

Analysis of the Multiplicity Matching Parameter in Suffix Trees

Mark Daniel Ward¹ and Wojciech Szpankowski^{2†}

¹*Department of Mathematics, Purdue University, West Lafayette, IN, USA. mward@math.purdue.edu*

²*Department of Computer Science, Purdue University, West Lafayette, IN, USA. spa@cs.purdue.edu*

In a suffix tree, the multiplicity matching parameter (MMP) M_n is the number of leaves in the subtree rooted at the branching point of the $(n + 1)$ st insertion. Equivalently, the MMP is the number of pointers into the database in the Lempel-Ziv '77 data compression algorithm. We prove that the MMP asymptotically follows the logarithmic series distribution plus some fluctuations. In the proof we compare the distribution of the MMP in suffix trees to its distribution in tries built over independent strings. Our results are derived by both probabilistic and analytic techniques of the analysis of algorithms. In particular, we utilize combinatorics on words, bivariate generating functions, pattern matching, recurrence relations, analytical poissonization and depoissonization, the Mellin transform, and complex analysis.

Keywords: suffix trees, combinatorics on words, pattern matching, autocorrelation polynomial, complex asymptotics, data compression

Contents

1	Introduction	2
2	Main Results	4
3	Proofs	7
3.1	BGF for the Multiplicity Matching Parameter of Independent Tries	7
3.2	BGF for the Multiplicity Matching Parameter of Suffix Trees	7
3.3	On the Autocorrelation Polynomial	10
3.4	Analytic Continuation	12
3.5	Singularity Analysis	13
3.6	Comparing Suffix Trees to Tries	14
3.7	Summary and Conclusion	18

[†]This research was supported by NSF Grant CCR-0208709, NIH grant R01 GM068959-01, and AFOSR Grant FA8655-04-1-3074.

1 Introduction

When transmitting data, the goal of *source coding (data compression)* is to represent the source with a minimum of symbols. On the other hand, the goal of *channel coding (error correction)* is to represent the source with a minimum of error probability in decoding. These goals are obviously in conflict. Traditionally, additional symbols are transmitted when performing error correction.

In [1], an algorithm for joint data compression and error correction is presented; the compression performance is not degraded because the algorithm requires no extra symbols for error correction. In this scheme, a Reed-Solomon error-correcting code is embedded into the Lempel-Ziv '77 data compression algorithm (see [2]). Lonardi and Szpankowski utilize the fact that the LZ'77 adaptive data compression algorithm is unable to remove all redundancy from the source. Our goal here is to precisely determine the number of redundant bits that are available to be utilized in the aforementioned scheme.

We recall the basic operation of the LZ'77 data compression algorithm. When n bits of the source have already been compressed, the LZ'77 encoder finds the longest prefix of the uncompressed data that also appears in the database (namely, the compressed portion of the data). The encoder performs the compression by storing a pointer into the database (and also the length of this prefix, as well as the next character of the source). Often, this longest prefix appears more than once in the database. Each of the database entries are *equally eligible* for use by the encoder; thus, *any of the analogous pointers* into the database is suitable. In practice, the choice of pointer among these candidates has no significance. On the other hand, by judiciously selecting the pointer, some error correction can be performed. For instance, if two pointers are available, the encoder could easily perform a parity check by choosing the first pointer for “0” and the second pointer for “1”. Lonardi and Szpankowski’s scheme for performing error correction is very elaborate. We refer the reader to their paper for more details.

We let M_n denote the *number of pointers into the database* when n bits have already been compressed (as described above). Throughout this paper, we are primarily interested in precisely determining the asymptotics of M_n . A thorough analysis of M_n yields a characterization of the degree to which error correction can be performed in the scheme discussed above. We note that $\lfloor \log_2 M_n \rfloor$ bits are available to be used for correcting errors.

Tries, especially suffix trees, provide a natural way to study M_n . We work here with strings of characters drawn independently from the binary alphabet $\mathcal{A} := \{0, 1\}$. We let p denote the probability of “0” and $q = 1 - p$ denotes the probability of “1”; without loss of generality, we assume that $q \leq p$ throughout the discussion.

We first recall the definition of a binary trie built over a set \mathcal{Y} of n strings. The construction is recursive. If $|\mathcal{Y}| = 0$, then the trie is empty. If $|\mathcal{Y}| = 1$, then $\text{trie}(\mathcal{Y})$ is a single node. Finally, if $|\mathcal{Y}| > 1$, then \mathcal{Y} is partitioned into two subsets, \mathcal{Y}_0 and \mathcal{Y}_1 , such that a string is in \mathcal{Y}_0 if its first symbol is 0, and a string is in \mathcal{Y}_1 if its first symbol is 1. Then $\text{trie}(\mathcal{Y}_0)$ and $\text{trie}(\mathcal{Y}_1)$ are each constructed in the same way, except that the splitting of sets at the k th step is based on the k th symbol of the string. This completes the definition of a binary trie.

Now we briefly recall the construction of a binary suffix tree built over a string $X = X_1X_2X_3 \dots$. The word $X^{(i)} = X_iX_{i+1}X_{i+2} \dots$ is the i th suffix of X , which begins at the i th position of X . Then a binary suffix tree is precisely a binary trie built over the first n suffixes of X , namely $X^{(1)}, X^{(2)}, \dots, X^{(n)}$.

In a suffix tree, M_n is exactly the number of leaves in the subtree rooted at the branching point of the $(n + 1)$ st insertion (cf. Figure 1). The strings in a suffix tree are highly dependent on each other, which apparently makes a precise analysis of M_n quite difficult; therefore, we also consider the analogous (but

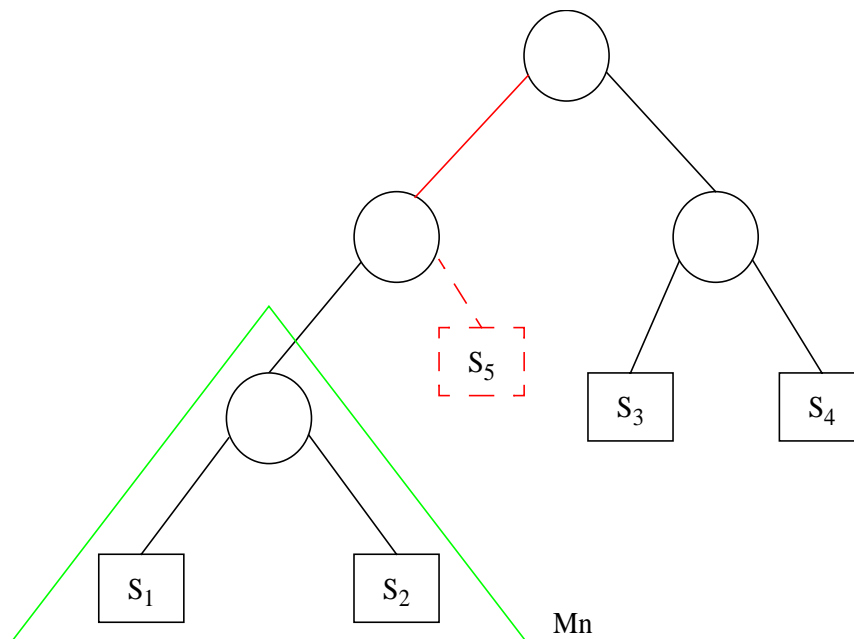


Fig. 1: Finding M_n^I in a trie.

simpler) situation in a trie built over *independent* strings; namely, we study M_n^I , which is the number of leaves in the subtree rooted at the branching point of the $(n + 1)$ st insertion in a trie built over $n + 1$ *independent* binary strings. For instance, in Figure 1, we have $n = 4$ and $M_n^I = 2$ because there are two leaves (namely, S_1 and S_2) in the subtree rooted at the branching point of the 5th insertion.

We are primarily concerned with comparing the distribution of M_n (a parameter of suffix trees) to the distribution of M_n^I (a parameter of tries built over independent strings). Our approach to the proof begins with the observation that a variety of parameters have the same asymptotic behavior regardless of whether they correspond to suffix trees or to tries built over independent strings. This was observed in ? and then made precise in ?, where the typical depth in a suffix tree is proven to be asymptotically the same as the typical depth in a trie built over independent strings when the underlying source is i.i.d. An extension of such results to an underlying Markovian model is presented in ?.

The limiting distribution of several trie parameters is given in ? and ?. More results about trie parameters are found in ?. The variance of the external path in a symmetric trie is given in ?. The depth of a digital trie with an underlying Markovian dependency is analyzed in ?. Many results about a variety of tree structures are collected in ?. Average-case studies of several parameters of suffix trees are found in ? and ?.

We briefly summarize the methodology of our proof. Our goal is to compare the distribution of M_n (the multiplicity matching parameter of a suffix tree) to the distribution of M_n^I (the MMP of a trie built over independent strings). Our proof that these two parameters have the same asymptotic distribution consists of several steps. We first derive bivariate generating functions for M_n and M_n^I , denoted as $M(z, u)$ and

$M^I(z, u)$, respectively. We noted above that a suffix tree is built over the suffixes $X^{(1)}, X^{(2)}, \dots, X^{(n)}$ of a string X . These suffixes are highly dependent on each other. Therefore, in deriving the bivariate generating function $M(z, u)$, an interesting obstacle arises: We need to determine the degree to which a suffix of X can overlap with itself. Fortunately, the autocorrelation polynomial $S_w(z)$ of a word w measures the amount of overlap of a word w with itself. The autocorrelation polynomial was introduced in [?] and was utilized extensively in [?] and [?]. The autocorrelation polynomial is defined as

$$S_w(z) = \sum_{k \in \mathcal{P}(w)} \mathbf{P}(w_{k+1}^m) z^{m-k} \quad (1)$$

where $m = |w|$ and where $\mathcal{P}(w)$ denotes the set of positions k of w satisfying $w_1 \dots w_k = w_{m-k+1} \dots w_m$, that is, w 's prefix of length k is equal to w 's suffix of length k . Using the autocorrelation polynomial, we can overcome the difficulties inherent in the fact that suffixes of a word X overlap with each other. By utilizing $S_w(z)$, we are able to obtain a succinct way of describing the bivariate generating function $M(z, u)$. Fortunately, the autocorrelation polynomial is well-understood. Note that the autocorrelation polynomial $S_w(z)$ has a $\mathbf{P}(w_{k+1}^m) z^{m-k}$ term if and only if w has an overlap with itself of length k . All words w overlap with themselves trivially, so all autocorrelation polynomials have a constant term (i.e., $z^{m-m} = z^0 = 1$ term). On the other hand, with high probability, w has very few large nontrivial overlaps with itself. Therefore, with high probability, all nontrivial overlaps of w with itself are small; such overlaps correspond to high-degree terms of $S_w(z)$.

In order to compare $M(z, u)$ and $M^I(z, u)$, we utilize complex analysis. Specifically, we take advantage of Cauchy's theorem, which allows us to analyze the poles of the generating functions $M(z, u)$ and $M^I(z, u)$ in order to obtain precise information about the distributions of M_n and M_n^I . During this residue analysis, it is necessary that the generating function for M_n is analytically continued from the unit disk to a larger disk.

Our ultimate conclusion is that the distribution of the multiplicity matching parameter M_n is asymptotically the same in suffix trees and independent tries, i.e., M_n and M_n^I have asymptotically the same distribution.

The asymptotics for the distribution and factorial moments of M_n^I were given in [?]. Specifically, M_n^I asymptotically follows the logarithmic series distribution (plus some fluctuations when $\ln p / \ln q$ is rational). Since we prove here that M_n and M_n^I have asymptotically the same distribution, then as a consequence, we see that M_n also asymptotically follows the logarithmic series distribution. One striking property of this distribution is the high concentration around the mean. We see that $E[M_n]$ is asymptotically $\frac{1}{h}$ (where h denotes the entropy of the source) and also M_n is highly concentrated around this average value; this property of M_n is very desirable for the error correction scheme described in [?].

This paper is a concise version of the first author's Ph.D. thesis; see [?].

2 Main Results

We consider the string $X = X_1 X_2 X_3 \dots$, where the X_i 's are i.i.d. random variables on $\mathcal{A} := \{0, 1\}$ with $\mathbf{P}(X_i = 0) = p$ and $\mathbf{P}(X_i = 1) = q$. (Without loss of generality, we assume throughout the discussion that $q \leq p$.) Let $X^{(i)}$ denote the i th suffix of X . In other words, $X^{(i)} = X_i X_{i+1} X_{i+2} \dots$. Consider the *longest* prefix w of $X^{(n+1)}$ such that $X^{(i)}$ also has w as a prefix, for some i with $1 \leq i \leq n$. Then M_n is defined as the number of $X^{(i)}$'s (with $1 \leq i \leq n$) that have w as a prefix. So

$$M_n = \#\{1 \leq i \leq n \mid X^{(i)} = X_i X_{i+1} X_{i+2} \dots \text{ has } w \text{ as a prefix}\}. \quad (2)$$

An alternate definition of M_n is available via *suffix trees*. First, consider a suffix tree built from the first $n + 1$ suffixes of X . Next, consider the *insertion point* of the $(n + 1)$ st suffix. Then M_n is exactly the *number of leaves* in the subtree rooted at the branching point of the $(n + 1)$ st insertion. For instance, suppose that the $(n + 1)$ st suffix starts with $w\beta$ for some $w \in \mathcal{A}^*$ and $\beta \in \mathcal{A}$. Then, examining the first n suffixes, if there are exactly k suffixes that begin with $w\alpha$ (where $\alpha = 1 - \beta$, i.e., $\{\alpha, \beta\} = \{0, 1\}$), and the other $n - k$ suffixes do not begin with w , we conclude that $M_n = k$.

Unfortunately, the strings in a suffix tree are highly dependent on each other; thus, a precise analysis of M_n is quite difficult. On the other hand, the asymptotic behavior of M_n^I , an analogous parameter of tries built over independent strings, is well-understood. Specifically, M_n^I asymptotically follows the logarithmic series distribution (plus some fluctuations when $\ln p / \ln q$ is rational). In ?, a precise analysis of M_n^I is given via the analysis of independent tries, using recurrence relations, analytical poissonization and deoissonization, the Mellin transform, and complex analysis.

To define M_n^I , we consider the situation described above, but we build a trie from $n + 1$ independent strings from \mathcal{A}^* . So we consider independent $X(i)$'s; specifically, we define $X(i) = X_1(i)X_2(i)X_3(i) \dots$, where $\{X_j(i) \mid i, j \in \mathbb{N}\}$ is a collection of i.i.d. random variables. We let w denote the *longest prefix* of $X(n + 1)$ such that $X(i)$ also has w as a prefix, for some i with $1 \leq i \leq n$. Then M_n^I is defined as the number of $X(i)$'s (with $1 \leq i \leq n$) that have w as a prefix. So

$$M_n^I = \#\{1 \leq i \leq n \mid X(i) = X_1(i)X_2(i)X_3(i) \dots \text{ has } w \text{ as a prefix}\}. \quad (3)$$

To define M_n^I via *tries*, first consider a trie built from the $n + 1$ independent strings from \mathcal{A}^* . Next, consider the *insertion point* of the $(n + 1)$ st string. Then M_n is exactly the *number of leaves* in the subtree rooted at the branching point of the $(n + 1)$ st insertion. As above, suppose that the $(n + 1)$ st string starts with $w\beta$. Then, examining the first n strings, if there are exactly k strings that begin with $w\alpha$ (again $\alpha = 1 - \beta$), and the other $n - k$ strings do not begin with w , we conclude that $M_n^I = k$.

Since we know from ? that M_n^I follows the logarithmic series distribution plus some fluctuations, then it suffices to prove that M_n has a similar asymptotic distribution. To accomplish this goal, we compare the distribution of M_n in suffix trees to the distribution of M_n^I in independent tries.

Briefly, our proof technique is the following: We let $M(z, u) = \sum_{1 \leq k, n \leq \infty} \mathbf{P}(M_n = k) u^k z^n$ and $M^I(z, u) = \sum_{1 \leq k, n \leq \infty} \mathbf{P}(M_n^I = k) u^k z^n$ denote the bivariate generating functions for M_n and M_n^I , respectively. To study these generating functions, we consider the w 's defined above. Specifically, for $M(z, u)$, we recall from (2) that if w denotes the longest prefix of $X^{(n+1)} = X_{n+1}X_{n+2}X_{n+3} \dots$ that appears as a prefix of any $X^{(i)} = X_iX_{i+1}X_{i+2} \dots$, then M_n enumerates the number of such occurrences of w . This approach to $M(z, u)$ allows us to sum over all $w \in \mathcal{A}^*$ instead of summing over $k, n \in \mathbb{N}$. Similarly, for $M^I(z, u)$, we utilize (3) to see that if w denotes the longest prefix of $X(n + 1) = X_1(n + 1)X_2(n + 1)X_3(n + 1) \dots$ that appears as a prefix of any $X_1(i)X_2(i)X_3(i) \dots$, then M_n^I is precisely the number of such occurrences of w . Therefore, to determine $M^I(z, u)$, we can sum over all $w \in \mathcal{A}^*$ instead of summing over the integers k and n .

We note that the $X^{(i)}$'s are highly dependent on each other. In fact, if $i \geq j$, then $X^{(i)} = X_iX_{i+1}X_{i+2} \dots$ is a substring of $X^{(j)} = X_jX_{j+1}X_{j+2} \dots$. This apparently makes the derivation of the bivariate generating function $M(z, u)$ quite difficult. We overcome this hurdle by succinctly describing the degree to which a suffix of X can overlap with itself. We accomplish this by utilizing the autocorrelation polynomial $S_w(z)$ of a word w , which measures the amount of overlap of a word w with itself. As mentioned above, the autocorrelation polynomial is defined as

$$S_w(z) = \sum_{k \in \mathcal{P}(w)} \mathbf{P}(w_{k+1}^m) z^{m-k} \quad (4)$$

where $\mathcal{P}(w)$ denotes the set of positions k of w satisfying $w_1 \dots w_k = w_{m-k+1} \dots w_m$, that is, w 's prefix of length k is equal to w 's suffix of length k . Via the autocorrelation polynomial, we are able to surmount the difficulties inherent in the overlapping suffixes. Thus, using $S_w(z)$, we obtain a succinct description of the bivariate generating function $M(z, u)$. The autocorrelation polynomial is well-understood; we utilize several results about $S_w(z)$ from [1] and [2]. In particular, when comparing $M(z, u)$ and $M^I(z, u)$, it is extremely useful to note that the autocorrelation polynomial $S_w(z)$ is close to 1 with high probability (for $|w|$ large).

In order to obtain information about the difference of the two BGFs as $Q(z, u) = M(z, u) - M^I(z, u)$, we utilize residue analysis. We make a comparison of the poles of $M(z, u)$ and $M^I(z, u)$ using Cauchy's theorem (integrating with respect to z). As a result, we prove that $Q_n(u) := [z^n]Q(z, u) = O(n^{-\varepsilon})$ uniformly for $|u| \leq p^{-1/2}$ as $n \rightarrow \infty$. Then we use another application of Cauchy's theorem (integrating with respect to u). Specifically, we extract the coefficient $\mathbf{P}(M_n = k) - \mathbf{P}(M_n^I = k) = [u^k z^n]Q(z, u)$ in order to prove our main result.

Theorem 2.1 *There exist $\varepsilon > 0$ and $b > 1$ such that*

$$\mathbf{P}(M_n = k) - \mathbf{P}(M_n^I = k) = O(n^{-\varepsilon} b^{-k}) \quad (5)$$

for large n .

Therefore, the distributions of M_n and M_n^I are asymptotically the same. We conclude that M_n also asymptotically follows the logarithmic series distribution (plus some fluctuations when $\ln p / \ln q$ is rational).

Theorem 2.2 *There exist $\varepsilon > 0$ and $\varepsilon_j > 0$ (for each $j \in \mathbb{N}$) depending on p such that the j th factorial moment of M_n is*

$$\mathbf{E}[(M_n)_j] = \Gamma(j) \frac{q(p/q)^j + p(q/p)^j}{h} + \gamma_j(\log_{1/p} n) + O(n^{-\varepsilon_j}) \quad (6)$$

where γ_j is a periodic function with mean 0 and small modulus if $\ln p / \ln q$ is rational, and otherwise $\gamma_j(x) \rightarrow 0$ as $x \rightarrow \infty$. Also $h = -p \log p - q \log q$ denotes the entropy of the source. The probability generating function of M_n is

$$\mathbf{E}[u^{M_n}] = -\frac{q \ln(1 - pu) + p \ln(1 - qu)}{h} + \gamma(\log_{1/p} n, u) + O(n^{-\varepsilon}), \quad (7)$$

for $|u| \leq p^{-1/2}$ where $\gamma(\cdot, u)$ is a periodic function with mean 0 and small modulus if $\ln p / \ln q$ is rational, and otherwise $\gamma_j(u, x) \rightarrow 0$ (uniformly for $|u| \leq p^{-1/2}$) as $x \rightarrow \infty$. More precisely,

$$\mathbf{E}[u^{M_n}] = \sum_{j=1}^{\infty} \left[\frac{p^j q + q^j p}{jh} + \sum_{k \in \mathbb{Z} \setminus \{0\}} -\frac{e^{2kr\pi i \log_{1/p} n} \Gamma(z_k) (p^j q + q^j p) (z_k)^{\bar{j}}}{j! (p^{-z_k+1} \ln p + q^{-z_k+1} \ln q)} \right] u^j + O(n^{-\varepsilon}) \quad (8)$$

when $\ln p / \ln q = r/t$ for some $r, t \in \mathbb{Z}$, we have $z_k = 2kr\pi i / \ln p$. Therefore, as $n \rightarrow \infty$, we conclude that M_n follows the logarithmic series distribution plus some fluctuations if $\ln p / \ln q = r/t$ is rational, i.e.,

$$\mathbf{P}(M_n = j) = \frac{p^j q + q^j p}{jh} + \sum_{k \neq 0} -\frac{e^{2kr\pi i \log_{1/p} n} \Gamma(z_k) (p^j q + q^j p) (z_k)^{\bar{j}}}{j! (p^{-z_k+1} \ln p + q^{-z_k+1} \ln q)} + O(n^{-\varepsilon}). \quad (9)$$

If $\ln p / \ln q$ is irrational, then M_n asymptotically follows the logarithmic series distribution, without fluctuations.

Note that the average value of M_n is asymptotically $\frac{1}{h}$, and also M_n is highly concentrated around the mean; this property of M_n is very desirable for the error correction scheme described in ?.

3 Proofs

We first derive the bivariate generating functions for M_n and M_n^I , denoted as $M(z, u)$ and $M^I(z, u)$, respectively. Then we prove a few useful lemmas concerning the autocorrelation polynomial. Next, we prove that $M(z, u)$ can be analytically continued from the unit disk to a larger disk. Afterwards, we determine the poles of $M(z, u)$ and $M^I(z, u)$. We write $Q(z, u) = M(z, u) - M^I(z, u)$; we use Cauchy's theorem to that $Q_n(u) := [z^n]Q(z, u) \rightarrow 0$ uniformly for $u \leq p^{-1/2}$ as $n \rightarrow \infty$. Then we apply Cauchy's theorem again to prove that $\mathbf{P}(M_n = k) - \mathbf{P}(M_n^I = k) = [u^k z^n]Q(z, u) = O(n^{-\varepsilon} b^{-k})$ for some $\varepsilon > 0$ and $b > 1$.

We conclude that the distribution of the multiplicity matching parameter M_n is asymptotically the same in suffix trees as in tries built over independent strings, i.e., M_n and M_n^I have asymptotically the same distribution. Therefore, M_n also follows the logarithmic series distribution plus some fluctuations.

3.1 BGF for the Multiplicity Matching Parameter of Independent Tries

First we obtain the bivariate generating function for M_n^I , which is the multiplicity matching parameter for a trie built over the *independent* strings $X(1), \dots, X(n+1)$, where $X(i) = X_1(i)X_2(i)X_3(i) \dots$ and $\{X_j(i) \mid i, j \in \mathbb{N}\}$ is a collection of i.i.d. random variables with $\mathbf{P}(X_j(i) = 0) = p$ and $\mathbf{P}(X_j(i) = 1) = q = 1 - p$. We let w denote the *longest prefix* of both $X(n+1)$ and at least one other string $X(i)$ for some $1 \leq i \leq n$. We write β to denote the $(|w| + 1)$ st character of $X(n+1)$. When $M_n^I = k$, we conclude that exactly k strings $X(i)$ have $w\alpha$ as a prefix, and the other $n - k$ strings $X(i)$ do not have w as a prefix at all. Thus the generating function for M_n^I is exactly

$$M^I(z, u) := \sum_{n=1}^{\infty} \sum_{k=1}^{\infty} \mathbf{P}(M_n^I = k) u^k z^n = \sum_{n=1}^{\infty} \sum_{k=1}^{\infty} \sum_{\substack{w \in \mathcal{A}^* \\ \alpha \in \mathcal{A}}} \mathbf{P}(w\beta) \binom{n}{k} (\mathbf{P}(w\alpha))^k (1 - \mathbf{P}(w))^{n-k} u^k z^n. \quad (10)$$

After simplifying, it follows immediately that

$$M^I(z, u) = \sum_{\substack{w \in \mathcal{A}^* \\ \alpha \in \mathcal{A}}} \frac{u \mathbf{P}(\beta) \mathbf{P}(w)}{1 - z(1 - \mathbf{P}(w))} \frac{z \mathbf{P}(w) \mathbf{P}(\alpha)}{1 - z(1 + u \mathbf{P}(w) \mathbf{P}(\alpha) - \mathbf{P}(w))}. \quad (11)$$

Our reasoning about $M^I(z, u)$ can be applied when we derive generating function $M(z, u)$ for M_n in the next section, but the situation will be more complicated, because the occurrences of w can overlap.

3.2 BGF for the Multiplicity Matching Parameter of Suffix Trees

Now we obtain the bivariate generating function for M_n , which is the multiplicity matching parameter for a suffix tree built over the first $n + 1$ *suffixes* $X^{(1)}, \dots, X^{(n+1)}$ of a string X (i.e., $X^{(i)} = X_i X_{i+1} X_{i+2} \dots$). The bivariate generating function for the multiplicity matching parameter is much more difficult to derive in the dependent (suffix tree) case than in the independent (trie) case, because the suffixes of X are dependent on each other. We let w denote the *longest prefix* of both $X^{(n+1)}$ and at least one $X^{(i)}$ for some $1 \leq i \leq n$. We write β to denote the $(|w| + 1)$ st character of $X^{(n+1)}$; when $M_n = k$, we conclude that exactly k suffixes $X^{(i)}$ have $w\alpha$ as a prefix, and the other $n - k$ strings $X^{(i)}$ do not have w as a prefix at all. Thus, we are

interested in finding strings with exactly k occurrences of $w\alpha$, ended on the right by an occurrence of $w\beta$, with no other occurrences of w at all. This set of words is exactly the language $\mathcal{R}_w\alpha(\mathcal{T}_w^{(\alpha)})^{k-1}\mathcal{T}_w^{(\alpha)}\beta$, where

$$\begin{aligned}\mathcal{R}_w &= \{v \mid v \text{ contains exactly one occurrence of } w, \text{ located at the right end}\} \\ \mathcal{T}_w^{(\alpha)} &= \{v \mid w\alpha v \text{ contains exactly two occurrences of } w, \text{ located at the left and right ends}\} \quad (12)\end{aligned}$$

So, the generating function for M_n is

$$M(z, u) = \sum_{k=1}^{\infty} \sum_{\substack{w \in \mathcal{A}^* \\ \alpha \in \mathcal{A}}} \sum_{s \in \mathcal{R}_w} \mathbf{P}(s\alpha) z^{|s|+1} u \left(\sum_{t \in \mathcal{T}_w^{(\alpha)}} \mathbf{P}(t\alpha) z^{|t|+1} u \right)^{k-1} \sum_{v \in \mathcal{T}_w^{(\alpha)}} \mathbf{P}(v\beta) z^{|v|+1-|w|-1}. \quad (13)$$

After simplifying the geometric sum, this yields

$$M(z, u) = \sum_{\substack{w \in \mathcal{A}^* \\ \alpha \in \mathcal{A}}} u \mathbf{P}(\beta) \frac{R_w(z)}{z^{|w|}} \frac{\mathbf{P}(\alpha) z T_w^{(\alpha)}(z)}{1 - \mathbf{P}(\alpha) z u T_w^{(\alpha)}(z)}. \quad (14)$$

We note that $R_w(z)/z^{|w|} = \mathbf{P}(w)/D_w(z)$ (?), where $D_w(z) = (1-z)S_w(z) + z^m \mathbf{P}(w)$ and where $S_w(z)$ denotes the autocorrelation polynomial for w . Recall that $S_w(z)$ measures the degree to which a word w overlaps with itself, and specifically

$$S_w(z) = \sum_{k \in \mathcal{P}(w)} \mathbf{P}(w_{k+1}^m) z^{m-k} \quad (15)$$

where $\mathcal{P}(w)$ denotes the set of positions k of w satisfying $w_1 \dots w_k = w_{m-k+1} \dots w_m$, that is, w 's prefix of length k is equal to w 's suffix of length k ; also, $m = |w|$. Returning to (14), it follows that

$$M(z, u) = \sum_{\substack{w \in \mathcal{A}^* \\ \alpha \in \mathcal{A}}} \frac{u \mathbf{P}(\beta) \mathbf{P}(w)}{D_w(z)} \frac{\mathbf{P}(\alpha) z T_w^{(\alpha)}(z)}{1 - \mathbf{P}(\alpha) z u T_w^{(\alpha)}(z)}. \quad (16)$$

In order to derive an explicit form of $M(z, u)$, we still need to find $T_w^{(\alpha)}(z)$. If we define

$$\mathcal{M}_w = \{v \mid wv \text{ contains exactly two occurrences of } w, \text{ located at the left and right ends}\} \quad (17)$$

then we observe that $\alpha \mathcal{T}_w^{(\alpha)}$ is exactly the subset of words of \mathcal{M}_w that begin with α ; We use $\mathcal{H}_w^{(\alpha)}$ to denote this subset (i.e., $\mathcal{H}_w^{(\alpha)} = \mathcal{M}_w \cap (\alpha \mathcal{A}^*)$), and thus $\alpha \mathcal{T}_w^{(\alpha)} = \mathcal{H}_w^{(\alpha)}$. So (16) simplifies to

$$M(z, u) = \sum_{\substack{w \in \mathcal{A}^* \\ \alpha \in \mathcal{A}}} \frac{u \mathbf{P}(\beta) \mathbf{P}(w)}{D_w(z)} \frac{H_w^{(\alpha)}(z)}{1 - u H_w^{(\alpha)}(z)}. \quad (18)$$

In order to compute $H_w^{(\alpha)}(z)$, we write $\mathcal{M}_w = \mathcal{H}_w^{(\alpha)} + \mathcal{H}_w^{(\beta)}$, where $\mathcal{H}_w^{(\beta)}$ is the subset of words from \mathcal{M}_w that start with β (i.e., $\mathcal{H}_w^{(\beta)} = \mathcal{M}_w \cap (\beta \mathcal{A}^*)$). (Note that every word of \mathcal{M}_w begins with either α or β , because the empty word $\varepsilon \notin \mathcal{M}_w$.) The following useful lemma is the last necessary ingredient to obtain an explicit formula for $M(z, u)$ from (18).

Lemma 3.1 Let $\mathcal{H}_w^{(\alpha)}$ denote the subset of words from \mathcal{M}_w that start with α . Then

$$H_w^{(\alpha)}(z) = \frac{D_{w\alpha}(z) - (1-z)}{D_w(z)}. \quad (19)$$

Proof We use the concepts and notation from ? and ? throughout. In particular, we define

$$\mathcal{U}_w = \{v \mid wv \text{ contains exactly one occurrence of } w \text{ (located at the left end)}\} \quad (20)$$

and we recall from (12) and (17) above that

$$\begin{aligned} \mathcal{R}_w &= \{v \mid v \text{ contains exactly one occurrence of } w, \text{ located at the right end}\} \\ \mathcal{M}_w &= \{v \mid wv \text{ contains exactly two occurrences of } w, \text{ located at the left and right ends}\} \end{aligned} \quad (21)$$

The following notation is similar but slightly adapted for our proof.

$$\mathcal{U}_w^{(\alpha)} = \{v \mid v \text{ starts with } \alpha, \text{ and } wv \text{ has exactly 1 occurrence of } w\alpha \text{ and no occurrences of } w\beta\}. \quad (22)$$

We note that the set of words with no occurrences of $w\beta$ is exactly $\mathcal{A}^* \setminus \mathcal{R}_{w\beta}(\mathcal{M}_{w\beta})^* \mathcal{U}_{w\beta}$, which has generating function

$$\frac{1}{1-z} - \frac{R_{w\beta}(z)U_{w\beta}(z)}{1-M_{w\beta}(z)}. \quad (23)$$

Now we describe the set of words with no occurrences of $w\beta$ in a different way. The set of words with no occurrences of $w\beta$ and at least one occurrence of $w\alpha$ is exactly $\mathcal{R}_w(\mathcal{H}_w^{(\alpha)})^* \mathcal{U}_w^{(\alpha)}$, which has generating function $R_w(z)U_w^{(\alpha)}(z)/(1-H_w^{(\alpha)}(z))$. The set of words with no occurrences of $w\beta$ and no occurrences of $w\alpha$ is exactly $\mathcal{R}_w + (\mathcal{A}^* \setminus \mathcal{R}_w(\mathcal{M}_w)^* \mathcal{U})$. (Note that the set of such words that end in w is exactly \mathcal{R}_w ; on the other hand, the set of such words that do not end in w is exactly $\mathcal{A}^* \setminus \mathcal{R}_w(\mathcal{M}_w)^* \mathcal{U}$.) So the set of words with no occurrences of $w\alpha$ and no occurrences of $w\beta$ has generating function $R_w(z) + 1/(1-z) - R_w(z)U_w(z)/(1-M_w(z))$. So the set of words with no occurrences of $w\beta$ has generating function

$$\frac{R_w(z)U_w^{(\alpha)}(z)}{1-H_w^{(\alpha)}(z)} + R_w(z) + \frac{1}{1-z} - \frac{R_w(z)U_w(z)}{1-M_w(z)}. \quad (24)$$

Combining (23) and (24), it follows that

$$\frac{1}{1-z} - \frac{R_{w\beta}(z)U_{w\beta}(z)}{1-M_{w\beta}(z)} = \frac{R_w(z)U_w^{(\alpha)}(z)}{1-H_w^{(\alpha)}(z)} + R_w(z) + \frac{1}{1-z} - \frac{R_w(z)U_w(z)}{1-M_w(z)}. \quad (25)$$

Now we find the generating function for $\mathcal{U}_w^{(\alpha)}$. For each word $v \in \mathcal{U}_w^{(\alpha)}$, either wv has exactly one or two occurrences of w . The subset of $\mathcal{U}_w^{(\alpha)}$ of the first type is exactly $V_w^{(\alpha)} := U_w \cap (\alpha\mathcal{A}^*)$, i.e., the subset of words from U_w that start with α . The subset of $\mathcal{U}_w^{(\alpha)}$ of the second type is exactly $\mathcal{H}_w^{(\alpha)}$. We observe that

$$\mathcal{V}_w^{(\alpha)} \cdot \mathcal{A} = (\mathcal{H}_w^{(\alpha)} + \mathcal{V}_w^{(\alpha)}) \setminus \{\alpha\} \quad (26)$$

(see ?), so $V_w^{(\alpha)}(z) = (H_w^{(\alpha)}(z) - \mathbf{P}(\alpha)z)/(z-1)$. Since $\mathcal{U}_w^{(\alpha)} = \mathcal{V}_w^{(\alpha)} + \mathcal{H}_w^{(\alpha)}$, it follows that

$$U_w^{(\alpha)}(z) = \frac{H_w^{(\alpha)}(z) - \mathbf{P}(\alpha)z}{z-1} + H_w^{(\alpha)}(z) = \frac{zH_w^{(\alpha)}(z) - \mathbf{P}(\alpha)z}{z-1}. \quad (27)$$

Recalling equation (25), we see that

$$\frac{1}{1-z} - \frac{R_{w\beta}(z)U_{w\beta}(z)}{1-M_{w\beta}(z)} = \frac{R_w(z)(zH_w^{(\alpha)}(z) - \mathbf{P}(\alpha)z)}{(1-H_w^{(\alpha)}(z))(z-1)} + R_w(z) + \frac{1}{1-z} - \frac{R_w(z)U_w(z)}{1-M_w(z)}. \quad (28)$$

Simplifying, and using $U_w(z) = (1-M_w(z))/(1-z)$ and $U_{w\beta}(z) = (1-M_{w\beta}(z))/(1-z)$ (see ?), it follows that

$$\frac{R_{w\beta}(z)}{R_w(z)} = \frac{z\mathbf{P}(\beta)}{1-H_w^{(\alpha)}(z)}. \quad (29)$$

Solving for $H_w^{(\alpha)}(z)$ and then using $R_w(z) = z^m\mathbf{P}(w)/D_w(z)$ and $R_{w\beta}(z) = z^{m+1}\mathbf{P}(w)\mathbf{P}(\beta)/D_{w\beta}(z)$ (see ?), it follows that

$$H_w^{(\alpha)}(z) = \frac{D_w(z) - D_{w\beta}(z)}{D_w(z)}. \quad (30)$$

Note $D_w(z) - D_{w\beta}(z) = (1-z)S_w(z) + z^m\mathbf{P}(w) - (1-z)S_{w\beta}(z) - z^{m+1}\mathbf{P}(w)\mathbf{P}(\beta) = (1-z)(S_{w\alpha}(z) - 1) + z^{m+1}\mathbf{P}(w)\mathbf{P}(\alpha) = D_{w\alpha}(z) - (1-z)$. Thus, (30) completes the proof of the lemma. \square

Using the lemma above, we finally observe a form of $M(z, u)$ that we summarize below.

Theorem 3.1 *Let $M(z, u) := \sum_{n=1}^{\infty} \sum_{k=1}^{\infty} \mathbf{P}(M_n = k)u^k z^n$ denote the bivariate generating function for M_n , the multiplicity matching parameter of a suffix tree built over the first $n+1$ suffixes $X^{(1)}, \dots, X^{(n+1)}$ of a string X . Then*

$$M(z, u) = \sum_{\substack{w \in \mathcal{A}^* \\ \alpha \in \mathcal{A}}} \frac{u\mathbf{P}(\beta)\mathbf{P}(w)}{D_w(z)} \frac{D_{w\alpha}(z) - (1-z)}{D_w(z) - u(D_{w\alpha}(z) - (1-z))} \quad (31)$$

for $|u| < 1$ and $|z| < 1$. Here $D_w(z) = (1-z)S_w(z) + z^m\mathbf{P}(w)$, and $S_w(z)$ denotes the autocorrelation polynomial for w , defined in (1).

3.3 On the Autocorrelation Polynomial

Throughout the rest of our analysis we assume that, without loss of generality, $p \geq q$. Note that $p \leq \sqrt{p} < 1$, so there exists $\rho > 1$ such that $\rho\sqrt{p} < 1$ (and thus $\rho p < 1$ too). Finally, define $\delta = \sqrt{p}$. We establish a few lemmas about the autocorrelation polynomial that will be important for our analysis. Recall that the autocorrelation polynomial is $S_w(z) = \sum_{k \in \mathcal{P}(w)} \mathbf{P}(w_{k+1}^m)z^{m-k}$, where $\mathcal{P}(w)$ denotes the set of positions k of w satisfying $w_1 \dots w_k = w_{m-k+1} \dots w_m$, that is, w 's prefix of length k is equal to w 's suffix of length k .

The autocorrelation polynomial $S_w(z)$ has a $\mathbf{P}(w_{k+1}^m)z^{m-k}$ term if and only if w has an overlap with itself of length k . Since each word w overlaps with itself trivially, then every autocorrelation polynomial has a constant term (i.e., $z^{m-m} = z^0 = 1$ term). With high probability, however, w has very few large nontrivial overlaps with itself. Therefore, with high probability, all nontrivial overlaps of w with itself are small; such overlaps correspond to high-degree terms of $S_w(z)$. Therefore, when w is a randomly chosen

long word, then $S_w(z)$ is very close to 1 with very high probability. The first lemma makes this notion mathematically precise.

Lemma 3.2 *If $\theta = (1 - \rho p)^{-1} > 1$, then*

$$\sum_{w \in \mathcal{A}^k} \llbracket |S_w(\rho) - 1| \leq (\rho\delta)^k \theta \rrbracket \mathbf{P}(w) \geq 1 - \delta^k \theta \quad (32)$$

where $\llbracket A \rrbracket = 1$ if A holds, and $\llbracket A \rrbracket = 0$ otherwise.

Proof Our proof is the one given in ?. Note that $S_w(z) - 1$ has a term of degree $i \leq j$ if and only if $m - i \in \mathcal{P}(w)$ with $1 \leq i \leq j$. Therefore, for each such i and each $w_1 \dots w_i$, there is *exactly one* word $w_{i+1} \dots w_k$ such that $S_w(z) - 1$ has a term of degree i . Therefore, for fixed j and k ,

$$\begin{aligned} & \sum_{w \in \mathcal{A}^k} \llbracket S_w(z) - 1 \text{ has a term of degree } \leq j \rrbracket \mathbf{P}(w) \\ & \leq \sum_{1 \leq i \leq j} \sum_{w_1, \dots, w_i \in \mathcal{A}^i} \mathbf{P}(w_1 \dots w_i) \sum_{w_{i+1}, \dots, w_k \in \mathcal{A}^{k-i}} \llbracket S_w(z) - 1 \text{ has a term of degree } i \rrbracket \mathbf{P}(w_{i+1} \dots w_k) \\ & \leq \sum_{1 \leq i \leq j} \sum_{w_1, \dots, w_i \in \mathcal{A}^i} \mathbf{P}(w_1 \dots w_i) p^{k-i} = \sum_{1 \leq i \leq j} p^{k-i} \leq \frac{p^{k-j}}{1-p} \end{aligned} \quad (33)$$

We use $j = \lfloor k/2 \rfloor$. Thus $\sum_{w \in \mathcal{A}^k} \llbracket \text{all terms of } S_w(z) - 1 \text{ have degree } > \lfloor k/2 \rfloor \rrbracket \mathbf{P}(w) \geq 1 - \delta^k \theta$.

Note that, if all terms of $S_w(z) - 1$ have degree $> \lfloor k/2 \rfloor$, then

$$|S_w(\rho) - 1| \leq \sum_{i > \lfloor k/2 \rfloor} (\rho p)^i = \frac{(\rho p)^{\lfloor k/2 \rfloor + 1}}{1 - \rho p} \leq \frac{(\rho p)^{k/2}}{1 - \rho p} \leq \frac{\rho^k p^{k/2}}{1 - \rho p} = (\rho\delta)^k \theta. \quad (34)$$

This completes the proof of the lemma. \square

Using this lemma, we can quickly obtain another result that is similar but slightly stronger.

First consider words w such that $|S_w(\rho) - 1| \leq (\rho\delta)^k \theta$. Write $S_w(z) = \sum_{i=0}^{k-1} a_i z^i$ and $S_{w\alpha}(z) = \sum_{i=0}^k b_i z^i$. Observe that either $b_i = 0$ or $b_i = a_i$. The following lemma follows immediately:

Lemma 3.3 *If $\theta = (1 - \rho p)^{-1} + 1$ and $\alpha \in \mathcal{A}$, then*

$$\sum_{w \in \mathcal{A}^k} \llbracket \max\{|S_w(\rho) - 1|, |S_{w\alpha}(\rho) - 1|\} \leq (\rho\delta)^k \theta \rrbracket \mathbf{P}(w) \geq 1 - \delta^k \theta. \quad (35)$$

Also, the autocorrelation polynomial is never too small. In fact

Lemma 3.4 *Define $c = 1 - \rho\sqrt{p} > 0$. Then there exists an integer $K \geq 1$ such that, for $|w| \geq K$ and $|z| \leq \rho$ and $|u| \leq \delta^{-1}$,*

$$|S_w(z) - u S_{w\alpha}(z) + u| \geq c. \quad (36)$$

Proof The proof consists of considering several cases. The only condition for K is $(1 + \delta^{-1}) \frac{(\rho p)^{K/2}}{1 - \rho p} \leq c/2$. The analysis is not difficult; all details are presented in ?. \square

3.4 Analytic Continuation

Our goal in this section is to prove the following:

Theorem 3.2 *The generating function $M(z, u)$ can be analytically continued for $|u| \leq \delta^{-1}$ and $|z| < 1$.*

The proof requires several lemmas and observations. We always assume $|u| \leq \delta^{-1}$.

Lemma 3.5 *If $0 < r < 1$, then there exists $C > 0$ and an integer K_1 (both depending on r) such that*

$$|D_w(z) - u(D_{w\alpha}(z) - (1-z))| \geq C \quad (37)$$

for $|w| \geq K_1$ and $|z| \leq r$ (and, as before, $|u| \leq \delta^{-1}$).

Proof Consider the K and c defined in Lemma 3.4, which tells us that, for all $|w| \geq K$, we have

$$|S_w(z) - uS_{w\alpha}(z) + u| \geq c \quad (38)$$

for $|z| \leq \rho$. So, for $|w| \geq K$, we have $|D_w(z) - u(D_{w\alpha}(z) - (1-z))| \geq (1-r)c - r^m p^m (1 - \delta^{-1}rp)$. Note that $r^m p^m (1 - \delta^{-1}rp) \rightarrow 0$ as $m \rightarrow \infty$. Therefore, replacing K by a larger K_1 if necessary, we can without loss of generality assume that $r^m p^m (1 - \delta^{-1}rp) \leq (1-r)c/2$. So we define $C = (1-r)c/2$, and the result follows immediately. \square

Now we can strengthen the previous lemma by dropping the “ K_1 ”, i.e., by not requiring w to be a long word:

Lemma 3.6 *If $0 < r < 1$, then there exists $C > 0$ (depending on r) such that*

$$|D_w(z) - u(D_{w\alpha}(z) - (1-z))| \geq C \quad (39)$$

for $|z| \leq r$ (and, as before, $|u| \leq \delta^{-1}$).

Proof Consider the K_1 defined in Lemma 3.5. Let C_0 denote the “ C ” from Lemma 3.5. There are only finitely many w 's with $|w| < K_1$, say w_1, \dots, w_i . For each such w_j (with $1 \leq j \leq i$), we note that $D_{w_j}(z) - u(D_{w_j\alpha}(z) - (1-z)) \neq 0$ for $|z| \leq r$ and $|u| \leq \delta^{-1}$, so there exists $C_j > 0$ such that $|D_{w_j}(z) - u(D_{w_j\alpha}(z) - (1-z))| \geq C_j$ for all $|z| \leq r$ and $|u| \leq \delta^{-1}$. Finally, we define $C = \min\{C_0, C_1, \dots, C_i\}$. \square

Finally, we prove Theorem 3.2.

Proof Consider $|z| \leq r < 1$. We proved in Lemma 3.6 there exists $C > 0$ depending on r such that, for all $|u| \leq \delta^{-1}$, we have $\frac{1}{|D_w(z) - u(D_{w\alpha}(z) - (1-z))|} \leq \frac{1}{C}$. Setting $u = 0$, we also have $\frac{1}{|D_w(z)|} \leq \frac{1}{C}$. Thus

$$|M(z, u)| \leq \frac{\mathbf{P}(\beta)\delta^{-1}}{C^2} \sum_{\alpha \in \mathcal{A}} \sum_{w \in \mathcal{A}^*} \mathbf{P}(w) |D_{w\alpha}(z) - (1-z)|. \quad (40)$$

Now we use Lemma 3.3. Consider w and α with $\max\{|S_w(\rho) - 1|, |S_{w\alpha}(\rho) - 1|\} \leq (\rho\delta)^m \theta$. It follows immediately that

$$|D_{w\alpha}(z) - (1-z)| = |(1-z)(S_{w\alpha}(z) - 1) + z^{m+1}\mathbf{P}(w)\mathbf{P}(\alpha)| \leq (1+r)(\rho\delta)^m \theta + r^{m+1}p^m p = O(s^m), \quad (41)$$

where $s = \max\{\rho\delta, r\rho\}$. Now consider the other w 's and α 's. We have

$$|D_{w\alpha}(z) - (1-z)| = |(1-z)(S_{w\alpha}(z) - 1) + z^{m+1}\mathbf{P}(w)\mathbf{P}(\alpha)| \leq \frac{(1+r)p\rho}{1-p\rho} + r^{m+1}p^m p \leq \frac{(1+r)p\rho}{1-p\rho} + 1, \quad (42)$$

so we define $C_1 = \frac{(1+r)p\rho}{1-p\rho} + 1$ to be a value which only depends on r (recall that r is fixed here). Thus

$$\begin{aligned} |M(z, u)| &\leq \frac{\mathbf{P}(\beta)\delta^{-1}}{C^2} \sum_{\alpha \in \mathcal{A}} \sum_{m \geq 0} \sum_{w \in \mathcal{A}^m} |\mathbf{P}(w)(D_{w\alpha}(z) - (1-z))| \\ &\leq \frac{\mathbf{P}(\beta)\delta^{-1}}{C^2} \sum_{\alpha \in \mathcal{A}} \sum_{m \geq 0} |(1 - \delta^m \theta)O(s^m) + \delta^m \theta C_1| \leq \frac{\mathbf{P}(\beta)\delta^{-1}}{C^2} \sum_{\alpha \in \mathcal{A}} \sum_{m \geq 0} O(s^m) = O(1) \end{aligned} \quad (43)$$

and this completes the proof of the theorem. \square

3.5 Singularity Analysis

We first determine (for $|u| \leq \delta^{-1}$) the zeroes of $D_w(z) - u(D_{w\alpha}(z) - (1-z))$ and (in particular) the zeroes of $D_w(z)$.

Lemma 3.7 *There exists an integer $K_2 \geq 1$ such that, for u fixed (with $|u| \leq \delta^{-1}$) and $|w| \geq K_2$, there is exactly one root of $D_w(z) - u(D_{w\alpha}(z) - (1-z))$ in the closed disk $\{z \mid |z| \leq \rho\}$.*

Proof Let K and c be defined as in Lemma 3.4. Without loss of generality (replacing K by a larger K_2 , if necessary), we can also assume that $2(p\rho)^{K_2} < c(\rho - 1)$ and $K_2 \geq K_1$ (where K_1 is defined in Lemma 3.5). Also, we can choose K_2 large enough (for use later) such that $\exists c_2 > 0$ with

$$\rho(1 - p^{K_2}(1 + \delta^{-1}p)) - 1 > c_2 \quad \text{and thus} \quad \rho(1 - p^{K_2}) - 1 > c_2. \quad (44)$$

We recall $0 < p\rho\delta^{-1} < 1$, and thus $0 < 1 - p\rho\delta^{-1} < 1$. Since $|u| < \delta^{-1}$ and $|z| \leq \rho$, then for $|w| \geq K_2$ we have $|\mathbf{P}(w)z^m(1 - uz\mathbf{P}(\alpha))| \leq (p\rho)^m(1 + \delta^{-1}p\rho) \leq 2(p\rho)^m < c(\rho - 1) \leq |(S_w(z) - uS_{w\alpha}(z) + u)(\rho - 1)|$. Therefore, for z on the circle $\{z \mid |z| = \rho\}$, we have $|\mathbf{P}(w)z^m(1 - uz\mathbf{P}(\alpha))| < |(S_w(z) - uS_{w\alpha}(z) + u)(z - 1)|$. Equivalently,

$$|(D_w(z) - u(D_{w\alpha}(z) - (1-z))) - ((S_w(z) - uS_{w\alpha}(z) + u)(z - 1))| < |(S_w(z) - uS_{w\alpha}(z) + u)(z - 1)|. \quad (45)$$

Therefore, by Rouché's Theorem, $D_w(z) - u(D_{w\alpha}(z) - (1-z))$ and $(S_w(z) - uS_{w\alpha}(z) + u)(z - 1)$ have the same number of zeroes inside the disk $\{z \mid |z| \leq \rho\}$. Since $|S_w(z) - uS_{w\alpha}(z) + u| \geq c$ inside this disk, we conclude that $(S_w(z) - uS_{w\alpha}(z) + u)(z - 1)$ has exactly one root in the disk. It follows that $D_w(z) - u(D_{w\alpha}(z) - (1-z))$ also has exactly one root in the disk. \square

When $u = 0$, this lemma implies (for $|w| \geq K_2$) that $D_w(z)$ has exactly one root in the disk $\{z \mid |z| \leq \rho\}$. Let A_w denote this root, and let $B_w = D'_w(A_w)$. Also let $C_w(u)$ denote the root of $D_w(z) - u(D_{w\alpha}(z) - (1-z))$ in the closed disk $\{z \mid |z| \leq \rho\}$. Finally, we define

$$E_w(u) := \left(\frac{d}{dz} (D_w(z) - u(D_{w\alpha}(z) - (1-z))) \right) \Big|_{z=C_w} = D'_w(C_w) - u(D'_{w\alpha}(C_w) + 1). \quad (46)$$

We have precisely determined the singularities of $M(z, u)$. Next, we make a comparison of $M(z, u)$ to $M^l(z, u)$, in order to show that M_n and M_n^l have asymptotically similar behaviors.

3.6 Comparing Suffix Trees to Tries

Now we define

$$Q(z, u) = M(z, u) - M^I(z, u). \quad (47)$$

Using the notation from (11) and (31), if we write

$$\begin{aligned} M_{w,\alpha}^I(z, u) &= \frac{u\mathbf{P}(\beta)\mathbf{P}(w)}{1-z(1-\mathbf{P}(w))} \frac{z\mathbf{P}(w)\mathbf{P}(\alpha)}{1-z(1+u\mathbf{P}(w)\mathbf{P}(\alpha)-\mathbf{P}(w))} \\ M_{w,\alpha}(z, u) &= \frac{u\mathbf{P}(\beta)\mathbf{P}(w)}{D_w(z)} \frac{D_{w\alpha}(z) - (1-z)}{D_w(z) - u(D_{w\alpha}(z) - (1-z))} \end{aligned} \quad (48)$$

then we have proven that

$$Q(z, u) = \sum_{\substack{w \in \mathcal{A}^* \\ \alpha \in \mathcal{A}}} (M_{w,\alpha}(z, u) - M_{w,\alpha}^I(z, u)). \quad (49)$$

We also define $Q_n(u) = [z^n]Q(z, u)$. We denote the contribution to $Q_n(u)$ from a specific w and α as $Q_n^{(w,\alpha)}(u) = [z^n](M_{w,\alpha}(z, u) - M_{w,\alpha}^I(z, u))$. Then we observe that

$$Q_n^{(w,\alpha)}(u) = \frac{1}{2\pi i} \oint (M_{w,\alpha}(z, u) - M_{w,\alpha}^I(z, u)) \frac{dz}{z^{n+1}} \quad (50)$$

where the path of integration is a circle about the origin with counterclockwise orientation.

We define

$$I_{w,\alpha}(\rho, u) = \frac{1}{2\pi i} \int_{|z|=\rho} (M_{w,\alpha}(z, u) - M_{w,\alpha}^I(z, u)) \frac{dz}{z^{n+1}}. \quad (51)$$

By Cauchy's theorem, we observe that the contribution to $Q_n(u)$ from a specific w and α is exactly

$$\begin{aligned} Q_n^{(w,\alpha)}(u) &= I_{w,\alpha}(\rho, u) - \operatorname{Res}_{z=A_w} \frac{M_{w,\alpha}(z, u)}{z^{n+1}} - \operatorname{Res}_{z=C_w(u)} \frac{M_{w,\alpha}(z, u)}{z^{n+1}} \\ &\quad + \operatorname{Res}_{z=1/(1-\mathbf{P}(w))} \frac{M_{w,\alpha}^I(z, u)}{z^{n+1}} + \operatorname{Res}_{z=1/(1+u\mathbf{P}(w)\mathbf{P}(\alpha)-\mathbf{P}(w))} \frac{M_{w,\alpha}^I(z, u)}{z^{n+1}}. \end{aligned} \quad (52)$$

To simplify this expression, note that

$$\begin{aligned} \operatorname{Res}_{z=A_w} \frac{M_{w,\alpha}(z, u)}{z^{n+1}} &= -\frac{\mathbf{P}(\beta)\mathbf{P}(w)}{B_w} \frac{1}{A_w^{n+1}} \\ \operatorname{Res}_{z=C_w(u)} \frac{M_{w,\alpha}(z, u)}{z^{n+1}} &= \frac{\mathbf{P}(\beta)\mathbf{P}(w)}{E_w(u)} \frac{1}{C_w(u)^{n+1}} \\ \operatorname{Res}_{z=1/(1-\mathbf{P}(w))} \frac{M_{w,\alpha}^I(z, u)}{z^{n+1}} &= \mathbf{P}(\beta)\mathbf{P}(w)(1-\mathbf{P}(w))^n \\ \operatorname{Res}_{z=1/(1+u\mathbf{P}(w)\mathbf{P}(\alpha)-\mathbf{P}(w))} \frac{M_{w,\alpha}^I(z, u)}{z^{n+1}} &= -\mathbf{P}(\beta)\mathbf{P}(w)(1+u\mathbf{P}(w)\mathbf{P}(\alpha)-\mathbf{P}(w))^n \end{aligned} \quad (53)$$

It follows from (52) that

$$Q_n^{(w,\alpha)}(u) = I_{w,\alpha}(\rho, u) + \frac{\mathbf{P}(\beta)\mathbf{P}(w)}{B_w} \frac{1}{A_w^{n+1}} - \frac{\mathbf{P}(\beta)\mathbf{P}(w)}{E_w(u)} \frac{1}{C_w(u)^{n+1}}$$

$$+ \mathbf{P}(\beta)\mathbf{P}(w)(1 - \mathbf{P}(w))^n - \mathbf{P}(\beta)\mathbf{P}(w)(1 + u\mathbf{P}(w)\mathbf{P}(\alpha) - \mathbf{P}(w))^n. \quad (54)$$

We next determine the contribution of the $z = A_w$ terms of $M(z, u)$ and the $z = 1/(1 - \mathbf{P}(w))$ terms of $M^I(z, u)$ to the difference $Q_n(u) = [z^n](M(z, u) - M^I(z, u))$.

Lemma 3.8 *The “ A_w terms” and the “ $1/(1 - \mathbf{P}(w))$ terms” (for $|w| \geq K_2$) altogether have only $O(n^{-\varepsilon})$ contribution to $Q_n(u)$, i.e.,*

$$\sum_{\substack{|w| \geq K_2 \\ \alpha \in \mathcal{A}}} \left(-\operatorname{Res}_{z=A_w} \frac{M_{w,\alpha}(z, u)}{z^{n+1}} + \operatorname{Res}_{z=1/(1-\mathbf{P}(w))} \frac{M_{w,\alpha}^I(z, u)}{z^{n+1}} \right) = O(n^{-\varepsilon}), \quad (55)$$

for some $\varepsilon > 0$.

Proof We define

$$f_w(x) = \frac{1}{A_w^{x+1} B_w} + (1 - \mathbf{P}(w))^x \quad (56)$$

for x real. So by (53) it suffices to prove that

$$\sum_{\substack{|w| \geq K_2 \\ \alpha \in \mathcal{A}}} \mathbf{P}(\beta)\mathbf{P}(w)f_w(x) = O(x^{-\varepsilon}). \quad (57)$$

Note that $\sum_{\substack{|w| \geq K_2 \\ \alpha \in \mathcal{A}}} \mathbf{P}(\beta)\mathbf{P}(w)f_w(x)$ is absolutely convergent for all x . Also $\bar{f}_w(x) = f_w(x) - f_w(0)e^{-x}$ is exponentially decreasing when $x \rightarrow +\infty$ and is $O(x)$ when $x \rightarrow 0$ (notice that we utilize the $f_w(0)e^{-x}$ term in order to make sure that $\bar{f}_w(x) = O(x)$ when $x \rightarrow 0$; this provides a fundamental strip for the Mellin transform in the next step). Therefore, its Mellin transform $\bar{f}_w^*(s) = \int_0^\infty \bar{f}_w(x)x^{s-1} dx$ is well-defined for $\Re(s) > -1$ (see ? and ?). We compute

$$\bar{f}_w^*(s) = \Gamma(s) \left(\frac{(\log A_w)^{-s} - 1}{A_w B_w} + (-\log(1 - \mathbf{P}(w)))^{-s} - 1 \right) \quad (58)$$

where Γ denotes the Euler gamma function, and we note that

$$\begin{aligned} (\log A_w)^{-s} &= \left(\frac{\mathbf{P}(w)}{S_w(1)} \right)^{-s} (1 + O(\mathbf{P}(w))) \\ (-\log(1 - \mathbf{P}(w)))^{-s} &= \mathbf{P}(w)^{-s} (1 + O(\mathbf{P}(w))) \end{aligned} \quad (59)$$

Also

$$\begin{aligned} A_w &= 1 + \frac{1}{S_w(1)} \mathbf{P}(w) + O(\mathbf{P}(w)^2) \\ B_w &= -S_w(1) + \left(-\frac{2S_w'(1)}{S_w(1)} + m \right) \mathbf{P}(w) + O(\mathbf{P}(w)^2) \end{aligned} \quad (60)$$

Therefore

$$\frac{1}{A_w B_w} = -\frac{1}{S_w(1)} + O(|w|\mathbf{P}(w)) \quad (61)$$

So $\bar{f}_w^*(s) = \Gamma(s) \left(\left(-\frac{1}{S_w(1)} + O(|w|\mathbf{P}(w)) \right) \left(\left(\frac{\mathbf{P}(w)}{S_w(1)} \right)^{-s} (1 + O(\mathbf{P}(w))) - 1 \right) + \mathbf{P}(w)^{-s} (1 + O(\mathbf{P}(w))) - 1 \right) = \Gamma(s) \left(\mathbf{P}(w)^{-s} (-S_w(1)^{s-1} + 1 + O(|w|\mathbf{P}(w))) + \frac{1}{S_w(1)} - 1 + O(|w|\mathbf{P}(w)) \right).$

We define $g^*(s) = \sum_{\substack{|w| \geq K_2 \\ \alpha \in \mathcal{A}}} \mathbf{P}(\beta) \mathbf{P}(w) \bar{f}_w^*(s)$. Then we compute

$$g^*(s) = \sum_{\alpha \in \mathcal{A}} \mathbf{P}(\beta) \sum_{|w| \geq K_2} \mathbf{P}(w) \bar{f}_w^*(s) = \sum_{\alpha \in \mathcal{A}} \mathbf{P}(\beta) \Gamma(s) \sum_{m=K_2}^{\infty} \left(\sup\{q^{-\Re(s)}, 1\} \delta \right)^m O(1), \quad (62)$$

where the last equality is true because $1 \geq p^{-\Re(s)} \geq q^{-\Re(s)}$ when $\Re(s)$ is negative, and also because $q^{-\Re(s)} \geq p^{-\Re(s)} \geq 1$ when $\Re(s)$ is positive. We always have $\delta < 1$. Also, there exists $c > 0$ such that $q^{-c} \delta < 1$. Therefore, $g^*(s)$ is analytic in $\Re(s) \in (-1, c)$. Working in this strip, we choose ε with $0 < \varepsilon < c$. Then we have

$$\sum_{\substack{|w| \geq K_2 \\ \alpha \in \mathcal{A}}} \mathbf{P}(\beta) \mathbf{P}(w) f_w(x) = \frac{1}{2\pi i} \int_{\varepsilon - i\infty}^{\varepsilon + i\infty} g^*(s) x^{-s} ds + \sum_{\substack{|w| \geq K_2 \\ \alpha \in \mathcal{A}}} \mathbf{P}(\beta) \mathbf{P}(w) f_w(0) e^{-x}. \quad (63)$$

Majorizing under the integral, we see that the first term is $O(x^{-\varepsilon})$ since $g^*(s)$ is analytic in the strip $\Re(s) \in (-1, c)$ (and $-1 < \varepsilon < c$). Also, the second term is $O(e^{-x})$. This completes the proof of the lemma. \square

Now we bound the contribution to $Q_n(u)$ from the $C_w(u)$ terms of $M(z, u)$ and the $z = 1/(1 + u\mathbf{P}(w)\mathbf{P}(\alpha) - \mathbf{P}(w))$ terms of $M^l(z, u)$.

Lemma 3.9 *The “ $C_w(u)$ terms” and the “ $1/(1 + u\mathbf{P}(w)\mathbf{P}(\alpha) - \mathbf{P}(w))$ terms” (for $|w| \geq K_2$) altogether have only $O(n^{-\varepsilon})$ contribution to $Q_n(u)$, for some $\varepsilon > 0$. More precisely,*

$$\sum_{\substack{|w| \geq K_2 \\ \alpha \in \mathcal{A}}} \left(-\operatorname{Res}_{z=C_w(u)} \frac{M_{w,\alpha}(z, u)}{z^{n+1}} + \operatorname{Res}_{z=1/(1+u\mathbf{P}(w)\mathbf{P}(\alpha)-\mathbf{P}(w))} \frac{M_{w,\alpha}^l(z, u)}{z^{n+1}} \right) = O(n^{-\varepsilon}). \quad (64)$$

Proof The proof technique is the same as the one for Lemma 3.8 above. \square

Next we prove that the $I_{w,\alpha}(\rho, u)$ terms in (54) have $O(n^{-\varepsilon})$ contribution to $Q_n(u)$.

Lemma 3.10 *The “ $I_{w,\alpha}(\rho, u)$ terms” (for $|w| \geq K_2$) altogether have only $O(n^{-\varepsilon})$ contribution to $Q_n(u)$, for some $\varepsilon > 0$. More precisely,*

$$\sum_{\substack{|w| \geq K_2 \\ \alpha \in \mathcal{A}}} I_{w,\alpha}(\rho, u) = O(n^{-\varepsilon}). \quad (65)$$

Proof Here we only sketch the proof. A rigorous proof is given in ?. Recall that

$$I_{w,\alpha}(\rho, u) = \frac{1}{2\pi i} \int_{|z|=\rho} u \mathbf{P}(\beta) \mathbf{P}(w) \left(\frac{1}{D_w(z)} \frac{D_{w\alpha}(z) - (1-z)}{D_w(z) - u(D_{w\alpha}(z) - (1-z))} - \frac{1}{1-z(1-\mathbf{P}(w))} \frac{z \mathbf{P}(w) \mathbf{P}(\alpha)}{1-z(1+u\mathbf{P}(w)\mathbf{P}(\alpha)-\mathbf{P}(w))} \right) \frac{dz}{z^{n+1}}. \quad (66)$$

By Lemma 3.7, K_2 was selected to be sufficiently large such that $(\rho p)^m(1 - \delta^{-1}\rho p) \leq (\rho - 1)c/2$. Thus, writing $C_1 = (\rho - 1)c/2$, we have $1/|D_w(z) - u(D_{w\alpha}(z) - (1 - z))| \leq 1/C_1$ and thus $1/|D_w(z)| \leq 1/C_1$. Also $1/|1 - z(1 - \mathbf{P}(w))| \leq 1/c_2$ and $1/|1 - z(1 + u\mathbf{P}(w)\mathbf{P}(\alpha) - \mathbf{P}(w))| \leq 1/c_2$ by (44). So we obtain

$$|I_{w,\alpha}(\rho, u)| = O(\rho^{-n})\mathbf{P}(w)(S_{w\alpha}(\rho) - 1) + O(\rho^{-n})\mathbf{P}(w)O((\rho p)^m). \quad (67)$$

Thus, by Lemma 3.3, $\sum_{\alpha \in \mathcal{A}} \sum_{|w|=m} |I_{w,\alpha}(\rho, u)| = O(\rho^{-n})O((\rho\delta)^m)$. We conclude $\sum_{\substack{|w| \geq K_2 \\ \alpha \in \mathcal{A}}} |I_{w,\alpha}(\rho, u)| = O(\rho^{-n})$, and the lemma follows. \square

Finally, we consider the contribution to $Q_n(u)$ from small words $|w|$. Basically, we prove that $|w|$ has a normal distribution with mean $\frac{1}{h} \log n$ and variance $\theta \log n$, where $h = -p \log p - q \log q$ denotes the entropy of the source, and θ is a constant. Therefore, $|w| \leq K_2$ is extremely unlikely, and as a result, the contribution to $Q_n(u)$ from words w with $|w| \leq K_2$ is very small.

Lemma 3.11 *The terms $\sum_{\substack{|w| < K_2 \\ \alpha \in \mathcal{A}}} (M_{w,\alpha}(z, u) - M_{w,\alpha}^I(z, u))$ altogether have only $O(n^{-\varepsilon})$ contribution to $Q_n(u)$.*

Proof Let D_n denote the depth of the $(n + 1)$ st insertion in a suffix tree, i.e., $D_n < k$ if and only if

$$X_{n+1} \dots X_{n+k} \neq X_{i+1} \dots X_{i+k} \quad \text{for all } 0 \leq i < n \quad (68)$$

i.e., $D_n = |w|$ in the notation of Section 3.2. Similarly, let D_n^I denote the depth of the $(n + 1)$ st insertion in a trie built over $n + 1$ independent strings, i.e., $D_n^I < k$ if and only if

$$X_1(n+1) \dots X_k(n+1) \neq X_1(i) \dots X_k(i) \quad \text{for all } 1 \leq i \leq n \quad (69)$$

i.e., $D_n^I = |w|$ in the notation of Section 3.1.

Therefore,

$$[z^n] \sum_{\substack{|w| < K_2 \\ \alpha \in \mathcal{A}}} (M_{w,\alpha}(z, u) - M_{w,\alpha}^I(z, u)) = \sum_{i < K_2} \sum_{k=1}^n (\mathbf{P}(M_n = k \ \& \ D_n = i) - \mathbf{P}(M_n^I = k \ \& \ D_n^I = i)) u^k. \quad (70)$$

Noting that $\mathbf{P}(M_n = k \ \& \ D_n = i) \leq \mathbf{P}(D_n = i)$ and $\mathbf{P}(M_n^I = k \ \& \ D_n^I = i) \leq \mathbf{P}(D_n^I = i)$, it follows that

$$[z^n] \sum_{\substack{|w| < K_2 \\ \alpha \in \mathcal{A}}} |M_{w,\alpha}(z, u) - M_{w,\alpha}^I(z, u)| \leq \sum_{i < K_2} \sum_{k=1}^n (\mathbf{P}(D_n = i) + \mathbf{P}(D_n^I = i)) |u|^k. \quad (71)$$

In ?, the typical depth D_{n+1}^T in a trie built over $n + 1$ independent strings was shown to be asymptotically normal with mean $\frac{1}{h} \log(n + 1)$ and variance $\theta \log(n + 1)$. We observe that D_n^I (defined in (69)) and D_{n+1}^T have the same distribution; to see this, observe that $\mathbf{P}(D_n^I < k) = \sum_{|w|=k} \mathbf{P}(w)(1 - \mathbf{P}(w))^n = \mathbf{P}(D_{n+1}^T < k)$. Therefore, D_n^I is also asymptotically normal with mean $\frac{1}{h} \log n$ and variance $\theta \log n$. In ?, we rigorously prove that D_n^I and D_n have asymptotically the same distribution, namely, a normal distribution with mean $\frac{1}{h} \log(n + 1)$ and variance $\theta \log(n + 1)$. Therefore, considering (71) (and noting that K_2 is a constant), it follows that

$$[z^n] \sum_{\substack{|w| < K_2 \\ \alpha \in \mathcal{A}}} |M_{w,\alpha}(z, u) - M_{w,\alpha}^I(z, u)| = O(n^{-\varepsilon}). \quad (72)$$

This completes the proof of the lemma. \square

All contributions to (54) have now been analyzed. We are finally prepared to summarize our results.

3.7 Summary and Conclusion

Combining the last four lemmas, we see that $Q_n(u) = O(n^{-\varepsilon})$ uniformly for $|u| \leq \delta^{-1}$, where $\delta^{-1} > 1$. For ease of notation, we define $b = \delta^{-1}$. Finally, we apply Cauchy's theorem again. We compute

$$\mathbf{P}(M_n = k) - \mathbf{P}(M_n^I = k) = [u^k z^n] Q(z, u) = [u^k] Q_n(u) = \frac{1}{2\pi i} \int_{|u|=b} \frac{Q_n(u)}{u^{k+1}} du. \quad (73)$$

Since $Q_n(u) = O(n^{-\varepsilon})$, it follows that

$$|\mathbf{P}(M_n = k) - \mathbf{P}(M_n^I = k)| \leq \frac{1}{|2\pi i|} (2\pi b) \frac{O(n^{-\varepsilon})}{b^{k+1}} = O(n^{-\varepsilon} b^{-k}). \quad (74)$$

So Theorem 2.1 holds. It follows that M_n and M_n^I have asymptotically the same distribution, and therefore M_n and M_n^I asymptotically have the same factorial moments. The main result of ? gives the asymptotic distribution and factorial moments of M_n^I . As a result, Theorem 2.2 follows immediately. Therefore, M_n follows the logarithmic series distribution, i.e., $\mathbf{P}(M_n = j) = \frac{p^j q + q^j p}{j! h}$ (plus some small fluctuations if $\ln p / \ln q$ is rational).

Acknowledgements

The second author warmly thanks Philippe Flajolet for his hospitality during his recent stay at INRIA.