# COUNTING FINITE LANGUAGES BY TOTAL WORD LENGTH

**Stefan Gerhold**[1]
*Vienna University of Technology, Wiedner Hauptstraße, Vienna, Austria*
`sgerhold at fam.tuwien.ac.at`

### Abstract

We investigate the number of sets of words that can be formed from a finite alphabet, counted by the total length of the words in the set. An explicit expression for the counting sequence is derived from the generating function, and asymptotics for large alphabet size and large total word length are discussed. Moreover, we derive a Gaussian limit law for the number of words in a random finite language.

## 1. Introduction and Basic Properties

The results of Chomsky and Schützenberger [1] on generating functions of formal languages are classical in combinatorics. This theory, and much of the related literature, is mainly concerned with problems of the kind: How many words of length $n$ are there in an infinite formal language, which is defined by some specification?

On the other hand, the following questions have received little attention: If the size of a *finite* language is fixed, in terms of the total length of all the words it contains, then how many words are there, on average? How many such finite languages are there? The answers are not immediate even without any restrictions on the words of the language (forbidden patterns etc.), so we focus on the unconstrained case in the present note. The analysis of these problems nicely illustrates the interplay between symbolic specifications and complex analytic properties of generating functions.

Not surprisingly, languages that consist of just one very long word dominate the count, in the sense that they contribute an exponential factor $m^n$. (Here and in what follows, we write $m$ for the size of the underlying alphabet, and $n$ for the total word length.) The subexponential factors can be found, in a rather straightforward way, from the behavior of the pertinent generating function at its dominating singularity.

Two noteworthy, and maybe surprising, features emerge: First, the asymptotic count is uniform in alphabet size $m$ and total word length $n$. Second, the average number of words in a finite language is concentrated around $\sqrt{n}$, independently of the alphabet size.

To fix notation, let $f_n = f_n(m)$ denote the number of formal languages (i.e., sets of words) with total word length $n$ over an alphabet with $m \geq 2$ symbols [3, I.37]. For instance, $f_2(2) = 5$ and $f_3(2) = 16$, as seen from the listings

$$\{a, b\}, \{aa\}, \{ab\}, \{ba\}, \{bb\}$$

respectively

$$\{a, aa\}, \{a, ab\}, \{a, ba\}, \{a, bb\}, \{b, aa\}, \{b, ab\}, \{b, ba\}, \{b, bb\}, \{aaa\},$$
$$\{aab\}, \{aba\}, \{abb\}, \{baa\}, \{bab\}, \{bba\}, \{bbb\}.$$

Another value is $f_2(3) = 12$, illustrated by

$$\{aa\}, \{ab\}, \{ac\}, \{ba\}, \{bb\}, \{bc\}, \{ca\}, \{cb\}, \{cc\}, \{a, b\}, \{a, c\}, \{b, c\}.$$

The sequence $f_n(2)$ is number **A102866** of Sloane's On-Line Encyclopedia of Integer Sequences.[2] The ordinary generating function (ogf) [3, I.37]

$$F(z) := \sum_{n=0}^{\infty} f_n z^n = \exp\left(\sum_{k=1}^{\infty} \frac{(-1)^{k-1}}{k} \frac{mz^k}{1 - mz^k}\right) \tag{1}$$

can be obtained by a standard procedure (the "power set construction" [3, I.2]; finite languages are sets of sequences built from alphabet elements). Its first terms are

$$F(z) = 1 + mz + \tfrac{1}{2}m(3m - 1)z^2 + m(\tfrac{13}{6}m^2 - \tfrac{1}{2}m + \tfrac{1}{3}) + \mathrm{O}(z^4). \tag{2}$$

Note that

$$F(z) = \exp\left(\frac{mz}{1 - mz}\right) \phi(z),$$

where

$$\phi(z; m) = \phi(z) := \exp\left(\sum_{k=2}^{\infty} \frac{(-1)^{k-1}}{k} \frac{mz^k}{1 - mz^k}\right) \tag{3}$$

is analytic for $|z| < 1/\sqrt{m}$. (Indeed, for $0 < \varepsilon < 1/\sqrt{m}$, $|z| \leq 1/\sqrt{m} - \varepsilon$, and $k \geq 2$, we have

$$m|z|^k \leq m|z|^2 \leq m(m^{-1/2} - \varepsilon)^2 = 1 - \varepsilon\sqrt{m}(2 - \varepsilon\sqrt{m}) =: 1 - \varepsilon',$$

whence

$$\left|\frac{mz^k}{1 - mz^k}\right| \leq \frac{m|z|^k}{\varepsilon'}.)$$

The dominating singularity of $F(z)$ is thus located at $z = 1/m$, leading to the rough approximation $f_n(m) \approx m^n$. Clearly, we have $f_n(m) > m^n$ for $m, n \geq 2$ (consider languages consisting only of one word). We will see in Theorem 3 below that the ratio $f_n(m)/m^n$ is $e^{2\sqrt{n}+O(\log n)}$.

Our first result is an explicit expression for $f_n(m)$, which can be obtained from (1). To state it, we write $\mathbf{i} \vdash n$, if the vector $\mathbf{i} = (i_1, \ldots, i_n) \in \mathbb{Z}_{\geq 0}^n$ represents a partition of $n$, in the sense that $i_1 + 2i_2 + \cdots + ni_n = n$.

**Theorem 1.** *For $m \geq 2$ and $n \geq 1$, we have*

$$f_n(m) = \sum_{\mathbf{i} \vdash n} \frac{A_1(m)^{i_1} \ldots A_n(m)^{i_n}}{i_1! \ldots i_n!}, \tag{4}$$

*where*

$$A_j(m) := \sum_{d|j} (-1)^{d-1} m^{j/d}/d, \qquad j \geq 1.$$

*Proof.* We expand the Lambert series [6] in the exponent of $F(z)$, using the geometric series formula, and then collect terms:

$$F(z) = \exp\left(\sum_{k=1}^{\infty} \frac{(-1)^{k-1}}{k} \sum_{j=1}^{\infty} m^j z^{kj}\right)$$

$$= \exp\left(\sum_{n=1}^{\infty} A_n(m)z^n\right)$$

$$= \prod_{n=1}^{\infty} \exp\left(A_n(m)z^n\right)$$

$$= \prod_{n=1}^{\infty} \sum_{i_n=0}^{\infty} \frac{A_n(m)^{i_n} z^{ni_n}}{i_n!}$$

$$= \sum_{n=0}^{\infty} z^n \sum_{\mathbf{i} \vdash n} \frac{A_1(m)^{i_1} \ldots A_n(m)^{i_n}}{i_1! \ldots i_n!}.$$

$\square$

## 2. Asymptotics for Large Alphabet Size

Next we derive the asymptotics of $f_n(m)$ as $m$, the cardinality of the alphabet, tends to infinity. Define $\kappa_n$ and $\mu_n = \mu_n(m)$ by

$$\sum_{n=0}^{\infty} \kappa_n z^n = \exp\left(\frac{z}{1-z}\right) \quad \text{and} \quad \sum_{n=0}^{\infty} \mu_n z^n = \phi(z).$$

Note that $n!\kappa_n$ is Sloane's **A000262** (several combinatorial interpretations are given on that web page), and that $\kappa_n$ has the representation

$$\kappa_n = \sum_{\mathbf{i} \vdash n} \frac{1}{i_1! \ldots i_n!}, \qquad n \geq 1. \tag{5}$$

Then we can write

$$f_n/m^n = [z^n] \exp\left(\frac{z}{1-z}\right) \phi(z/m)$$
$$= \kappa_n + \kappa_{n-1}\mu_1/m + \cdots + \kappa_0\mu_n/m^n.$$

If the dependence of $\mu_n$ on $m$ is not too strong, the first term on the right-hand side should dominate when $m \to \infty$. This is indeed the case:

**Theorem 2.** *If $n \geq 1$ is fixed and $m \to \infty$, we have*

$$f_n(m) \sim \kappa_n m^n. \tag{6}$$

*Proof.* Since, as $m \to \infty$,

$$A_j(m) = m^j + \mathrm{O}(m^{j/2}), \qquad j \geq 1,$$

we have

$$A_j(m)^{i_j} = m^{j i_j}(1 + \mathrm{O}(m^{-j/2})),$$

whence, for $\mathbf{i} \vdash n$,

$$A_1(m)^{i_1} \ldots A_n(m)^{i_n} = m^n(1 + \mathrm{O}(m^{-1/2})).$$

The result thus follows from (4) and (5).                              $\square$

Note that $\kappa_1 = 1$, $\kappa_2 = \frac{3}{2}$, and $\kappa_3 = \frac{13}{6}$, in line with (2).

## 3. Asymptotics for Large Total Word Length

**Theorem 3.** *For large total word length $n$, the sequence $f_n = f_n(m)$ has the asymptotics*

$$f_n \sim \frac{\phi(1/m)}{2\sqrt{e\pi}} \times \frac{m^n e^{2\sqrt{n}}}{n^{3/4}}, \qquad n \to \infty, \tag{7}$$

*where $\phi$ is defined in (3). More precisely, there is a full asymptotic expansion of the form*

$$f_n \sim \frac{\phi(1/m)}{2\sqrt{e\pi}} \times \frac{m^n e^{2\sqrt{n}}}{n^{3/4}} \left(1 + \sum_{j \geq 1} c_j n^{-j/2}\right), \qquad n \to \infty. \tag{8}$$

The expansion (8) is a special case of a result by Wright [7], who studied Taylor coefficients of functions of the form "analytic factor times exponential of a pole" (he also allowed for a logarithmic factor). Alternatively, the first order asymptotics (7) can be also be obtained from the Hayman-admissibility [3, 4] of the ogf $F(z)$.

Nevertheless, we sketch the proof of (7) and (8), because we will revisit it in Sections 4 and 5. The steps are very similar to the saddle point analysis [3, Example VIII.7] of $\exp(z/(1-z))$, the ogf of $\kappa_n$, slightly perturbed by the presence of the analytic factor $\phi(z)$. Let us shift the dominating singularity from $z = 1/m$ to $z = 1$. Then the integrand in Cauchy's formula

$$f_n = f_n(m) = \frac{m^n}{2i\pi} \oint \frac{F(z/m)}{z^{n+1}} dz \tag{9}$$

has an approximate saddle point at $z = \hat{z} := 1 - 1/\sqrt{n}$. We write $z = \hat{z}e^{i\theta}$, where $\theta = \arg(z)$ is constrained by

$$|\theta| < n^{-\alpha}, \qquad \tfrac{2}{3} < \alpha < \tfrac{3}{4}, \tag{10}$$

so that $z$ lies in a small arc around the saddle point. In this range we have the uniform expansions

$$z^{-n-1} = \exp\left(\sqrt{n} + \tfrac{1}{2} - in\theta + O(n^{-1/2})\right), \qquad n \to \infty, \tag{11}$$

and

$$\frac{1}{1-z} = \sqrt{n} + i\theta n - n^{3/2}\theta^2 + O(n^{1/2-\alpha}). \tag{12}$$

Recall that $\phi(z/m)$ is analytic at $z = 1$. Therefore, by (11) and (12), the local expansion of the integrand in (9) at the saddle point $\hat{z}$ is

$$\frac{F(z/m)}{z^{n+1}} = \phi(1/m) \exp\left(-\tfrac{1}{2} + 2\sqrt{n} - n^{3/2}\theta^2\right) \times \left(1 + O(n^{1/2-\alpha})\right), \tag{13}$$

valid as $n \to \infty$, uniformly w.r.t. $\theta$ in the range (10). Note that

$$\int_{-n^{-\alpha}}^{n^{-\alpha}} e^{-n^{3/2}\theta^2} d\theta \sim \sqrt{\pi} n^{-3/4},$$

so that integrating (13) from $-n^{-\alpha}$ to $n^{-\alpha}$ yields the right-hand side of (7). To prove (7), it remains to show that the integral from $n^{-\alpha}$ to $\pi$ grows slower (the other half of the tail is handled by symmetry). There is a $C > 0$ such that

$$\left|\frac{F(z/m)}{z^{n+1}}\right| \le C|z|^{-n} \exp \Re \left(\frac{1}{1-z}\right), \qquad |z| < 1. \tag{14}$$

If $z = \hat{z}e^{i\theta}$ lies on the integration contour, then the factor $|z|^{-n}$ in (14) is $O(e^{\sqrt{n}})$. The remaining factor $\exp \Re(1/(1-z))$ decreases if $|\theta| = |\arg(z)|$ increases, hence

$$\int_{n^{-\alpha}}^{\pi} \exp \Re \left(\frac{1}{1 - \hat{z}e^{i\theta}}\right) d\theta \le \pi \exp \Re \left(\frac{1}{1 - \hat{z}e^{i\theta}}\right)\Bigg|_{\theta=n^{-\alpha}}$$
$$= \exp\left(\sqrt{n} - n^{3/2-2\alpha} + O(1)\right).$$

(The last line is obtained by recapitulating the derivation of (13), with $\theta = n^{-\alpha}$.)
Hence

$$\left| \oint_{n^{-\alpha} < |\theta| < \pi} \frac{F(z/m)}{z^{n+1}} \mathrm{d}z \right| \leq \exp\left(2\sqrt{n} - n^{3/2 - 2\alpha} + \mathrm{O}(1)\right). \qquad (15)$$

This indeed grows slower than $\mathrm{e}^{2\sqrt{n}}/n^{3/4}$, so that the proof of (7) is complete.

The availability of a full asymptotic expansion is a typical feature of the saddle point method [3, Example VIII.4]. To justify (8), first note that it suffices to check that the central part $\int_{-n^{-\alpha}}^{n^{-\alpha}} \mathrm{d}\theta$ of the Cauchy integral has such an expansion, as the tail estimate (15) lies asymptotically below (8). This expansion is found by tedious, but straightforward calculations: Take more terms in (11), (12), and the Taylor series of $\phi(z/m)$ at the saddle point, and integrate from $\theta = -n^{-\alpha}$ to $n^{-\alpha}$.

## 4. Joint Asymptotics

Note that the limits $m \to \infty$ and $n \to \infty$ commute in the following sense: Since we have $\kappa_n \sim 1/(2\sqrt{e\pi})\mathrm{e}^{2\sqrt{n}}/n^{3/4}$ [3, Prop. 8.4], the right-hand side of (6) has, as $n \to \infty$, the same asymptotics as the right-hand side of (7) for $m \to \infty$. We will now show that letting $m$ and $n$ tend to infinity simultaneously yields the same result, regardless of their respective speeds.

**Theorem 4.** *If both the word length and the alphabet size tend to infinity, we have*

$$f_n(m) \sim \frac{1}{2\sqrt{e\pi}} \times \frac{m^n \mathrm{e}^{2\sqrt{n}}}{n^{3/4}}, \qquad m, n \to \infty.$$

*Proof.* The result can be obtained by an adaption of the proof of Theorem 3. Again we use Cauchy's formula, with the same saddle point contour as before:

$$f_n(m) = \frac{m^n}{2\pi} \hat{z}^{-n} \int_{-\pi}^{\pi} F(\hat{z}\mathrm{e}^{\mathrm{i}\theta}/m)\mathrm{e}^{-\mathrm{i}(n+1)\theta} \mathrm{d}\theta$$

$$= \frac{m^n}{2\pi} \hat{z}^{-n} \int_{-\pi}^{\pi} \exp\left(\frac{\hat{z}\mathrm{e}^{\mathrm{i}\theta}}{1 - \hat{z}\mathrm{e}^{\mathrm{i}\theta}}\right) \phi\left(\frac{\hat{z}\mathrm{e}^{\mathrm{i}\theta}}{m}; m\right) \mathrm{e}^{-\mathrm{i}(n+1)\theta} \mathrm{d}\theta. \qquad (16)$$

We will show that

$$\phi\left(\frac{\hat{z}\mathrm{e}^{\mathrm{i}\theta}}{m}; m\right) \to 1, \qquad m, n \to \infty, \text{ uniformly w.r.t. } \theta \in [-\pi, \pi]. \qquad (17)$$

Assuming this we are done. Indeed, assertion (17) shows at the same time the validity of the local expansion (13), with $\phi(1/m)$ replaced by 1, and the persistence of the tail estimate (15).

To prove (17), notice that

$$\phi\Big(\frac{\hat{z}\mathrm{e}^{\mathrm{i}\theta}}{m};m\Big) = \exp\left(\sum_{k=2}^{\infty}\frac{(-1)^{k-1}}{k}\frac{m^{1-k}\hat{z}^k\mathrm{e}^{ki\theta}}{1-m^{1-k}\hat{z}^k\mathrm{e}^{ki\theta}}\right). \tag{18}$$

We have $|m^{1-k}\hat{z}^k\mathrm{e}^{ki\theta}| < \frac{1}{2}$ for $m \geq 2$, and hence

$$\left|\sum_{k=2}^{\infty}\frac{(-1)^{k-1}}{k}\frac{m^{1-k}\hat{z}^k\mathrm{e}^{ki\theta}}{1-m^{1-k}\hat{z}^k\mathrm{e}^{ki\theta}}\right| \leq \sum_{k=2}^{\infty}m^{1-k}\hat{z}^k$$

$$= \sum_{k=2}^{\infty}m^{1-k}\left(1-\frac{1}{\sqrt{n}}\right)^k$$

$$= \frac{(1-1/\sqrt{n})^2}{m(1-1/m+1/(m\sqrt{n}))}.$$

Thus the exponent in (18) is uniformly o(1), which establishes (17).                    □

## 5. The Distribution of the Number of Words

A natural parameter to consider is the number $W_n$ of words in a random finite language of total word length $n$. (The alphabet size $m \geq 2$ is fixed throughout this section.) The appropriate bivariate ogf, with $z$ marking total word length and $u$ marking number of words, is given by

$$F(z,u) := \exp\left(\sum_{k=1}^{\infty}\frac{(-1)^{k-1}}{k}\frac{mz^ku^k}{1-mz^k}\right).$$

The expected number of words is then

$$\mathbf{E}[W_n] = f_n^{-1}[z^n]\partial_u F(z,u)|_{u=1}. \tag{19}$$

Notice that

$$\partial_u F(z,u)|_{u=1} = F(z)\sum_{k=1}^{\infty}\frac{(-1)^{k-1}mz^k}{1-mz^k}, \tag{20}$$

so that the asymptotic analysis of $[z^n]\partial_u F(z,u)|_{u=1}$ is an easy extension of the one of $f_n = [z^n]F(z)$ in Section 3: Close to the saddle point, the new factor resulting from the right-hand side of (20) is

$$\frac{1}{1-z} = \sqrt{n}(1+\mathrm{o}(1)).$$

Hence $[z^n]\partial_u F(z,u)|_{u=1} \sim \sqrt{n}f_n$, so that, by (19), the expectation of $W_n$ satisfies

$$\mathbf{E}[W_n] \sim \sqrt{n}, \qquad n \to \infty.$$

Similarly, one can obtain the asymptotics $\sigma(W_n) \sim n^{1/4}/\sqrt{2}$ for the standard deviation.

Finally, we show that the law of the normalized number of words is asymptotically normal. To do so, we appeal to a result by Sačkov [3, 5]; alternatively, Drmota et al.'s concept of extended Hayman-admissibility [2] could have been used.

**Theorem 5.** *The number of words $W_n$ in a random finite language admits a Gaussian limit law:*

$$\frac{W_n - a_n}{b_n} \to \mathcal{N}(0,1), \qquad n \to \infty,$$

*in distribution, where the scaling constants satisfy $a_n \sim \sqrt{n}$ and $b_n \sim n^{1/4}/\sqrt{2}$.*

*Proof.* As is well known, combinatorial limit laws can often be obtained by an asymptotic analysis of the probability generating function

$$\mathbf{E}[u^{W_n}] = f_n^{-1}[z^n]F(z,u). \tag{21}$$

Again, we adapt the proof of Theorem 3. If $u$ ranges in a fixed small neighborhood of $u = 1$, the expansion (13) generalizes to the uniform local expansion

$$\frac{F(z/m,u)}{z^{n+1}} = \phi(1/m,u;m)\exp\big(-\tfrac{1}{2}u + 2\sqrt{un} - u^{-1/2}n^{3/2}\theta^2\big) \times \big(1 + \mathrm{O}(n^{1/2-\alpha})\big),$$

where

$$\phi(z,u;m) := \exp\left(\sum_{k=2}^{\infty} \frac{(-1)^{k-1}}{k} \frac{mz^k u^k}{1 - mz^k}\right).$$

Integrating from $\theta = -n^{-\alpha}$ to $n^{-\alpha}$, and taking into account (7), we infer that (21) has the uniform asymptotics

$$\mathbf{E}[u^{W_n}] \sim \exp(h_n(u)), \qquad n \to \infty,$$

with

$$h_n(u) := 2(\sqrt{u} - 1)\sqrt{n} + \tfrac{1}{4}\log u + \log\frac{\phi(1/m,u;m)}{\phi(1/m;m)}.$$

Note that, for $n \to \infty$,

$$h_n'(1) = \sqrt{n} + \mathrm{O}(1),$$
$$h_n''(1) = -\tfrac{1}{2}\sqrt{n} + \mathrm{O}(1),$$
$$h_n'''(1) = \tfrac{3}{4}\sqrt{n} + \mathrm{O}(1),$$

so that the function $h_n(u)$ satisfies the conditions of Theorem 9.13 in Flajolet and Sedgewick's monograph [3], itself taken from Sačkov [5]. We conclude that

$$\frac{W_n - h_n'(1)}{(h_n'(1) + h_n''(1))^{1/2}}$$

converges in distribution to a standard normal random variable. $\qquad\square$

## 6. Possible Extensions

Let $\mathcal{L}$ be some (infinite) formal language containing $b_n$ words of length $n$, with ogf

$$B(z) = b_1 z + b_2 z^2 + \dots$$

Then

$$F_{\mathcal{L}}(z) = \sum_{k=0}^{\infty} f_{\mathcal{L},n} z^n = \exp\left( \sum_{k=1}^{\infty} \frac{(-1)^{k-1}}{k} B(z^k) \right) \tag{22}$$

is the ogf of finite languages built from words of $\mathcal{L}$. (So far, we had $\mathcal{L} = \bigcup_{n \geq 1} \mathcal{A}^n$, the language of *all* words over a finite alphabet $\mathcal{A} = \{\mathsf{a}_1, \dots, \mathsf{a}_m\}$, and $b_n = m^n$.)

A large class of languages admits rational generating functions [1, 3]. The explicit expression from Theorem 1 can be adapted, mutatis mutandis, to any $\mathcal{L}$ featuring a rational ogf $B(z)$. Moreover, the ogf (22) is amenable to Hayman's method for any rational $B(z)$. The dominating factor of $f_{\mathcal{L},n}$ follows from the location of the dominating pole of $B(z)$, while the order of the pole determines the subexponential factor. On the other hand, the results in Sections 4 and 5 require uniformity properties that are not immediately obvious for general $B(z)$. As a natural question for future research, we ask for conditions on $\mathcal{L}$ that ensure uniform asymptotics and/or a Gaussian limit law for $f_{\mathcal{L},n}$.

## References

[1] N. Chomsky and M. P. Schützenberger, *The algebraic theory of context-free languages*, in Computer programming and formal systems, North-Holland, Amsterdam, 1963, pp. 118–161.

[2] M. Drmota, B. Gittenberger, and T. Klausner, *Extended admissible functions and Gaussian limiting distributions*, Math. Comp., 74 (2005), pp. 1953–1966.

[3] P. Flajolet and R. Sedgewick, *Analytic Combinatorics*, Cambridge University Press, 2009.

[4] W. K. Hayman, *A generalisation of Stirling's formula*, J. Reine Angew. Math., 196 (1956), pp. 67–95.

[5] V. N. Sačkov, *Veroyatnostnye metody v kombinatornom analize*, "Nauka", Moscow, 1978.

[6] H. S. Wilf, *generatingfunctionology*, A K Peters Ltd., Wellesley, MA, third ed., 2006.

[7] E. M. Wright, *The coefficients of a certain power series*, J. London Math. Soc., 7 (1932), pp. 256–262.