_2 = \frac{1}{\pi_{h_1}} [(\mathbf{P}^2 - \Pi)\mathbf{Z}^2]_{h_m, h_1},$$

where $\mathbf{Z} = (\mathbf{I} - (\mathbf{P} - \Pi))^{-1}$ is the fundamental matrix of the underlying Markov chain. In the **Bernoulli model**, $E_1 = E_2 = 0$ since $\mathbf{P} = \Pi$, and (14) reduces to an equality for $n \geq 2m - 1$. Thus

$$(19) \quad \text{Var } O_n(H) = nc_1 + c_2,$$

with

$$\begin{aligned} c_1 &= P(H)(2A_H(1) - 1 - (2m - 1)P(H)), \\ c_2 &= P(H)((m - 1)(3m - 1)P(H) - (m - 1)(2A_H(1) - 1) - 2A'_H(1)). \end{aligned}$$

(ii) Distribution: Case $r = O(1)$. Let ρ_H be the root of $D_H(z) = 0$ of smallest modulus and multiplicity one. Then ρ_H is real positive and lies outside the unit circle $|z| < 1$, and there exists $\rho > \rho_H$ such that

$$(20) \quad \Pr\{O_n(H) = r\} = \sum_{j=1}^{r+1} (-1)^j a_j \binom{n}{j-1} \rho_H^{-(n+j)} + O(\rho^{-n}),$$

where

$$(21) \quad a_{r+1} = \frac{\rho_H^m P(H) (\rho_H - 1)^{r-1}}{(D'_H(\rho_H))^{r+1}},$$

and the remaining coefficients can be computed according to

$$(22) \quad a_j = \frac{1}{(r + 1 - j)!} \lim_{z \rightarrow \rho_H} \frac{d^{r+1-j}}{dz^{r+1-j}} (T^{(r)}(z)(z - \rho_H)^{r+1})$$

with $j = 1, 2, \dots, r$.

(iii) Central Limit Regime: Case $r = EO_n + x\sqrt{\text{Var } O_n}$. For $x = O(1)$ we have, as $n \rightarrow \infty$,

$$(23) \quad \Pr\{O_n(H) = r\} = \frac{1}{\sqrt{2\pi c_1 n}} e^{-x^2/2} \left(1 + O\left(\frac{1}{\sqrt{n}}\right) \right),$$

where c_1 is defined in (15) above.

(iv) Large Deviations: Case $r = (1 + \delta)EO_n$. Let $a = (1 + \delta)P(H)$ with $\delta > 0$. For complex t , define $\rho(t)$ to be the root of

$$(24) \quad 1 - e^t M_H(e^\rho) = 0,$$

while ω_a and σ_a are defined as

$$(25) \quad -\rho'(\omega_a) = a,$$

$$(26) \quad -\rho''(\omega_a) = \sigma_a^2.$$

Then

$$(27) \quad \Pr\{O_n(H) = (1 + \delta)EO_n\} = \frac{1}{\sigma_a \sqrt{2\pi(n - m + 1)}} e^{-(n-m+1)I(a)} \left(1 + O\left(\frac{1}{n}\right) \right),$$

where $I(a) = a\omega_a + \rho(\omega_a)$.

As mentioned before, the above results find applications in information theory and molecular biology. For example, *relative frequency* is an important concept in information theory, and it is defined as

$$\Delta_n(H) = \frac{O_n(H)}{n - m + 1}.$$

Relative frequency appears in the definition of types and typical types (see [8]), and is often used to estimate information source statistics. As an easy corollary to Theorem 2.2, we obtain the following second-order characterization of $\Delta_n(H)$:

COROLLARY 2.1. *Under the hypotheses of Theorem 2.2, the following hold:*

(i) [Central Limit Regime] For $x = O(1)$,

$$(28) \Pr \left\{ \Delta_n(H) = P(H) + x \sqrt{\frac{c_1}{n - m + 1}} \right\} = \frac{1}{\sqrt{2\pi c_1 n}} e^{-x^2/2} \left(1 + O\left(\frac{1}{\sqrt{n}}\right) \right).$$

(ii) [Large Deviations] For $a = (1 + \delta)P(H)$ with $\delta > 0$,

$$(29) \Pr\{\Delta_n(H) \geq (1 + \delta)P(H)\} \\ = \frac{1}{\sigma_a \sqrt{2\pi(n - m + 1)}(1 - e^{-I(a)})} e^{-(n-m+1)I(a)} \left(1 + O\left(\frac{1}{n}\right) \right),$$

where ω_a and $I(a)$ are as defined in Theorem 2.2(iii).

3. Analysis. The key element of our analysis is a derivation of the generating function $T(z, u)$ presented in Theorem 2.1. The first part of the discussion below is quite general and works uniformly for both the Bernoulli model and the Markovian model. It is based on constructing some special languages and finding relationships among them. Later, in Section 3.2 we translate these relations into formulae for generating functions.

3.1. Combinatorial Relationships on Certain Languages. A collection of words sharing a given property is commonly called a *language*. This section is devoted to presenting combinatorial relationships between some languages that help to derive some results in a uniform manner. In this section we do not make any probabilistic assumption.

We start with some definitions:

DEFINITION 3. Given a pattern H :

- (i) Let \mathcal{T} be a language of words containing at least one occurrence of H , and, for any integer $r \geq 1$, let \mathcal{T}_r be the language of words containing exactly r occurrences of H .
- (ii) We define \mathcal{R} as the set of words containing only one occurrence of H , located at the right end. We also define \mathcal{U} as

$$(30) \quad \mathcal{U} = \{u : H \cdot u \in \mathcal{T}_1\},$$

where the operation \cdot means concatenation of words. In other words, a word $u \in \mathcal{U}$ if $H \cdot u$ has exactly one occurrence of H at the left end of $H \cdot u$.

(iii) Let \mathcal{M} be the language:

$$\mathcal{M} = \{w : H \cdot w \in \mathcal{T}_2 \text{ and } H \text{ occurs at the right end of } H \cdot w\},$$

that is, \mathcal{M} is a language such that $H \cdot \mathcal{M}$ has exactly two occurrences of H at the left and right ends of a word from \mathcal{M} .

We can now describe languages \mathcal{T} and \mathcal{T}_r in terms of \mathcal{R} , \mathcal{M} , and \mathcal{U} . This will further lead to a simple formula for the generating function of $O_n(H)$.

THEOREM 3.1. *Language \mathcal{T} satisfies the fundamental equation*

$$(31) \quad \mathcal{T} = \mathcal{R} \cdot \mathcal{M}^* \cdot \mathcal{U}.$$

Notably, language \mathcal{T}_r can be represented for any $r \geq 1$ as follows:

$$(32) \quad \mathcal{T}_r = \mathcal{R} \cdot \mathcal{M}^{r-1} \cdot \mathcal{U}.$$

Here, by definition, $\mathcal{M}^0 := \{\varepsilon\}$ and $\mathcal{M}^ := \bigcup_{r=0}^{\infty} \mathcal{M}^r$.*

PROOF. We first prove (32) and obtain our decomposition of \mathcal{T}_r as follows: The first occurrence of H in a word belonging to \mathcal{T}_r determines a prefix p that is in \mathcal{R} . Then a nonempty word w that creates the second occurrence of H is concatenated. Hence, w is in \mathcal{M} . This process is repeated $r - 1$ times. Finally, after the last H occurrence a suffix u that does not create a new occurrence of H is added. Equivalently, Hu is such that u is in \mathcal{U} , and w is a proper subword of Hu . Finally, a word belongs to \mathcal{T} if, for some $1 \leq r < \infty$, it belongs to \mathcal{T}_r . The set union $\bigcup_{r=1}^{\infty} \mathcal{M}^{r-1}$ yields precisely \mathcal{M}^* . \square

We now prove the following result that summarizes relationships between the languages \mathcal{R} , \mathcal{M} and \mathcal{U} .

THEOREM 3.2. *The languages \mathcal{M} , \mathcal{U} , and \mathcal{R} satisfy*

$$(33) \quad \bigcup_{k \geq 1} \mathcal{M}^k = \mathcal{W} \cdot H + \mathcal{A} - \{\varepsilon\},$$

$$(34) \quad \mathcal{U} \cdot \mathcal{S} = \mathcal{M} + \mathcal{U} - \{\varepsilon\},$$

$$(35) \quad H \cdot \mathcal{M} = \mathcal{S} \cdot \mathcal{R} - (\mathcal{R} - H),$$

where \mathcal{W} is the set, of all words, \mathcal{S} is the alphabet set and $+$ and $-$ are disjoint union and subtraction of languages. In particular, a combination of (34) and (35) gives

$$(36) \quad H \cdot \mathcal{U} \cdot \mathcal{S} - H \cdot \mathcal{U} = (\mathcal{S} - \varepsilon)\mathcal{R}.$$

Additionally, we have

$$(37) \quad \mathcal{T}_0 \cdot H = \mathcal{R} \cdot \mathcal{A}.$$

PROOF. All the relations above are proved in a similar fashion. We first deal with (33). Let k be the number of H occurrences in $\mathcal{W} \cdot H$. By definition, $k \geq 1$ and the last occurrence is on the right: this implies that $\mathcal{W} \cdot H \subseteq \bigcup_{k \geq 1} \mathcal{M}^k$. Furthermore, a word w in $\bigcup_{k \geq 1} \mathcal{M}^k$ is not in $\mathcal{W} \cdot H$ iff its size $|w|$ is smaller than $|H|$. Then the second H occurrence in Hw overlaps with H , which means that w is in $\mathcal{A} - \varepsilon$.

We now turn to (34). When a character s is added immediately after a word u from \mathcal{U} , two cases may occur: either Hus still does not contain a second occurrence of H , which means that us is a nonempty word of \mathcal{U} , or a new H appears, clearly at the right end. Then us is in \mathcal{M} . Furthermore, the whole set $\mathcal{M} + (\mathcal{U} - \varepsilon)$ is attained, i.e., a strict prefix of \mathcal{M} cannot contain a new H occurrence. Hence, it is in \mathcal{U} , and a strict prefix of a \mathcal{U} -word is in \mathcal{U} .

We now prove (35). Let $x = sw$ be a word in $H \cdot \mathcal{M}$ where s is a symbol from \mathcal{S} . As x contains exactly two occurrences of H located at its left and right ends, w is in \mathcal{R} and x is in $\mathcal{S} \cdot \mathcal{R} - \mathcal{R}$. Reciprocally, if a word swH from $\mathcal{S} \cdot \mathcal{R}$ is not in \mathcal{R} , then swH contains a second H occurrence starting in sw . As wH is in \mathcal{R} , the only possible position is at the left end, and then x is in $H \cdot \mathcal{M}$. We now rewrite:

$$\mathcal{S} \cdot \mathcal{R} - \mathcal{R} = \mathcal{S} \cdot \mathcal{R} - (\mathcal{R} \cap \mathcal{S} \cdot \mathcal{R}) = \mathcal{S} \cdot \mathcal{R} - (\mathcal{R} - H),$$

which yields $H \cdot \mathcal{M} - H = (\mathcal{S} - \varepsilon) \cdot \mathcal{R}$.

Deriving (37) is only a little more intricate. Let t be some word in \mathcal{T}_0 . We consider the factorization $t = w_1w_2$ such that w_2 is the largest suffix that is also an $(m - k)$ -prefix of H , with $k \in HH$ and $m = |H|$. In other words, w_2 is the largest suffix satisfying the equation $w_2 \cdot H = H \cdot a$, where a is in \mathcal{A} . If w_1H were not in \mathcal{R} , a second occurrence of H would occur in w_1H starting in w_1 . As $w_1Ha = w_1w_2H$, this contradicts the maximal property of w_2 . Therefore, $\mathcal{T}_0 \cdot H \subseteq \mathcal{R} \cdot \mathcal{A}$. Finally, we consider a word w_1Ha in $\mathcal{R} \cdot \mathcal{A}$. We may rewrite it as $H \cdot a = w_2 \cdot H$. It suffices now to show that $w_1w_2 \in \mathcal{T}_0$. Indeed, since $|w_2| < |H|$, any occurrence of H would go across w_1 and w_1H would contain two occurrences of H , which contradicts the definition of \mathcal{R} . This proves $\mathcal{R} \cdot \mathcal{A} \subseteq \mathcal{T}_0 \cdot H$, and completes the proof of Theorem 3.2. \square

3.2. *Associated Generating Functions.* In the previous section we did not make any probabilistic assumption. Thus, Theorem 3.2 is true for any model, including the Bernoulli and Markovian ones. In this section we translate the language relationships into generating functions. Therefore, we need to return to our probabilistic assumptions. Most of our derivations deal with the Markovian model.

To transfer our language relations into generating functions, we need a few rules associated with two operations on languages: namely, disjoint union $+$ and concatenation \cdot become the sum operation and the multiplication operation on generating functions. We start with the following simply property that is true for both probabilistic models:

- (P1) Let \mathcal{L}_1 and \mathcal{L}_2 be two arbitrary languages with generating functions $L_1(z)$ and $L_2(z)$, respectively. Then the language $\mathcal{L} = \mathcal{L}_1 + \mathcal{L}_2$ is transferred into the generating function $L(z)$ such that

$$L(z) = L_1(z) + L_2(z).$$

To translate the concatenation operation, it is necessary to consider the Bernoulli and the Markovian models separately. We start with the **Bernoulli model**:

(P2) We now consider a new language \mathcal{L} which is constructed from the concatenation of two other languages, say \mathcal{L}_1 and \mathcal{L}_2 , that is, $\mathcal{L} = \mathcal{L}_1 \cdot \mathcal{L}_2$. In the *Bernoulli model* the generating function $L(z)$ of \mathcal{L} becomes

$$L(z) = L_1(z)L_2(z)$$

since $P(wv) = P(w)P(v)$ for $w \in \mathcal{L}_1$ and $v \in \mathcal{L}_2$. In particular, the generating function $L(z)$ of $\mathcal{L} = \mathcal{S} \cdot \mathcal{L}_1$ is $L(z) = zL_1(z)$, where \mathcal{S} is the alphabet set, since $S(z) = \sum_{p \in \mathcal{S}} P(s)z = z$.

In the **Markovian model** $P(wv) \neq P(w)P(v)$, thus property (P2) is no longer true. We have to replace it by a more sophisticated one. We have to condition \mathcal{L}_2 on symbols preceding a word from \mathcal{L}_2 (i.e., belonging to \mathcal{L}_1). In general, for a K -order Markov chain, one must distinguish V^K ending states for \mathcal{L}_1 and V^K initial states for \mathcal{L}_2 . For simplicity of presentation, we only consider first-order Markov chains (i.e., $K = 1$), and we write $\ell(w)$ for the last symbol of a word w . In particular, to rewrite property (P2) we must introduce the following conditional generating function for a language \mathcal{L} :

$$L_i^j(z) = \sum_{w \in \mathcal{L}} P(w, \ell(w) = j | w_1 = i) z^{|w|}.$$

Then for the Markovian model property (P2) becomes:

(P2) Let $\mathcal{L} = \mathcal{W} \cdot \mathcal{V}$. Then

$$(38) \quad L_k^l(z) = \sum_{i,j \in \mathcal{S}} p_{ji} W_k^j(z) V_i^l(z),$$

where $W_k^j(z)$ and $V_i^l(z)$ are conditional generating functions for \mathcal{W} and \mathcal{V} , respectively. To prove this, let $w \in \mathcal{W}$ and $v \in \mathcal{V}$. Observe that

$$\begin{aligned} P(wv) &= \sum_{j \in \mathcal{S}} P(wv, \ell(w) = j) \\ &= \sum_{j \in \mathcal{S}} P(w, \ell(w) = j) P(v | \ell(w) = j) \\ &= \sum_{j \in \mathcal{S}} \sum_{i \in \mathcal{S}} P(w, \ell(w) = j) p_{ji} P(v | v_1 = i). \end{aligned}$$

After conditioning on the first symbol of \mathcal{W} and the last symbol of \mathcal{V} , we prove (38).

In passing, we observe that in the Markovian model of our problem one actually must deal only with two kinds of words: We have words w for which no assumption is made on the preceding words (e.g., these are the words in language \mathcal{R} with generating function $R(z)$); and we deal with words for which the preceding word admits H as a suffix (e.g., words in \mathcal{U} and \mathcal{M} whose H -conditional generating functions are $U_H(z)$ and $M_H(z)$, respectively).

The lemma below together with Theorem 3.1 proves (6) and (7) of Theorem 2.1.

LEMMA 3.1. *The generating functions associated with languages \mathcal{M}, \mathcal{U} , and \mathcal{R} satisfy*

$$(39) \quad \frac{1}{1 - M_H(z)} = A_H(z) + P(H)z^m \left(\frac{1}{1 - z} + F(z) \right),$$

$$(40) \quad U_H(z) = \frac{M_H(z) - 1}{z - 1},$$

$$(41) \quad R(z) = P(H)z^m \cdot U_H(z),$$

provided the underlying Markov chain is stationary.

PROOF. We first prove (40). Interestingly, it does not need the stationarity assumption. We consider the language relationship (34) from Theorem 3.2 which we rewrite as $\mathcal{U} \cdot \mathcal{S} - \mathcal{U} = \mathcal{M} - \varepsilon$. Observe that $\sum_{j \in \mathcal{S}} P_{i,j}z = z$. Hence, set $\mathcal{U} \cdot \mathcal{S}$ yields (conditioning on the left occurrence of H)

$$\sum_{w \in \mathcal{U}} \sum_{j \in \mathcal{S}} P(w_j|H)z^{|w_j|} = \sum_{i \in \mathcal{S}} \sum_{w \in \mathcal{U}, \ell(w)=i} P(w|H)z^{|w|} \sum_{j \in \mathcal{S}} P_{i,j}z = U_H(z) \cdot z.$$

Of course, $\mathcal{M} - \varepsilon$ and \mathcal{U} translate into $M_H(z) - 1$ and $U_H(z)$, and (40) is proved.

We now turn our attention to (41), and we use relationship (35) of Theorem 3.2. Observe that $\mathcal{S} \cdot \mathcal{R}$ can be rewritten as

$$\sum_{j,i \in \mathcal{S}^2} \sum_{i w \in \mathcal{R}} P(ji w)z^{|ji w|} = z^2 \sum_{j \in \mathcal{S}} \sum_{i \in \mathcal{S}} \pi_j P_{j,i} \sum_{i w \in \mathcal{R}} P(w|w_{-1} = i)z^{|w|}.$$

However, due to the stationarity of the underlying Markov chain $\sum_j \pi_j P_{j,i} = \pi_i$. As $\pi_i P(w|w_{-1} = i) = P(iw)$, we get $zR(z)$. Furthermore, in (35) $H \cdot \mathcal{M} - H$ translates into $P(H)z^m \cdot (M_H(z) - 1)$. Nonetheless, by (40), this becomes $P(H)z^m \cdot U_H(z)(z - 1)$, and, after a simplification, we prove (41).

Finally, we deal with (39), and prove it using (33) from Theorem 3.2. The left-hand side of (33) involves language \mathcal{M} , hence we must condition on the left occurrence of H . In particular, $\bigcup_{r \geq 1} \mathcal{M}^r + \varepsilon$ of (33) translates into $1/(1 - M_H(z))$. Now we deal with $\mathcal{W} \cdot H$ of the right-hand side of (33). *Conditioning* on the left occurrence of H , the generating function $W(z)H(z)$ of $\mathcal{W} \cdot H$ becomes

$$\begin{aligned} W_H(z)H(z) &= \sum_{n \geq 0} \sum_{|w|=n} z^{n+m} P(wH|w_{-1} = \ell(H)) \\ &= \sum_{n \geq 0} \sum_{|w|=n} z^n P(wh_1|w_{-1} = \ell(H))P(v = h_2 \cdots h_m|v_{-1} = h_1)z^m. \end{aligned}$$

We have $P(v = h_2 \cdots h_m|v_{-1} = h_1)z^m = (1/\pi_{h_1})z^m P(H)$, and, for $n \geq 0$,

$$\sum_{|w|=n} P(wh_1|w_{-1} = \ell(H)) = [\mathbf{P}^{n+1}]_{\ell(H), h_1},$$

where, we recall, $\ell(H) = h_m$ is the last character of H . In summary: language $\mathcal{W} \cdot H$ contributes $P(H)z^m [(1/\pi_{h_1}) \sum_{n \geq 0} \mathbf{P}^{n+1} z^n]_{\ell(H), h_1}$, while language $\mathcal{A} - \{\varepsilon\}$ introduces $A_H(z) - 1$. We now observe that, for any symbols i and j ,

$$\left[\frac{1}{\pi_j} \sum_{n \geq 0} \Pi z^n \right]_{i,j} = \sum_{n \geq 0} z^n = \frac{1}{1 - z}.$$

Using the equality $\mathbf{P}^{n+1} - \Pi = (\mathbf{P} - \Pi)^{n+1}$ (which follows from a consecutive application of the identity: $\Pi\mathbf{P} = \Pi$), we finally obtain the sum in (39). This completes the proof of the theorem. \square

REMARK. The generating function of language \mathcal{T}_0^j (i.e., no H occurrence with the last symbol of a word from \mathcal{T}_0^j being j) in the Markov case was previously derived by Chrysaphinou and Papastavridis in [5]. We observe that the generating function $T^{(0)}(z)$ of \mathcal{T}_0 easily follows from (6) and the equation $T^{(0)}(z) = 1/(1 - z) - \sum_{r \geq 1} T^{(r)}(z)$.

3.3. *Moments and Limiting Distribution.* In this final subsection we derive the first two moments of O_n as well as asymptotics for $\Pr\{O_n = r\}$ for different ranges of r , that is, we prove Theorem 2.2. Actually, we should mention that using general results on Markov chains and renewal theory one immediately guesses that the limiting distribution must be normal for $r = EO_n + O(\sqrt{n})$. However, here the challenge is to estimate precisely the variance. Our approach offers an easy, uniform, and precise derivation of all moments, including the variance, as well as local limit distributions (including the convergence rate) for the central and large deviations regimes. We use an analytic approach (see [2], [13], [26], and [31]).

A. *Moments.* First, from Theorem 2.1 we conclude that

$$T_u(z, 1) = \frac{z^m P(H)}{(1 - z)^2},$$

$$T_{uu}(z, 1) = \frac{2z^m P(H)M_H(z)D_H(z)}{(1 - z)^3},$$

where $T_u(z, 1)$ and $T_{uu}(z, 1)$ are first and second derivatives of $T(z, u)$ at $u = 1$. Now, we observe that both expressions admit as a numerator a function that is analytic beyond the unit circle. This allows for a very simple computation of the expectation and variance based on the following basic formula:

$$(42) \quad [z^n](1 - z)^{-p} = \frac{\Gamma(n + p)}{\Gamma(p)\Gamma(n + 1)},$$

where $[z^n]$ means the coefficient of z^n . To obtain EO_n we proceed as follows, for $n \geq m$:

$$EO_n = [z^n]T_u(z, 1) = P(H)[z^{n-m}](1 - z)^{-2} = (n - m + 1)P(H).$$

We denote

$$\Phi(z) = 2z^m P(H)M_H(z)D_H(z),$$

which is a polynomial in the Bernoulli case. We use the Taylor expansion

$$\Phi(z) = \Phi(1) + (z - 1)\Phi'(1) + \frac{(z - 1)^2}{2}\Phi''(1) + (z - 1)^3 f(z),$$

where $f(z)$ is a polynomial of degree $2m - 2$. It follows that $[z^n](z - 1)f(z)$ is 0 for $n \geq 2m - 1$ and, using formula (42), we get

$$EO_n(O_n - 1) = [z^n]T_{uu}(z, 1) = \Phi(1)\frac{(n + 2)(n + 1)}{2} - \Phi'(1)(n + 1) + \frac{1}{2}\Phi''(1).$$

Observing that $M_H(z)D_H(z) = D_H(z) + (1 - z)$, we use MAPLE to obtain a precise formula for the variance (see (14) of Theorem 2.2). In the Markov case, we have to compute the additional term

$$[z^n] \frac{2(z^{2m} P(H)^2 F(z))}{(1 - z)^2},$$

where $F(z)$ is analytic beyond the unit circle for $|z| \leq R$, with $R > 1$. The Taylor expansion of $F(z)$ is $E_1 + (1 - z)E_2$ and applying (42) again yields the result. In a similar manner, we can compute all the moments of O_n .

B. Distribution: Case $r = O(1)$. Now, we prove part (ii) of Theorem 2.2, that is, we establish an asymptotic expression for $\Pr\{O_n = r\}$ for $r = O(1)$. We first rewrite the formula on $T^{(r)}(z)$ as follows:

$$(43) \quad T^{(r)}(z) = \frac{z^m P(H)(D_H(z) + z - 1)^{r-1}}{D_H^{r+1}(z)}.$$

Observe that $\Pr\{O_n = r\}$ is the coefficient at z^n of $T^{(r)}(z)$. By Hadamard’s theorem (see [31] and [35]), the asymptotics of the coefficients of $T^{(r)}(z)$ depend on the singularities of $T^{(r)}(z)$. In our case, the generating function is a rational function, thus we can only expect poles (for which the denominator $D_H(z)$ vanishes). The next lemma establishes the existence and properties of such a pole.

LEMMA 3.2. *The equation $D_H(z) = 0$ has at least one root, and all its roots are of modulus greater than 1.*

PROOF. A root of $D_H(z) = (1 - z)/(1 - M_H(z))$ is clearly a pole of $1/(1 - M_H(z))$. As $1/(1 - M_H(z))$ is the generating function of a language, it converges for $|z| < 1$ and has no pole of modulus smaller than 1. Since $D_H(1) \neq 0$, $z = 1$ is a simple pole of $1/(1 - M_H(z))$. As all its coefficients are real and positive, there is no other pole of modulus $|z| = 1$. It follows that all roots of $D_H(z)$ are of modulus greater than 1. The existence of a root is guaranteed since $D_H(z)$ is either a polynomial (Bernoulli model) or a ratio of polynomials (Markov model). □

In view of the above, the generating function $T^{(r)}(z)$ can be expanded around its root of smallest modulus, say ρ_H , as Laurent’s series (see [26], [35], and [40]):

$$(44) \quad T^{(r)}(z) = \sum_{j=1}^{r+1} \frac{a_j}{(z - \rho_H)^j} + \tilde{T}^{(r)}(z),$$

where $\tilde{T}^{(r)}(z)$ is analytical in $|z| < \rho'$ and ρ' is defined as $\rho' = \inf\{|\rho| : \rho > \rho_H \text{ and } D_H(\rho) = 0\}$. The constants a_j satisfy formulae (22). This formula simplifies into (21) for the leading constant a_{-r-1} . As a consequence of analyticity [40] we have, for $1 < \rho_H < \rho < \rho'$, $[z^n]\tilde{T}^{(r)}(z) = O(\rho^{-n})$. Hence, the term $\tilde{T}^{(r)}(z)$ contributes only to the lower terms in the asymptotic expansion of $T^{(r)}(z)$.

We need an asymptotic expansion for the first terms in (43). This is rather a standard computation (see [31] and [40]), but for completeness we provide a short proof. The following chain of identities is easy to justify for any $\rho > 0$:

$$\begin{aligned} \sum_{j=1}^{r+1} \frac{a_j}{(z - \rho)^j} &= \sum_{j=1}^{r+1} \frac{a_j(-1)^j}{\rho^j(1 - (z/\rho)^j)} \\ &= \sum_{j=1}^{r+1} (-1)^j a_j \rho^{-j} \sum_{n=0}^{\infty} \binom{n+j-1}{n} \left(\frac{z}{\rho}\right)^n \\ &= \sum_{n=1}^{\infty} z^n \sum_{j=1}^{\min\{r+1, n\}} (-1)^j a_j \binom{n}{j-1} \rho^{-(n+j)}. \end{aligned}$$

After some algebra, we prove part (ii) of Theorem 2.2.

C. Central Limit Theorem: Case $r = EO_n + xO(\sqrt{n})$. We now establish part (iii) of Theorem 2.2, that is, we compute $\Pr\{O_n = r\}$ for $r = EO_n + x\sqrt{\text{Var } O_n}$ when $x = O(1)$. Let $\mu_n = EO_n(H) = (n - m + 1)P(H)$ and $\sigma_n^2 = \text{Var } O_n(H) \sim c_1n$. To establish asymptotic normality of $(O_n(H) - \mu_n)/\sigma_n$, it suffices, according to Lévy’s theorem, to prove the following (see also [2]):

$$(45) \quad \lim_{n \rightarrow \infty} e^{-\tau\mu_n/\sigma_n} T_n(e^{\tau/\sigma_n}) = e^{\tau^2/2}$$

for complex τ . Again, by Cauchy’s theorem

$$T_n(u) = \frac{1}{2\pi i} \oint \frac{T(z, u)}{z^{n+1}} dz = \frac{1}{2\pi i} \oint \frac{uP(H)}{D_H^2(z)(1 - uM_H(z))z^{n+1-m}} dz,$$

where integration is along a circle around the origin. The evaluation of this integral is standard and it appeals to the Cauchy residue theorem. Namely, we enlarge the circle of integration to a bigger one, say $R > 1$, such that the bigger circle contains the dominating pole of the integrand function. Observe that the Cauchy integral over the bigger circle is $O(R^{-n})$. We now substitute (for simplicity of further derivations) $u = e^t$ and $z = e^\rho$. Then the poles of the integrand are the roots of the equation

$$(46) \quad 1 - e^t M_H(e^\rho) = 0.$$

This equation implicitly defines in some neighborhood of $t = 0$ a unique C^∞ function $\rho(t)$, satisfying $\rho(0) = 0$. Notably, all other roots ρ satisfy $\inf|\rho| = \rho' > 0$. Then the residue theorem with $e^{\rho'} > R > e^\rho > 1$ leads to

$$(47) \quad T_n(e^t) = C(t)e^{-(n+1-m)\rho(t)} + O(R^{-n}),$$

where

$$C(t) = \frac{P(H)}{D_H^2(\rho(t))M_H'(\rho(t))}.$$

To study some properties of $\rho(t)$, we observe that the cumulant formula implies $EO_n(H) = [t] \log T_n(e^t)$ and $\sigma_n^2 = [t^2] \log T_n(e^t)$ where, we recall, $[t^r]f(t)$ denotes

the coefficient of $f(t)$ at t^r . In our case, $\mu_n \sim -n\rho'(0)$ as well as $\sigma_n^2 \sim -n\rho''(0)$. In (47), now set $t = \tau/\sigma_n \rightarrow 0$ for some complex τ . Since uniformly in t we have $\rho(t) = t\rho'(0) + \rho''(0)t^2/2 + O(t^3)$ for $t \rightarrow 0$, our estimate (47) leads to

$$\begin{aligned} e^{-\tau\mu_n/\sigma_n}T_n(e^{\tau/\sigma_n}) &= \exp\left(\frac{\tau^2}{2} + O\left(\frac{n\tau^3}{\sigma_n^3}\right)\right) \\ &= e^{\tau^2/2}\left(1 + O\left(\frac{1}{\sqrt{n}}\right)\right), \end{aligned}$$

which completes the proof of Theorem 2.2(iii).

Actually, we can proceed as in [17] or [21] to obtain a much more refined local limit result. For example, a direct application of results from [17] (see Chapter 4.3.3) leads to the following for $x = o(n^{1/6})$:

$$(48) \quad \begin{aligned} \Pr\{O_n = EO_n + x\sqrt{nc_1}\} \\ = \frac{1}{\sqrt{2\pi nc_1}}e^{-x^2/2}\left(1 - \frac{\kappa_3}{2c_1^{3/2}\sqrt{n}}\left(x - \frac{x^3}{3}\right)\right) + O(n^{-3/2}), \end{aligned}$$

where κ_3 is a constant (i.e., the third cumulant).

D. Large Deviations: Case $r = (1 + \delta)EO_n$. Finally, we consider the large deviations result. From (47) we conclude that

$$\lim_{n \rightarrow \infty} \frac{\log T_n(e^t)}{n} = -\rho(t).$$

Thus, directly from the Gärtner–Ellis theorem [4], [9] we prove that

$$\lim_{n \rightarrow \infty} \frac{\log \Pr\{O_n > na\}}{n} = -I(a),$$

where

$$I(a) = a\omega_a + \rho(\omega_a)$$

with ω_a being a solution of $-\rho'(t) = a$. A stronger version of the above follows directly from Theorem 3.1 of [4]. To derive our result of Theorem 2.2, we use (49) and the “shift of mean” technique as discussed below (see [4], [17], [21], and [31]).

As in the central limit regime, we could use Cauchy’s formula to compute the probability $\Pr\{O_n = r\}$ for $r = EO_n + xO(\sqrt{n})$. However, formula (49) is only good for $x = O(1)$. To expand its validity, we shift the mean of the generating function $T_n(u)$ to a new value, say $m = an = (1 + \delta)P(H)(n - m + 1)$, so we can again apply the central limit formula (49) around the new mean. To accomplish this, we rewrite (47) as

$$T_n(e^t) = C(t)[g(t)]^{n-m+1},$$

where $g(t) = e^{-\rho(t)}$, and for simplicity of this discussion we dropped the $O(R^{-n})$ term. The above suggests that $T_n(e^t)$ is the moment generating function of a sum S_n of $n - m + 1$ “almost” independent random variables X_1, \dots, X_{n-m+1} and Y whose

moment generating functions are $g(t)$ and $C(t)$, respectively. Observe that $ES_n = (n - m + 1)P(H)$ while we need to estimate the tail of S_n around $(1 + \delta)(n - m + 1)P(H)$. To achieve it, we introduce a new random variable \tilde{X}_i whose moment generating function $\tilde{g}(t)$ is

$$\tilde{g}(t) = \frac{g(t + \omega)}{g(\omega)},$$

where ω will be chosen later. Then the mean and the variance of the new variable \tilde{X} is

$$E\tilde{X} = \frac{g'(\omega)}{g(\omega)} = -\rho'(\omega),$$

$$\text{Var}\tilde{X} = \frac{g''(\omega)}{g(\omega)} - \left(\frac{g'(\omega)}{g(\omega)}\right)^2 = -\rho''(\omega).$$

We now choose ω_a such that

$$-\rho'(\omega_a) = \frac{g'(\omega_a)}{g(\omega_a)} = a = P(H)(1 + \delta).$$

Then the new sum $\tilde{S}_n = Y + \tilde{X}_1 + \dots + \tilde{X}_{n-m+1}$ has a new mean $(1 + \delta)P(H)(n - m + 1) = a(n - m + 1)$, and hence we can apply the central limit result (49) to \tilde{S}_n . To translate from \tilde{S}_n to S_n we use the following simple formula:

$$(49) \quad [e^{tN}](C(t)g^n(t)) = \frac{g^n(\omega)}{e^{\omega N}} [e^{tN}] \left(\frac{C(t)g^N(t + \omega)}{g^N(\omega)} \right),$$

where $N = a(n - m + 1)$ and $[e^{tN}]g(t)$ denotes the coefficient of $g(t)$ at e^{tN} . Now we can apply (49) to the right-hand side of the above to obtain

$$[e^{tN}] \left(\frac{C(t)g^N(t + \omega)}{g^N(\omega)} \right) = \frac{1}{\sigma_a \sqrt{2\pi(n - m + 1)}} (1 + O(n^{-1})) + O(n^{-5/2}).$$

Finally, using (49) and the above, we prove Theorem 2.2(iv).

Acknowledgments. It is our pleasure to acknowledge several discussions with A. Dembo, A. Odlyzko, P. Pevzner, E. Coward and E. Sutinen on the topic of this paper. Our deepest appreciation goes to Philippe Flajolet for his inspirations and support over the last decade.

References

- [1] D. Barbara and T. Imielinski, Sleepers and Workoholics—Caching in Mobile Wireless Environments, *Proc. ACM SIGMOD*, pp. 1–15, Minneapolis, 1994.
- [2] E. Bender, Central and Local Limit Theorems Applied to Asymptotic Enumeration, *J. Combin. Theory Ser. A*, **15**, 91–111, 1973.
- [3] S. Breen, M. Waterman, and N. Zhang, Renewal Theory for Several Patterns, *J. Appl. Probab.*, **22**, 228–234, 1985.

- [4] J. Bucklew and J. Sadowsky, A Contribution to the Theory of Chernoff Bounds, *IEEE Trans. Inform. Theory*, **39**, 249–254, 1993.
- [5] C. Chrysaphinou and S. Papastavridis, The Occurrence of Sequence of Patterns in Repeated Dependent Experiments, *Theory Probab. Appl.*, **79**, 167–173, 1990.
- [6] E. Coward, Exact Calculation of Word Occurrence Probabilities and Its Applications to the Search of Repeats, *Proc. Mathematical Analysis of Biological Sequences*, Rouen, 1997.
- [7] M. Crochemore and W. Rytter, *Text Algorithms*, Oxford University Press, New York, 1995.
- [8] I. Csiszár and J. Körner, *Information Theory: Coding Theorems for Discrete Memoryless Systems*, Academic Press, New York, 1981.
- [9] R. Ellis, Large Deviations for a General Class of Random Vectors, *Ann. Probab.*, **12**, 1–12, 1984.
- [10] W. Feller, *An Introduction to Probability and Its Applications*, Vol. 1, Wiley, New York, 1968.
- [11] J. Fickett, Recognition of Protein Coding Regions in DNA Sequences, *Nucleic Acids Res.*, **10**, 5303–5318, 1982.
- [12] P. Flajolet, X. Gourdon, and C. Martinez, Patterns in Random Binary Search Trees, *Random Structures Algorithms*, **11**, 223–244, 1997.
- [13] P. Flajolet and M. Soria, General Combinatorial Schemas: Gaussian Limit Distributions and Exponential Tails, *Discrete Math.*, **114**, 159–180, 1993.
- [14] I. Fudos, E. Pitoura, and W. Szpankowski, On Pattern Occurrences in a Random Text, *Inform. Process. Lett.*, **57**, 307–312, 1996.
- [15] M. S. Gelfand, Prediction of Function in DNA Sequence Analysis, *J. Comput. Biol.*, **2**, 87–117, 1995.
- [16] M. Geske, A. Godbole, A. Schafner, A. Skolnick, and G. Wallstrom, Compound Poisson Approximations for World Patterns Under Markovian Hypotheses, *J. Appl. Probab.*, **32**, 877–892, 1995.
- [17] D. Greene and D. E. Knuth, *Mathematics for the Analysis of Algorithms*, Birkhäuser, Boston, 1990.
- [18] L. Guibas and A. Odlyzko, Maximal Prefix-Synchronized Codes, *SIAM J. Appl. Math.*, **35**, 401–418, 1978.
- [19] L. Guibas and A. Odlyzko, Periods in Strings, *J. Combin. Theory Ser. A*, **30**, 19–43, 1981.
- [20] L. Guibas and A. W. Odlyzko, String Overlaps, Pattern Matching, and Nontransitive Games, *J. Combin. Theory Ser. A*, **30**, 183–208, 1981.
- [21] H.-K. Hwang, Théorèmes Limites pour les Structures Combinatoires et les Fonctions Arithmétiques, Thèse de Doctorat, l'École Polytechnique, 1994.
- [22] P. Jacquet and W. Szpankowski, Autocorrelation on Words and Its Applications. Analysis of Suffix Trees by String-Ruler Approach, *J. Combin. Theory Ser. A*, **66**, 237–269, 1994.
- [23] P. Jokinen and E. Ukkonen, Two Algorithms for Approximate String Matching in Static Texts, *Proc. MFCS 91, Lecture Notes in Computer Science*, vol. 520, pp. 240–248, Springer-Verlag, Berlin, 1991.
- [24] S. Karlin, C. Bruge, and A. Campbell, Statistical Analysis of Counts and Distributions of Restriction Sites in DNA Sequences, *Nucl. Acids Res.*, **20**, 1363–1370, 1992.
- [25] S. Karlin and F. Ost, Counts of Long Aligned Word Matches Among Random Letter Sequences, *Ann. Probab.*, **19**, 293–351, 1987.
- [26] D. E. Knuth, *The Art of Computer Programming: Fundamental Algorithms*, vol. 1., Addison-Wesley, Reading, MA, 1973.
- [27] S. R. Li, A Martingale Approach to the Study of Occurrences of Sequence Patterns in Repeated Experiments, *Ann. Probab.*, **8**, 1171–1176, 1980.
- [28] T. Luczak and W. Szpankowski, A Suboptimal Lossy Data Compression Based on Approximate Pattern Matching, *IEEE Trans. Inform. Theory*, **43**, 1439–1451, 1997.
- [29] K. Marton and P. Shields, The Positive-Divergence and Blowing-up Properties, *Israel J. Math.*, **80**, 331–348, 1994.
- [30] P. T. Nielsen, On the Expected Duration of a Search for Fixed Pattern in Random Data, *IEEE Trans. Inform. Theory*, **19**, 702–704, 1973.
- [31] A. Odlyzko, Asymptotic Enumeration, in *Handbook of Combinatorics*, vol. II, pp. 1063–1229 (Eds. R. Graham, M. Götschel, and L. Lovász), Elsevier Science, Amsterdam, 1995.
- [32] P. Pevzner, M. Borodovsky, and A. Mironov, Linguistic of Nucleotide Sequences: The Significance of Deviations from Mean: Statistical Characteristics and Prediction of the Frequency of Occurrence of Words, *J. Biomol. Struct. Dynamics*, **6**, 1013–1026, 1991.
- [33] B. Prum, F. Rodolphe, and E. Turckheim, Finding Words with Unexpected Frequencies in Deoxyribonucleic Acid Sequence, *J. Roy. Statist. Soc. Ser. B*, **57**, 205–220, 1995.

- [34] M. Régnier and W. Szpankowski, On the Approximate Pattern Occurrence in a Text, *Proc. SEQUENCE '97*, Positano, 1997.
- [35] R. Remmert, *Theory of Complex Functions*, Springer-Verlag, New York, 1991.
- [36] S. Schbath, Etude Asymptotique du Nombre d'Occurrences d'un Mot dans une Chaîne de Markov et Application à la Recherche de Mots de Fréquence Exceptionnelle dans les Séquences d'ADN, Thèse, Université René Descartes Paris V, 1995.
- [37] E. Sutinen and W. Szpankowski, On the Collapse of q -Gram Filtration, "Fun with Algorithms," Elba, 1998.
- [38] W. Szpankowski, Asymptotic Properties of Data Compression and Suffix Trees, *IEEE Trans. Inform. Theory*, **39**, 1647–1659, 1993.
- [39] M. Waterman, *Introduction to Computational Biology*, Chapman & Hall, New York, 1995.
- [40] H. Wilf, *generatingfunctionology*, Academic Press, Boston, 1990.
- [41] Z. Zhang and E. Yang, An On-Line Universal Lossy Data Compression Algorithm via Continuous Codebook Refinement—Part II: Optimality for Phi-Mixing Source Models, *IEEE Trans. Inform. Theory*, **42**, 822–836, 1996.