

Joint String Complexity for Markov Sources

DEDICATED TO OUR FRIEND AND MENTOR PHILIPPE FLAJOLET

Philippe Jacquet
Bell Labs
Alcatel-Lucent
91620 Nozay
France
philippe.jacquet@alcatel-lucent.com

Wojciech Szpankowski¹
Department of Computer Science
Purdue University
W. Lafayette, IN 47907
U.S.A.
spa@cs.purdue.edu

Abstract

String complexity is defined as the cardinality of a set of all distinct words (factors) of a given string. For two strings, we define *joint string complexity* as the set of words that are common to both strings. We also relax this definition and introduce *joint semi-complexity* restricted to the common words appearing at least twice in both strings. In this paper we analyze joint complexity and joint semi-complexity when both strings are generated by a Markov source. The problem turns out to be quite challenging requiring subtle singularity analysis and saddle point method over infinity many saddle points leading to novel oscillatory phenomena with single and double periodicities.

1 Introduction

In the last decades, several attempts have been made to capture mathematically the concept of “complexity” of a sequence. The notion is connected with quite deep mathematical properties, including the rather elusive concept of randomness in a string (see e.g., [3, 10, 11]).

We are interested in studying some measures of complexity of a set of strings. We recall that the *string complexity* for a single string is defined as the cardinality of the set of distinct words (factors) of a given string [9]. For two strings the *joint complexity* is the cardinality of the set of common words to both strings, while the *joint semi-complexity* is the set of words that occur at least twice in both strings. Hereafter, we mostly analyze the average joint semi-complexity when both strings are generated by non-identical Markov sources. In fact, in this conference version to avoid cumbersome notation and much longer derivations, we only present (asymptotics) results when one of the source is uniform.

String complexity has a number of applications. It captures the “richness of the language” used in a sequence. For example, sequences with low complexity contain a large number of repeated substrings and they eventually become periodic (e.g., tandem repeats in a DNA sequence). In order to identify unusually low- or high-complexity strings one needs to determine how the complexities of the strings under study deviate from the average or maximum string complexity. On the other hand, joint string complexity is a very efficient way of evaluating similarity degree of two sequences. For example, genome sequences of two dogs will contain more common words than genome sequences of a dog and a cat. Similarly, the set of common words of one author’s texts is larger than the set of common words between two texts from two different authors. Finally, as our analysis shows the joint complexity or the joint semi-complexity can be used to discriminate between identically distributed sources and non-identical sources. Indeed, we shall prove that with non identical sources, the string joint complexity is of order

¹This work was supported in part by NSF Science and Technology Center on Science of Information Grant CCF-0939370, NSF Grants DMS-0800568 and CCF-0830140, NSA Grant H98230-08-1-0092, AFOSR Grant FA8655-11-1-3076, and the MNSW grant N206 369739.

$O(n^\kappa)$ for some $\kappa < 1$ when both strings are of length n . When the sources are identical the complexity is $O(n)$ (not treated in the paper). Furthermore, if both strings are identical (*i.e.* we deal with a single string complexity), then the string complexity is $O(n^2)$ [9].

Single string complexity was studied extensively in the past. The literature is reviewed in [9] where precise analysis of string complexity is discussed for strings generated by unbiased memoryless sources. Another analysis of the same situation was also proposed in [4] where for the first time the joint string complexity for memoryless sources is presented. It was evident from [4] that precise analysis of the joint complexity is quite challenging due to intricate singularity analysis and infinite number of saddle points. In this paper we deal with the joint string complexity and semi-complexity for Markov sources. To the best of our knowledge this problem was never tackled before. As expected, its analysis is very sophisticated but at the same time quite rewarding. It requires not only generalized (two-dimensional) depoissonization and generalized (two-dimensional) Mellin transforms but also subtle singularity analysis and unusual saddle point over infinity many saddle points. Furthermore, unlike other similar analyses [12] the non-binary case leads to new oscillatory phenomenon with single and double periodicities. In the long version of this paper, we shall also show that traditional dichotomy between rational and irrational cases must be extended to *commensurability* and *non-commensurability*.

2 Main Results

In this section we first define precisely the joint string complexity and semi-complexity. Hereafter, we mostly concentrate on the semi-complexity. Then we derive the basic functional equation describing the average joint semi-complexity for two strings generated by Markov sources. Finally, in this conference version we present asymptotic results only for a simplified version when one of the string is generated by an unbiased source while the other source is Markovian.

We begin by introducing some general notation. Let ω and σ be two sequences over alphabet \mathcal{A} . We denote by $|\omega|_\sigma$ the number of times σ occurs in ω . For example, $|abbba|_{bb} = 2$. By convention $|\omega|_\epsilon = |\omega|$, where ϵ is an empty string.

Throughout we denote by X a string (text) whose complexity we plan to study. We also assume that its length $|X|$ is equal to n . We consider string semi-complexity of a single string X , that is, the set of *distinct* substrings of X that occur twice in X , excluding the empty string ϵ . We denote this set as S_X : $S_X := \{\omega : \omega \neq \epsilon \ \& \ |X|_\omega \geq 2\}$. For example, if $X = aabaa$, then $S_X = \{a, b, aa\}$. We also have

$$|S_X| = \sum_{\sigma \in \mathcal{A}^*} 1_{|X|_\sigma > 1},$$

where 1_A is the indicator function of A . It is also related to the number of nodes in the associated suffix tree of X [13] (see also [5]). In passing we point out that the average string complexity was studied in Janson *at al.* [9] and Jacquet [4].

Now, let X and Y be two sequences (not necessarily of the same length). We define the *joint semi-complexity* as the cardinality of the set $S_{X,Y} = S_X \cap S_Y$ of common words that appear at least twice in both strings, excluding the empty string. For example, if $X = aabaa$ and $Y = abbba$, then $S_{X,Y} = \{a\}$. In fact, the joint semi-complexity corresponds to the number of common nodes in two suffix trees built from X and Y (excluding the root), and

$$|S_{X,Y}| = \sum_{\sigma \in \mathcal{A}^*} 1_{|X|_\sigma > 1} \times 1_{|Y|_\sigma > 1}.$$

In this conference paper, we only analyze the *average joint semi-complexity*. Extension to the joint complexity is rather straightforward and it is left for the full paper. We now

assume that both strings X and Y are generated by two *independent Markov sources* with the transition probabilities $P_i(a|b)$ for source i , $i \in \{1, 2\}$, where $(a, b) \in \mathcal{A}^2$. We denote by \mathbf{P}_1 (resp. \mathbf{P}_2) the transition matrix of Markov source 1 (resp. source 2).

We denote by $S_{n,m}$ the conditional average string semi-complexity, that is, the average number of common words occurring at least twice in two strings generated by stationary Markov sources. As in [5], we can prove that

$$S_{n,m} = T_{n,m} + O(n^{-\epsilon}), \quad (1)$$

where $T_{m,n}$ is the average number of common internal nodes between two tries T_1 and T_2 built over n and m independent strings, respectively, generated by Markov sources. We should point out that a single trie built over a Markov source was already analyzed in [8].

Let $a \in \mathcal{A}$. We denote $T_{a,m,n}$ the average number of common internal nodes between two tries when all the $n + m$ independent strings start with symbol a . Observe that $T_{a,n,m}$ for $n, m \geq 2$ satisfies the following recurrence for all $b \in \mathcal{A}$

$$T_{b,n,m} = 1 + \sum_{a \in \mathcal{A}} \sum_{n_a, m_a} \binom{n}{n_a} \binom{m}{m_a} (P_1(a|b))^{n_a} (1 - P_1(a|b))^{n - n_a} (P_2(a|b))^{m_a} (1 - P_2(a|b))^{m - m_a} T_{a, n_a, m_a}.$$

The double Poisson transform of $T_{a,n,m}$

$$T_a(z_1, z_2) = \sum_{n,m} T_{a,n,m} \frac{z_1^n z_2^m}{n!m!} e^{-z_1 - z_2} \quad (2)$$

translates the above recurrence into the following functional equation:

$$T_b(z_1, z_2) = (1 - (1 + z_1)e^{-z_1})(1 - (1 + z_2)e^{-z_2}) + \sum_{a \in \mathcal{A}} T_a(P_1(a|b)z_1, P_2(a|b)z_2). \quad (3)$$

Furthermore, the cumulative double Poisson transform

$$T(z_1, z_2) = \sum_{n,m} T_{n,m} \frac{z_1^n z_2^m}{n!m!} e^{-z_1 - z_2} \quad (4)$$

of the unconditional average joint semi-complexity $T_{n,m}$ satisfies

$$T(z_1, z_2) = (1 - (1 + z_1)e^{-z_1})(1 - (1 + z_2)e^{-z_2}) + \sum_{a \in \mathcal{A}} T_a(\pi_1(a)z_1, \pi_2(a)z_2). \quad (5)$$

Here $\pi_i(a)$ is the stationary probability of $a \in \mathcal{A}$ for the source model $i \in \{1, 2\}$. The identity (5) means that every strings in tries T_1 and T_2 have their first symbol distributed according to the stationary distribution.

In the next section we shall first first analyze $T(z_1, z_2)$ for large $z_1, z_2 \rightarrow \infty$ in a cone around the real axes, and then recall two-dimensional (double) analytic depoissonization (cf. Section 3.4) to translate $T(z_1, z_2)$ into $T_{n,m}$ by observing that $T_{n,m} \sim T(n, m)$. In order to analyze asymptotically (as $z_1, z_2 \rightarrow \infty$) the above system of functional equations (3)&(5), we resort (no surprises here!) to the double Mellin transform [1, 13]. We define the double Mellin transform of $T_a(z_1, z_2)$ as

$$\tilde{T}_a(s_1, s_2) = \int_0^\infty \int_0^\infty z_1^{s_1-1} z_2^{s_2-1} T_a(z_1, z_2) dz_1 dz_2. \quad (6)$$

Notice that for any $a \in \mathcal{A}$: $T_a(z_1, z_2) = O(z_1^2 z_2^2)$ when $z_1, z_2 \rightarrow 0$ and $T_a(z_1, z_2) = O(|z_1| + |z_2|)$ when $z_1, z_2 \rightarrow \infty$. Thus the Mellin transform is defined for $-2 < \Re(s_1), \Re(s_2) < -1$.

Applying basic properties of the Mellin transform [1, 13] we find for all $b \in \mathcal{A}$

$$\tilde{T}_b(s_1, s_2) = (s_1 + 1)\Gamma(s_1)(s_2 + 1)\Gamma(s_2) + \sum_{a \in \mathcal{A}} (P_1(a|b))^{-s_1} (P_2(a|b))^{-s_2} \tilde{T}_a(s_1, s_2). \quad (7)$$

Let now $\mathbf{T}(s_1, s_2)$ be the vector consisting of $\tilde{T}_a(s_1, s_2)$ for $a \in \mathcal{A}$. Then

$$\mathbf{T}(s_1, s_2) = (s_1 + 1)\Gamma(s_1)(s_2 + 1)\Gamma(s_2)\mathbf{1} + \mathbf{P}(s_1, s_2)\mathbf{T}(s_1, s_2) \quad (8)$$

where $\mathbf{1}$ the unit vector composed of 1's, and $\mathbf{P}(s_1, s_2)$ is the matrix whose (a, b) elements are $P_1(a|b))^{-s_1} (P_2(a|b))^{-s_2}$. In other words, $\mathbf{P}(s_1, s_2) = \mathbf{P}_1^{-s_1*} * \mathbf{P}_2^{-s_2*}$, where $*$ indicates the Schur product (*i.e.* element-wise) and \mathbf{P}^{x*} is the Schur power obtained by raising all elements of \mathbf{P} to the power x .

Putting all together we arrive at the following matrix equation

$$\mathbf{T}(s_1, s_2) = (s_1 + 1)\Gamma(s_1)(s_2 + 1)\Gamma(s_2)(\mathbf{I} - \mathbf{P}(s_1, s_2))^{-1}\mathbf{1}. \quad (9)$$

Moreover, if $\tilde{T}(s_1, s_2)$ is the double Mellin transform of $T(z_1, z_2)$, then from (5) we find

$$\tilde{T}(s_1, s_2) = (s_1 + 1)\Gamma(s_1)(s_2 + 1)\Gamma(s_2) \left(1 + \langle \boldsymbol{\pi}(s_1, s_2) | (\mathbf{I} - \mathbf{P}(s_1, s_2))^{-1}\mathbf{1} \rangle\right), \quad (10)$$

where $\boldsymbol{\pi}(s_1, s_2)$ is the vector $[(\pi_1(a))^{-s_1} (\pi_2(a))^{-s_2}]_{a \in \mathcal{A}}$ and $\langle \cdot | \cdot \rangle$ is the scalar product.

As pointed out in the introduction, we focus in this conference paper on asymptotics of a simple model in which one of the source is *uniform and memoryless*, that is, $P_1(a|b) = \frac{1}{|\mathcal{A}|}$, or $\mathbf{P}_1 = \frac{1}{|\mathcal{A}|}\mathbf{1} \otimes \mathbf{1}$, and when $\mathbf{P}_2 \neq \mathbf{P}_1$. To simplify our notation for all $(a, b) \in \mathcal{A}^2$ we shall write $P_2(a|b) = P(a|b)$ and $\mathbf{P}_2 = \mathbf{P}$. Therefore

$$\mathbf{P}(s_1, s_2) = |\mathcal{A}|^{s_1} \mathbf{P}(s) \quad (11)$$

with $\mathbf{P}(s) = \mathbf{P}^{-s_2*}$. We also write $\pi(a) = \pi_2(a)$ and the vector $\boldsymbol{\pi}(s) = [\pi(a)^{-s}]_{a \in \mathcal{A}}$, thus

$$\boldsymbol{\pi}(s_1, s_2) = |\mathcal{A}|^{s_1} \boldsymbol{\pi}(s). \quad (12)$$

Let $\lambda(s_1, s_2)$ be the main (largest) eigenvalue of $\mathbf{P}(s_1, s_2)$. We have

$$\lambda(s_1, s_2) = |\mathcal{A}|^{s_1} \lambda(s). \quad (13)$$

where $\lambda(s)$ is the main eigenvalue of matrix $\mathbf{P}(s)$. We also define $\mathbf{u}(s)$ as the right eigenvector of \mathbf{P}^{-s*} and $\boldsymbol{\zeta}(s)$ as the left eigenvector. We shall normalize $\langle \boldsymbol{\zeta}(s) | \mathbf{u}(s) \rangle = 1$.

Define $L(s) = \log_{|\mathcal{A}|} \lambda(s)$. In the Appendix we prove the following lemma.

Lemma 1 *The function $L(s)$ is convex when s is real.*

Let c_2 be the unique real solution of $L'(s) = 1$, that is,

$$\frac{\lambda'(c_2)}{\lambda(c_2)} = \log |\mathcal{A}|. \quad (14)$$

Let also $c_1 = -\log_{|\mathcal{A}|} \lambda(c_2)$ and $\kappa = -c_1 - c_2$. The next lemma is proved in the Appendix.

Lemma 2 *We always have $c_2 > -1$. When all $(a, b) \in \mathcal{A}^2$, $P(a|b) > 0$, then $c_2 < 0$.*

We now can formulate our main result.

Theorem 1 Let $f(s) = \langle \boldsymbol{\pi}(s) | \mathbf{u}(s) \rangle$ and $g(s) = \langle \boldsymbol{\zeta}(s) | \mathbf{1} \rangle$. Furthermore, with $\Psi(s)$ being the Euler psi function, define

$$\alpha_2 = L''(c_2) \quad (15)$$

$$\begin{aligned} \beta_2(s_1, s_2) = & -\alpha_2 \left(\Psi(s_1) + \frac{1}{1+s_1} + \log |\mathcal{A}| \right) \\ & + \Psi'(s_1) - \frac{1}{(s_1+1)^2} + \Psi'(s_2) - \frac{1}{(s_2+1)^2} \\ & + \frac{f''(s_2)}{f(s_2)} - \left(\frac{f'(s_2)}{f(s_2)} \right)^2 + \frac{g''(s_2)}{g(s_2)} - \left(\frac{g'(s_2)}{g(s_2)} \right)^2. \end{aligned} \quad (16)$$

(i) If $c_2 < 0$, then

$$T_{n,n} = n^\kappa \frac{f(c_2)g(c_2)(c_1+1)\Gamma(c_1)(c_2+1)\Gamma(c_2)}{\lambda(c_2) \log |\mathcal{A}| \sqrt{2\pi(\alpha_2 \log n + \beta_2(c_1, c_2))}} + n^\kappa Q(\log n) + o\left(\frac{n^\kappa}{\sqrt{\log n}}\right), \quad (17)$$

where depending on some properties of $\mathbf{P}(c_2)$ the quantity $\sqrt{x}Q(x)$ tends to a periodic or a double periodic function when $x \rightarrow \infty$ (with amplitude of order 10^{-6}), as detailed in Theorem 2.

(ii) If $c_2 > 0$, then

$$T_{n,n} = \sum_{k \in \mathbb{Z}} \frac{f(0)g(0)}{\lambda(0)} \frac{(1-L(0) - \frac{2ik\pi}{\log |\mathcal{A}|})\Gamma(-L(0) - \frac{2ik\pi}{\log |\mathcal{A}|})}{\log |\mathcal{A}|} n^{L(0)+2ik\pi/\log |\mathcal{A}|} + O(n^{L(0)-\epsilon}). \quad (18)$$

In order to present succinctly properties of function $Q(x)$, we need to introduce more notations. We say that matrix \mathbf{P} is *rationally balanced* if there exists $c \in \mathcal{A}$ and $\nu \in \mathbb{R}$ such that

$$\nu(P(a|b) + P(c|a) - P(c|b)) \in \mathbb{Z}. \quad (19)$$

Notice that this property does not depend on the pivot symbol c . In fact, we apply this to a matrix $\log^*(\mathbf{M})$ whose elements are logarithm of the elements of a given matrix \mathbf{M} .

Finally, let $\partial\mathcal{K}$ be the set of complex tuples (s_1, s_2) such that $\Re(s_1) = c_1$ and $\Re(s_2) = c_2$, and $|\mathcal{A}|^{s_1}\lambda(s_2) = 1$. Notice that $(c_1, c_2) \in \partial\mathcal{K}$. We define $\partial\mathcal{K}^* = \partial\mathcal{K} - \{(c_1, c_2)\}$.

Theorem 2 The function $Q(x)$ introduced in Theorem 1 can be expressed as

$$Q(x) = \sum_{(s_1, s_2) \in \partial\mathcal{K}^*} e^{ix\Im(s_1+s_2)} \frac{f(s_2)g(s_2)(s_1+1)(s_2+1)\Gamma(s_1)\Gamma(s_2)}{\lambda(s_2) \log |\mathcal{A}| \sqrt{2\pi(\alpha_2 x + \beta_2(s_1, s_2))}}. \quad (20)$$

If the matrix $\log^*\left(\frac{1}{P(x|x)}\mathbf{P}\right)$ is rationally balanced (i.e., rational case), then let ν be the smallest non negative real such that

$$\nu \log\left(\frac{1}{P(c|c)}\mathbf{P}\right) \in \mathbb{Z}|\mathcal{A}|^2.$$

Then

$$\partial\mathcal{K} = \left\{ \left(c_1 + \frac{2ik\pi}{\log |\mathcal{A}|} + 2i\pi\ell\nu \frac{\log P(c|c)}{\log |\mathcal{A}|}, c_2 + 2i\pi\ell\nu \right), (k, \ell) \in \mathbb{Z}^2 \right\}, \quad (21)$$

and $\sqrt{x}Q(x)$ is asymptotically double periodic. Otherwise (i.e., irrational case),

$$\partial\mathcal{K} = \left\{ \left(c_1 + \frac{2ik\pi}{\log 2}, c_2 \right), k \in \mathbb{Z} \right\}$$

and $\sqrt{x}Q(x)$ is asymptotically simply periodic.

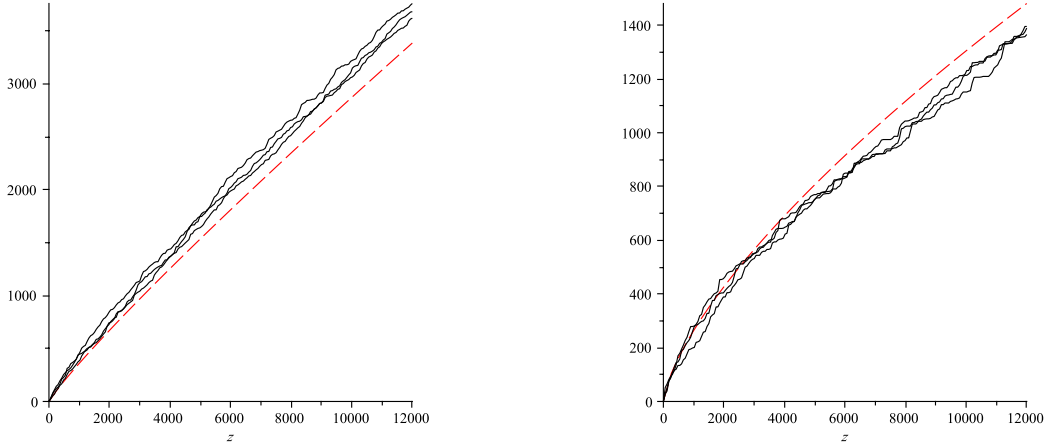


Figure 1: Joint semi-complexity: three simulated trajectories (black) versus asymptotic average (dashed red); case $c_2 > 0$ (left) and case $c_2 < 0$ (right).

In Figure 1 we plot the joint semi-complexity for several pairs of strings X and Y emitted by the same sources. String X is generated by a Markov source with the transition matrix \mathbf{P} , and string Y is generated by a uniform memoryless source. We consider two Markov sources for X with the following transition matrix:

$$\mathbf{P} = \begin{bmatrix} 0 & 0.5 \\ 1 & 0.5 \end{bmatrix} \quad \mathbf{P} = \begin{bmatrix} 0.2 & 0.8 \\ 0.8 & 0.2 \end{bmatrix}. \quad (22)$$

For the first \mathbf{P} (left plot in Figure 1) we have $c_2 > 0$ (cf. Theorem 1(ii)) while for the second \mathbf{P} (right plot in Figure 1) we have $c_2 < 0$ (cf. Theorem 1(i)).

3 Analysis

In this section we present the proof of our main result. We start with the double Mellin transform, followed by a brief discussion of properties of the kernel set $\partial\mathcal{K}$ of the solutions of $|\mathcal{A}|^{s_1}\lambda(s_2) = 1$. Then we discuss the saddle point method and finally the double dePoissonization. Some technical lemmas are in the Appendix.

3.1 Double Inverse Mellin Transform

In (10) of previous section, we computed the double Mellin transform $\tilde{T}(s_1, s_2)$ of the double Poisson transform $T(z_1, z_2)$. We now recover $T(z_1, z_2)$ for large $z_1, z_2 \rightarrow \infty$ around the real axes. We now set $z_1 = z_2 = z$. By the inverse double Mellin transform [4] we have for all $-2 < \rho < -1$:

$$\begin{aligned} T(z, z) &= \frac{1}{(2i\pi)^2} \iint_{\Re(s_1)=\Re(s_2)=\rho} \tilde{T}(s_1, s_2) z^{-s_1-s_2} ds_1 ds_2 \\ &= \frac{1}{(2i\pi)^2} \iint_{\Re(s_1)=\Re(s_2)=\rho} \Xi(s_1)\Xi(s_2) \left(1 + \langle \boldsymbol{\pi}(s_1, s_2) | (\mathbf{I} - \mathbf{P}(s_1, s_2))^{-1} \mathbf{1} \rangle\right) z^{-s_1-s_2} ds_1 ds_2, \end{aligned} \quad (23)$$

where $\Xi(s) := (s+1)\Gamma(s)$. We now denote by $\lambda_1(s), \lambda_2(s), \dots, \lambda_{|\mathcal{A}|}(s)$ the eigenvalues of the matrix $\mathbf{P}(s)$ in the non-increasing order of their modulus. Thus $\lambda(s) := \lambda_1(s)$. We also denote by $\mathbf{u}_i(s)$ (resp. $\boldsymbol{\zeta}_i(s)$) the right (resp. left) eigenvector of $\mathbf{P}(s)$ associated with $\lambda_i(s)$ normalizing

$\langle \zeta_i(s) | \mathbf{u}_i(s) \rangle = 1$. We write $(\zeta_1(s), \mathbf{u}_1(s))$ for $(\zeta(s), \mathbf{u}(s))$. By the spectral representation of matrices [13], we arrive at

$$(\mathbf{I} - \mathbf{P}(s_1, s_2))^{-1} = \sum_{i=1}^{|\mathcal{A}|} \frac{1}{1 - |\mathcal{A}|^{s_1} \lambda_i(s_2)} \mathbf{u}_i(s_2) \otimes \zeta_i(s_2) \quad (24)$$

where \otimes denotes the tensor product.

The eigenvalues $\lambda_i(s)$ are analytic functions of s as long as they are strictly decreasing (*i.e.* $|\lambda_{i-1}(s)| > |\lambda_i(s)| > |\lambda_{i+1}(s)|$ for all i). To simplify our analysis, we also postulate that none of the eigenvalue is identically equal to zero, that is, we assume $\log \lambda_i(s)$ exists except on a countable set $\mathcal{R} = \{s : \exists i : \lambda_i(s) = 0\}$. It should be, however, pointed out that there are cases when some eigenvalues are identically equal to zero. For example, for memoryless sources we have for all $i \geq 2$: $\lambda_i(s) \equiv 0$. But these cases are easy to handle by just excluding these null eigenvalues; we will not address it here (see [4]).

Let $L_{i,k}(s_2) = \frac{1}{-\log |\mathcal{A}|} (\log \lambda_i(s_2) + 2ik\pi)$ and $f_i(s) = \langle \boldsymbol{\pi}(s) | \mathbf{u}_i(s) \rangle$ and $g_i(s) = \langle \zeta_i(s) | \mathbf{1} \rangle$. Define

$$I(z, \rho) = \frac{1}{2i\pi} \int_{\Re(s)=\rho} \sum_{k \in \mathbb{Z}} \sum_{i=1}^{|\mathcal{A}|} \frac{f_i(s) g_i(s_2) \Xi(-L_{i,k}(s)) \Xi(s)}{\lambda_i(s) \log |\mathcal{A}|} z^{L_{i,k}(s)-s} ds. \quad (25)$$

Lemma 3 For any $M > 0$ and $\rho \in]-2, -1[$, we have

$$T(z, z) = I(z, \rho) + 1 + O(z^{1-M}). \quad (26)$$

Proof: To evaluate asymptotically (23), we move the line of integration with respect to s_1 from ρ to M for some $M > 1$. There are simple poles at $L_{i,k}(s_2)$ for $i \in \{1, \dots, |\mathcal{A}|\}$ and $k \in \mathbb{Z}$ where $(\mathbf{I} - \mathbf{P}(s_1, s_2))^{-1}$ ceases to exist (see next subsection for further discussion). The residues of these poles are exactly equal to $I(z, \rho)$. Furthermore, there is a pole at $s_1 = 0$ of $\Gamma(s_1)$ with residues 1. Thus by the Cauchy residue theorem the estimate of $T(z, z)$ becomes

$$T(z, z) = I(z, \rho) + \frac{1}{2i\pi} \int_{\Re(s_2)=\rho} \Xi(s_2) z^{-s_2} ds_2 + O(z^{1-M}). \quad (27)$$

But the above integral is equal to $1 + O(z^{1-M})$ (from the pole at $s_2 = 0$). ■

Our task is to estimate the integral $I(z, \rho)$ by finding a value of ρ contributing the most to $I(z, \rho)$. It turns out that this value is c_2 (when $c_2 < 0$), and it has an interesting property of being a *saddle* point with respect to function $z^{L(s)-s}$. The saddle point is at $L'(s) = 1$ which coincides with (14). But since $c_2 \notin]-2, -1[$ and $I(z, \rho)$ may have additional singularities when ρ moves towards c_2 , we must take this into account.

When $c_2 > 0$ the important contribution is at $\rho = 0$ due to the simple pole of function $\Xi(s)$. This case will be treated separately. Nevertheless, in the Appendix we prove the following lemma. The proof is rather technical since we have to cope with potential singularities between ρ and c_2 .

Lemma 4 If $c_2 < 0$, then for all $M > 0$ we have $T(z, z) = I(z, c_2) + 1 + O(z^{1-M})$.

3.2 Kernel Analysis

In this section we assume that $c_2 < 0$. Here, we establish some properties of the set $\partial\mathcal{K}$ where $\tilde{T}(s_1, s_2)$ ceases to exist. In particular, we show that if $(s_1, s_2) \in \partial\mathcal{K}$, then $(s_1 + \frac{2ik\pi}{\log|\mathcal{A}|}, s_2) \in \partial\mathcal{K}$ for all $k \in \mathbb{Z}$. We also show that the main eigenvalue of \mathbf{P} is the dominant singularity.

Let \mathbf{P} be a matrix on $\mathcal{A} \times \mathcal{A}$ of complex coefficients p_{ab} for all $(a, b) \in \mathcal{A}^2$. Let \mathbf{Q} be a matrix q_{ab} . In the following we say \mathbf{P} and \mathbf{Q} are *conjugate* if there exists a non-zero complex vector $(x_a)_{a \in \mathcal{A}}$ such that $q_{ab} = \frac{x_a}{x_b} p_{ab}$. We say that such matrices are *imaginary conjugate* if $|x_a| = 1$ for all $a \in \mathcal{A}$.

Observe that: (i) two conjugate matrices have the same eigenvalue set; (ii) if $\mathbf{u} = (u_a)_{a \in \mathcal{A}}$ is right eigenvector of \mathbf{P} , then $(x_a u_a)_{a \in \mathcal{A}}$ is right eigenvector of \mathbf{Q} . Similarly, if $(\zeta_a)_{a \in \mathcal{A}}$ is left eigenvector of \mathbf{P} , then $(\frac{1}{x_a} \zeta_a)_{a \in \mathcal{A}}$ is the left eigenvector of \mathbf{Q} .

The following lemma is essential and adapted from [8]. We prove it in the Appendix.

Lemma 5 *Let $\mathbf{M} = [m_{ab}]_{(a,b) \in \mathcal{A}^2}$ be a matrix such that $m_{ab} \geq 0$. We assume that 1 is the largest eigenvalue of \mathbf{M} . Let \mathbf{Q} be a matrix with coefficients $q_{ab} = e^{i\theta_{ab}} m_{ab}$ where θ_{ab} is real. The matrix \mathbf{Q} has eigenvalue 1 if and only if it is imaginary conjugate to matrix \mathbf{M} .*

Corollary 1 *Let $c \in \mathcal{A}$. The matrix \mathbf{Q} has eigenvalue 1 if and only if for all $(a, b) \in \mathcal{A}^2$:*

$$\frac{1}{2\pi} (\theta_{ab} + \theta_{ca} - \theta_{cb}) \in \mathbb{Z} . \quad (28)$$

Proof: We have $e^{i(\theta_a - \theta_b)} = \frac{e^{i\theta_{cb}}}{e^{i\theta_{ca}}}$, thus $e^{i(\theta_{cb} - \theta_{ca})} = e^{i\theta_{ab}}$. ■

Lemma 6 *Let $c \in \mathcal{A}$. A tuple (s_1, s_2) belongs to $\partial\mathcal{K}$ iff for all $(a, b) \in \mathcal{A}^2$ we have*

$$\frac{\Im(s_1)}{2\pi} \log |\mathcal{A}| - \frac{\Im(s_2)}{2\pi} \log \frac{P(a|b)P(c|a)}{P(c|b)} \in \mathbb{Z} . \quad (29)$$

Proof: Set $\mathbf{M} = \mathbf{P}(c_1, c_2)$ and $\mathbf{Q} = \mathbf{P}(s_1, s_2)$. Then, it follows directly from Corollary 1 with $e^{i\theta_{ab}} = |\mathcal{A}|^{i\Im(s_1)} (P(a|b))^{-i\Im(s_2)}$. ■

Lemma 7 *For all s_2 such that $\Re(s_2) = c_2$, we have $\exists s_1 : (s_1, s_2) \in \partial\mathcal{K} \iff |\lambda(s_2)| < \lambda(c_2)$.*

Proof: By the Perron-Frobenius, we have $|\lambda(s_2)| \leq \lambda(c_2)$ since $\Re(s_2) = c_2$ and $|\mathbf{P}(s_2)| = \mathbf{P}(c_2)$ (by taking the modulus element-wise). If $|\lambda(s_2)| = \lambda(c_2)$, then there will be t_1 such that $|\mathcal{A}|^{it_1} \lambda(s_2) = \lambda(c_2)$, and therefore $(c_1 + it_1, s_2) \in \partial\mathcal{K}$. ■

Lemma 8 *We have $\lambda(c_2) > \lambda_2(c_2)$.*

Proof: It follows from Perron-Frobenius that the main eigenvalue is unique. ■

Let \mathcal{U} be a complex neighborhood of 0 such that $\forall s \in \mathcal{U} : |\lambda_1(c_2 + s)| > |\lambda_2(c_2 + s)|$. In the Appendix we prove the following lemma.

Lemma 9 *Let s_k be a sequence such that $\Re(s_k) = c_2$ and $|\lambda(s_k)| \rightarrow \lambda(c_2)$. Then for all $s \in \mathcal{U}$ we have*

$$\forall i : \lim_{k \rightarrow \infty} \frac{\lambda_i(s_k + s)}{\lambda(s_k + s)} = \frac{\lambda_i(c_2 + s)}{\lambda(c_2 + s)}, \quad (30)$$

and

$$\lim_{k \rightarrow \infty} L(s_k + s) - L(s_k) = L(c_2 + s) - L(c_2) \quad (31)$$

$$\lim_{k \rightarrow \infty} L'(s_k + s) = L'(c_2 + s). \quad (32)$$

where the convergence also holds for any further derivative of function $L'(s)$.

In passing, we have $L'(s_k) \rightarrow 1$ and $L''(s_k) \rightarrow \alpha_2$. Finally, we prove in the Appendix our main lemma of this section.

Lemma 10 *There exists $\epsilon > 0$ such that for all $i \neq 1$ and for all s such that $\Re(s) = c_2$:*

$$|\lambda_i(s)| < \lambda(c_2) - \epsilon. \quad (33)$$

Proof: This is a consequence of previous lemmas. Suppose that there exists s_k such that $|\lambda_2(s_k)| \rightarrow \lambda(c_2)$. This implies that $|\lambda(s_k)| \rightarrow \lambda(c_2)$, but by previous lemma $|\lambda_2(s_k)| \rightarrow \lambda_2(c_2) = \lambda(c_2) - \epsilon$. ■

3.3 Saddle Points Analysis

We shall complete the proof by the saddle points analysis. We recall that for all $M > 0$: $T(z, z) = I(z, c_2) + 1 + O(z^{1-M})$ where $I(z, c)$ is given in (25). We first assume $c_2 < 0$.

Lemma 11 *There exists $\epsilon > 0$ such that*

$$T(z, z) = \sum_{k \in \mathbb{Z}} \frac{1}{2i\pi} \int_{\Re(s)=c_2} \frac{f_1(s)g_1(s)\Xi(-L_{i,k}(s))\Xi(s)}{\lambda_1(s) \log |\mathcal{A}|} z^{L_{1,k}(s)-s} ds + O(z^{\kappa-\epsilon}) \quad (34)$$

where $\kappa = -c_1 - c_2$.

Proof: By Lemma 10 for all $i \neq 1$ we have $\log_{|\mathcal{A}|} |\lambda_i(s)| < \log_{|\mathcal{A}|} \lambda(c_2) - \epsilon$ for some $\epsilon > 0$, and we can ignore other eigenvalues. Let

$$H_i(s, z) = \sum_{k \in \mathbb{Z}} \frac{f_i(s)g_i(s)}{\lambda_i(s)} \Xi(-L_{i,k}(s)) z^{L_{i,k}(s)}. \quad (35)$$

The contribution of $\int_{\Re(s)=c_2} H_i(s, z)\Xi(s)z^{-s}ds$ is of order $\int_{\Re(s)=c_2} |\Xi(s)|z^{\Re(L_i(s)-s)}ds$ which in turn is of order $z^{L(c_2)-c_2-\epsilon} = z^{\kappa-\epsilon}$. ■

3.3.1 The Rational Case

We assume now that the matrix $\log^* \left(\frac{1}{P(c|c)} \mathbf{P} \right)$ is *rationally balanced*. The matrix $\mathbf{P}(s + 2i\pi\nu)$ is then imaginary conjugate with the matrix $P(c|c)^{2i\pi\nu} \mathbf{P}(s)$ and $L(s + 2i\pi\nu) = L(s) + 2i\pi\nu \log P(c|c)$. Thus $\Re(L(c_2 + it))$ is periodic in t with period $2\pi\nu$. Furthermore, $L'(s)$ is also periodic with period $2\pi\nu$. Thus, $s_\ell = c_2 + 2i\pi\ell\nu$ for $\ell \in \mathbb{Z}$ are saddle points of $z^{L(s)-s}$.

We concentrate on the term $k = 0$ of the right-hand side of (34). Define

$$b_2(s) = \frac{d^2}{ds^2} \log \left(\frac{f(s)g(s)}{\lambda(s)} \Xi(-L(s)) \Xi(s) \right). \quad (36)$$

Notice that $b(s) = \beta_2(-L(s), s)$. Since the function

$$\log \left(\frac{f(s)g(s)}{\lambda(s)} \Xi(-L(s)) \Xi(s) \right)$$

has bounded variations, we have the classic saddle point result [2, 13]

$$\begin{aligned} & \frac{1}{2i\pi} \int_{\Re(s)=c_2} \frac{f(s)g(s)}{\lambda(s)} \Xi(-L(s)) \Xi(s) z^{L(s)-s} ds = \\ & = \sum_{\ell} \frac{f(s_{\ell})g(s_{\ell})}{\lambda(s_{\ell})} \Xi(-L(s_{\ell})) \Xi(s_{\ell}) \frac{z^{L(s_{\ell})-s_{\ell}}}{\sqrt{2\pi(\alpha_2 \log z + b_2(s_{\ell}))}} (1 + o(1)). \end{aligned} \quad (37)$$

Notice that $\Re(L(s_{\ell}) - s_{\ell}) = \kappa$. When adding the contribution from the $L(s) + \frac{2ik\pi}{\log|\mathcal{A}|}$ we obtain the expression for $Q(\log z)$ with $\partial\mathcal{K} = \{(-L(s_{\ell}) - \frac{2ik\pi}{\log|\mathcal{A}|}, s_{\ell}), (k, \ell) \in \mathbb{Z}^2\}$. The double periodicity comes from the fact that $\sqrt{x}Q(x) = \sum_{k,\ell} q_{k,\ell} e^{i(k\alpha + \ell\beta)x} + o(1)$ when $x \rightarrow \infty$ for some *incommensurable*² pair of real numbers (α, β) and complex numbers $\{q_{k,\ell}\}_{(k,\ell) \in \mathbb{Z}^2}$.

3.3.2 The Irrational Case

We now turn to the irrational case. Let $A > 0$ be a number such that for all $|s| \leq A$ we have $|\lambda(c_2 + s)| > |\lambda_2(c_2 + s)|$; thus by Lemma 9 $L(c_2 + s)$ is analytic. We assume that $c_2 < 0$ is the only saddle point on $\Re(s) = c_2$ for $|\Im(s)| \leq A$. We also postulate that there is $\alpha_3 > 0$ such that

$$|t| \leq A \Rightarrow \Re(L(c_2 + it) - L(c_2)) \leq -\alpha_3 t^2. \quad (38)$$

From the previous analysis we know that

$$\begin{aligned} & \frac{1}{2i\pi} \sum_{k \in \mathbb{Z}} \int_{\Re(s)=c_2, |\Im(s)| \leq A} \frac{f(s)g(s)}{\lambda(s)} \Xi \left(-L(s) - \frac{2ik\pi}{\log|\mathcal{A}|} \right) \Xi(s) z^{L(s)-s+2ik\pi/\log|\mathcal{A}|} ds = \\ & = Q(\log z)(1 + o(1)). \end{aligned} \quad (39)$$

Assume now (38) and define

$$\xi(s) = \sum_{k \in \mathbb{Z}} \left| \Xi \left(s - \frac{2ik\pi}{\log|\mathcal{A}|} \right) \right|. \quad (40)$$

The function $\xi(s)$ is continuous and bounded as long as $\Re(s)$ is bounded. Our aim is to prove that

$$\frac{1}{2i\pi} \int_{\Re(s)=c_2, |\Im(s)| > A} \left| \frac{f(s)g(s)}{\lambda(s)} \right| \xi(-L(s)) |\Xi(s)| z^{\Re(L(s))-c_2} ds = o\left(\frac{z^{\kappa}}{\sqrt{\log z}}\right), \quad (41)$$

which completes the proof of Theorem 1.

We know that $|f(s)g(s)| \leq f(c_2)g(c_2)$. In addition, we know that for $\Re(s) = c_2$ we have $\Re(L(s)) < L(c_2)$ as long as $\Im(s) \neq 0$. We also have $|\lambda(s)| > \epsilon'$ for some $\epsilon' > 0$ since the matrix $\mathbf{P}(s)$ stays away from the null matrix. Therefore, we need to estimate

$$\int_{\Re(s)=c_2, |\Im(s)| > A} |\Xi(s)| z^{\Re(L(s))-c_2} ds. \quad (42)$$

²A pair of numbers (α, β) is *commensurable* if there exists a real number ν such that the vector $(\nu\alpha, \nu\beta) \in \mathbb{Z}^2$; otherwise the pair is *incommensurable*.

For any $\epsilon > 0$, the portion of the line $\Re(s) = c_2$, where $\Re(L(s)) < L(c_2) - \epsilon$, contributes $z^{\kappa - \epsilon}$ to $T(z, z)$. Our attention must turn to the values of s on this line such that $\Re(L(s))$ is arbitrary close to $L(c_2)$. In particular, we are interested in the local maxima of $\Re(L(s))$ that are arbitrary close to $L(c_2)$. Indeed, these local maxima play a role in the saddle point method.

Let us consider the sequence of those maxima denoted by s_ℓ for $\ell \in \mathbb{N}$ such that $\Re(L(s_\ell)) \rightarrow L(c_2)$. By Lemma 9 we know that for all real t $L(s_\ell + it) - L(s_\ell) \rightarrow L(c_2 + it) - L(c_2)$ and that $L'(s_\ell + it) \rightarrow L'(c_2 + it)$. Therefore for all real t such $|t| \leq A$

$$\limsup_{\ell \rightarrow \infty} (\Re(L(s_\ell + it)) - \Re(L(s_\ell))) \leq -\alpha_3 t^2 \quad (43)$$

We define $I(A)$ to be the set of complex numbers s such that $\Re(s) = c_2$ and $\min_\ell \{|s - s_\ell|\} > A$.

Lemma 12 *There exists ϵ such that for all $s \in I(A)$: $\Re(L(s)) < L(c_2) - \epsilon$.*

Proof: Assume $s \in I(A)$. Since s is not a local maxima, we study the variation of $\Re(L(s))$ around the local maxima s_ℓ . Without loss of generality we assume that $s_\ell - A$ is between s and s_ℓ and therefore $\Re(L(s_\ell - A)) > \Re(s)$. Since $\limsup \Re(L(s_\ell - A)) < L(c_2) - \alpha_3 A^2 < L(c_2) - \epsilon$ the lemma is proven. \blacksquare

In view of the above, we conclude that

$$\int_{\Re(s)=c_2, |\Im(s)|>A} |\Xi(s)| z^{\Re(L(s))-c_2} ds \leq \sum_\ell \int_{|t|\leq A} |\Xi(s_\ell + it)| z^{\Re(L(s_\ell+it))-c_2} dt + O(z^{\kappa-\epsilon}). \quad (44)$$

Since $\Xi(s) = (s+1)\Gamma(s)$ on the line $\Re(s) = c_2$, there exists a real $B > 0$ such that $\forall s$:

$$\Re(s) = c_2 \Rightarrow \max_{|t|\leq A} \{|\Xi(s+it)|\} \leq B|\Xi(s)|. \quad (45)$$

Therefore, our analysis can be limited to

$$\sum_\ell \int_{|t|\leq A} |\Xi(s_\ell)| z^{\Re(L(s_\ell+it))-c_2} dt. \quad (46)$$

Lemma 13 *For ℓ tending to infinity, the s_ℓ are separated by a distance at least equal to A .*

Proof: First, let us assume that $\ell, \ell' \rightarrow \infty$ and $|s_\ell - s_{\ell'}| \rightarrow 0$, then we have

$$L'(s_{\ell'}) = L'(s_\ell) + (s_{\ell'} - s_\ell)L''(s_\ell) + O(|s_\ell - s_{\ell'}|^2). \quad (47)$$

Since $L''(s_\ell) \rightarrow \alpha_2 \neq 0$, then we cannot have $L'(s_{\ell'}) = 1$, thus $s_{\ell'}$ cannot be a local maximum of $\Re(L(s))$. Second, if $\liminf |s_\ell - s_{\ell'}| > \epsilon$ for some $\epsilon > 0$ with $|s_\ell - s_{\ell'}| < A$, then using the inequality

$$\limsup \Re(L(s_{\ell'})) - \Re(L(s_\ell)) \leq -\alpha_3 |s_\ell - s_{\ell'}|^2 < -\alpha_3 \epsilon^2 \quad (48)$$

we cannot have $\Re(L(s_{\ell'})) \rightarrow L(c_2)$. \blacksquare

The consequence of the previous lemma and the properties of function $\Xi(s)$ is that $\sum_\ell |\Xi(s_\ell)| < \infty$. Therefore,

$$\sum_\ell \int_{|t|\leq A} |\Xi(s_\ell)| z^{\Re(L(s_\ell+it))-c_2} dt = z^\kappa \sum_\ell |\Xi(s_\ell)| z^{\Re(L(s_\ell))-L(c_2)} \int_{|t|\leq A} z^{\Re(L(s_\ell+it))-\Re(L(s_\ell))} dt. \quad (49)$$

Since $\limsup_{\ell \rightarrow \infty} \Re(L(s_\ell + it)) - \Re(L(s_\ell)) \leq -\alpha_3 t^2$, we have [13]

$$\limsup_{\ell \rightarrow \infty} \int_{|t| \leq A} z^{\Re(L(s_\ell + it)) - \Re(L(s_\ell))} dt \leq \frac{1}{\sqrt{\pi \alpha_3 \log z}}, \quad (50)$$

and since $\lim_{z \rightarrow \infty} z^{\Re(s_\ell) - L(c_2)} = 0$, by the dominating convergence theorem, we arrive at

$$\sum_{\ell} |\Xi(s_\ell)| z^{\Re(s_\ell) - L(c_2)} \int_{|t| \leq A} z^{\Re(L(s_\ell + it)) - \Re(L(s_\ell))} dt = o\left(\frac{1}{\sqrt{\log z}}\right). \quad (51)$$

The case $c_2 > 0$ is discussed in the Appendix. It relies on the classic contour and singularities methods.

3.4 Double Depoissonization

To complete the proof of Theorem 1 we need to translate $T(z_1, z_2)$ into the original $T_{n,m}$. We accomplish it through a two-dimensional (double) depoissonization that we discuss next.

Let $a_{n,m}$ be a two-dimensional (double) sequence of complex numbers. We define the double Poisson transform $f(z_1, z_2)$ of $a_{n,m}$ as

$$f(z_1, z_2) = \sum_{n,m \geq 0} \frac{z_1^n z_2^m}{n! m!} e^{-z_1 - z_2}.$$

It is relatively straightforward to extend the one-dimensional depoissonization result of [6] to the two-dimensional case [4, 7]. In the Appendix we prove the following.

Lemma 14 *Let \mathcal{S}_θ be a cone of angle θ around the real axis. Assume that there exist $B > 0$, $D > 0$, $\alpha < 1$ and β such that for $|z_1|, |z_2| \rightarrow \infty$:*

- if $z_1, z_2 \in \mathcal{S}_\theta$: $|f(z_1, z_2)| = B(|z_1|^\beta + |z_2|^\beta)$;
- if $z_1, z_2 \notin \mathcal{S}_\theta$: $|f(z_1, z_2)e^{z_1 + z_2}| = D e^{\alpha|z_1| + \alpha|z_2|}$;
- if $z_i \in \mathcal{S}_\theta$ and $z_j \notin \mathcal{S}_\theta$ for $\{i, j\} = \{1, 2\}$: $|f(z_1, z_2)e^{z_j}| < D|z_i|^\beta e^{\alpha|z_j|}$.

Then

$$a_{n,m} = f(n, m) + O\left(\frac{n^\beta}{m} + \frac{m^\beta}{n}\right).$$

This leads to our final result of this section.

Lemma 15 *Generating function $T(z_1, z_2)$ satisfies the condition of Lemma 14 with $\beta = 1$. Therefore,*

$$T_{n,m} = T(n, m) + O\left(\frac{n}{m} + \frac{m}{n}\right).$$

References

- [1] P. Flajolet, X. Gourdon, and P. Dumas, Mellin Transforms and Asymptotics: Harmonic sums, *Theoretical Computer Science*, 144, 3–58, 1995.
- [2] P. Flajolet and R. Sedgewick, *Analytic Combinatorics*, Cambridge University Press, Cambridge, 2008.
- [3] Ilie, L., Yu, S., and Zhang, K. Repetition Complexity of Words In *Proc. COCOON* 320–329, 2002.
- [4] P. Jacquet, Common words between two random strings, *IEEE Intl. Symposium on Information Theory*, 1495-1499, 2007.
- [5] P. Jacquet, and W. Szpankowski, Autocorrelation on Words and Its Applications. Analysis of Suffix Trees by String-Ruler Approach, *J. Combinatorial Theory Ser. A*, 66, 237–269, 1994.
- [6] P. Jacquet, and W. Szpankowski, Analytical Depoissonization and Its Applications, *Theoretical Computer Science*, 201, 1–62, 1998.
- [7] P. Jacquet, and W. Szpankowski, *Analytic Pattern Matching: DNA, Lempel-Ziv, and Trees*, Cambridge Press, 2012.
- [8] . P. Jacquet, W. Szpankowski, and J. Tang, Average Profile of the Lempel-Ziv Parsing Scheme for a Markovian Source, *Algorithmica*, 31, 318-360, 2001.
- [9] S. Janson, S. Lonardi and W. Szpankowski, On Average Sequence Complexity, *Theoretical Computer Science*, 326, 213-227, 2004.
- [10] Li, M., and Vitanyi, P. *Introduction to Kolmogorov Complexity and its Applications*. Springer-Verlag, Berlin, Aug. 1993.
- [11] Niederreiter, H., Some computable complexity measures for binary sequences, In *Sequences and Their Applications*, Eds. C. Ding, T. Hellseth and H. Niederreiter Springer Verlag, 67-78, 1999.
- [12] G. Park, H.K. Hwang, P. Nicodeme, and W. Szpankowski, Profile of Tries, *SIAM J. Computing*, 8, 1821-1880, 2009.
- [13] W. Szpankowski, *An analysis of Algorithms on Sequences*, John Wiley, New York, 2001.

Appendix

Proof of Lemma 1: We have

$$\begin{cases} \lambda(s) &= \langle \boldsymbol{\zeta}(s) | \mathbf{P}(s) \mathbf{u}(s) \rangle = \sum_{a,b} \zeta_a(s) u_b(s) P(a|b)^{-s}, \\ \lambda'(s) &= \langle \boldsymbol{\zeta}(s) | \mathbf{P}'(s) \mathbf{u}(s) \rangle = \sum_{a,b} \zeta_a(s) u_b(s) P(a|b)^{-s} (-\log P(a|b)), \\ \lambda''(s) &= \langle \boldsymbol{\zeta}(s) | \mathbf{P}''(s) \mathbf{u}(s) \rangle = \sum_{a,b} \zeta_a(s) u_b(s) P(a|b)^{-s} (\log P(a|b))^2. \end{cases} \quad (52)$$

Let $x_{a,b} = \frac{1}{\lambda(s)} \zeta_a(s) u_b(s) P(a|b)^{-s}$, we have $\sum_{a,b} x_{a,b} = 1$. Also

$$L''(s) = \sum_{a,b} x_{a,b} (\log P(a|b))^2 - \left(\sum_{a,b} x_{a,b} \log P(a|b) \right)^2. \quad (53)$$

By the Perron Frobenius theorem $\boldsymbol{\zeta}(s)$ and $\mathbf{u}(s)$ have positive coefficients. Hence, by the convexity of the quadratic function we must have $L''(s) \geq 0$. \blacksquare

Proof of Lemma 2: For $c_2 > -1$ it suffices to show $L'(-1) < 1$. Since $\lambda(-1) = 1$, we we only need to prove that $\lambda'(-1) < \log |\mathcal{A}|$. We also have $\boldsymbol{\zeta}(-1) = \mathbf{1}$. Observe

$$\lambda'(-1) = \sum_{a,b} \zeta_a(-1) u_b(-1) (-P(a|b) \log P(a|b)). \quad (54)$$

Since $\sum_{a,b} \zeta_a(-1) u_b(-1) = \langle \boldsymbol{\zeta}(-1) | \mathbf{1} \rangle \langle \mathbf{1} | \mathbf{u}(-1) \rangle = |\mathcal{A}|$, and by concavity of the function $-x \log x$ on the interval $[0, 1]$ we have

$$\lambda'(-1) < - \left(\sum_{a,b} \zeta_a(-1) u_b(-1) P(a|b) \right) \log \left(\sum_{a,b} \frac{1}{|\mathcal{A}|} \zeta_a(-1) u_b(-1) P(a|b) \right) = \log |\mathcal{A}|; \quad (55)$$

the inequality is strict because the $P(a|b)$'s are not all identical.

For $c_2 < 0$, when all $P(a|b) > 0$, it suffices to prove that $L'(0) > 1$. Then by the increasing nature of $L'(s)$ we will necessarily have $c_2 < 0$. We have $\mathbf{P}(0)$ entirely made of ones, thus $\lambda(0) = |\mathcal{A}|$ and $\mathbf{u}(0) = \frac{1}{|\mathcal{A}|} \mathbf{1}$ and $\boldsymbol{\zeta}(0) = \mathbf{1}$. By convexity of $-\log$ function, and since the $P(a|b)$'s are not all identical

$$\frac{\lambda'(0)}{\lambda(0)} = - \sum_{a,b} \frac{1}{|\mathcal{A}|^2} \log P(a|b) > - \log \left(\sum_{a,b} \frac{1}{|\mathcal{A}|^2} P(a|b) \right). \quad (56)$$

Since $\sum_{a,b} P(a|b) = |\mathcal{A}|$ we get $\frac{\lambda'(0)}{\lambda(0)} > \log |\mathcal{A}|$. \blacksquare

Proof of Lemma 4: We move the integration line in $I(z, \rho)$ from the vertical line $\Re(s_2) = \rho$ to the vertical line $\Re(s_2) = c_2$. The function

$$J_k(s, z) = \sum_i \frac{f_i(s) g_i(s)}{\lambda_i(s)} \Xi(-L_{i,k}(s)) z^{L_{i,k}(s)} \quad (57)$$

satisfies the (somewhat complicated) identity

$$J_k(s, z) = \langle \boldsymbol{\pi}(s) | (\mathbf{P}(s))^{-1} \exp \left(-\frac{\log z}{\log |\mathcal{A}|} (\log \mathbf{P}(s) + 2ik\pi \mathbf{I}) \right) \Xi \left(-\frac{1}{\log |\mathcal{A}|} (\log \mathbf{P}(s) + 2ik\pi \mathbf{I}) \right) \mathbf{1} \rangle, \quad (58)$$

knowing that any analytical function f can be applied to matrix $\mathbf{P}(s)$ as long its eigenvalues do not correspond to a singularity of f . Therefore the only singularities that we meet when we move the integration line of $I(z, \rho)$ are the elements of $\mathcal{R} = \{s : \lambda_i(s) = 0, \text{ for some } i\}$.

We have

$$I(z, \rho) = \frac{1}{2i\pi} \int_{\Re(s)=\rho} \sum_i H_i(s, z) \Xi(s) z^{-s} ds . \quad (59)$$

If $\theta \in \mathcal{R}$, and we have $\lambda_i(\theta) = 0$, then the function $L_{i,k}(s) = \frac{1}{-\log|\mathcal{A}|} (\log \lambda_i(s) + 2ik\pi)$ is meromorphic around θ . However if θ is a simple root of $\lambda_i(s)$, then moving around θ would be equivalent to add 1 to the integer k : $\log \lambda_i(s) \rightarrow \log \lambda_i(s) + 2i\pi$. If the root is of multiplicity ℓ it is equivalent to add ℓ to the integer k . In any case the function $H_i(s, z)$ being invariant when ℓ is added to k , turns out to be fully analytic around θ , and the integration path in $I(z, \rho)$ can be moved over θ .

However, the function $H_i(s, z)$ is still singular on $s = \theta$ therefore there will be a contribution coming from the integration of $H_i(s, z) \Xi(s) z^{-s}$ on an arbitrary small loop around θ . Since $\Re(L_{i,k}(s)) \rightarrow -\infty$ when $s \rightarrow \theta$, having $\Re(L_{i,k}(s)) < -M$ will guarantee that the contribution is in $O(z^{1-M})$ and can therefore be included in the error term. ■

Proof of Lemma 5: Let $\mathbf{u} = (u_a)_{a \in \mathcal{A}}$ be the right eigenvector of \mathbf{M} and $(v_a)_{a \in \mathcal{A}}$ be the right eigenvector of \mathbf{Q} . Let also $v_a = x_a u_a$. If 1 is the eigenvalue, we have for all $c \in \mathcal{A}$:

$$(1 - e^{i\theta_{cc}} m_{cc}) u_c = \sum_{b \neq c} m_{cb} u_b e^{i\theta_{cb}} \frac{x_b}{x_c} . \quad (60)$$

If $e^{i\theta_{cc}} \neq 1$, then

$$|(1 - e^{i\theta_{cc}} m_{cc}) u_c| > (1 - m_{cc}) u_c \quad (61)$$

By the Perron-Frobenius theorem all u_a are real non negative. Suppose that $|x_c| = \max_{a \in \mathcal{A}} \{|x_a|\}$. If $\exists d \in \mathcal{A}$: $\frac{|x_d|}{|x_c|} < 1$ or if $(b, b') \in (\mathcal{A} - \{c\})^2$: $e^{i\theta_{cb}} \frac{x_b}{x_c} \neq e^{i\theta_{cb'}} \frac{x_{b'}}{x_c}$. Then

$$\left| \sum_{b \neq c} m_{cb} u_b e^{i\theta_{cb}} \frac{x_b}{x_c} \right| < \sum_{b \neq c} m_{cb} u_b . \quad (62)$$

But we also know that

$$(1 - m_{cc}) u_c = \sum_{b \neq c} m_{cb} u_b . \quad (63)$$

Therefore, we have $e^{i\theta_{cc}} = 1$ and for all $b \in \mathcal{A}$: $|x_b| = |x_c|$, and for all $(b, b') \in (\mathcal{A} - \{c\})^2$: $e^{i\theta_{cb}} \frac{x_b}{x_c} = e^{i\theta_{cb'}} \frac{x_{b'}}{x_c}$. But since for all $b \in \mathcal{A}$ $|x_b| = |x_c|$ every symbol in \mathcal{A} can play the role of c . Since for all $c \in \mathcal{A}$

$$(1 - m_{cc}) = \sum_{b \neq c} m_{cb} u_b e^{i\theta_{cb}} \frac{x_b}{x_c} = \sum_{b \neq c} m_{cb} u_b , \quad (64)$$

we simply have $\forall (a, b) \in \mathcal{A}$: $e^{i\theta_{ab}} \frac{x_b}{x_a} = 1$. Denoting $x_a = e^{i\theta_a}$ we get the expected result. The inverse proposition is immediate. ■

Proof of Lemma 9: Consider the matrix $\frac{1}{\lambda(s_k)} \mathbf{P}(s_k)$. Since the coefficients of this matrix are bounded, there is no loss in generality to consider the sequence of matrices converging to a matrix \mathbf{M} . The matrix \mathbf{M} and matrix $\mathbf{Q} = \frac{1}{\lambda(c_2)} \mathbf{P}(c_2)$ defined in Lemma 5 are imaginary conjugate *i.e.* the coefficients of \mathbf{M} are of the form

$$e^{i(\theta_a - \theta_b)} \frac{1}{\lambda(c_2)} P(a|b)^{-c_2} \quad (65)$$

for some vector of real numbers θ_a . Therefore, \mathbf{M} and $\frac{1}{\lambda(c_2)}\mathbf{P}(c_2)$ have the same spectrum. The spectrum of $\frac{1}{\lambda(s_k)}\mathbf{P}(s_k)$ converges to the spectrum of \mathbf{M} . Furthermore, the right eigenvector $\mathbf{u}(s_k)$ converges to the vector $e^{i\theta_a}u_a(c_2)$ and the left eigenvector $\zeta(s_k)$ converges to $e^{-i\theta_a}\zeta_a(c_2)$.

For any complex number s we have the identity

$$\frac{1}{\lambda(s_k)}\mathbf{P}(s_k + s) = \frac{1}{\lambda(s_k)}\mathbf{P}(s_k) * \mathbf{P}(s). \quad (66)$$

Thus $\frac{1}{\lambda(s_k)}\mathbf{P}(s_k + s)$ converges to $\mathbf{M} * \mathbf{P}(s)$ and is conjugate to $\frac{1}{\lambda(c_2)}\mathbf{P}(c_2 + s)$. Hence, the logarithm of the eigenvalue $L(s_k + s) - L(s_k)$ converges to $L(c_2 + s) - L(c_2)$. The property $|\lambda(c_2 + s)| > \lambda_2(c_2 + s)$ for all $s \in \mathcal{U}$ implies the analyticity of $L(c_2 + s)$, and therefore $L'(c_2 + s)$. Since the eigen spectrum of $\frac{1}{\lambda(s_k)}\mathbf{P}(s_k + s)$ converges to the eigen spectrum of $\frac{1}{\lambda(c_2)}\mathbf{P}(c_2 + s)$, we also have $\lambda(s_k + s) > \lambda_2(s_k + s)$ when k is large enough which implies the analyticity of $L(s_k + s)$. Thus by Ascoli theorem the derivatives converge, too. \blacksquare

Proof of Theorem 1 for the case $c_2 > 0$: While moving the integration path from $\Re(s) = \rho$ to $\Re(s) = c_2$, as in Lemma 4 we meet the pole of $\Xi(s)$ at $s = 0$. There, the function $\sum_i H_i(s, n)\Xi(s)n^{-s}$ has residue $\sum_i H_i(0, n)$. Thus for all $M > 0$:

$$T(n, n) = \sum_i H_i(0, n) + I(c_2, n) + O(n^{1-M}). \quad (67)$$

We notice that $I(c_2, n) = O(n^\kappa)$ which is $O(n^{L(0)-\epsilon})$. Furthermore,

$$\sum_i H_i(0, n) = \sum_i \frac{f_i(0)g_i(0)}{\lambda_i(0)} \sum_{k \in \mathbb{Z}} \Xi \left(-L_i(s) - \frac{2ik\pi}{\log |\mathcal{A}|} \right) n^{L_i(0) + 2ik\pi / \log |\mathcal{A}|}, \quad (68)$$

and for $i > 1$, $L_i(0) < L(0) - \epsilon$, thus $T(n, n) = H_1(0, n) + O(n^{L(0)-\epsilon})$. \blacksquare

Proof of Lemma 14: Let

$$f_n(z_2) = \sum_m a_{n,m} \frac{z_2^m}{m!} e^{-z_2}.$$

We notice that $f(z_1, z_2)$ is the Poisson transform of the sequence $f_n(z_2)$ with respect to the variable z_1 .

First depoissonization. For $z_2 \in \mathcal{S}_\theta$ we have the estimates

$$\begin{aligned} z_1 \in \mathcal{S}_\theta : |f(z_1, z_2)| &< B(|z_1|^\beta + |z_2|^\beta) \\ z_1 \notin \mathcal{S}_\theta : |f(z_1, z_2)e^{z_1}| &< D|z_2|^\beta e^{\alpha|z_1|}. \end{aligned}$$

Therefore, from the one-dimensional analytic depoissonization of [6, 13] for $z_2 \in \mathcal{S}_\theta$, we have for all integers $k > 0$

$$f_n(z_2) = f(n, z_2) + O \left(n^{\beta-1} + \frac{|z_2|^\beta}{n} \right) + O(|z_2|^\beta n^{\beta-k}).$$

Similarly, when $z_2 \notin \mathcal{S}_\theta$ we have

$$\begin{aligned} z_1 \in \mathcal{S}_\theta : |f(z_1, z_2)e^{z_2}| &< D|z_1|^\beta e^{\alpha|z_2|} \\ z_1 \notin \mathcal{S}_\theta : |f(z_1, z_2)e^{z_1+z_2}| &< D e^{\alpha|z_1| + \alpha|z_2|}. \end{aligned}$$

Thus for all integer k and $\forall z_2 \notin \mathcal{S}_\theta$

$$f_n(z_2)e^{z_2} = f(n, z_2)e^{z_2} + O(n^{\beta-1}e^{\alpha|z_2|}) + O(n^{\beta-k}e^{\alpha|z_2|}).$$

Second depoissonization The two results on $f_n(z_2)$, respectively for $z_2 \in \mathcal{S}_\theta$ and $z_2 \notin \mathcal{S}_\theta$, allow us to depoissonize $f_n(z_2)$. For all $k > \beta$:

- for $z_2 \in \mathcal{S}_\theta$: $f_n(z_2) = O(n^\beta + |z_2|^\beta)$;
- for $z_2 \notin \mathcal{S}_\theta$: $f_n(z_2)e^{z_2} = O(n^\beta e^{\alpha|z_2|})$.

These estimates are uniform. Therefore,

$$a_{n,m} = f_n(m) + O\left(\frac{n^\beta}{m} + \frac{m^\beta}{n}\right) + O\left(n^\beta m^{\beta-k}\right).$$

Since

$$f_n(m) = f(n, m) + O\left(n^{\beta-1} + \frac{m^\beta}{n}\right)$$

and setting $k > \beta + 1$, we prove the desired estimate. \blacksquare

Proof of Lemma 15: To prove the lemma, we need to establish three conditions (i)-(iii) of Lemma 14. We accomplish it through a generalization of the so called *increasing domain* approach discussed in [6, 13].

We first prove the lemma for the generating functions $T_a(z_1, z_2)$ for every $a \in \mathcal{A}$. Let $r = \max_{(a,b) \in \mathcal{A}^2, i \in \{1,2\}} \{P_i(a|b)\}$. We denote \mathcal{S}_k part of the cone \mathcal{S}_θ that contains points such that $|z| < \rho^{-k}$. Notice that $\mathcal{S}_k \subset \mathcal{S}_{k+1}$ for all integer k . We also notice $C(z_1, z_2) = O((|z_1| + |z_2|)^2)$ when $z_1, z_2 \rightarrow 0$, therefore we can define

$$B_k = \max_{a \in \mathcal{A}, (z_1, z_2) \in \mathcal{S}_k \times \mathcal{S}_k} \frac{|T_a(z_1, z_2)|}{|z_1| + |z_2|} < \infty.$$

We use the functional equation (3):

$$T_b(z_1, z_2) = (1 - (1 + z_1)e^{-z_1})(1 - (1 + z_2)e^{-z_2}) + \sum_{a \in \mathcal{A}} T_a(P_1(a|b)z_1, P_2(a|b)z_2). \quad (69)$$

In the above equation, we notice that if $(z_1, z_2) \in \mathcal{S}_{k+1} \times \mathcal{S}_{k+1} - \mathcal{S}_k \times \mathcal{S}_k$, then for all $(a, b) \in \mathcal{A}^2$ $(P_1(a|b)z_1, P_2(a|b)z_2)$ are in $\mathcal{S}_k \times \mathcal{S}_k$ and therefore we have for some fixed $\beta > 0$ and for all $b \in \mathcal{A}$:

$$|T_b(z_1, z_2)| \leq B_k \left(\sum_{a \in \mathcal{A}} P_1(a|b)|z_1| + P_2(a|b)|z_2| \right) + \beta = B_k(|z_1| + |z_2|) + \beta \quad (70)$$

since $|1 - (1 + z_i)e^{-z_i}|$ is uniformly bounded for all integers k by some $\sqrt{\beta}$ for both $i \in \{1, 2\}$ when $(z_1, z_2) \in \mathcal{S}_k$. Thus, we can derive the following recurrent inequality:

$$B_{k+1} \leq B_k + \beta \max_{(z_1, z_2) \in \mathcal{S}_{k+1} \times \mathcal{S}_{k+1} - \mathcal{S}_k \times \mathcal{S}_k} \left\{ \frac{1}{|z_1| + |z_2|} \right\} = B_k + \beta \rho^k. \quad (71)$$

We should notice that

$$\min_{(z_1, z_2) \in \mathcal{S}_{k+1} \times \mathcal{S}_{k+1} - \mathcal{S}_k \times \mathcal{S}_k} \{|z_1| + |z_2|\} = \rho^{-k} \quad (72)$$

because maybe only one of the number z_i has modulus greater than ρ^{-k} . It turns out that $\lim_{k \rightarrow \infty} B_k < \infty$, establishing condition (i) of Lemma 14.

Now we are going to establish condition (iii). For this end we define \mathcal{G} as the complementary cone of \mathcal{S}_θ and \mathcal{G}_k as the portion made of the point of modulus smaller than ρ^{-k} . We will use $\cos \theta < \alpha < 1$, therefore $\forall z \in \mathcal{G}$: $|e^z| < e^{\alpha|z|}$. We define D_k as

$$D_k = \max_{a \in \mathcal{A}, (z_1, z_2) \in \mathcal{G}_k \times \mathcal{G}_k} \frac{|T_a(z_1, z_2)e^{z_1+z_2}|}{\exp(\alpha|z_1| + \alpha|z_2|)}. \quad (73)$$

We define $G_a(z_1, z_2) = T_a(z_1, z_2)e^{z_1+z_2}$, we have the equation

$$G_b(z_1, z_2) = (e^{z_1} - 1 - z_1)(e^{z_2} - 1 - z_2) + \sum_{a \in \mathcal{A}} T_a(P_1(a|b)z_1, P_2(a|b)z_2) e^{1-P_1(a|b)z_1+(1-P_2(a|b))z_2} . \quad (74)$$

We notice that if $(z_1, z_2) \in \mathcal{G}_{k+1} \times \mathcal{G}_{k+1} - \mathcal{G}_k \times \mathcal{G}_k$, then all $(P_1(a|b)z_1, P_2(a|b)z_2)$ are in $\mathcal{G}_k \times \mathcal{G}_k$ and therefore we have for all $b \in \mathcal{A}$:

$$|G_b(z_1, z_2)| \leq D_k \left(\sum_{a \in \mathcal{A}} \exp((P_1(a|b)\alpha + (1 - P_1(a|b)) \cos \theta)|z_1| + (P_2(a|b)\alpha + (1 - P_2(a|b)) \cos \theta)|z_2|) \right) + (e^{\cos \theta |z_1|} + 1 + |z_1|)(e^{\cos \theta |z_2|} + 1 + |z_2|).$$

We notice that $\forall (a, b) \in \mathcal{A}^2$ and $\forall i \in \{1, 2\}$:

$$P_i(a|b)\alpha + (1 - P_i(a|b)) \cos \theta - \alpha \leq -(1 - \rho)(\alpha - \cos \theta) , \quad (75)$$

We also have $e^{\cos \theta |z_i|} + 1 + |z_i| \leq e^{\cos \theta |z_i|} (2 + \frac{1}{e \cos \theta})$, therefore

$$\frac{|G_b(z_1, z_2)|}{\exp(\alpha(|z_1| + |z_2|))} \leq D_k |\mathcal{A}| e^{-(1-\rho)(\alpha - \cos \theta)(|z_1| + |z_2|)} + (2 + \frac{1}{e \cos \theta})^2 e^{-(\alpha - \cos \theta)(|z_1| + |z_2|)} . \quad (76)$$

Since $(z_1, z_2) \in \mathcal{G}_{k+1} \times \mathcal{G}_{k+1} - \mathcal{G}_k \times \mathcal{G}_k$ implies $|z_1| + |z_2| \geq \rho^{-k}$ it follows

$$D_{k+1} \leq \max \left\{ D_k, |\mathcal{A}| D_k e^{-(1-\rho)(\alpha - \cos \theta)\rho^{-k}} + (2 + \frac{1}{e \cos \theta})^2 e^{-(\alpha - \cos \theta)\rho^{-k}} \right\} . \quad (77)$$

We clearly have $\lim_{k \rightarrow \infty} D_k < \infty$ and condition (iii) is established.

The proof of condition (ii) for z_1 and z_2 being in \mathcal{S}_θ and \mathcal{G} is a mixture of the above proofs.

The proof on the unconditional generating function $T(z_1, z_2)$ is a trivial extension based on (5). ■