

1 Hidden Words Statistics for Large Patterns

2 Svante Janson 

3 Department of Mathematics, Uppsala University, PO Box 480, SE-751 06 Uppsala, Sweden

4 svante.janson@math.uu.se

5 Wojciech Szpankowski 

6 Center for Science of Information, Department of Computer Science, Purdue University, West

7 Lafayette, IN, USA

8 spa@cs.purdue.edu

9 Abstract

10 We study here the so called *subsequence pattern matching* also known as *hidden pattern matching* in
11 which one searches for a given pattern w of length m as a *subsequence* in a random text of length
12 n . The quantity of interest is the number of occurrences of w as a subsequence (i.e., occurring in
13 *not* necessarily consecutive text locations). This problem finds many applications from intrusion
14 detection, to trace reconstruction, to deletion channel, and to DNA-based storage systems. In all of
15 these applications, the pattern w is of variable length. To the best of our knowledge this problem
16 was only tackled for a fixed length $m = O(1)$ [8]. In our main result Theorem 5 we prove that
17 for $m = o(n^{1/3})$ the number of subsequence occurrences is normally distributed. In addition, in
18 Theorem 6 we show that under some constrains on the structure of w the asymptotic normality can
19 be extended to $m = o(\sqrt{n})$. For a special pattern w consisting of the same symbol, we indicate
20 that for $m = o(n)$ the distribution of number of subsequences is either asymptotically normal or
21 asymptotically log normal. We conjecture that this dichotomy is true for all patterns. We use
22 Hoeffding's projection method for U -statistics to prove our findings.

23 **2012 ACM Subject Classification** General and reference → General literature; General and reference

24 **Keywords and phrases** Hidden pattern matching, subsequences, probability, U-statistics, projection
25 method

26 **Digital Object Identifier** 10.4230/LIPICs...

27 **Funding** *Svante Janson*: Supported by the Knut and Alice Wallenberg Foundation.

28 *Wojciech Szpankowski*: This work was supported by NSF Center for Science of Information (CSoI)
29 Grant CCF-0939370, and in addition by NSF Grant CCF-1524312.



© Svante Janson and Wojciech Szpankowski;
licensed under Creative Commons License CC-BY

Leibniz International Proceedings in Informatics

LIPICs Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

1 Introduction and Motivation

One of the most interesting and least studied problem in pattern matching is known as the *subsequence string matching* or the *hidden pattern matching* [12]. In this case, we search for a pattern $w = w_1 w_2 \cdots w_m$ of length m in the text $\Xi^n = \xi_1 \dots \xi_n$ of length n as *subsequence*, that is, we are looking for indices $1 \leq i_1 < i_2 < \cdots < i_m \leq n$ such that $\xi_{i_1} = w_1, \xi_{i_2} = w_2, \dots, \xi_{i_m} = w_m$. We say that w is *hidden* in the text Ξ^n . We do not put any constraints on the gaps $i_{j+1} - i_j$, so in language of [8] this is known as the *unconstrained hidden pattern matching*. The most interesting quantity of such a problem is the number of subsequence occurrences in the text generated by a random source. In this paper, we study the limiting distribution of this quantity when m , the length of the pattern, grows with n .

Hereafter, we assume that a memoryless source generates the text Ξ , that is, all symbols are generated independently with probability p_a for symbol $a \in \mathcal{A}$, where the alphabet \mathcal{A} is assumed to be finite. We denote by $p_w = \prod_j p_{w_j}$ the probability of the pattern w . Our goal is to understand the probabilistic behavior, in particular, the limiting distribution of the number of subsequence occurrences that we denote by $Z := Z_\Xi(w)$. It is known that behavior of Z depends on the order of magnitude of the pattern length m . For example, for the *exact pattern matching* (i.e., the pattern w must occur as a *string* in consecutive positions of the text), the limiting distribution is normal for $m = O(1)$ (more precisely, when $np_w \rightarrow \infty$, hence up to $m = O(\log n)$), but it becomes a Pólya–Aeppli distribution when $np_w \rightarrow \lambda > 0$ for some constant λ , and finally (conditioned on being non-zero) it turns into a geometric distribution when $np_w \rightarrow 0$ [12] (see also [2]). We might expect a similar behaviour for the subsequence pattern matching. In [8] it was proved by analytic combinatoric methods that the number of subsequence occurrences, $Z_\Xi(w)$, is asymptotically normal when $m = O(1)$, and not much is known beyond this regime. (See also [3]. Asymptotic normality for fixed m follows also by general results for U -statistics [10].) However, in many applications – as discussed below – we need to consider patterns w whose lengths grow with n . In this paper, we prove two main results. In Theorem 5 we establish that for $m = o(n^{1/3})$ the number of subsequence occurrences is normally distributed. Furthermore, in Theorem 6 we show that under some constraints on the structure of w , the asymptotic normality can be extended to $m = o(\sqrt{n})$. Moreover, for the special pattern $w = a^m$ consisting of the same symbol repeated, we show in Theorem 4 that for $m = o(\sqrt{n})$, the distribution of number of occurrences is asymptotically normal, while for larger m (up to cn for some $c > 0$) it is asymptotically log-normal. We conjecture that this dichotomy is true for a large class of patterns.

Regarding methodology, unlike [8] we use here probabilistic tools. We first observe that Z can be represented as a U -statistic (see (1)). This suggests to apply the Hoeffding [10] projection method to prove asymptotic normality of Z for some large patterns. Indeed, we first decompose Z into a sum of orthogonal random variables with variances of decreasing order in n (for m not too large), and show that the variable of the largest variance converges to a normal distribution, proving our main results Theorems 5 and 6.

The hidden pattern matching problem, especially for large patterns, finds many applications from intrusion detection, to trace reconstruction, to deletion channel, to DNA-based storage systems [1; 4; 5; 6; 12; 17]. Here we discuss below in some detail two of them, namely the deletion channel and the trace reconstruction problem.

A deletion channel [5; 6; 7; 14; 17; 19] with parameter d takes a binary sequence $\Xi^n = \xi_1 \cdots \xi_n$ where $\xi_i \in \mathcal{A}$ as input and deletes each symbol in the sequence independently with probability d . The output of such a channel is then a *subsequence* $\zeta = \zeta(x) = \xi_{i_1} \dots \xi_{i_M}$ of

Ξ , where M follows the binomial distribution $\text{Binom}(n, (1-d))$, and the indices i_1, \dots, i_M correspond to the bits that are *not* deleted. Despite significant effort [6; 14; 15; 17; 19] the mutual information between the input and output of the deletion channel and its capacity are still unknown. We hope to provide a more detailed characterization of the mutual information for memoryless sources using results of this and forthcoming papers. Indeed, it turns out that the mutual information $I(\Xi^n; \zeta(\Xi^n))$ can be exactly formulated as the problem of the subsequence pattern matching. In [5] it was proved that

$$I(\Xi^n; \zeta(\Xi^n)) = \sum_w d^{n-|w|} (1-d)^{|w|} (\mathbb{E}[Z_{\Xi^n}(w) \log Z_{\Xi^n}(w)] - \mathbb{E}[Z_{\Xi^n}(w)] \log \mathbb{E}[Z_{\Xi^n}(w)]),$$

74 where the sum is over all binary sequences of length smaller than n and $Z_{\Xi^n}(w)$ is the number
75 of subsequence occurrences w in the text Ξ^n . As one can see, to find precise asymptotics of
76 the mutual information we need to understand probabilistic behavior of Z for $m = O(n)$ and
77 typical w , which is our long term goal. The trace reconstruction problem [4; 11; 16; 18] is
78 related to the deletion channel problem since we are asking how many copies of the output
79 deletion channel we need to see until we can reconstruct the input sequence with high
80 probability.

81 2 Main Results

82 In this section we formulate precisely our problem and present our main results. Proofs are
83 delayed till the next section.

84 2.1 Problem formulation and notation

85 We consider a random string $\Xi^n = \xi_1 \dots \xi_n$ of length n . We assume that ξ_1, ξ_2, \dots are i.i.d.
86 random letters from a finite alphabet \mathcal{A} ; each letter ξ_i has the distribution $\mathbb{P}(\xi_i = a) = p_a$
87 where $a \in \mathcal{A}$, for some given vector $\mathbf{p} = (p_a)_{a \in \mathcal{A}}$; we assume $p_a > 0$, $a \in \mathcal{A}$.

88 Let $w = w_1 \dots w_m$ be a fixed string of length m over the same alphabet \mathcal{A} . We assume
89 $n \geq m$. Let $p_w := \prod_{j=1}^m p_{w_j}$, which is the probability that $\xi_1 \dots \xi_m$ equals w .

90 Let $Z = Z_{n,w}(\xi_1 \dots \xi_n)$ be the number of occurrences of w as a substring of $\xi_1 \dots \xi_n$. For
91 a set \mathcal{S} (in our case $[n]$ or $[m]$) and $k \geq 0$, let $\binom{\mathcal{S}}{k}$ be the collection of sets $\alpha \subseteq \mathcal{S}$ with $|\alpha| = k$.
92 Thus, $|\binom{\mathcal{S}}{k}| = \binom{|\mathcal{S}|}{k}$. For $k = 0$, $\binom{\mathcal{S}}{0}$ contains just the empty set \emptyset . For $k = 1$, we identify
93 $\binom{\mathcal{S}}{1}$ and \mathcal{S} in the obvious way. We write $\alpha \in \binom{[n]}{k}$ as $\{\alpha_1, \dots, \alpha_k\}$, where we assume that
94 $\alpha_1 < \dots < \alpha_k$. Then

$$95 \quad Z = \sum_{\alpha \in \binom{[n]}{m}} I_\alpha, \quad \text{where} \quad I_\alpha = \prod_{j=1}^m \mathbf{1}\{\xi_{\alpha_j} = w_j\}, \quad \alpha_1 < \dots < \alpha_m. \quad (1)$$

96 ► Remark 1. In the limit theorems, we are studying the asymptotic distribution of Z . We
97 then assume that $n \rightarrow \infty$ and (usually) $m \rightarrow \infty$; we thus implicitly consider a sequence of
98 words $w^{(n)}$ of lengths $m_n = |w^{(n)}|$. But for simplicity we do not show this in the notation.

99 We have $\mathbb{E} I_\alpha = p_w$ for every α . Hence,

$$100 \quad \mathbb{E} Z = \sum_{\alpha \in \binom{[n]}{m}} \mathbb{E} I_\alpha = \binom{n}{m} p_w. \quad (2)$$

101 Further, let $Y_\alpha := p_w^{-1} I_\alpha$, so $\mathbb{E} Y_\alpha = 1$, and

$$102 \quad Z^* := p_w^{-1} Z = \sum_{\alpha \in \binom{[n]}{m}} Y_\alpha, \quad (3)$$

103 so $\mathbb{E} Z^* = \binom{n}{m}$ and

$$104 \quad Z^* - \mathbb{E} Z^* = p_w^{-1} Z - \binom{n}{m} = \sum_{\alpha \in \binom{[n]}{m}} (Y_\alpha - 1). \quad (4)$$

105 We also write $\|Y\|_p := (\mathbb{E} |Y|^p)^{1/p}$ for the L^p norm of a random variable Y , while $\|\mathbf{x}\|$
 106 is the usual Euclidean norm of a vector \mathbf{x} in some \mathbb{R}^m . C denotes constants that may
 107 be different at different occurrences; they may depend on the alphabet \mathcal{A} and $(p_a)_a$, but
 108 not on n, m or w . Finally, \xrightarrow{d} and \xrightarrow{p} mean convergence in distribution and probability,
 109 respectively.

110 We are now ready to present our main results regarding the limiting distribution of Z ,
 111 the number of subsequence $w = a_1 \dots a_m$ occurrences when $m \rightarrow \infty$. We start with a simple
 112 example, namely, $w = a^m = a \cdot \dots \cdot a$ for some $a \in \mathcal{A}$, and show that depending on whether
 113 $m = o(\sqrt{n})$ or not the number of subsequences will follow asymptotically either the normal
 114 distribution or the log-normal distribution.

115 Before we present our results we consider asymptotically normal and log-normal distribu-
 116 tions in general, and discuss their relation.

117 2.2 Asymptotic normality and log-normality

118 If X_n is a sequence of random variables and a_n and b_n are sequences of real numbers, with
 119 $b_n > 0$, then $X_n \sim \text{AsN}(a_n, b_n)$ means that

$$120 \quad \frac{X_n - a_n}{\sqrt{b_n}} \xrightarrow{d} N(0, 1). \quad (5)$$

122 We say that X_n is *asymptotically normal* if $X_n \sim \text{AsN}(a_n, b_n)$ for some a_n and b_n , and
 123 *asymptotically log-normal* if $\ln X_n \sim \text{AsN}(a_n, b_n)$ for some a_n and b_n (this assumes $X_n \geq 0$).
 124 Note that these notions are equivalent when the asymptotic variance b_n is small, as made
 125 precise by the following lemma.

126 ► **Lemma 2.** *If $b_n \rightarrow 0$, and a_n are arbitrary, then*

$$127 \quad \ln X_n \sim \text{AsN}(a_n, b_n) \iff X_n \sim \text{AsN}(e^{a_n}, b_n e^{2a_n}). \quad (6)$$

129 **Proof.** By replacing X_n by X_n/e^{a_n} , we may assume that $a_n = 0$. If $\ln X_n \sim \text{AsN}(0, b_n)$
 130 with $b_n \rightarrow 0$, then $\ln X_n \xrightarrow{p} 0$, and thus $X_n \xrightarrow{p} 1$. It follows that $\ln X_n/(X_n - 1) \xrightarrow{p} 1$
 131 (with $0/0 := 1$), and thus

$$132 \quad \frac{X_n - 1}{b_n^{1/2}} = \frac{X_n - 1}{\ln X_n} \frac{\ln X_n}{b_n^{1/2}} \xrightarrow{d} N(0, 1), \quad (7)$$

134 and thus $X_n \sim \text{AsN}(1, b_n)$. The converse is proved by the same argument. ◀

135 ► **Remark 3.** Lemma 2 is best possible. Suppose that $\ln X_n \sim \text{AsN}(a_n, b_n)$. If $b_n \rightarrow b > 0$,
 136 then $\ln(X_n/e^{a_n}) = \ln X_n - a_n \xrightarrow{d} N(0, b)$, and thus

$$137 \quad X_n/e^{a_n} \xrightarrow{d} e^{\zeta_b}, \quad \zeta_b \sim N(0, b). \quad (8)$$

139 In this case (and only in this case), X_n thus converges in distribution, after scaling, to a
 140 log-normal distribution. If $b_n \rightarrow \infty$, then no linear scaling of X_n can converge in distribution
 141 to a non-degenerate limit, as is easily seen.

2.3 A simple example

We consider first a simple example where the asymptotic distribution can be found easily by explicit calculations. Fix $a \in \mathcal{A}$ and let $w = a^m = a \cdots a$, a string with m identical letters. Then, if $N = N_a$ is the number of occurrences of a in $\xi_1 \cdots \xi_n$, then

$$Z = \binom{N_a}{m}. \quad (9)$$

We will show that Z is asymptotically normal if m is small, and log-normal for larger m .

► **Theorem 4.** *Suppose that $m < np_a$, with $np_a - m \gg n^{1/2}$.*

(i) *Then*

$$\ln Z \sim \text{AsN}\left(\ln \binom{np_a}{m}, n \left| \ln \left(1 - \frac{m}{np_a}\right) \right|^2 p_a(1 - p_a)\right) \quad (10)$$

(ii) *In particular, if $m = o(n)$, then*

$$\ln Z \sim \text{AsN}\left(\ln \binom{np_a}{m}, (p_a^{-1} - 1) \frac{m^2}{n}\right) \quad (11)$$

(iii) *If $m = o(n^{1/2})$, then this implies*

$$Z / \mathbb{E} Z \sim \text{AsN}\left(1, (p_a^{-1} - 1) \frac{m^2}{n}\right) \quad (12)$$

and thus

$$Z \sim \text{AsN}\left(\mathbb{E} Z, (p_a^{-1} - 1) \frac{m^2}{n} (\mathbb{E} Z)^2\right). \quad (13)$$

Proof. (i) We have $N_a \sim \text{Bin}(n, p_a)$. Define $Y := N_a - np_a$. Then, by the Central Limit Theorem,

$$Y \sim \text{AsN}(0, np_a(1 - p_a)). \quad (14)$$

By (9), we have

$$\begin{aligned} \ln Z - \ln \binom{np_a}{m} &= \ln \binom{np_a + Y}{m} - \ln \binom{np_a}{m} \\ &= \ln \Gamma(np_a + Y + 1) - \ln \Gamma(np_a + Y - m + 1) - \ln m! \\ &\quad - (\ln \Gamma(np_a + 1) - \ln \Gamma(np_a - m + 1) - \ln m!) \\ &= \int_{y=0}^Y \int_{x=-m}^0 (\ln \Gamma)''(np_a + x + y + 1) dx dy. \end{aligned} \quad (15)$$

We fix a sequence $\omega_n \rightarrow \infty$ such that $np_a - m \gg \omega_n \gg n^{1/2}$; this is possible by the assumption. Note that (14) implies that $Y/\omega_n \xrightarrow{P} 0$, and thus $\mathbb{P}(|Y| \leq \omega_n) \rightarrow 1$. We may thus in the sequel assume $|Y| \leq \omega_n$. We assume also that n is so large that $np_a - m \geq 2\omega_n > 0$.

Stirling's formula implies, by taking the logarithm and differentiating twice (in the complex half-plane $\text{Re } z > \frac{1}{2}$, say)

$$(\ln \Gamma)''(x) = \frac{1}{x} + O\left(\frac{1}{x^2}\right) = \frac{1}{x} \left(1 + O\left(\frac{1}{x}\right)\right), \quad x \geq 1. \quad (16)$$

179 Consequently, (15) yields, noting the assumptions just made imply $|Y| \leq \omega_n \leq \frac{1}{2}(np_a - m)$,

$$\begin{aligned}
 180 \quad \ln Z - \ln \binom{np_a}{m} &= \int_{y=0}^Y \int_{x=-m}^0 \frac{1}{np_a + x + y + 1} \left(1 + O\left(\frac{1}{np_a - m}\right)\right) dx dy \\
 181 &= \int_{y=0}^Y \int_{x=-m}^0 \frac{1}{np_a + x} \left(1 + O\left(\frac{\omega_n}{np_a - m}\right)\right) dx dy \\
 182 &= \left(1 + O\left(\frac{\omega_n}{np_a - m}\right)\right) Y \int_{x=-m}^0 \frac{1}{np_a + x} dx \\
 183 &= (1 + o(1)) Y \ln \frac{np_a}{np_a - m}. \tag{17}
 \end{aligned}$$

185 Consequently, using also (14), we obtain

$$186 \quad \frac{\ln Z - \ln \binom{np_a}{m}}{n^{1/2} \left| \ln \left(1 - \frac{m}{np_a}\right) \right|} = (1 + o_p(1)) \frac{Y}{n^{1/2}} \xrightarrow{d} N(0, p_a(1 - p_a)), \tag{18}$$

188 which is equivalent to (10).

189 (ii) If $m = o(n)$, then $\left| \ln \left(1 - \frac{m}{np_a}\right) \right| \sim \frac{m}{np_a}$, and (11) follows.

190 (iii) If $m = o(n^{1/2})$, then (ii) applies, so (11) holds; hence Lemma 2 implies

$$191 \quad Z / \binom{np_a}{m} \sim \text{AsN}\left(1, (p_a^{-1} - 1) \frac{m^2}{n}\right). \tag{19}$$

193 Furthermore,

$$194 \quad \mathbb{E} Z = \binom{n}{m} p_a^m = \frac{n^m e^{O(m^2/n)}}{m!} p_a^m \sim \frac{n^m}{m!} p_a^m \tag{20}$$

196 and, similarly, $\binom{np_a}{m} \sim \frac{n^m p_a^m}{m!}$. Hence, $\mathbb{E} Z \sim \binom{np_a}{m}$ and (12) follows from (19); (13) is an
 197 immediate consequence. ◀

198 2.4 General results

199 We now present our main results. However, first we discuss the road map of our approach.
 200 First, we observe that the representation (1) shows that Z can be viewed as a U -statistic.
 201 For convenience, we consider Z^* , which differs from Z by a constant factor only, and show in
 202 (40) that $Z^* - \mathbb{E} Z^*$ can be decomposed into a sum $\sum_{\ell=1}^m V_\ell$ of orthogonal random variables
 203 V_ℓ such that, when m is not too large, $\text{Var}(\sum_{\ell=2}^m V_\ell) = o(\text{Var} V_1)$. Next, in Lemma 11 we
 204 prove that V_1 appropriately normalized converges to the standard normal distribution. This
 205 will allow us to conclude the asymptotic normality of Z .

206 In this paper, we only consider the region $m = o(n^{1/2})$. First, for $m = o(n^{1/3})$ we claim
 207 that the number of subsequence occurrences always is asymptotically normal.

208 ▶ **Theorem 5.** *If $m = o(n^{1/3})$, then*

$$209 \quad Z \sim \text{AsN}\left(\binom{n}{m} p_w, \sigma_1^2 p_w^2\right), \tag{21}$$

210 where

$$212 \quad \sigma_1^2 = \sum_{i=1}^n \sum_{a \in \mathcal{A}} p_a^{-1} \left(\sum_{j: w_j=a} \binom{i-1}{j-1} \binom{n-i}{m-j} \right)^2 - n \binom{n-1}{m-1}^2. \tag{22}$$

214 Furthermore, $\mathbb{E} Z = \binom{n}{m} p_w$ and $\text{Var} Z \sim p_w^2 \sigma_1^2$.

215 In the second main result, we restrict the patterns w to such that are not typical for the
216 random text; however, we will allow $m = o(n^{1/2})$.

217 ► **Theorem 6.** *Let $\mathbf{q} = (q_a)_{a \in \mathcal{A}}$ be the proportions of the letters in w , i.e., $q_a :=$
218 $\frac{1}{m} \sum_{j=1}^m \mathbf{1}\{w_j = a\}$. Suppose that $\liminf_{n \rightarrow \infty} \|\mathbf{q} - \mathbf{p}\| > 0$. If further $m = o(n^{1/2})$, then the
219 asymptotic normality (21) holds.*

220 3 Analysis and Proofs

221 In this section we will prove our main results. We start with some preliminaries.

222 3.1 Preliminaries and more notation

223 Let, for $a \in \mathcal{A}$,

$$224 \quad \varphi_a(x) := p_a^{-1} \mathbf{1}\{x = a\} - 1. \quad (23)$$

225 Thus, letting ξ be any random variable with the distribution of ξ_i ,

$$226 \quad \mathbb{E} \varphi_a(\xi) = 0, \quad a \in \mathcal{A}. \quad (24)$$

228 Let $p_* := \min_a p_a$ and

$$229 \quad B := p_*^{-1} - 1. \quad (25)$$

230 ► **Lemma 7.** *Let φ_a and B be as above.*

231 (i) *For every $a \in \mathcal{A}$,*

$$232 \quad \mathbb{E}[\varphi_a(\xi)^2] = p_a^{-1} - 1 \leq B. \quad (26)$$

234 (ii) *For some $c_1 > 0$ and every $a \in \mathcal{A}$,*

$$235 \quad \|\varphi_a(\xi)\|_2 = (p_a^{-1} - 1)^{1/2} \geq c_1. \quad (27)$$

237 (iii) *For any vector $\mathbf{r} = (r_a)_{a \in \mathcal{A}}$ with $\sum_a r_a = 1$,*

$$238 \quad \left\| \sum_{a \in \mathcal{A}} r_a \varphi_a(\xi) \right\|_2 \geq \|\mathbf{r} - \mathbf{p}\| := \left(\sum_{a \in \mathcal{A}} |r_a - p_a|^2 \right)^{1/2}. \quad (28)$$

240 **Proof.** The definition (23) yields

$$241 \quad \mathbb{E}[\varphi_a(\xi)^2] = p_a^{-2} \text{Var}[\mathbf{1}\{\xi = a\}] = p_a^{-2} p_a (1 - p_a) = p_a^{-1} - 1. \quad (29)$$

242 Hence, (26) and (27) follow, with B given by (25).

243 Finally, for every $x \in \mathcal{A}$, by (23) again,

$$244 \quad \sum_{a \in \mathcal{A}} r_a \varphi_a(x) = r_x p_x^{-1} - \sum_{a \in \mathcal{A}} r_a = r_x / p_x - 1 \quad (30)$$

246 and thus

$$247 \quad \mathbb{E} \left(\sum_{a \in \mathcal{A}} r_a \varphi_a(\xi) \right)^2 = \sum_{a \in \mathcal{A}} p_a (r_a / p_a - 1)^2 = \sum_{a \in \mathcal{A}} p_a^{-1} (r_a - p_a)^2 \quad (31)$$

249 and (28) follows. ◀

250 **3.2 A decomposition**

251 The representation (1) shows that Z is a special case of a U -statistic. For fixed m , the
 252 general theory of Hoeffding [10] applies and yields asymptotic normality. (Cf. [13, Section 4]
 253 for a related problem.) For increasing m (our main interest), we can still use the orthogonal
 254 decomposition of [10], which in our case takes the following form.

255 By the definitions in Section 2.1 and (23),

$$256 \quad Y_\alpha = \prod_{j=1}^m (p_{w_j}^{-1} \mathbf{1}\{\xi_{\alpha_j} = w_j\}) = \prod_{j=1}^m (\varphi_{w_j}(\xi_{\alpha_j}) + 1). \quad (32)$$

257 By multiplying out this product, we obtain

$$258 \quad Y_\alpha = \sum_{\gamma \subseteq [m]} \prod_{j \in \gamma} \varphi_{w_j}(\xi_{\alpha_j}). \quad (33)$$

259 Hence,

$$260 \quad Z^* = \sum_{\alpha \in \binom{[n]}{m}} Y_\alpha = \sum_{\alpha \in \binom{[n]}{m}} \sum_{\gamma \subseteq [m]} \prod_{j \in \gamma} \varphi_{w_j}(\xi_{\alpha_j}) = \sum_{\alpha \in \binom{[n]}{m}} \sum_{\gamma \subseteq [m]} \prod_{k=1}^{|\gamma|} \varphi_{w_{\gamma_k}}(\xi_{\alpha_{\gamma_k}}). \quad (34)$$

261 We rearrange this sum. First, let $\ell := |\gamma| \in [m]$, and consider all terms with a given ℓ . For
 262 each α and γ , with $|\gamma| = \ell$, let

$$263 \quad \alpha_\gamma := \{\alpha_{\gamma_1}, \dots, \alpha_{\gamma_\ell}\} \in \binom{[n]}{\ell}. \quad (35)$$

264 For given $\gamma \in \binom{[m]}{\ell}$ and $\beta \in \binom{[n]}{\ell}$, the number of $\alpha \in \binom{[n]}{m}$ such that $\alpha_\gamma = \beta$ equals the
 265 number of ways to choose, for each $k \in [\ell + 1]$, $\gamma_k - \gamma_{k-1} - 1$ elements of α in a gap of length
 266 $\beta_k - \beta_{k-1} - 1$, where we define $\beta_0 = \gamma_0 = 0$ and $\beta_{\ell+1} = n + 1$, $\gamma_{\ell+1} = m + 1$; this number is

$$267 \quad c(\beta, \gamma) := \prod_{k=1}^{\ell+1} \binom{\beta_k - \beta_{k-1} - 1}{\gamma_k - \gamma_{k-1} - 1}. \quad (36)$$

268 Consequently, combining the terms in (34) with the same α_γ ,

$$269 \quad Z^* = \sum_{\ell=0}^m \sum_{\gamma \in \binom{[m]}{\ell}} \sum_{\beta \in \binom{[n]}{\ell}} c(\beta, \gamma) \prod_{k=1}^{\ell} \varphi_{w_{\gamma_k}}(\xi_{\beta_k}). \quad (37)$$

270 We define, for $0 \leq \ell \leq m$ and $\beta \in \binom{[n]}{\ell}$,

$$271 \quad V_{\ell, \beta} := \sum_{\gamma \in \binom{[m]}{\ell}} c(\beta, \gamma) \prod_{k=1}^{\ell} \varphi_{w_{\gamma_k}}(\xi_{\beta_k}) \quad (38)$$

272 and

$$273 \quad V_\ell := \sum_{\beta \in \binom{[n]}{\ell}} V_{\ell, \beta}. \quad (39)$$

274 Thus (37) yields the decomposition

$$275 \quad Z^* = \sum_{\ell=0}^m V_\ell. \quad (40)$$

276

277 For $\ell = 0$, $\binom{[n]}{0}$ contains only the empty string \emptyset , and

$$278 \quad V_0 = V_{0,\emptyset} = \binom{n}{m} = \mathbb{E} Z^*. \quad (41)$$

280 Furthermore, note that two summands in (37) with different β are orthogonal, as a consequence of (24) and independence of different ξ_i . Consequently, the variables $V_{\ell,\beta}$ ($\ell \in [m]$, $\beta \in \binom{[m]}{\ell}$) are orthogonal, and hence the variables V_ℓ ($\ell = 0, \dots, m$) are orthogonal.

283 Let

$$284 \quad \sigma_\ell^2 := \text{Var}(V_\ell) = \mathbb{E} V_\ell^2 = \sum_{\beta \in \binom{[m]}{\ell}} \mathbb{E} V_{\ell,\beta}^2, \quad 1 \leq \ell \leq m. \quad (42)$$

285 Note also that by the combinatorial definition of $c(\beta, \gamma)$ given before (36), we see that

$$286 \quad \sum_{\beta \in \binom{[n]}{\ell}} c(\beta, \gamma) = \binom{n}{m}, \quad (43)$$

288 since this is just the number of $\alpha \in \binom{[n]}{m}$, and

$$289 \quad \sum_{\gamma \in \binom{[m]}{\ell}} c(\beta, \gamma) = \binom{n-\ell}{m-\ell}, \quad (44)$$

291 since this sum is the total number of ways to choose $m - \ell$ elements of the $n - \ell$ elements of α in the gaps.

293 3.3 The projection method

294 We use the projection method used by Hoeffding [10] to prove asymptotic normality for U -statistics. Translated to the present setting, the idea of the projection method is to approximate $Z^* - \mathbb{E} Z^* = Z^* - V_0$ by V_1 , thus ignoring all terms with $\ell \geq 2$ in the sum in (40). In order to do this, we estimate variances.

298 First, by (26) and the independence of the ξ_i ,

$$299 \quad \left\| \prod_{k=1}^{\ell} \varphi_{w_{\gamma_k}}(\xi_{\beta_k}) \right\|_2 = \left(\prod_{k=1}^{\ell} \mathbb{E} |\varphi_{w_{\gamma_k}}(\xi_{\beta_k})|^2 \right)^{1/2} \leq B^{\ell/2}. \quad (45)$$

300 By Minkowski's inequality, (38), (45) and (44),

$$301 \quad \|V_{\ell,\beta}\|_2 \leq \sum_{\gamma \in \binom{[m]}{\ell}} c(\beta, \gamma) B^{\ell/2} = B^{\ell/2} \binom{n-\ell}{m-\ell} \quad (46)$$

303 or, equivalently,

$$304 \quad \mathbb{E} V_{\ell,\beta}^2 \leq B^\ell \binom{n-\ell}{m-\ell}^2. \quad (47)$$

306 This leads to the following estimates.

307 ► **Lemma 8.** For $1 \leq \ell \leq m$,

$$308 \quad \sigma_\ell^2 := \mathbb{E} V_\ell^2 \leq \widehat{\sigma}_\ell^2 := B^\ell \binom{n}{\ell} \binom{n-\ell}{m-\ell}^2. \quad (48)$$

309

XX:10 Hidden Words

310 **Proof.** The definition of V_ℓ in (39) and (47) yield, since the summands $V_{\ell,\beta}$ are orthogonal,

$$311 \quad \sigma_\ell^2 := \mathbb{E} V_\ell^2 = \sum_{\beta \in \binom{[n]}{\ell}} \mathbb{E} V_{\ell,\beta}^2 \leq \binom{n}{\ell} B^\ell \binom{n-\ell}{m-\ell}^2, \quad (49)$$

312 as needed. ◀

314 Note that, for $1 \leq \ell < m$,

$$315 \quad \frac{\widehat{\sigma}_{\ell+1}^2}{\widehat{\sigma}_\ell^2} = B \frac{\binom{n}{\ell+1} \binom{n-\ell-1}{m-\ell-1}^2}{\binom{n}{\ell} \binom{n-\ell}{m-\ell}^2} = B \frac{n-\ell}{\ell+1} \left(\frac{m-\ell}{n-\ell} \right)^2 \leq B \frac{m^2}{(\ell+1)n}. \quad (50)$$

317 ▶ **Lemma 9.** *If $m \leq B^{-1/2} n^{1/2}$, then*

$$318 \quad \text{Var}(Z^* - V_1) \leq B^2 m^2 \binom{n-1}{m-1}^2. \quad (51)$$

320 **Proof.** By (50) and the assumption, for $1 \leq \ell < m$,

$$321 \quad \frac{\widehat{\sigma}_{\ell+1}^2}{\widehat{\sigma}_\ell^2} \leq \frac{1}{\ell+1} \leq \frac{1}{2}, \quad (52)$$

323 and thus, summing a geometric series,

$$324 \quad \begin{aligned} \text{Var}(Z^* - V_1) &= \sum_{\ell=2}^m \text{Var}(V_\ell) \leq \sum_{\ell=2}^m \widehat{\sigma}_\ell^2 \leq \sum_{\ell=2}^m 2^{2-\ell} \widehat{\sigma}_2^2 \leq 2 \widehat{\sigma}_2^2 \\ &= B^2 n(n-1) \binom{n-2}{m-2}^2 \leq B^2 m^2 \binom{n-1}{m-1}^2. \end{aligned} \quad (53)$$

327 ◀

328 3.4 The first term V_1

329 For $\ell = 1$, we identify $\binom{[n]}{\ell}$ and $[n]$, and we write $V_{1,i} := V_{1,\{i\}}$. Note that, by (36),

$$330 \quad c(i, j) := c(\{i\}, \{j\}) = \binom{i-1}{j-1} \binom{n-i}{m-j}. \quad (54)$$

332 Thus (39) and (38) become

$$333 \quad V_1 = \sum_{i=1}^n V_{1,i} \quad (55)$$

335 with, using (54),

$$336 \quad V_{1,i} = \sum_{j=1}^m c(i, j) \varphi_{w_j}(\xi_i) = \sum_{j=1}^m \binom{i-1}{j-1} \binom{n-i}{m-j} \varphi_{w_j}(\xi_i). \quad (56)$$

338 Note that $V_{1,i}$ is a function of ξ_i , and thus the random variables $V_{1,i}$ are independent.

339 Furthermore, (24) implies $\mathbb{E} V_{1,i} = 0$. Let $\tau_i^2 := \text{Var} V_{1,i} = \mathbb{E} V_{1,i}^2$. Then, see (42),

$$340 \quad \sigma_1^2 = \text{Var} V_1 = \sum_{i=1}^n \text{Var} V_{1,i} = \sum_{i=1}^n \tau_i^2. \quad (57)$$

341

342 Observe that it follows from (56) and (23) that

$$343 \quad \tau_i^2 = \sum_{a \in \mathcal{A}} p_a^{-1} \left(\sum_{j: w_j = a} \binom{i-1}{j-1} \binom{n-i}{m-j} \right)^2 - \binom{n-1}{m-1}^2. \quad (58)$$

344 Taking $\ell = 1$ in (47) yields the upper bound

$$345 \quad \tau_i^2 = \mathbb{E} V_{1,i}^2 \leq B \binom{n-1}{m-1}^2, \quad i \in [n]. \quad (59)$$

347 Summing over i , or using (48), we obtain

$$348 \quad \sigma_1^2 := \mathbb{E} V_1^2 \leq \hat{\sigma}_1^2 := Bn \binom{n-1}{m-1}^2. \quad (60)$$

350 We notice that the upper bound is achievable. Indeed, for $w = a \cdots a$, by (58) and (57),

$$351 \quad \tau_i^2 = (p_a^{-1} - 1) \binom{n-1}{m-1}^2, \quad \sigma_1^2 = n(p_a^{-1} - 1) \binom{n-1}{m-1}^2. \quad (61)$$

353 We show also a general lower bound.

354 ► **Lemma 10.** *There exists $c, c' > 0$ such that*

$$355 \quad \sigma_1^2 \geq \frac{c}{m} \hat{\sigma}_1^2 = c' \frac{n}{m} \binom{n-1}{m-1}^2. \quad (62)$$

357 **Proof.** We consider the first term in the sum in (56) separately, and write

$$358 \quad V_{1,i} = c(i, 1) \varphi_{w_1}(\xi_i) + V'_{1,i}, \quad (63)$$

360 where

$$361 \quad V'_{1,i} := \sum_{j=2}^m c(i, j) \varphi_{w_j}(\xi_i). \quad (64)$$

363 We have, by (54), $c(i, 1) = \binom{n-i}{m-1}$. Consequently, for any $i \in [n]$,

$$364 \quad \frac{c(i, 1)}{c(1, 1)} = \frac{\binom{n-i}{m-1}}{\binom{n-1}{m-1}} = \frac{\prod_{k=0}^{m-2} (n-i-k)}{\prod_{k=0}^{m-2} (n-1-k)} = \prod_{k=0}^{m-2} \left(1 - \frac{i-1}{n-1-k} \right) \\ 365 \quad \geq 1 - \sum_{k=0}^{m-2} \frac{i-1}{n-1-k} \geq 1 - \frac{m(i-1)}{n-m+1}. \quad (65)$$

367 Let $\delta \leq 1/4$ be a fixed small positive number, chosen later. Assume that $i \leq 1 + \delta n/m$.

368 In particular, either $i = 1$ or $m \leq m(i-1) \leq \delta n < n/2$, and thus (65) implies

$$369 \quad \frac{c(i, 1)}{c(1, 1)} \geq 1 - \frac{m(i-1)}{n-m} \geq 1 - \frac{\delta n}{n/2} = 1 - 2\delta. \quad (66)$$

371 By (44), (66) implies

$$372 \quad \sum_{j=2}^m c(i, j) = \binom{n-1}{m-1} - c(i, 1) = c(1, 1) - c(i, 1) \leq 2\delta c(1, 1). \quad (67)$$

373

XX:12 Hidden Words

374 Hence, by (64), Minkowski's inequality and (26), cf. (46),

$$375 \quad \|V'_{1,i}\|_2 \leq \sum_{j=2}^m c(i,j) \|\varphi_{w_j}(\xi_i)\|_2 \leq \sum_{j=2}^m c(i,j) B^{1/2} \leq 2\delta B^{1/2} c(1,1). \quad (68)$$

377 Furthermore, (27) and (66) yield

$$378 \quad \|c(i,1)\varphi_{w_1}(\xi_i)\|_2 \geq c(i,1)c_1 \geq c_1(1-2\delta)c(1,1) \geq \frac{1}{2}c_1c(1,1). \quad (69)$$

380 Finally, (63) and the triangle inequality yield, using (69) and (68),

$$381 \quad \|V_{1,i}\|_2 \geq \|c(i,1)\varphi_{w_1}(\xi_i)\|_2 - \|V'_{1,i}\|_2 \geq (\frac{1}{2}c_1 - 2\delta B^{1/2})c(1,1). \quad (70)$$

383 We now choose $\delta := c_1/(8B^{1/2})$, and find that for some $c_2 > 0$,

$$384 \quad \tau_i^2 := \|V_{1,i}\|_2^2 \geq c_2c(1,1)^2, \quad i \leq 1 + \delta n/m. \quad (71)$$

386 Consequently, by (57),

$$387 \quad \sigma_1^2 = \sum_{i=1}^n \tau_i^2 \geq \frac{\delta n}{m} c_2 c(1,1)^2 = c_3 \frac{n}{m} \binom{n-1}{m-1}^2. \quad (72)$$

389 This proves (62), with $c' := c_3$ and $c = c'/B$. \blacktriangleleft

390 The next lemma is proved in the Appendix in which we verify Lyapunov's condition to
391 prove asymptotic normality of V_1 .

392 **► Lemma 11.** *Suppose that $m = o(n)$. Then V_1 is asymptotically normal:*

$$393 \quad V_1/\sigma_1 \xrightarrow{d} N(0,1). \quad (73)$$

395 3.5 Proofs of Theorem 5 and 6

396 We next prove a general theorem showing asymptotic normality under some conditions.

397 **► Theorem 12.** *Suppose that $n \rightarrow \infty$ and that*

$$398 \quad m^2 \binom{n-1}{m-1}^2 = o(\sigma_1^2). \quad (74)$$

400 *Then*

$$401 \quad \text{Var } Z = p_w^2 \text{Var } Z^* \sim p_w^2 \sigma_1^2 \quad (75)$$

403 *and*

$$404 \quad \frac{Z^* - \mathbb{E} Z^*}{\sigma_1} \xrightarrow{d} N(0,1), \quad (76)$$

$$405 \quad \frac{Z - \mathbb{E} Z}{(\text{Var } Z)^{1/2}} = \frac{Z^* - \mathbb{E} Z^*}{(\text{Var } Z^*)^{1/2}} \xrightarrow{d} N(0,1). \quad (77)$$

407 **Proof.** By Lemma 9 and (74),

$$408 \quad \text{Var} \left(\frac{Z^* - V_1}{\sigma_1} \right) = \frac{\text{Var}(Z^* - V_1)}{\sigma_1^2} \leq B^2 \frac{m^2 \binom{n-1}{m-1}^2}{\sigma_1^2} = o(1). \quad (78)$$

410 Hence, recalling $\mathbb{E} V_1 = 0$,

$$411 \frac{Z^* - \mathbb{E} Z^* - V_1}{\sigma_1} \xrightarrow{\mathbb{P}} 0. \quad (79)$$

412
413 Combining (73) and (79), we obtain (76).

414 Furthermore, by (78), and since the terms in (40) are orthogonal,

$$415 \text{Var } Z^* = \text{Var } V_1 + \text{Var}(Z^* - V_1) = \sigma_1^2 + o(\sigma_1^2) \sim \sigma_1^2, \quad (80)$$

416
417 which yields (75), and also shows that we may replace σ_1 by $(\text{Var } Z^*)^{1/2}$ in (76), which yields
418 (77); the equality in (77) is a trivial consequence of (3). ◀

419 Now we are ready to prove our main results.

420 **Proof of Theorem 5.** By Lemma 10,

$$421 \frac{m^2 \binom{n-1}{m-1}^2}{\sigma_1^2} \leq C \frac{m^3}{n} = o(1). \quad (81)$$

422
423 Thus (74) holds, and the result follows by Theorem 12 together with (2) and (3). ◀

424 Recall that in Theorem 6, the range of m is improved, assuming that w is *not* typical for
425 the random source with probabilities $\mathbf{p} = (p_a)_{a \in \mathcal{A}}$ that we consider.

426 **Proof of Theorem 6.** By Theorem 12, with (74) verified by Lemma 13 below. ◀

427 ▶ **Lemma 13.** Let $\mathbf{q} = (q_a)_{a \in \mathcal{A}}$ be the proportions of the letters in w . Then

$$428 \sigma_1^2 \geq \frac{m^2}{n} \binom{n}{m}^2 \|\mathbf{q} - \mathbf{p}\|^2 = n \binom{n-1}{m-1}^2 \|\mathbf{q} - \mathbf{p}\|^2. \quad (82)$$

429
430 **Proof.** Let

$$431 \psi_i(x) := \sum_{j=1}^m c(i, j) \varphi_{w_j}(x). \quad (83)$$

432
433 Thus (56) is $V_{1,i} = \psi_i(\xi_i)$, and (57) is, since $\mathbb{E} \psi_i(\xi) = 0$,

$$434 \sigma_1^2 = \text{Var } V_1 = \sum_{i=1}^n \mathbb{E}[\psi_i(\xi_i)^2] = \mathbb{E} \sum_{i=1}^n \psi_i(\xi)^2. \quad (84)$$

435
436 Hence, by the Cauchy–Schwarz inequality,

$$437 n\sigma_1^2 = n \mathbb{E} \sum_{i=1}^n \psi_i(\xi)^2 \geq \mathbb{E} \left(\sum_{i=1}^n \psi_i(\xi) \right)^2. \quad (85)$$

438
439 Furthermore, by (83) and (43)

$$440 \sum_{i=1}^n \psi_i(x) = \sum_{i=1}^n \sum_{j=1}^m c(i, j) \varphi_{w_j}(x) = \sum_{j=1}^m \binom{n}{m} \varphi_{w_j}(x) = \binom{n}{m} \sum_{a \in \mathcal{A}} m q_a \varphi_a(x). \quad (86)$$

441
442 Hence, (28) yields

$$443 \left\| \sum_{i=1}^n \psi_i(\xi) \right\|_2 = m \binom{n}{m} \left\| \sum_{a \in \mathcal{A}} q_a \varphi_a(\xi) \right\|_2 \geq m \binom{n}{m} \|\mathbf{q} - \mathbf{p}\|. \quad (87)$$

444
445 Combining (85) and (87) yields (82). ◀

446 **References**

- 447 1 M. Atallah, R. Gwadera, and W. Szpankowski. Reliable Detection of Episodes in Event
448 Sequences, *Third IEEE International Conference on Data Mining (ICDM-03)*, 67–74,
449 Melbourne, Florida, 2003.
- 450 2 E. A. Bender and F. Kochman. The distribution of subword counts is usually normal,
451 *European Journal of Combinatorics* 14, 265–275, 1993.
- 452 3 J. Bourdon and B. Vallée. Generalized Pattern Matching Statistics, *Mathematics and*
453 *computer science* (Colloquium Proceedings, Versailles, 2002), B. Chauvin et al. Editors,
454 Birkhäuser Verlag, 229–245, 2002.
- 455 4 M. Cheraghchi, R. Gabrys, O. Milenkovic, and J. Ribeiro. Coded trace reconstruction,
456 [arXiv:1903.09992](https://arxiv.org/abs/1903.09992), 2019.
- 457 5 M. Drmota, K. Viswanathan, and W. Szpankowski. Mutual Information for a Deletion
458 Channel, *ISIT 2012*, Boston, 2012.
- 459 6 S. Diggavi and M. Grossglauser. Information transmission over finite buffer channels,
460 *IEEE Trans. Information Theory*, 52, 1226–1237, 2006.
- 461 7 R.L. Dobrushin. Shannon’s theorem for channels with synchronization errors, *Prob. Info.*
462 *Trans.*, 18–36, 1967.
- 463 8 P. Flajolet, W. Szpankowski, and B. Vallée. Hidden word statistics, *Journal of the ACM*,
464 53, 1–37, 2006.
- 465 9 Allan Gut. *Probability: A Graduate Course*, 2nd ed., Springer, New York, 2013.
- 466 10 W. Hoeffding. A class of statistics with asymptotically normal distribution, *Ann. Math.*
467 *Statistics* 19, 293–325, 1984.
- 468 11 N. Holden and R. Lyones. Lower Bounds for Trace Reconstruction, [arXiv:1808.02336](https://arxiv.org/abs/1808.02336),
469 2018.
- 470 12 P. Jacquet and W. Szpankowski. *Analytic Pattern Matching: From DNA to Twitter*,
471 Cambridge University Press, 2015.
- 472 13 S. Janson, B. Nakamura, and D. Zeilberger. On the asymptotic statistics of the number
473 of occurrences of multiple permutation patterns, *J. Comb.*, 6, 117–143, 2015.
- 474 14 A. Kalai, M. Mitzenmacher, and M. Sudan. Tight asymptotic bounds for the deletion
475 channel with small deletion probabilities, *ISIT*, Austin, 2010.
- 476 15 Y. Kanoria and A. Montanari, On the deletion channel with small deletion probability,
477 *ISIT*, Austin, 2010; see [arXiv:1104.5546](https://arxiv.org/abs/1104.5546) for an extension.
- 478 16 A. McGregor, E. Price, and S. Vorotnikova. Trace Reconstruction Revisited, *European*
479 *Symposium on Algorithms*, 689–700, 2014.
- 480 17 M. Mitzenmacher. A survey of results for deletion channels and related synchronization
481 channels, *Probab. Surveys*, 1–33, 2009.
- 482 18 Y. Peres and A. Zhai. Average-case reconstruction for the deletion channel: subpolynomi-
483 ally many traces suffice, *FOCS*, 2017.
- 484 19 R. Venkataramanan, S. Tatikonda, and K. Ramchandran. Achievable rates for channels
485 with deletions and insertions, *ISIT*, St. Petersburg, Russia, 2011.

486 **Appendix**487 **3.6 Proof of Lemma 11**

488 We show that the central limit theorem applies to the sum $V_1 = \sum_i V_{1,i}$ in (55). The terms
489 $V_{1,i}$ are independent and have means $\mathbb{E} V_{1,i} = 0$. We verify Lyapunov's condition.

490 Since the random variable ξ takes values in the finite set \mathcal{A} , the linear space \mathcal{V} of functions
491 of ξ has finite dimension $|\mathcal{A}|$. Moreover, every function in \mathcal{V} is bounded. The L^2 and L^3
492 norms $\|\cdot\|_2$ and $\|\cdot\|_3$ are thus finite on \mathcal{V} , and are thus both norms on the finite-dimensional
493 vector space \mathcal{V} ; hence there exists a constant C such that for any function f ,

$$494 \quad \|f(\xi)\|_3 \leq C\|f(\xi)\|_2. \quad (88)$$

496 In particular, since the definition (56) shows that $V_{1,i}$ is a function of $\xi_i \stackrel{d}{=} \xi$,

$$497 \quad \|V_{1,i}\|_3 \leq C\|V_{1,i}\|_2 = C\tau_i, \quad 1 \leq i \leq n. \quad (89)$$

499 Furthermore, by (59) and (62),

$$500 \quad \frac{\max_i \tau_i^2}{\sigma_1^2} \leq \frac{B \binom{n-1}{m-1}^2}{c' \frac{n}{m} \binom{n-1}{m-1}^2} = C \frac{m}{n} = o(1). \quad (90)$$

502 Consequently, using (89), (57) and (90),

$$503 \quad \frac{\sum_{i=1}^n \mathbb{E} |V_{1,i}|^3}{\sigma_1^3} = \frac{\sum_{i=1}^n \|V_{1,i}\|_3^3}{\sigma_1^3} \leq \frac{C \sum_{i=1}^n \tau_i^3}{\sigma_1^3} \leq C \frac{\max_i \tau_i \sum_{i=1}^n \tau_i^2}{\sigma_1^3}$$

$$504 \quad = C \frac{\max_i \tau_i}{\sigma_1} = o(1). \quad (91)$$

506 This shows the Lyapunov condition, and thus a standard form of the central limit theorem,
507 [9, Theorem 7.2.4 or 7.6.2], yields (73).