



UNIVERSIDAD DE BUENOS AIRES

FACULTAD DE CIENCIAS EXACTAS Y NATURALES

DEPARTAMENTO DE COMPUTACIÓN

# **Secuencias completamente equidistribuidas basadas en secuencias de De Bruijn**

Tesis de Licenciatura en Ciencias de la Computación

Emilio Guido Almansi  
ealmansi@gmail.com, LU: 674/12

Directora: Verónica Becher

Buenos Aires, 2019

## SECUENCIAS COMPLETAMENTE EQUIDISTRIBUIDAS BASADAS EN SECUENCIAS DE DE BRUIJN

En este trabajo estudiamos una secuencia de números reales completamente equidistribuidos publicada por Donald Knuth en 1965. La noción de equidistribución completa se utiliza en su sentido clásico; es decir, que para una secuencia dada, todas sus subsecuencias finitas y contiguas de cualquier longitud presentan una distribución uniforme. En un artículo de 1963, Joel Franklin considera esta propiedad como un primer requerimiento de pseudoaleatoriedad en secuencias determinísticas, y prueba que la equidistribución completa implica muchas otras propiedades importantes de las secuencias aleatorias. El trabajo de Knuth se basa en secuencias de De Bruijn, las cuales tienen también una relación cercana con la noción de equidistribución y pueden ser generadas en tiempo constante amortizado por el algoritmo FKM (Fredricksen, Kessler, Maiorana, 1978). Presentamos una variante de la secuencia de Knuth mediante una construcción similar, aunque más sencilla, y damos una prueba elemental de que la secuencia generada también es completamente equidistribuida.

**Palabras claves:** secuencia aleatoria, equidistribución completa, secuencia de De Bruijn, algoritmo FKM, secuencia de Ford, secuencia de Knuth.

# COMPLETELY EQUIDISTRIBUTED SEQUENCES BASED ON DE BRUIJN SEQUENCES

In this work, we study a construction published by Donald Knuth in 1965 yielding a completely equidistributed sequence of real numbers. Complete equidistribution is interpreted in its classic sense; namely, that finite contiguous subsequences of any length have a uniform distribution within a given sequence. Joel Franklin in a paper from 1963 suggests this as a first requirement for pseudorandomness in deterministic sequences, and proves that complete equidistribution implies many other important statistical properties shared by all random sequences. Knuth's work is based on De Bruijn sequences, which are also closely related to equidistribution and can be generated by the FKM algorithm (Fredricksen, Kessler, Maiorana, 1978) in amortized constant time. We provide a variant of Knuth's sequence via a similar, albeit simpler, construction and give an elementary proof showing that the sequence it yields is also completely equidistributed.

**Keywords:** Random Sequence, Complete Equidistribution, De Bruijn Sequence, FKM Algorithm, Ford Sequence, Knuth Sequence.

## CONTENTS

1. Preliminaries . . . . .	1
1.1 Random Sequences . . . . .	1
1.2 Notational Conventions . . . . .	1
1.3 Complete Equidistribution . . . . .	2
1.4 Weyl's Criterion . . . . .	4
1.5 De Bruijn Sequences and Ford Sequences . . . . .	4
2. Completely Equidistributed Sequences Based on De Bruijn Sequences . . . . .	6
2.1 Knuth's Sequence . . . . .	6
2.2 Linearly Increasing Alphabet Sizes . . . . .	8
2.2.1 Proof of Theorem 1 . . . . .	9
2.2.2 Alternative Proof of Theorem 1 . . . . .	18

## 1. PRELIMINARIES

### 1.1 Random Sequences

In engineering, computer science, and other branches of science we are often required to simulate random processes. Simulations usually involve the generation of non-random, deterministic sequences of numbers which approximate a sequence of independent, random samples from a given probability distribution. While a deterministic sequence can never be random in the sense of not following any patterns or being entirely unpredictable, nothing prevents one from identifying properties common to all random sequences and constructing pseudo-random sequences satisfying these properties.

Significant work has been done in this direction. In a paper by Franklin [4], the notion of equidistribution is discussed and presented as a first requirement for randomness. Franklin proves that the sequence of powers  $\alpha, \alpha^2, \alpha^3, \dots$  is completely equidistributed modulo 1 for *almost all*  $\alpha > 1$ . However, no specific value of  $\alpha$  is provided that verifies this property. A concrete construction yielding a completely equidistributed sequence is given by Knuth [7], and based on De Bruijn sequences of increasing order and alphabet size. In this work, we provide a similar albeit simpler construction and prove that the sequence it yields is also completely equidistributed.

Since any real computer has a finite word-length and a finite amount of memory, it can only produce numbers of limited precision and sequences that are ultimately periodic. Thus, in practice a pseudo-random number generator (PRNG) is unable to produce a truly completely equidistributed sequence of real numbers. Instead, practical PRNGs are evaluated by empirical randomness tests to show that they possess suitable equidistribution properties (see [8]), and are often equidistributed only up to a finite number of dimensions. For example, numbers produced by the Mersenne Twister, introduced by Matsumoto and Nishimura in [10], show 623-dimensional equidistribution.

In subsequent sections, we ignore any practical limitations by considering a computational model with infinite memory and with infinite word-length, where real numbers can be stored and computed with perfect precision.

### 1.2 Notational Conventions

Since notation on sequences can differ across the literature, we state the conventions used throughout this work.

When discussing a sequence  $X = x_1, x_2, \dots$ , the expression  $X_i$  denotes the  $i$ -th element

of the sequence, with the first element having an index of 1.

The expression  $X_{1:n}$  denotes a prefix of length  $n$  from the sequence  $X$ . Namely,  $X_{1:n} = x_1, x_2, \dots, x_n$ .

Given two sequences  $X$  and  $Y$ , the expression  $\langle X; Y \rangle$  denotes the sequence obtained by concatenating  $Y$  after  $X$ , when this operation is well-defined. We also extend this notation to more than two input sequences; for example, if  $A = 1, 2, 3$ ,  $B = 3$ , and  $C = 4, 5$ , then  $\langle A; B; C \rangle = 1, 2, 3, 3, 4, 5$ .

### 1.3 Complete Equidistribution

In order to define complete equidistribution, we first need to define a notion of "probability" for deterministic sequences. Let  $P = p_1, p_2, \dots$  be an infinite sequence of predicates and  $\sigma$  a function such that  $\sigma(p_i) = 1$  if  $p_i$  is *true*, and  $\sigma(p_i) = 0$  otherwise. We define:

$$Pr(P) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \sigma(p_i)$$

when this limit exists.

For example, given a sequence  $x_1, x_2, \dots$ , we can define:

$$Pr(x_i < x_{i+1}) = \lim_{N \rightarrow \infty} \frac{1}{N} \sum_{i=1}^N \sigma(x_i < x_{i+1})$$

if the limit exists. Note that  $Pr(x_i < x_{i+1})$  denotes  $Pr(P)$  where  $P = (x_i < x_{i+1})_{i=1}^{\infty}$ .

We say that an infinite sequence  $X = x_1, x_2, \dots$  of real numbers is *equidistributed in the unit interval* if the probability of finding  $x_i$  in any subinterval is proportional to the length of the subinterval. Formally, if for every interval  $I = [u, v) \subseteq [0, 1)$  the following is true:

$$Pr(x_i \in I) = |I| = v - u.$$

Analogously, we say that an infinite sequence  $\bar{X} = \bar{x}_1, \bar{x}_2, \dots$  of  $k$ -dimensional vectors of real numbers is *equidistributed in the unit cube* if for every set  $I = [u_1, v_1) \times \dots \times [u_k, v_k) \subseteq [0, 1)^k$  the following is true:

$$Pr(\bar{x}_i \in I) = |I| = \prod_{d=1}^k v_d - u_d.$$

For every positive integer  $k$  we define the *windows sequence of  $X$  of order  $k$* , denoted  $W_k(X)$ , as the sequence of  $k$ -dimensional vectors containing every possible window (or contiguous subsequence) of  $X$  from left to right:

$$\begin{aligned} W_k(X) &= (x_1, x_2, \dots, x_k), (x_2, x_3, \dots, x_{k+1}), \dots \\ &= ((x_i, x_{i+1}, \dots, x_{i+k-1}))_{i=1}^{\infty}. \end{aligned} \tag{1.1}$$

We say that  $X$  is  *$k$ -distributed* in the unit interval if its windows sequence of order  $k$  is equidistributed in the unit cube. Note that, as per the definition above,  $W_k(X)$  includes windows with superposition. If one instead considers windows without superposition, then the derived notion of  $k$ -distribution is not equivalent. See Theorem 19 in [4] for an example of a sequence which is 2-distributed in the former sense, but not in the latter.

If  $X$  is  $k$ -distributed for every positive integer  $k$ , then we additionally say that  $X$  is *completely equidistributed*. In [4], Franklin proves that if  $X$  is completely equidistributed, then it also satisfies many other statistical properties common to all random sequences. For example, for every fixed  $k$  the lag  $k$  autocorrelation of  $X$  is zero, and the probability of  $k$  consecutive terms having any specific relative order is  $1/k!$ .

For convenience, we extend the notion of windows sequences to finite and cyclic sequences. In the case of a finite sequence  $Y = y_1, y_2, \dots, y_l$  of length  $l$ , we define  $W_k(Y)$  similarly to how we did in 1.1:

$$\begin{aligned} W_k(Y) &= (y_1, y_2, \dots, y_k), \dots, (y_{l-k+1}, y_{l-k+2}, \dots, y_l) \\ &= ((y_i, y_{i+1}, \dots, y_{i+k-1}))_{i=1}^{\max(0, l-k+1)}. \end{aligned}$$

Note that the windows sequence of  $Y$  of order  $k$  is empty whenever  $l < k$ , having  $\max(0, l - k + 1)$  terms in general.

In order to avoid ambiguity with the previous definition, for a cyclic sequence  $Z = z_1, z_2, \dots, z_l$  of size  $l$  we denote its windows sequence of order  $k$  as  $W_k^c(Z)$  instead. In this case, the derived sequence has exactly  $l$  terms:

$$\begin{aligned} W_k^c(Z) &= (z_1, z_2, \dots, z_k), \dots, (z_l, z_1, \dots, z_{k-1}) \\ &= ((z_i, z_{i+1}, \dots, z_{i+k-1}))_{i=1}^l \end{aligned}$$

where the indices are taken modulo  $l$ .

## 1.4 Weyl's Criterion

First formulated by Hermann Weyl in a paper from 1916 [14], Weyl's Criterion states that a sequence  $X = x_1, x_2, \dots$  of real numbers is equidistributed in the unit interval if and only if for all non-zero integers  $l$ :

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N e^{2\pi i l x_n} = 0,$$

which allows questions about equidistribution to be reduced to bounds on exponential sums.

The criterion can also be generalized into higher dimensions in the following way. If  $\bar{X} = \bar{x}_1, \bar{x}_2, \dots$  is a sequence of  $k$ -dimensional vectors of real numbers, then  $\bar{X}$  is equidistributed in the unit cube if and only if for all non-zero  $k$ -dimensional vectors of integers  $\bar{\ell} = (l_1, \dots, l_k)$ :

$$\lim_{N \rightarrow \infty} \frac{1}{N} \sum_{n=1}^N e^{2\pi i \bar{\ell} \cdot \bar{x}_n} = 0,$$

where  $\bar{\ell} \cdot \bar{x}_n$  denotes the dot product of  $\bar{\ell}$  and  $\bar{x}_n$ .

See [13], pages 7 and 48, for proof of these statements.

## 1.5 De Bruijn Sequences and Ford Sequences

De Bruijn sequences are well-studied sequences from the field of combinatorics on words; see [1] for a historical introduction. Additionally, these sequences show a close link to the property of equidistribution; namely, that any De Bruijn sequence  $X$  of order  $k$  can be extended into a  $k$ -distributed sequence  $X' = \langle X; X; \dots \rangle$ . Throughout this work, we use the terms De Bruijn sequence and Ford sequence according to the following definitions.

**Definition.** A  $b$ -ary De Bruijn sequence of order  $k$  is any sequence of length  $b^k$  which, when viewed as a cycle, contains every possible  $b$ -ary sequence of length  $k$  exactly once as a contiguous subsequence.

**Example.** Listed next are two distinct binary De Bruijn sequences of order 3:

$$\begin{aligned} &0, 0, 0, 1, 0, 1, 1, 1; \\ &0, 0, 0, 1, 1, 1, 0, 1. \end{aligned}$$



Note that all possible binary sequences of length 3 appear exactly once as a contiguous subsequence in each example. This includes those instances such as 1, 0, 0 which wrap around the right-hand end of the sequences.

**Example.** Next, we list two distinct 4-ary De Bruijn sequences of order 2:

0, 0, 1, 0, 2, 0, 3, 1, 1, 2, 1, 3, 2, 2, 3, 3;

0, 0, 1, 1, 2, 1, 3, 1, 0, 2, 2, 3, 2, 0, 3, 3.

There exists a total of  $\frac{(b!)^{b^k-1}}{b^k}$  different  $b$ -ary De Bruijn sequences of order  $k$ . This result first became well-known due to De Bruijn [2], although a later paper by the same author credits C. Flye Sainte-Marie with having previously proven the result for the case  $b = 2$  in 1894.

A simple algorithm for producing a  $b$ -ary De Bruijn sequence of order  $k$  is as follows. First, set  $x_1 = x_2 = \dots = x_k = 0$ . Then, pick  $x_{k+i} \in \{0, 1, \dots, b-1\}$  for  $i = 1 \dots b^k - k$  according to the following rules: i)  $x_{1+i}, \dots, x_{k+i}$  does not duplicate any previous contiguous subsequence of size  $k$ , and ii)  $x_{k+i} = 0$  only when no other options are available that satisfy rule i).

The algorithm as stated above is exhibited by Knuth in [7] and attributed to a manuscript from 1957 by Ford [3]. However, according to a survey by Fredricksen in [5], multiple variants of the same algorithm can be found in the literature tracing back to 1934 in a paper by Martin [9]. In Chapter 2, we use lexicographically least De Bruijn sequences, to which we refer as Ford sequences following the work by Fredricksen, for constructing completely equidistributed sequences.

**Definition.** A  $b$ -ary Ford sequence of order  $k$ , denoted  $F^{(b,k)}$ , is the lexicographically least  $b$ -ary De Bruijn sequence of order  $k$ .

**Example.** The first sequences in each of the examples above are, respectively, the binary Ford sequence of order 3, and the 4-ary Ford sequence of order 2.

In [6], Fredricksen and Maiorana introduce an algorithm to generate a  $b$ -ary Ford sequence of order  $k$  extending previous work by Fredricksen and Kessler. Ruskey, Savage, and Wang later refer to this algorithm as the FKM (Fredricksen, Kessler, Maiorana) algorithm, and prove that it has a constant amortized running time (see [12]).

## 2. COMPLETELY EQUIDISTRIBUTED SEQUENCES BASED ON DE BRUIJN SEQUENCES

In this chapter, we first present the completely equidistributed sequence given by Knuth in [7]. Next, we discuss the given construction and motivate the variant introduced in Section 2.2. We then state the main result of this work about the complete equidistribution of the newly introduced sequence, and provide an elementary proof of this fact based only on the defining property of De Bruijn sequences. Finally, we present a simpler alternative proof based on Weyl's Criterion and a proposition from the theory of linear modular congruences.

### 2.1 Knuth's Sequence

We define Knuth's sequence, denoted as  $K$ , following the construction presented in [7]. In its original formulation, the sequence  $K$  can be constructed using any given family of De Bruijn sequences. Here, we exhibit the construction using Ford sequences since they are easily defined in a univocal manner and can be generated efficiently, but any other fixed family of De Bruijn sequences would yield the same result.

Given a natural number  $n$ , we define:

i) an  $A$  sequence of order  $n$ , denoted  $A^{(n)}$ , as the finite sequence of rational numbers obtained from dividing by  $2^n$  each of the terms in a  $2^n$ -ary Ford sequence of order  $n$ :

$$A^{(n)} = \frac{f_1}{2^n}, \frac{f_2}{2^n}, \dots, \frac{f_{2^{n^2}}}{2^n} = \left( \frac{f_i}{2^n} \right)_{i=1}^{2^{n^2}}$$

where  $F^{(2^n, n)} = f_1, \dots, f_{2^{n^2}}$

and, ii) a  $B$  sequence of order  $n$ , denoted  $B^{(n)}$ , as  $n2^{2n}$  consecutive copies of  $A^{(n)}$ :

$$B^{(n)} = \left\langle \underbrace{A^{(n)}; A^{(n)}; \dots; A^{(n)}}_{n2^{2n} \text{ times}} \right\rangle.$$

By construction, the size of  $A^{(n)}$  is  $|A^{(n)}| = |F^{(2^n, n)}| = 2^{n^2}$ , and the size of  $B^{(n)}$  is  $|B^{(n)}| = n2^{2n}|A^{(n)}| = n2^{2n}2^{n^2}$ . Note as well that, for any given  $n$ , all terms in  $A^{(n)}$  and in  $B^{(n)}$  are numbers in the set  $\left\{0, \frac{1}{2^n}, \frac{2}{2^n}, \dots, \frac{2^n-1}{2^n}\right\} \subset [0, 1)$ .

For example, when  $n = 2$ :

$$\begin{aligned}
 F^{(4,2)} &= 0, 0, 1, 0, 2, 0, 3, 1, 1, 2, 1, 3, 2, 2, 3, 3 \\
 A^{(2)} &= \frac{0}{4}, \frac{0}{4}, \frac{1}{4}, \frac{0}{4}, \frac{2}{4}, \frac{0}{4}, \frac{3}{4}, \frac{1}{4}, \frac{1}{4}, \frac{2}{4}, \frac{1}{4}, \frac{3}{4}, \frac{2}{4}, \frac{2}{4}, \frac{3}{4}, \frac{3}{4} \\
 B^{(2)} &= \left\langle \underbrace{A^{(2)}; \dots; A^{(2)}}_{32 \text{ times}} \right\rangle = \underbrace{\frac{0}{4}, \frac{0}{4}, \dots, \frac{3}{4}, \frac{3}{4}}_{A^{(2)}}, \dots, \underbrace{\frac{0}{4}, \frac{0}{4}, \dots, \frac{3}{4}, \frac{3}{4}}_{A^{(2)}}
 \end{aligned}$$

and  $|A^{(2)}| = 16$ ,  $|B^{(2)}| = 512$ .

We now define Knuth's sequence, denoted as  $K$ , as the infinite sequence of real numbers resulting from the concatenation of all possible  $B$  sequences in increasing order:

$$K = \langle B^{(1)}; B^{(2)}; B^{(3)}; \dots \rangle.$$

**Theorem** (Knuth 1965, [7], page 268). *The sequence  $K$  is completely equidistributed.*

Knuth provides an elementary proof of the theorem stated above. Two choices in the construction yielding  $K$  may seem arbitrary at first glance, but play an important role in this proof. Namely, that the number of repetitions of each  $A$  sequence within a  $B$  sequence is  $n2^{2n}$ , and the fact that the alphabet sizes of the composing Ford sequences grow exponentially as  $2^n$ .

On account of the first choice, one can adapt Knuth's proof in a rather straightforward manner to show that a sufficient condition for the complete equidistribution of  $K$  is for the number of repetitions to grow asymptotically faster than  $2^{2n}$ . A proof of this fact is omitted here since it follows from Knuth's work, together with the technique used in the following section to obtain an analogous result. Knuth's choice of  $n2^{2n}$  repetitions is sufficient to achieve complete equidistribution, but a more reduced number of repetitions such as  $\lceil \log(n+1) \rceil 2^{2n}$  would suffice as well.

In regard to the second choice, Knuth's use of alphabet sizes which grow exponentially as powers of 2 allows one to reason more easily about the rational numbers comprised in the sequence  $K$  in terms of their binary representations. In particular, Knuth uses this device to derive properties of the distribution of the  $m$  most-significant bits of the terms in a  $2^n$ -ary Ford sequence, where  $m \leq n$ . See [7, Lemma 2].

In the next section, we show that it is possible to construct the sequence using linearly increasing alphabet sizes instead, while preserving the property of complete equidistribution. In turn, this allows for a much more reduced number of repetitions of each  $A$  sequence within a  $B$  sequence.

## 2.2 Linearly Increasing Alphabet Sizes

We now introduce the main contribution of this work, which is a variant of Knuth's sequence based on Ford sequences with linearly increasing alphabet sizes.

Within the current section, consider  $t : \mathbb{N} \mapsto \mathbb{N}$  to be an arbitrary but fixed function. Later, we study which conditions  $t$  must satisfy in order for the generated sequence to be completely equidistributed. Similar to the previous section, we define for any given natural number  $n$ :

i) a  $C$  sequence of order  $n$ , denoted  $C^{(n)}$ , as the finite sequence of rational numbers obtained from dividing by  $n$  each of the terms in an  $n$ -ary Ford sequence of order  $n$ :

$$C^{(n)} = \frac{f_1}{n}, \frac{f_2}{n}, \dots, \frac{f_{n^n}}{n} = \left( \frac{f_i}{n} \right)_{i=1}^{n^n}$$

where  $F^{(n,n)} = f_1, \dots, f_{n^n}$

and, ii) a  $D$  sequence of order  $n$ , denoted  $D^{(n)}$ , as  $t(n)$  consecutive copies of  $C^{(n)}$ :

$$D^{(n)} = \left\langle \underbrace{C^{(n)}; C^{(n)}; \dots; C^{(n)}}_{t(n) \text{ times}} \right\rangle.$$

Note again that the size of  $C^{(n)}$  is  $|C^{(n)}| = |F^{(n,n)}| = n^n$ , and the size of  $D^{(n)}$  is  $|D^{(n)}| = t(n)|C^{(n)}| = t(n)n^n$ . In this case, for any given  $n$  all terms in  $C^{(n)}$  and in  $D^{(n)}$  are numbers in the set  $\left\{0, \frac{1}{n}, \frac{2}{n}, \dots, \frac{n-1}{n}\right\} \subset [0, 1)$ .

The key difference between the way  $C$  sequences are constructed when compared to  $A$  sequences from the previous section is that, as the order of the sequence grows, the alphabet size for the underlying Ford sequence grows linearly  $(1, 2, 3, 4, \dots)$  rather than exponentially  $(2, 4, 8, 16, \dots)$ .

For example, when  $n = 3$  and  $t$  is equal to the identity function:

$$\begin{aligned} F^{(3,3)} &= 0, 0, 0, 1, 0, 0, 2, 0, 1, 1, 0, 1, 2, 0, 2, 1, 0, 2, 2, 1, 1, 1, 2, 1, 2, 2, 2 \\ C^{(3)} &= \frac{0}{3}, \frac{0}{3}, \frac{0}{3}, \frac{1}{3}, \frac{0}{3}, \frac{0}{3}, \frac{2}{3}, \frac{0}{3}, \frac{1}{3}, \frac{1}{3}, \frac{0}{3}, \frac{1}{3}, \frac{2}{3}, \frac{0}{3}, \frac{2}{3}, \frac{1}{3}, \frac{0}{3}, \frac{2}{3}, \frac{2}{3}, \frac{1}{3}, \frac{1}{3}, \frac{1}{3}, \frac{2}{3}, \frac{1}{3}, \frac{2}{3}, \frac{2}{3}, \frac{2}{3} \\ D^{(3)} &= \left\langle C^{(3)}; C^{(3)}; C^{(3)} \right\rangle = \underbrace{\frac{0}{3}, \frac{0}{3}, \dots, \frac{2}{3}, \frac{2}{3}}_{C^{(3)}} \underbrace{\frac{0}{3}, \frac{0}{3}, \dots, \frac{2}{3}, \frac{2}{3}}_{C^{(3)}} \underbrace{\frac{0}{3}, \frac{0}{3}, \dots, \frac{2}{3}, \frac{2}{3}}_{C^{(3)}} \end{aligned}$$

and  $|C^{(3)}| = 27$ ,  $|D^{(3)}| = 81$ .

We now define the sequence  $L$  as the infinite sequence of real numbers resulting from the concatenation of all possible  $D$  sequences in increasing order:

$$L = \langle D^{(1)}; D^{(2)}; D^{(3)}; \dots \rangle.$$

**Theorem 1.** *If  $t : \mathbb{N} \mapsto \mathbb{N}$  is a non-decreasing function and  $\lim_{n \rightarrow \infty} n/t(n) = 0$ , then the sequence  $L$  is completely equidistributed.*

**Example.** If  $t(n) = n^2$ , then:

$$L = \left\langle C^{(1)}; \underbrace{C^{(2)}; \dots; C^{(2)}}_{4 \text{ copies}}; \underbrace{C^{(3)}; \dots; C^{(3)}}_{9 \text{ copies}}; \dots \right\rangle,$$

and  $L$  is completely equidistributed.

### 2.2.1 Proof of Theorem 1

In order to present our proof of Theorem 1, we first establish some preliminary definitions.

Consider a prefix of  $L$  of length  $N$ , denoted  $L_{1:N}$ . It is always possible to find numbers  $p, q, r \in \mathbb{N}$  such that:

$$L_{1:N} = \left\langle D^{(1)}; \dots; D^{(r-1)}; \underbrace{C^{(r)}; \dots; C^{(r)}}_{q \text{ times}}; C_{1:p}^{(r)} \right\rangle$$

where  $0 \leq q < t(r)$  and  $1 \leq p \leq r^r$ . Here,  $r$  is the order of the rightmost, possibly incomplete  $D$  sequence present in  $L_{1:N}$ . The number  $q$  is the amount of complete  $C$  sequences of order  $r$  appearing before the rightmost, possibly incomplete  $C$  sequence, while  $p$  is the amount of terms present in said sequence. Note that the values of  $p$ ,  $q$  and  $r$  are uniquely determined by the value of  $N$ .

By considering the length of the sequence on each side of the previous equation, we obtain a functional relationship between  $N$ ,  $p$ ,  $q$ , and  $r$ :

$$\begin{aligned} N &= \sum_{s=1}^{r-1} |D^{(s)}| + q|C^{(r)}| + p \\ &= \sum_{s=1}^{r-1} t(s)s^s + qr^r + p. \end{aligned} \tag{2.1}$$

Let  $k$  be a positive integer and  $I = [u_1, v_1) \times \cdots \times [u_k, v_k)$  a set such that  $I \subseteq [0, 1)^k$ , where both  $k$  and  $I$  have arbitrary but fixed values. Let  $N$  range freely over the natural numbers, and the quantity  $\nu_N$  denote the number of windows of  $L$  of size  $k$  starting at indices  $i = 1 \dots N$  that belong to the set  $I$ :

$$\nu_N = \sum_{i=1}^N \sigma\left((W_k(L))_i \in I\right).$$

We can now express the probability, in the sense defined in Chapter 1, of any given window of  $L$  of size  $k$  belonging to the set  $I$  as:

$$Pr\left((W_k(L))_i \in I\right) = \lim_{N \rightarrow \infty} \frac{\nu_N}{N}.$$

Consider sufficiently large values of  $N$  such that  $k < r$ . This is always possible since  $r$  is an unbounded, non-decreasing function of  $N$ . We can decompose  $L_{1:N}$  into four consecutive sections; namely, sequences  $S^{(1)}$ ,  $S^{(2)}$ ,  $S^{(3)}$  and  $S^{(4)}$ :

$$\begin{aligned} L_{1:N} &= \left\langle S^{(1)}; S^{(2)}; S^{(3)}; S^{(4)} \right\rangle, \quad \text{where} \\ S^{(1)} &= \left\langle D^{(1)}; D^{(2)}; \dots; D^{(k-1)} \right\rangle \\ S^{(2)} &= \left\langle D^{(k)}; D^{(k+1)}; \dots; D^{(r-1)} \right\rangle \\ S^{(3)} &= \left\langle \underbrace{C^{(r)}; \dots; C^{(r)}}_{q \text{ times}} \right\rangle \\ S^{(4)} &= C_{1:p}^{(r)}. \end{aligned}$$

Note that  $S^{(1)}$  and  $S^{(3)}$  can potentially be empty, such as when  $k = 1$  or  $q = 0$ , respectively.

We denote the cumulative sums of the sizes of the sequences defined above as  $n_0 = 0$ , and  $n_j = n_{j-1} + |S^{(j)}|$  for  $j = 1, 2, 3, 4$ . Now, we can similarly decompose  $\nu_N$  into five parts:

$$\begin{aligned} \nu_N &= \nu_N^{(1)} + \nu_N^{(2)} + \nu_N^{(3)} + \nu_N^{(4)} + \varepsilon_b, \quad \text{where} \\ \nu_N^{(j)} &= \sum_{i=1+n_{j-1}}^{n_j-k+1} \sigma\left((W_k(L))_i \in I\right) \quad j = 1, 2, 3, 4 \\ &= \sum_{i=1}^{|S^{(j)}|-k+1} \sigma\left((W_k(S^{(j)}))_i \in I\right) \end{aligned} \tag{2.2}$$

for some  $\varepsilon_b \leq 3(k-1)$ .

For each  $j = 1, 2, 3, 4$ , the quantity  $\nu_N^{(j)}$  accounts for windows contained entirely within the sequence  $S^{(j)}$ , and  $\varepsilon_b$  accounts for all windows crossing any of the three borders between the four sections. This is enough to account for all possible windows, since any given window is either entirely contained in some section, or it starts at a given section and ends at a subsequent one, thereby crossing a border.

Before obtaining more precise expressions for these quantities, we first state the following three technical propositions.

**Proposition 2.** *If  $n \in \mathbb{N}$  and  $x, y \in \mathbb{R}$  such that  $[x, y) \subseteq [0, n)$ , then the number of integers from the set  $\{0, 1, \dots, n-1\}$  contained in  $[x, y)$  is equal to  $y - x + \varepsilon$  for some  $\varepsilon \in (-1, 1)$ .*

*Proof.* Since  $0 \leq y$ , there are exactly  $\lceil y \rceil = y + \varepsilon_y$  non-negative integers in the set  $[0, y)$  for some  $\varepsilon_y \in [0, 1)$ . Similarly for  $x$ , there are exactly  $\lceil x \rceil = x + \varepsilon_x$  non-negative integers in the set  $[0, x)$  for some  $\varepsilon_x \in [0, 1)$ . The difference between these two quantities is equal to the number of non-negative integers contained in the set  $[x, y)$ , which is  $y - x + (\varepsilon_y - \varepsilon_x)$ . Observing that  $(\varepsilon_y - \varepsilon_x) \in (-1, 1)$ , and that all non-negative integers between  $x$  and  $y$  belong to the set  $\{0, 1, \dots, n-1\}$ , the proof is complete.  $\square$

**Proposition 3.** *If  $k$  is a positive integer and  $a_1, a_2, \dots, a_k, b_1, b_2, \dots, b_k$  are sequences of real numbers of length  $k$ , then the product of their element-by-element sums can be expanded in the following way:*

$$\prod_{d=1}^k a_d + b_d = \prod_{d=1}^k a_d + \sum_{j=1}^{2^k-1} \left[ \prod_{d=1}^k \begin{cases} a_d & \left\lfloor \frac{j}{2^{d-1}} \right\rfloor \text{ is even} \\ b_d & \text{otherwise} \end{cases} \right].$$

*Proof.* By induction on  $k$ . First, note that the property holds for  $k = 1$ :

$$\begin{aligned} \prod_{d=1}^1 a_d + b_d &= a_1 + b_1, \text{ and} \\ \prod_{d=1}^1 a_d + \sum_{j=1}^1 \left[ \prod_{d=1}^1 \begin{cases} a_d & \left\lfloor \frac{j}{2^{d-1}} \right\rfloor \text{ is even} \\ b_d & \text{otherwise} \end{cases} \right] &= a_1 + b_1. \end{aligned}$$

Next, we see that the inductive step holds for any  $k$ . First,

$$\begin{aligned}
\prod_{d=1}^{k+1} a_d + b_d &= (a_{k+1} + b_{k+1}) \prod_{d=1}^k a_d + b_d, \quad \text{and by I. H.} \\
&= (a_{k+1} + b_{k+1}) \left[ \prod_{d=1}^k a_d + \sum_{j=1}^{2^k-1} \left[ \prod_{d=1}^k \left\{ \begin{array}{cc} a_d & \lfloor \frac{j}{2^{d-1}} \rfloor \text{ is even} \\ b_d & \text{otherwise} \end{array} \right\} \right] \right] \\
&= \prod_{d=1}^{k+1} a_d + a_{k+1} \sum_{j=1}^{2^k-1} \left[ \prod_{d=1}^k \left\{ \begin{array}{cc} a_d & \lfloor \frac{j}{2^{d-1}} \rfloor \text{ is even} \\ b_d & \text{otherwise} \end{array} \right\} \right] \\
&\quad + b_{k+1} \prod_{d=1}^k a_d + b_{k+1} \sum_{j=1}^{2^k-1} \left[ \prod_{d=1}^k \left\{ \begin{array}{cc} a_d & \lfloor \frac{j}{2^{d-1}} \rfloor \text{ is even} \\ b_d & \text{otherwise} \end{array} \right\} \right].
\end{aligned}$$

Since for every  $j = 1 \dots 2^k - 1$  the value  $\lfloor \frac{j}{2^k} \rfloor = 0$  and is therefore even, we can add the factor  $a_{k+1}$  to the product in the second term simply by raising the upper limit to  $k+1$ . Similarly, in the fourth term we can add the factor  $b_{k+1}$  to the product by raising the upper limit to  $k+1$  and changing the limits in the sum to  $j = (2^k + 1) \dots (2^k + 2^k - 1)$ . This is true because adding  $2^k$  to  $j$  does not change the value of  $\lfloor \frac{j}{2^{d-1}} \rfloor$  for any  $d \leq k$ , but when  $d = k+1$  the value  $\lfloor \frac{j}{2^k} \rfloor = 1$  and is therefore odd. The third term can be rewritten as a similar product for a value of  $j = 2^k$ , and substituting into the equation above:

$$\begin{aligned}
\prod_{d=1}^{k+1} a_d + b_d &= \prod_{d=1}^{k+1} a_d + \sum_{j=1}^{2^k-1} \left[ \prod_{d=1}^{k+1} \left\{ \begin{array}{cc} a_d & \lfloor \frac{j}{2^{d-1}} \rfloor \text{ is even} \\ b_d & \text{otherwise} \end{array} \right\} \right] \\
&\quad + \sum_{j=2^k}^{2^k-1} \left[ \prod_{d=1}^{k+1} \left\{ \begin{array}{cc} a_d & \lfloor \frac{j}{2^{d-1}} \rfloor \text{ is even} \\ b_d & \text{otherwise} \end{array} \right\} \right] \\
&\quad + \sum_{j=2^k+1}^{2^k+2^k-1} \left[ \prod_{d=1}^k \left\{ \begin{array}{cc} a_d & \lfloor \frac{j}{2^{d-1}} \rfloor \text{ is even} \\ b_d & \text{otherwise} \end{array} \right\} \right] \\
&= \prod_{d=1}^{k+1} a_d + \sum_{j=1}^{2^{k+1}-1} \left[ \prod_{d=1}^{k+1} \left\{ \begin{array}{cc} a_d & \lfloor \frac{j}{2^{d-1}} \rfloor \text{ is even} \\ b_d & \text{otherwise} \end{array} \right\} \right],
\end{aligned}$$

which completes the proof. □

**Proposition 4.** *Given  $n \in \mathbb{N}$ , the following holds:*

$$\sum_{i=1}^n i^{i-1} \leq 2n^{n-1}.$$

*Proof.* By induction on  $n$ . The property holds for  $n = 1$  and  $n = 2$ :



$$\sum_{i=1}^1 i^{i-1} \leq 2, \quad \sum_{i=1}^2 i^{i-1} \leq 4$$

and the inductive step holds for  $n \geq 2$ :

$$\begin{aligned} \sum_{i=1}^{n+1} i^{i-1} &= \underbrace{\sum_{i=1}^n i^{i-1}}_{\leq 2n^{n-1}} + (n+1)^n \leq nn^{n-1} + (n+1)^n \leq 2(n+1)^n. \\ &\text{by I. H.} \end{aligned}$$

Therefore, the property holds for all  $n \in \mathbb{N}$ .  $\square$

We now obtain an expression for the number of windows of a  $C$  sequence which are contained in the set  $I$ . This is useful for evaluating  $\nu_N$ , as seen later on.

**Lemma 5.** *Given a positive integer  $k$  and a set  $I = [u_1, v_1) \times \cdots \times [u_k, v_k)$  where  $I \subseteq [0, 1)^k$ , let  $n \in \mathbb{N}$  such that  $k \leq n$  and consider the sequence  $C^{(n)}$  as a cyclic sequence. Then, for some  $\varepsilon \in (-1, 1)$ :*

$$\sum_{i=1}^{n^n} \sigma\left((W_k^c(C^{(n)}))_i \in I\right) = n^n |I| + n^{n-1}(2^k - 1)\varepsilon.$$

*Proof.* The expression on the left-hand side counts the number of windows of size  $k$  in  $C^{(n)}$  that are contained in  $I$ . First, note that any given window is contained in the set  $I$  if and only if the following is true:

$$\begin{aligned} (W_k^c(C^{(n)}))_i \in I &\iff u_1 \leq C_i^{(n)} < v_1 \\ &\quad \vdots \\ u_k &\leq C_{i+k-1}^{(n)} < v_k \end{aligned}$$

where  $i = 1 \dots n^n$  and indices are taken modulo  $n^n$ .

Since all terms in  $C^{(n)}$  are numbers in the set  $\{0, \frac{1}{n}, \dots, \frac{n-1}{n}\}$ , we multiply both sides of each inequality by  $n$ , allowing us to reason about integers belonging to a Ford sequence instead of rational numbers. We obtain the following:

$$\begin{aligned} (W_k^c(C^{(n)}))_i \in I &\iff nu_1 \leq F_i^{(n,n)} < nv_1 \\ &\quad \vdots \\ nu_k &\leq F_{i+k-1}^{(n,n)} < nv_k. \end{aligned}$$

As per Proposition 2, for each inequality above with  $d = 1 \dots k$  there are exactly  $nv_d - nu_d + \varepsilon_d$  possible solutions in the set  $\{0, 1, \dots, n-1\}$  for some value  $\varepsilon_d \in (-1, 1)$ . This yields a total of  $\prod_{d=1}^k [n(v_d - u_d) + \varepsilon_d]$  possible solutions to the system of inequalities. Each solution, when seen as an  $n$ -ary sequence of length  $k$ , appears exactly  $n^{n-k}$  times in  $F^{(n,n)}$ . This is true because there are  $n^{n-k}$  ways of extending an  $n$ -ary sequence of length  $k$  to one of length  $n$  and, by construction, each of these appears exactly once in  $F^{(n,n)}$  when viewed as a cycle. Since  $i$  ranges exactly once over each possible window of  $F^{(n,n)}$ , then:

$$\begin{aligned} \sum_{i=1}^{n^n} \sigma\left((W_k^c(C^{(n)}))_i \in I\right) &= n^{n-k} \prod_{d=1}^k [n(v_d - u_d) + \varepsilon_d] \\ &= n^n \prod_{d=1}^k [(v_d - u_d) + \varepsilon_d/n]. \end{aligned} \quad (2.3)$$

Using Proposition 3, we can expand this into the following:

$$\begin{aligned} n^n \prod_{d=1}^k [(v_d - u_d) + \varepsilon_d/n] &= n^n \prod_{d=1}^k (v_d - u_d) \\ &\quad + n^n \sum_{j=1}^{2^k-1} \left[ \prod_{d=1}^k \left\{ \begin{array}{ll} (v_d - u_d) & \lfloor \frac{j}{2^{d-1}} \rfloor \text{ is even} \\ \varepsilon_d/n & \text{otherwise} \end{array} \right\} \right]. \end{aligned} \quad (2.4)$$

If we define  $\varepsilon'_j$  for  $j = 1 \dots 2^k - 1$  as:

$$\varepsilon'_j/n = \prod_{d=1}^k \left\{ \begin{array}{ll} (v_d - u_d) & \lfloor \frac{j}{2^{d-1}} \rfloor \text{ is even} \\ \varepsilon_d/n & \text{otherwise} \end{array} \right\}$$

then, for each  $j$ , the value  $\varepsilon'_j \in (-1, 1)$ . This is true because the product on the right-hand side is composed of terms  $(v_d - u_d) \in (-1, 1)$  and  $\varepsilon_d/n \in (-1/n, 1/n)$  and, since  $j > 0$ , there is always at least one term of the second kind. Given that  $|I| = \prod_{d=1}^k (v_d - u_d)$ , we can further simplify equation 2.4 to get:

$$n^n \prod_{d=1}^k [(v_d - u_d) + \varepsilon_d/n] = n^n |I| + n^n \sum_{j=1}^{2^k-1} \varepsilon'_j/n. \quad (2.5)$$

Finally, since  $-(2^k - 1) < \sum_{j=1}^{2^k-1} \varepsilon'_j < (2^k - 1)$ , there exists some  $\varepsilon \in (-1, 1)$  such that:

$$\sum_{j=1}^{2^k-1} \varepsilon'_j = (2^k - 1)\varepsilon$$

and hence, by 2.3 and 2.5:

$$\sum_{i=1}^{n^n} \sigma\left((W_k^c(C^{(n)}))_i \in I\right) = n^n |I| + n^{n-1}(2^k - 1)\varepsilon.$$

□

*Proof of Theorem 1.* Using Lemma 5, we now prove the main result of this work.

Remember that  $k$  is a positive integer and  $I = [u_1, v_1) \times \cdots \times [u_k, v_k)$  is a set such that  $I \subseteq [0, 1)^k$ , where both  $k$  and  $I$  have arbitrary but fixed values. Let  $N$  range freely over the natural numbers. Next, we obtain an expression for  $\nu_N/N$  and compute its limit when  $N \rightarrow \infty$ . Recall the following definitions:

$$\begin{aligned} \nu_N^{(2)} &= \sum_{i=1}^{|S^{(2)}|-k+1} \sigma\left((W_k(S^{(2)}))_i \in I\right) \\ \nu_N^{(3)} &= \sum_{i=1}^{|S^{(3)}|-k+1} \sigma\left((W_k(S^{(3)}))_i \in I\right) \\ S^{(2)} &= \langle D^{(k)}; D^{(k+1)}; \dots; D^{(r-1)} \rangle \\ S^{(3)} &= \left\langle \underbrace{C^{(r)}; \dots; C^{(r)}}_{q \text{ times}} \right\rangle. \end{aligned}$$

Note that the sequences  $S^{(2)}$  and  $S^{(3)}$  are entirely composed of complete  $C$  sequences of increasing orders which are larger than or equal to  $k$ . Moreover, with the exception of the last, rightmost instance in each of  $S^{(2)}$  and  $S^{(3)}$ , every single  $C$  sequence is immediately succeeded by another  $C$  sequence of the same or the following order, including those which are part of a  $D$  sequence. Additionally, any window starting at the right-hand end of a  $C$  sequence necessarily finishes within the first  $k - 1$  elements of the following  $C$  sequence, all of which are guaranteed to be 0.

Therefore, the amount of windows of size  $k$  contained in  $I$  ranging over  $S^{(2)}$  and  $S^{(3)}$  is equal to the sum over each composing  $C$  sequence *viewed as a cycle*, with an error of at most  $k - 1$  due to the fact that we are counting only windows entirely contained within each sequence:

$$\begin{aligned} \nu_N^{(2)} &= \sum_{s=k}^{r-1} \underbrace{\left[ t(s) \sum_{i=1}^{s^s} \sigma\left((W_k^c(C^{(s)}))_i \in I\right) \right]}_{C \text{ sequences contained in } D^{(s)}} + \varepsilon_{\nu_N^{(2)}} \\ \nu_N^{(3)} &= q \sum_{i=1}^{r^r} \sigma\left((W_k^c(C^{(r)}))_i \in I\right) + \varepsilon_{\nu_N^{(3)}} \end{aligned}$$

for some values  $\varepsilon_{\nu_N^{(2)}} \leq k - 1$ , and  $\varepsilon_{\nu_N^{(3)}} \leq k - 1$ .

From Lemma 5:

$$\begin{aligned}\nu_N^{(2)} &= \sum_{s=k}^{r-1} \left[ t(s) \left( s^s |I| + s^{s-1} (2^k - 1) \varepsilon_s \right) \right] + \varepsilon_{\nu_N^{(2)}} \\ \nu_N^{(3)} &= q \left( r^r |I| + r^{r-1} (2^k - 1) \varepsilon_r \right) + \varepsilon_{\nu_N^{(3)}}\end{aligned}$$

for some values of  $\varepsilon_i \in (-1, 1)$ ,  $i = k \dots r$ .

Substituting back into  $\nu_N$  from equation 2.2:

$$\begin{aligned}\nu_N &= \nu_N^{(1)} \\ &+ \sum_{s=k}^{r-1} \left[ t(s) \left( s^s |I| + s^{s-1} (2^k - 1) \varepsilon_s \right) \right] + \varepsilon_{\nu_N^{(2)}} \\ &+ q \left( r^r |I| + r^{r-1} (2^k - 1) \varepsilon_r \right) + \varepsilon_{\nu_N^{(3)}} \\ &+ \nu_N^{(4)} + \varepsilon_b\end{aligned}$$

and factoring out terms multiplied by  $|I|$ , we get:

$$\begin{aligned}\nu_N &= |I| \left[ \sum_{s=k}^{r-1} t(s) s^s + q r^r \right] \\ &+ \nu_N^{(1)} \\ &+ \sum_{s=k}^{r-1} \left[ t(s) s^{s-1} (2^k - 1) \varepsilon_s \right] + \varepsilon_{\nu_N^{(2)}} \\ &+ q r^{r-1} (2^k - 1) \varepsilon_r + \varepsilon_{\nu_N^{(3)}} \\ &+ \nu_N^{(4)} + \varepsilon_b.\end{aligned}$$

We now rewrite the first term using the relationship between  $p$ ,  $r$ ,  $q$ , and  $N$  from equation 2.1:

$$\begin{aligned}
\nu_N &= |I| \left[ N - \sum_{s=1}^{k-1} t(s) s^s - p \right] \\
&\quad + \nu_N^{(1)} \\
&\quad + \sum_{s=k}^{r-1} \left[ t(s) s^{s-1} (2^k - 1) \varepsilon_s \right] + \varepsilon_{\nu_N^{(2)}} \\
&\quad + q r^{r-1} (2^k - 1) \varepsilon_r + \varepsilon_{\nu_N^{(3)}} \\
&\quad + \nu_N^{(4)} + \varepsilon_b
\end{aligned}$$

and after dividing both sides by  $N$  and rearranging terms we obtain:

$$\begin{aligned}
\frac{\nu_N}{N} - |I| &= \frac{p}{N} \left[ \frac{\nu_N^{(4)}}{p} - |I| \right] \\
&\quad + \frac{2^k - 1}{N} \left[ \sum_{s=k}^{r-1} t(s) s^{s-1} \varepsilon_s + q r^{r-1} \varepsilon_r \right] \\
&\quad + \frac{1}{N} \left[ \nu_N^{(1)} - |I| \sum_{s=1}^{k-1} t(s) s^s + \varepsilon_{\nu_N^{(2)}} + \varepsilon_{\nu_N^{(3)}} + \varepsilon_b \right].
\end{aligned}$$

Taking limits on both sides as  $N \rightarrow \infty$ , the third term on the right-hand side approaches 0 since the contents of the brackets are dependent on  $k$  and bounded as a function of  $N$ . In regard to the second term, using the fact that  $t$  is non-decreasing together with Proposition 4 we can see that:

$$\begin{aligned}
\frac{\sum_{s=k}^{r-1} t(s) s^{s-1} \varepsilon_s}{N} &\leq \frac{t(r-1) \sum_{s=1}^{r-1} s^{s-1}}{t(r-1)(r-1)^{(r-1)}} \leq \frac{2(r-1)^{(r-2)}}{(r-1)^{(r-1)}} = \frac{2}{r-1}, \text{ and} \\
\frac{q r^{r-1} \varepsilon_r}{N} &\leq \frac{q r^{r-1}}{q r^r} = \frac{1}{r},
\end{aligned}$$

and since  $r$  is an unbounded, non-decreasing function of  $N$ , this term approaches 0 as well.

Finally, consider the first term on the right-hand side and note that  $\left[ \frac{\nu_N^{(4)}}{p} - |I| \right] \in [-1, 1]$ , since both  $\frac{\nu_N^{(4)}}{p}, |I| \in [0, 1]$ . Moreover, since  $p \leq r^r$  and using the identity:

$$\frac{(x+1)^{(x+1)}}{x^x} = (x+1) \left( 1 + \frac{1}{x} \right)^x$$

for  $x = r - 1$ , we can see that for large values of  $r$ :

$$\frac{p}{N} \leq \frac{r^r}{t(r-1)(r-1)^{(r-1)}} = \frac{r}{t(r-1)} \left(1 + \frac{1}{r-1}\right)^{r-1} \leq \frac{r}{t(r-1)} e$$

which by hypothesis also approaches 0 as  $N \rightarrow \infty$ . Hence,

$$\lim_{N \rightarrow \infty} \frac{\nu_N}{N} = |I|$$

and, since  $k$  and  $I$  were chosen arbitrarily, the sequence  $L$  is completely equidistributed and the proof of Theorem 1 is complete. □

### 2.2.2 Alternative Proof of Theorem 1

We provide an alternative proof of Theorem 1 based on Weyl's Criterion, as stated in Section 1.4, using a proposition from the theory of linear modular congruences and a well-known fact from the study of the roots of unity.

Before applying Weyl's Criterion to the sequence  $L$ , we first prove the following lemma.

**Lemma 6.** *Given a positive integer  $k$  and a non-zero  $k$ -dimensional vector of integers  $\bar{l} = (l_1, \dots, l_k)$ , let  $n \in \mathbb{N}$  such that  $n > \max(k, \min(|l_1|, \dots, |l_k|))$  and consider the sequence  $C^{(n)}$  as a cyclic sequence. Then:*

$$\sum_{j=1}^{n^n} e^{2\pi i \bar{l} \cdot \bar{w}_j} = 0, \tag{2.6}$$

where  $W_k^c(C^{(n)}) = \bar{w}_1, \bar{w}_2, \dots, \bar{w}_{n^n}$ .

*Proof.* If we let:

$$W_k^c(F^{(n,n)}) = \bar{f}_1, \bar{f}_2, \dots, \bar{f}_{n^n}$$

then from the definition of a  $C$  sequence, it follows that  $\bar{w}_j = \frac{1}{n} \bar{f}_j$  for  $j = 1 \dots n^n$ .

Then, if we define  $\Gamma = \{0, 1, \dots, n-1\}$ , every  $\bar{\gamma} \in \Gamma^k$  appears exactly  $n^{n-k}$  times in the sequence  $W_k^c(F^{(n,n)})$ . This is true because there are  $n^{n-k}$  ways of extending an  $n$ -ary sequence of length  $k$  to one of length  $n$  and, by construction, each of these appears exactly once in  $F^{(n,n)}$  when viewed as a cycle.

Substituting into the left-hand side of equation 2.6:

$$\sum_{j=1}^{n^n} e^{2\pi i \bar{\ell} \cdot (\frac{1}{n} \bar{f}_j)} = \sum_{j=1}^{n^n} e^{\frac{2\pi i}{n} \bar{\ell} \cdot \bar{f}_j} = n^{n-k} \sum_{\bar{\gamma} \in \Gamma^k} e^{\frac{2\pi i}{n} \bar{\ell} \cdot \bar{\gamma}}.$$

Since  $\bar{\ell} \cdot \bar{\gamma} \in \mathbb{Z}$  and the function  $\exp : \mathbb{Z} \rightarrow \mathbb{C}$ ,  $\exp(m) = e^{\frac{2\pi i}{n} m}$  is periodic with a period equal to  $n$ , then:

$$\sum_{\bar{\gamma} \in \Gamma^k} e^{\frac{2\pi i}{n} \bar{\ell} \cdot \bar{\gamma}} = \sum_{r=0}^{n-1} \sum_{\substack{\bar{\gamma} \in \Gamma^k \\ \bar{\ell} \cdot \bar{\gamma} \equiv r}} e^{\frac{2\pi i}{n} r}$$

where the congruence  $\bar{\ell} \cdot \bar{\gamma} \equiv r$  is taken modulo  $n$ .

The conditions for the existence of solutions to equations of the form  $\bar{\ell} \cdot \bar{\gamma} \equiv r \pmod{n}$ ,  $\bar{\gamma} \in \Gamma^k$ , are well understood (see [11], page 114). In particular, this equation only has solutions when  $\gcd(l_1, \dots, l_k, n) = g$  divides  $r$  and, in such case, the total number of solutions is equal to  $gn^{k-1}$ . Therefore:

$$\begin{aligned} \sum_{r=0}^{n-1} \sum_{\substack{\bar{\gamma} \in \Gamma^k \\ \bar{\ell} \cdot \bar{\gamma} \equiv r}} e^{\frac{2\pi i}{n} r} &= \sum_{\substack{r=0 \\ g|r}}^{n-1} gn^{k-1} e^{\frac{2\pi i}{n} r} \\ &= gn^{k-1} \sum_{r'=0}^{\lfloor \frac{n-1}{g} \rfloor} e^{\frac{2\pi i}{n} gr'} \end{aligned}$$

where the last step comes from substituting  $r = gr'$ . If we also substitute  $n = gn'$  and observe that  $\lfloor \frac{n-1}{g} \rfloor = n' - 1$ :

$$\sum_{r'=0}^{\lfloor \frac{n-1}{g} \rfloor} e^{\frac{2\pi i}{n} gr'} = \sum_{r'=0}^{n'-1} e^{\frac{2\pi i}{n'} r'}.$$

The term on the right-hand side is the sum of all the roots of unity of order  $n'$ . It is a well-known fact that this sum is equal to 0 whenever  $n' > 1$ . Finally, since  $n > \min(|l_1|, \dots, |l_k|) \geq g$ , then  $n' = n/g > 1$ , and the proof is complete.

□

*Proof of Theorem 1.* Remember that  $k$  is a positive integer and  $\bar{\ell} = (l_1, \dots, l_k)$  is a non-zero  $k$ -dimensional vector of integers, where both  $k$  and  $\bar{\ell}$  have arbitrary but fixed values.

Let  $N$  range freely over the natural numbers. Similarly to Section 2.2.1, we consider a prefix of the sequence  $L$  of length  $N$ , denoted  $L_{1:N}$ , and we define the values  $p$ ,  $q$ , and  $r$  equivalently.

Letting  $W_k(L) = \bar{w}_1, \bar{w}_2, \dots$ , we analogously define the complex value  $\nu_N$  as the Weyl sum over the first  $N$  windows of size  $k$  of  $L$ :

$$\nu_N = \sum_{j=1}^N e^{2\pi i \bar{\ell} \cdot \bar{w}_j}.$$

Let  $m = \max(k, \min(|l_1|, \dots, |l_k|))$ , and consider sufficiently large values of  $N$  such that  $m < r$ . Like before, this is always possible since  $r$  is an unbounded, non-decreasing function of  $N$ . We can decompose  $L_{1:N}$  into four consecutive sections; namely, sequences  $S^{(1)}$ ,  $S^{(2)}$ ,  $S^{(3)}$  and  $S^{(4)}$ :

$$\begin{aligned} L_{1:N} &= \langle S^{(1)}; S^{(2)}; S^{(3)}; S^{(4)} \rangle, \quad \text{where} \\ S^{(1)} &= \langle D^{(1)}; D^{(2)}; \dots; D^{(m-1)} \rangle \\ S^{(2)} &= \langle D^{(m)}; D^{(m+1)}; \dots; D^{(r-1)} \rangle \\ S^{(3)} &= \left\langle \underbrace{C^{(r)}; \dots; C^{(r)}}_{q \text{ times}} \right\rangle \\ S^{(4)} &= C_{1:p}^{(r)} \end{aligned}$$

and define  $\nu_N^{(1)}$ ,  $\nu_N^{(2)}$ ,  $\nu_N^{(3)}$  and  $\nu_N^{(4)}$  as the Weyl sums over windows entirely contained within each respective section, such that:

$$\nu_N = \nu_N^{(1)} + \nu_N^{(2)} + \nu_N^{(3)} + \nu_N^{(4)} + \varepsilon_b$$

for some complex number  $\varepsilon_b$  with  $|\varepsilon_b| \leq 3(k-1)$  which accounts for all windows crossing over any border.

Following the same reasoning from Section 2.2.1, the values of  $\nu_N^{(2)}$  and  $\nu_N^{(3)}$  can be computed as the Weyl sums over the  $C$  sequences composing  $S^{(2)}$  and  $S^{(3)}$  viewed as cycles, plus two error terms accounting for right borders. However, in this case, due to Lemma 6 and the fact that all  $C$  sequences in  $S^{(2)}$  and  $S^{(3)}$  have orders greater than  $m$ , these sums vanish to zero. Therefore,  $\nu_N = \nu_N^{(1)} + \varepsilon_{\nu_N^{(2)}} + \varepsilon_{\nu_N^{(3)}} + \nu_N^{(4)} + \varepsilon_b$ , for some complex values  $\varepsilon_{\nu_N^{(2)}}, \varepsilon_{\nu_N^{(3)}}$  with  $|\varepsilon_{\nu_N^{(2)}}|, |\varepsilon_{\nu_N^{(3)}}| \leq k-1$ .

We now consider the limit of  $\nu_N/N$  as  $N \rightarrow \infty$ . Note that:



$$\begin{aligned}
\left| \frac{\nu_N}{N} \right| &\leq \frac{1}{N} \left| \nu_N^{(1)} + \varepsilon_{\nu_N^{(2)}} + \varepsilon_{\nu_N^{(3)}} + \varepsilon_b \right| + \frac{1}{N} \left| \nu_N^{(4)} \right| \\
&\leq \frac{1}{N} \left[ \left| \nu_N^{(1)} \right| + |\varepsilon_{\nu_N^{(2)}}| + |\varepsilon_{\nu_N^{(3)}}| + |\varepsilon_b| \right] + \frac{p}{N} \\
&\leq \frac{1}{N} \left[ \left| \nu_N^{(1)} \right| + (k-1) + (k-1) + 3(k-1) \right] + \frac{p}{N}.
\end{aligned}$$

The numerator in the first term is dependent on  $k$  and  $m$ , and constant as a function of  $N$ . As shown before, the second term approaches 0 as  $r \rightarrow \infty$ . Therefore, the sum of both terms approaches 0 as  $N \rightarrow \infty$ , which in turn implies that  $\nu_N/N$  vanishes as well.

Given that  $k$  and  $\bar{\ell}$  were chosen arbitrarily, Weyl's Criterion is satisfied for all values of  $k$  and the sequence  $L$  is completely equidistributed.

□

## BIBLIOGRAPHY

- [1] J. Berstel and D. Perrin. The origins of combinatorics on words. *European Journal of Combinatorics*, 28(3):996 – 1022, 2007.
- [2] N. G. de Bruijn. A combinatorial problem. *Proc. Koninklijke Nederlandse Academie van Wetenschappen*, 49:758–764, 1946.
- [3] L. R. Ford. A cyclic arrangement of  $m$ -tuples. *Journal of Combinatorial Theory*, (Report P-1071), 1957.
- [4] J. N. Franklin. Deterministic simulation of random processes. *Mathematics of Computation*, 17(81):28–59, 1963.
- [5] H. Fredricksen. A survey of full length nonlinear shift register cycle algorithms, 1982.
- [6] H. Fredricksen and J. Maiorana. Necklaces of beads in  $k$  colors and  $k$ -ary de bruijn sequences. *Discrete Mathematics*, 23(3):207 – 210, 1978.
- [7] D. E. Knuth. Construction of a random sequence. *BIT Numerical Mathematics*, 5(4):246–250, 1965.
- [8] Pierre L’Ecuyer and Richard Simard. TestU01: A C library for empirical testing of random number generators. *ACM Trans. Math. Softw.*, 33(4):22:1–22:40, 2007.
- [9] M. H. Martin. A problem in arrangements. *Bull. Amer. Math. Soc.*, 40(12):859–864, 1934.
- [10] M. Matsumoto and T. Nishimura. Mersenne Twister: A 623-dimensionally equidistributed uniform pseudo-random number generator. *ACM Trans. Model. Comput. Simul.*, 8(1):3–30, 1998.
- [11] P.J. McCarthy. *Introduction to Arithmetical Functions*. Universitext Series. Springer-Verlag, 1986.
- [12] F. Ruskey, C. Savage, and T. M. Y. Wang. Generating necklaces. *Journal of Algorithms*, 13(3):414 – 430, 1992.
- [13] R. Tijdeman. Review: L. Kuipers and H. Niederreiter. Uniform distribution of sequences. *Bull. Amer. Math. Soc.*, 81(4):672–675, 1975.
- [14] H. Weyl. Über die Gleichverteilung von Zahlen mod. Eins. *Mathematische Annalen*, 77(3):313–352, 1916.