

Enumerating De Bruijn Sequences

Vladimir Raphael Rosenfeld

Institute of Evolution, University of Haifa, Mount Carmel, Haifa 31905, Israel

E-mail: vladimir@research.haifa.ac.il

Abstract

A cycle is a sequence $a_1a_2 \cdots a_r$ taken in a circular order—that is, a_1 follows a_r , and $a_2 \cdots a_r a_1, \dots, a_r a_1 \cdots a_{r-1}$ are all the same cycle as $a_1a_2 \cdots a_r$. Given natural numbers $q \geq 1$ and $s \geq 2$, a cycle of s^q letters is called a *complete cycle* [1, 2], or *De Bruijn sequence*, if subsequences $a_i a_{i+1} \cdots a_{i+q-1}$ ($1 \leq i \leq s^q$) consist of all possible s^q ordered sequences $b_1 b_2 \cdots b_q$ over the alphabet A ($|A| = s$). In 1946, De Bruijn proved [1] (see [2]) that the number of complete cycles, under $s = 2$, is equal to $2^{2^{q-1}-q}$.

We propose the overall proof for $s \geq 2$, which determines the number of the De Bruijn sequences to be equal to $(s!)^{s^{q-1}}/s^q$. The demonstration is based on our recent results concerning the characteristic polynomial and permanent of the arc-graph [17], applied herein to some auxiliary digraphs.

Wherever possible, the main subject is discussed in the wider context of related combinatorial problems, which first include counting the *linear De Bruijn sequences*.

Obtained results can be used for calculating the number of monocyclic and linear compounds, formed from s sorts of atoms, obeying the specified combinatorial restrictions. The former is equivalent to finding the number of respective *necklaces with s kinds of beads*.

1 Introduction

This paper was provoked by the study of complex sequences being carried on by the research group under the supervision of Profs. Edward Trifonov and Alexander Bolshoy, in the Genome Diversity Center of the University of Haifa (see [24–29]). In particular, Dr. Valery Kirzhner defined a minimal generating sequence in DNA as the sequence of minimal length that produces all possible amino acids; thus, it should contain all triplets of nucleotides, taking into account the table of identity of some triplets. Such a minimal

sequence is, in some sense, the most complex; and the mathematical formalization of it leads to De Bruijn sequences.

As the first stage of work, our main objective herein is counting the generalized De Bruijn sequences, or else minimal generating sequences containing all words of a given length over a given alphabet, which disregard the equivalency of some triplets, specified by their identity table. Additionally to this topic, we supplement a brief discussion of its general background [1-12] that touches similar mathematical problems and their possible applications.

More specifically, one can consider all s^q possible words, or s -ary sequences, consisting of q characters over the alphabet A ($|A| = s$). Further, it can be proposed to construct the shortest (of length s^q) circular sequence that contains exactly once every possible q -character sequence, as its subword. In this closed cycle, every two adjacent q -subwords have exactly $q-1$ letters in common while each q -character subword has exactly 2 adjacent neighbors. A circular sequence possessing these properties is called a complete cycle (see [1, 2]), or De Bruijn sequence, due to N. G. De Bruijn who published the first paper on the subject [1]. (Note that his surname, in the literature, is also spelled by himself and other authors as de Bruijn and DeBruijn.)

Herein, we can already utilize, as a ready fact (see [1-3]), that such sequences do exist; our task remains to calculate their number for any given pair q and s of positive integers ($q \geq 1; s \geq 2$). De Bruijn proved, in 1946, that the number of complete cycles, under $s=2$, is $2^{2^{q-1}-q}$. The graph-theoretical ideas of his proof (see [1,2]) hold good for the general case as well. In a few words, the proof includes the construction of some auxiliary digraphs (also called De Bruijn graphs [3]) and the subsequent count of all Eulerian circuits in these graphs. The latter task can easily be performed by merging well-known methods [13-16] with our latest results on the spectrum of the arc-graph [17]. The targeted quantity equals $(s!)^{s^{q-1}}/s^q$ and resembles the partial De Bruijn result [1], in notation.

To return to the above genetical problem, one should realize that the mathematical biologist needs, first of all, to count linear De Bruijn sequences, which are the shortest unclosed sequences consisting of the same sets of q -character s -ary subwords. Because a complete cycle is circularly asymmetric, by definition, cutting it in any of s^q possible positions results in s^q distinct s^q -words, if one reads them only from the beginning. However, these linear words must be curtailed because one can find only $s^q - q + 1$ original q -subwords of the respective complete cycle, in any of them. But this situation can be corrected by adding the first $q-1$ characters of an obtained unclosed word to its end, which results in a longer $(s^q + q - 1)$ -character word that already contains every q -character subword of the respective complete cycle exactly once. The obtained word is a linear De Bruijn sequence. Since cutting a complete s^q -cycle in all possible ways gives rise to s^q linear De Bruijn sequences, the number of these sequences is equal to $(s!)^{s^{q-1}}$. Under this, together with each complete cycle, there should independently be considered its mate, wherein the same sequence of letters is read in the opposite direction.

Aside from biological objects, the properties of closed and unclosed De Bruijn sequences can be utilized in the synthetical chemistry of cyclic and linear molecules, respectively. Cases in point are engineering and design of new reagents for Analytical Chemistry or drugs that employ the principles of Combinatorial Chemistry. At the first stage of synthesis, when the general prognosis should be done, the researcher is much interested in devising "the most concentrated" all-inclusive molecule which allows one to simultaneously incorporate, in one reagent, all spatial compositions of reactive groups to be attested.

Moreover, such a substance should enable every mentioned composition of groups (in our case, displayed by a different segment of a De Bruijn sequence) to contest for the best credits under equal starting conditions. Then, when the optimal molecular substructures are already determined, one can turn to the synthesis of rather simple molecules that exclude "badly behaved" parts of the first "supermolecule". Clearly, such a tack could economize syntheticist's time.

The last chemical example, even though it was described briefly, puts forward the idea of replacing an intact De Bruijn s^q -sequence with all possible sets of shorter sequences (collectively comprising the same set of s^q q -subwords); certainly, none of these shorter ones can be a De Bruijn sequence, by definition. Here, the solution for distributing a complete cycle immediately comes from our recent finding for the permanent of the arc-graph [17]. Interestingly, the number of all such sets, including an intact complete cycle, as a one-element set, is exactly equal to the number of linear De Bruijn sequences: $(s!)^{s^{q-1}}$. In our opinion, the above problems and their solutions can better be discussed in the wider context of similar combinatorial questions. However, planning to consider some additional problems in the subsequent sections, we have no intention whatever to make a detailed survey in this paper. For this reason, all references will be given in minimal numbers. We would like only to stress that other trends also exist and are all interesting as well. Wherever possible, we shall also propose problems that the reader can try to solve. Our general goal is to enhance the interest of chemists in Mathematics and, conversely, attract mathematicians to the wider range of problems that come from Chemistry, Biology and other sciences.

Now we must supply mathematical requisites that will be used by us later, in the main section.

2 Preliminaries

This section culls just all known facts from Combinatorics and (Spectral) Theory of Graphs that will be needed for proving our targeted results; all information concerning allied areas will be given in Miscellaneous.

2.1 De Bruijn sequences

A cycle is a sequence $a_1 a_2 \cdots a_r$ taken in a circular order—that is, a_1 follows a_r , and $a_2 \cdots a_r a_1, \dots, a_r a_1 \cdots a_{r-1}$ are all the same cycle as $a_1 a_2 \cdots a_r$. Given natural numbers $q \geq 1$ and $s \geq 2$, a cycle of s^q letters is called a *complete cycle* [1, 2], or *De Bruijn sequence*, if subsequences $a_i a_{i+1} \cdots a_{i+q-1}$ ($1 \leq i \leq s^q$) consist of all possible s^q ordered sequences $b_1 b_2 \cdots b_q$ over the alphabet A ($|A| = s$).

In 1946, De Bruijn [1] (see [2]) proved his famous theorem:

Theorem 1. *For $s = 2$ and each positive integer q there are exactly $2^{2^{q-1}-q}$ complete cycles of length 2^q .*

In particular, for $q = 1, 2, 3$, there exist the following complete cycles:

$$\begin{array}{ll} q = 1, & 01, \\ q = 2, & 0011, \\ q = 3, & 00010111, \\ & 00011101. \end{array}$$

Apparently, cutting a complete s^q -cycle ($q \geq 2$) in all s^q positions generates s^q distinct words since any such cycle is circularly asymmetric, by definition. However, every s^q -word obtained in this fashion contains only $s^q - q + 1$ basic subwords of length q , out of those belonging to the complete cycle. A minimal word of length $s^q + q - 1$ that incorporates just the same set of s^q basic q -subwords as an intact complete cycle is called a *linear De Bruijn sequence*. Obviously, a linear De Bruijn sequence can be obtained by adding the first $q - 1$ letters of any s^q -word, obtained by cutting a complete cycle, to the end of this word.

The following result can be regarded as a corollary of De Bruijn's theorem:

Corollary 1.1. *For $s = 2$ and each positive integer q there are exactly $2^{2^q - 1}$ linear De Bruijn sequences of length $2^q + q - 1$.*

As a brief illustration, we shall consider the cases $q = 1$ and 2, as these follow from the above example for circular De Bruijn sequences:

$$\begin{aligned} q = 1, & \quad 01, \\ & \quad 10, \\ q = 2, & \quad 00110, \\ & \quad 01100, \\ & \quad 11001, \\ & \quad 10011. \end{aligned}$$

Another generalization of complete cycle is a *De bruijn s^q -set* of sequences which are not De Bruin sequences on their own, except for the case when a De Bruijn set consists of exactly one De Bruijn sequence, but collectively have the same aggregated length s^q and also produce the same set of all s -ary words of length q ; see Theorem 11 and Corollary 11.1, in Section 3.

In order to proceed, we need to introduce some graph-theoretical notions (see [13–21]). A *directed graph*, or *digraph*, D of order n consists of a finite nonempty set V of different objects that are called *vertices*, or *points*, together with a given set E containing m ordered pairs of different vertices of the set V . A pair (u, v) , or uv , of vertices from V is called an *arc* of a digraph D that emanates from a vertex u and enters a vertex v ; under $u = v$, an arc uu (vv) is called a *self-loop* lying in the point u (v). If an arc uv exists, in D , we say that a vertex u is adjacent to a vertex v ; and a vertex u and an arc uv are *incident* to each other, as well as an arc uv and a vertex v are. The *out-degree* $d^+(v)$ of a vertex v is the number of arcs that go out of it, including self-loops; symmetrically, the *in-degree* $d^-(v)$ of v is the number of arcs (and self-loops) that come into it. In lieu of the term *degree*, we also use its synonym *valency*, which may seem preferable while describing chemical objects.

Following [1–3], we need to define the series $\mathcal{G}_s = \{G_{s,q}\}_{q=1}^\infty$ ($s \geq 2$) of special digraphs that will be used by us in the further proof; here, the numbers s and q have the same interpretation as above. Initially, we set $G_{s,1}$ to be a one-vertex graph possessing s self-loops. The set $V_{s,q}$ of vertices of a digraph $G_{s,q}$ ($q \geq 2$) consists of all s^{q-1} ordered sequences, or words, of $q - 1$ letters over the alphabet A while the set E of arcs (and self-loops) is in one-one correspondence with all s^q words of q letters over A . Under this, the arc uv labeled by a word $a_1 a_2 \cdots a_{q-1} a_q$ emanates from a vertex $u = a_1 a_2 \cdots a_{q-1}$ and enters a vertex $v = a_2 \cdots a_{q-1} a_q$. In other words, arcs $a_1 a_2 \cdots a_{q-1} a_q$ and $a_2 a_3 \cdots a_q a_{q+1}$ share a common incident vertex $a_2 a_3 \cdots a_{q-1} a_q$. It is easy to see that the arc set $E_{s,q}$ of a

digraph $G_{s,q}$ is simultaneously the vertex set $V_{s,q+1}$ of the next digraph $G_{s,q+1}$, in \mathcal{G}_s (see [1–3]). But what is rather more important, $G_{s,q+1}$ ($q \geq 1$) can be obtained from $G_{s,q}$ by the process that can locally be called taking the arc-graph $\Gamma(G_{s,q})$ of a digraph $G_{s,q}$ (see [17] or Subsection 2.3, below); under this, $G_{s,q+1} = \Gamma(G_{s,q})$. The members of the series \mathcal{G}_s were called in [3] (see [6]) *De Bruijn graphs*. Herein, we shall adapt the methods applied in [1–3], wherein estimating the number of complete s^q -cycles was reduced to calculating the number of certain circular walks in the respective De Bruijn graph $G_{s,q}$.

2.2 Counting Eulerian circuits in digraphs

A digraph D is called *Eulerian* if there exists a closed spanning walk W traversing every arc, in D , exactly once and consistently with its orientation; under this, the number of arcs entering any vertex of D equals the number of arcs emanating from it. The mentioned closed walk W , in D , is called an *Eulerian circuit*. The circular order of arcs in an Eulerian circuit is of value because one and the same Eulerian digraph D admits more than one Eulerian circuit whenever the order of circularly touring its arcs may be varied. The last circumstance plays a crucial role when Eulerian circuits formalize the cyclical motion of particles in the respective models of statistical physics, where every possible closed walk of a particle must necessarily be taken into account [17]. All the above can readily be adapted to undirected graphs if one considers every edge as a pair of opposite darts. In the last sense, any connected undirected graph G admits at least one Eulerian circuit passing along every edge strictly twice and just in opposite directions.

The adjacency matrix of an unweighted digraph D with n vertices is an $n \times n$ matrix $C = C(D) = \{c_{ij}\}_{i,j=1}^n$ of zeros and ones, wherein an entry $c_{ij} = 1$ iff (if and only if) there is an arc ij (or a self-loop ii , if $i = j$) that goes out of a vertex i and enters a vertex j of D (see [13–17; 20, 21]). Another matrix pertaining to D is its *Laplace*, *Kirchhoff*, or *admittance*, matrix $T = T(D) = \{t_{ij}\}_{i,j=1}^n$, whose entries are defined as follows (see [13–16]):

$$t_{ij} = \begin{cases} c_{ij}, & \text{if } i \neq j; \text{ and} \\ c_{ii} - d^+(i), & \text{if } i = j. \end{cases}$$

Thus, the sum of entries in each column of T equals 0. Here, we do not consider an equivalent version T^* of T , wherein similar manipulations involve the columns of the original matrix C , instead. The reader can consider T^* on his/her own, as an exercise, substituting the respective in-degrees $d^-(j)$ for the out-degrees $d^+(i)$, in the definition of T above.

Every Laplace matrix $T(D)$ (or $T^*(D)$) of an Eulerian digraph D has the property that all its cofactors T_{ij} (or T_{ij}^*) are equal; moreover, here, $T_{ij} = T_{ij}^*$ as well (see [13–16]). Just in case, we recall that a cofactor T_{ij} is the respective minor, of T , multiplied by $(-1)^{i+j}$, where the mentioned minor is in turn the determinant $\det M_{ij}$ of an $(n-1) \times (n-1)$ matrix M_{ij} , obtained by scoring out the i^{th} row and j^{th} column in T .

The common cofactor $c(D) = T_{ij} = \text{const}$ of the Laplace matrix of an Eulerian digraph D is equal to the number of oriented spanning trees that go out of (or come into) any vertex i of D (see [13–16]).

At this point, we shall cite the famous matrix-tree theorem for graphs (see [13–16]), which was first proven by De Bruijn and van Aardenne-Ehrenfest [18], viz.:

Theorem 2. *The number $\varepsilon(D)$ of Eulerian circuits in a labeled Eulerian digraph D is*

equal to

$$\varepsilon(D) = c \prod_{i=1}^n (d_i - 1)!, \quad (1)$$

where c is the common value of cofactors T_{ij} in T ; and $d_i = d^+(i) = d^-(i)$.

Theorem 2 plays a very important role herein due to the following statement that comes hand in hand with it (see [1–3]):

Proposition 3. *The number of complete cycles of length s^q over the alphabet A ($|A| = s \geq 2; q \geq 1$) is equal to the number $\varepsilon(G_{s,q})$ of Eulerian circuits in the respective De Bruijn graph $G_{s,q}$.*

Proof. (Sketch.) By the definition of a De Bruijn digraph $G_{s,q}$, every arc of it corresponds to a distinct word of length q over the alphabet A ; and all these arcs together exactly comprise all s^q possible s -ary words of q letters. Since each Eulerian circuit, in $G_{s,q}$, traverses each of its arcs exactly once, it is in one-one correspondence with one complete cycle. Hence, we at once arrive at the proof. \square

Some facts from the Spectral Theory of Graphs [16] are needed for us right now, before beginning the next subsection. Let I denote the identity matrix, that is, a diagonal matrix, whose diagonal entries are all 1s while the other entries are all 0s. The *characteristic polynomial* $P(D; x)$ of a (di-)graph D is the characteristic polynomial of its adjacency matrix $C(D)$ (see [16]); that is,

$$P(D; x) = P(C(D); x) = \det[xI - C(D)].$$

Similarly, the *Laplacian polynomial* of D is defined (see [16]):

$$L(D; x) = P(T(D); x) = \det[xI - T(D)].$$

Herein, we need to employ the spectral method [16] of calculating the common cofactor $c = c(D)$. Since all cofactors of T are equal to c , one can deduce, in particular, that the principal $(n - 1) \times (n - 1)$ minors of T are all equal to c . From the Spectral Theory of Graphs (or Matrices) [16], it immediately follows that

$$c = c(D) = \frac{1}{n} L'(D; x) |_{x=0}, \quad (2)$$

where $L'(D; x) = \frac{d}{dx} L(D; x)$.

However, for all regular digraphs (with $d^+(i) = d^-(i) = d = \text{const}$, as we have for De Bruijn graphs) the Laplacian polynomial $L(D; x)$ can readily be calculated through the respective characteristic polynomial as follows:

$$L(D; x) = P(D; x + d). \quad (3)$$

Therefore, we arrive at an equivalent result, earlier derived for multigraphs by Hutschenreuther [19] (see p. 39 in [16]), viz.:

Proposition 4. *The common value c of the cofactors T_{ij} in T can be calculated as*

$$c = c(D) = \frac{1}{n} P'(D; x) |_{x=d} \quad (4)$$

We shall use this result in the next subsection.

2.3 Spectral properties of the arc-graph

Part of the information about the properties of the arc-graph will be borrowed by us from our previous paper [17]; other properties will be proven directly in this subsection. Let $D = D(V, E)$ be a digraph with the set V of vertices and set E of arcs (self-loops, if any, are considered as self-adjacent arcs whose head and tail coincide); $|V| = n$, $|E| = m$. The arc-graph $\Gamma(D) = \Gamma(E, U)$ of a digraph D is a derivative digraph whose vertex set $V(\Gamma)$ is the set E of arcs of D ; each ordered pair ij and kl of arcs, of D , is a pair of adjacent vertices in Γ iff the head j of ij coincides with the tail k of kl ($j = k$), whether the remaining tail i and head l coincide or not.

For the sake of completeness, note that the arc-graph $\Gamma(H)$ of an undirected graph $H = H(V, E)$ can also be constructed if we initially replace each edge ij with a pair of opposite darts ($1 \leq i, j \leq |V| = n; |E| = m$), which results in the so-called symmetric digraph $S = S(H) = S(V, E^*)$ ($|E^*| = 2|E| + \text{number of self-loops, if any}$), and then revert to the above pattern.

Rosenfeld [17] obtained the following general result:

Theorem 5. *Let $P(G; x)$ and $P(\Gamma(G); x)$ be the characteristic polynomial of an arbitrary weighted (di-)graph G and that of its arc-graph $\Gamma(G)$, respectively. Then*

$$P(\Gamma(G); x) = x^{m-n} P(G; x), \quad (5)$$

where m and n are the numbers of vertices in $\Gamma(G)$ and G , respectively.

In other words, the spectra of $\Gamma(G)$ and G may differ only in the number of zero eigenvalues and this difference in the multiplicities is $|m - n|$.

The Greek character " Γ " in " $\Gamma(G)$ " can be considered as an operator Γ transforming a graph G into another one $\Gamma(G)$. This operator has some remarkable properties. In particular, it can give for any Eulerian digraph G with not less than 2 proper arcs ij ($i \neq j$) outgoing from each of its vertices i , and an arbitrary number $\# \geq 0$ of self-loops, an infinite series of such digraphs: $\Gamma^0(G) := G, \Gamma^1(G) = \Gamma(G), \Gamma^2(G) = \Gamma(\Gamma(G)), \dots, \Gamma^{q+1}(G) = \Gamma(\Gamma^q(G))$ ($q \geq 0$), whose spectra differ only in the number of zero eigenvalues.

The reader familiar with [1-3] can immediately see that an instance of the last series $\{G_q\}_{q=1}^{\infty}$ of digraphs is the series \mathcal{G}_s of the De Bruijn graphs, whose original definition obeys the same Γ -constructive property (see above). In other words, this is tantamount to the following statement:

Proposition 6. *The series $\mathcal{G}_s = \{G_{s,q}\}_{q=1}^{\infty}$ of De Bruijn graphs is a recurrent sequence of digraphs, wherein $G_{s,1}$ is a one-vertex digraph with s self-loops and $G_{s,q+1} = \Gamma(G_{s,q})$.*

Proof. To prove it, one should compare the criteria of the adjacency of arcs in $G_{s,q}$, reconsidered as vertices of $G_{s,q+1}$, given in [1-3] and in [17]. Since the two criteria coincide for constructing all the graphs $G_{s,q+1}$ in \mathcal{G}_s , the proof is immediate. \square

Now we can readily calculate the characteristic polynomial of a digraph $G_{s,q}$ ($s \geq 2; q \geq 1$); the solution will be stated as

Lemma 7. *The characteristic polynomial of $G_{s,q}$ is*

$$P(G_{s,q}) = x^{s^{q-1}-1} (x - s) \quad (s \geq 2; q \geq 1). \quad (6)$$

Proof. By virtue of Proposition 6, the repetitive application of Theorem 5 demonstrates that every digraph $G_{s,q}$ ($s \geq 2; q \geq 1$) possesses only the nonzero eigenvalue $\lambda = s$ (namely, that of the one-vertex digraph $G_{s,q}$, with s self-loops). Since the number of vertices in a digraph $G_{s,q}$ is equal to s^{q-1} , it possesses exactly $s^{q-1} - 1$ zero eigenvalues. Considering all s^{q-1} eigenvalues together, we at once arrive at the proof. \square

Proposition 4 and Lemma 7 immediately afford, as their common corollary, the following

Lemma 8. *The common value c of cofactors $T_{ij}(G_{s,q})$ in a Laplace matrix $T(G_{s,q})$ of a digraph $G_{s,q}$ is equal to*

$$c = c(G_{s,q}) = \frac{s^{s^{q-1}}}{s^q} \quad (s \geq 2; q \geq 1). \quad (7)$$

Proof. First, calculate $P'(G_{s,q}; x)$, using the R.H.S. of (6) for it:

$$[x^{s^{q-1}-1}(x-s)]' = (s^{q-1}-1)x^{s^{q-1}-2}(x-s) + x^{s^{q-1}-1}.$$

Hence, under $x = s$, Proposition 3 gives

$$\frac{1}{s^{q-1}} P'(G_{s,q}; x)|_{x=s} = \frac{1}{s^{q-1}} [0s^{s^{q-1}-2}(s^{q-1}-1) + s^{s^{q-1}-1}] = \frac{s^{s^{q-1}-1}}{s^{q-1}} = \frac{s^{s^{q-1}}}{s^q},$$

which is the proof. \square

Another important property of the operator Γ is that Γ "unties" every Eulerian circuit θ of a graph $\Gamma^q(G)$, transferring it into an oriented cycle $\Gamma(\theta)$ of Γ^{q+1} ($q \geq 0$) with the same weight $w(\Gamma(\theta)) = w(\theta)$. Here, we recall that the weight $w(\sigma)$ of any cycle σ , in an arbitrary digraph D , is the product of the weights of arcs comprising σ . Moreover, Γ assures one-to-one correspondence between the set of all Eulerian circuits on G and the set of all oriented cycles of $\Gamma(G)$.

We shall also present a partial result for the tail coefficient of the permanent polynomial $P^+(\Gamma(G); x)$ of the arc-graph $\Gamma(G)$ of an Eulerian digraph G . In particular, G may be the above symmetric digraph $S(H)$ and, consequently, the arc-graph $\Gamma(H)$ of an undirected graph H can also be considered below in place of $\Gamma(G)$. The reader interested in calculating the tail coefficient for all sorts of weighted Eulerian (di-)graphs can see [17], where this problem was completely resolved.

First, it is worth recalling that the permanent polynomial $P^+(H; x)$ of a weighted digraph H is the permanent polynomial $P^+(C(H); x)$ of its adjacency matrix $C(H)$; herein, $P^+(C(H); x) = \text{per}[xI + C(H)]$, where I is a diagonal identity matrix (see p. 34 in [16]). Thus, the tail coefficient of $P^+(C(H); x)$ is simply $\text{per}C(H)$ of the adjacency matrix $C(H)$. Below, we shall derive a corollary of the general weighted version that was proven by Rosenfeld [17], viz.:

Proposition 9. *Let $C(\Gamma(G))$ be the adjacency matrix of the arc-graph $\Gamma(G)$ of an unweighted Eulerian digraph G . Then*

$$\text{per}[C(\Gamma(G))] = \prod_{i=1}^n d_i!, \quad (8)$$

where d_i stands for the out-degree of a vertex i in G ; and the product of factorials $d_i!$ is taken over all (indices of) vertices of G .

We want to specially introduce the definition of Eulerian subcircuit because it may otherwise seem ambiguous. Namely, an *Eulerian subcircuit* of a digraph D is the Eulerian circuit of its Eulerian subgraph $D_1 \subseteq D$ that takes into account exactly one circular order in which all arcs of D_1 can be traversed. In general, there may be more than one circular order for passing all arcs of D_1 ; therefore, the number of Eulerian subcircuits corresponding to D_1 may be more than 1.

Graph-theoretically, $\text{per}[C(\Gamma(G))]$ is the number of ways in which all arcs of G can be covered by its arc-disjoint Eulerian subcircuits (see [17]). To facilitate referring to this fact in the subsequent text, we shall derive the following working corollary of the last proposition:

Corollary 9.1. *Let $G_{s,q}$ ($s \geq 2; q \geq 1$) be a De Bruijn digraph. Then the number of ways in which all s^q arcs of $G_{s,q}$ can be covered by Eulerian subcircuits is $(s!)^{s^{q-1}}$.*

Proof. Setting the values $d_i = s$ and $n = s^{q-1}$ in (8) at once affords the proof. \square
Also, due to the above "untying" properties of the operator Γ , the permanent $\text{per}[C(\Gamma(G))]$ is the number of spanning cycle covers of Γ (that collectively cover all vertices of Γ). Therefore, we can end this subsection by formulating another corollary, viz.:

Corollary 9.2. *Let $G_{s,q+1}$ ($s \geq 2; q \geq 1$) be a De Bruijn digraph. Then the number of ways in which all s^q vertices of $G_{s,q+1}$ can be covered by oriented cycles is $(s!)^{s^{q-1}}$.*

Proof. Recalling that all s^q vertices of $G_{s,q+1}$ are exactly all arcs of $G_{s,q}$ and applying Corollary 9.1 to the last digraph, we immediately arrive at the proof. \square

At this point, it is time to summarize the tack which will be followed by us, in the next section.

2.4 Our tack

We shall keep the general ideas expounded in [1-3], according to which the enumeration of complete s^q -cycles, can be reduced to counting the number of Eulerian circuits in the respective DeBruin graph $G_{s,q}$ ($s \geq 2; q \geq 1$). Under this, we shall employ our recent results concerning the spectral properties of iterated arc-graphs [17], which are exemplified herein by the De Bruijn graphs. It will enable us to obtain the overall solution for all $s \geq 2$ and $q \geq 1$. We also plan to discuss some related combinatorial problems, in Miscellaneous.

3 Main results

We at once begin this section with its master theorem:

Theorem 10. *For positive integers $s \geq 2$ and $q \geq 1$ there are exactly $(s!)^{s^{q-1}-q}$ complete cycles of length s^q .*

Proof. By virtue of Lemma 8, c on the R.H.S. of (1) is equal to $s^{s^{q-1}}/s^q$ (see the R.H.S. of (7)); and, by definition of the De Bruijn graphs $G_{s,q}$, the degree $d_i = s$ ($1 \leq i \leq s^{q-1}$). With these specific values on the R.H.S. of (1), Theorem 2 gives

$$\varepsilon(G_{s,q}) = \frac{s^{s^{q-1}}}{s^q} [(s-1)!]^{s^{q-1}} = s^{s^{q-1}-q}.$$

But, by virtue of Proposition 3, $\varepsilon(G_{s,q})$ is also the number of complete s^q -cycles, whence the proof is immediate. \square

Theorem 10 gives, as its elementary corollaries, De Bruijn's Theorem (Theorem 1, herein) and Corollary 1.1. Moreover, we can formulate here "the generalized Corollary 1.1", viz.:

Corollary 10.1. *For positive integers $s \geq 2$ and $q \geq 1$ there exist exactly $(s!)^{s^{q-1}}$ linear De Bruijn sequences of length $s^q + q - 1$.*

Proof. It immediately follows from Theorem 10 and the definition of a linear De Bruijn sequence. \square

We can also calculate the number of De Bruijn s^q -sets due to the following theorem:

Theorem 11. *For positive integers $s \geq 2$ and $q \geq 1$, the number of De Bruijn s^q -sets is equal to $(s!)^{s^{q-1}}$.*

Proof. This generalizes the Proof of Theorem 10, where the number of complete s^q -cycle is calculated as the number $\varepsilon(G_{s,q})$ of Eulerian circuits of a De Bruijn graph $G_{s,q}$. Now, in lieu of that, we should consider the number of all possible covers of all arcs of $G_{s,q}$ by its Eulerian subcircuits. But the last number is given by Corollary 8.1 as $(s!)^{s^{q-1}}$. Hence, we immediately arrive at the proof. \square

Thus, one can come to the following common corollary of Corollary 10.1 and Theorem 11:

Corollary 11.1. *For positive integers $s \geq 2$ and $q \geq 1$ the number of linear De Bruijn sequences of length $s^q + q - 1$ equals the number of De Bruijn s^q -sets: $(s!)^{s^{q-1}}$.*

In our opinion, such a coincidence may lead to new like findings concerning De Bruijn sequences and/or their generalizations. But, at this point, we must stop our consideration of this topic and turn to discussing other combinatorial problems that, however, resemble by their appearance the above ones.

4 Miscellaneous

This section is a small compilation that seems to be close to the main text, done at the author's choice. It is a mere discussion of known results and methods [4–12; 20–23] but contains, at the end, some open problems that can be proposed to the reader.

4.1 Kautz s -ary closed sequences

A *Kautz s -ary closed sequence* is a circular sequence of l s -ary digits $0, 1, \dots, s-1$ such that consecutive digits are distinct and all subsequences of length q are distinct, too [6]. Thus, Kautz sequences are non-DeBruijn sequences included in the respective De Bruijn s^q -sets, with an additional proviso that equal digits may never be adjacent therein. Kautz sequences can also be represented by the series $\mathcal{H} = \{H_{s,q}\}_{q=s}^{\infty}$ ($s \geq 2$) of Kautz digraphs [6] that resemble De Bruijn graphs [1–3]. Namely, $H_{s,1} = K_s$, where K_s is a complete s -vertex digraph without self-loops, and $H_{s,q+1} = \Gamma(H_{s,q}) = \Gamma^q K_s$. Villar [6] proved that Kautz sequences exist for all lengths l except for 1 and $r(r-1)^{q-1} - 1$, where $q \geq 2$ and $r = s(s-1)^q$ is the number of arcs in a digraph $H_{s,q}$. To the best of our knowledge, counting of Kautz sequences (of admissible lengths) is hitherto an unsolved problem.

4.2 Other sequences with adjacency restrictions

The enumeration of s -ary circular sequences of length q is tantamount to the enumeration of q -bead necklaces with s kinds of beads (and other combinatorial restrictions, if any). Namely, the latter interpretation was adopted by Lloyd [8] for enumerating s -ary q -sequences with any given restrictions put on the adjacency of different ciphers (which can be adjacent or self-adjacent and which not). However, it should be noted that his calculation always considers two circular sequences as equal if one can be obtained from the other by reading the original sequence in the opposite direction. The instances of De Bruijn and Kautz sequences, however, do not admit such reversing of the circular order. Nevertheless, the work of Lloyd [8] is of paramount importance for chemists whenever they want to know the number of cyclic substitutional isomers with a given number of sorts of substituents and specified adjacency restrictions put on them.

In Graph Theory, it is known that the number of colorations of a labeled l -cycle ($l \geq 2$) with s colors provided that no two adjacent vertices are colored the same color is equal to $(s-1)^l + (-1)^l(s-1)$; see Theorem IX.23 in [15]. The respective result for a labeled path spanning l vertices is $s(s-1)^{l-1}$; see Theorem IX.24 in [15]. Here, we recall that labeled graphs take into account no symmetry whatever, even if they possess it. Nevertheless, using the so-called inclusion/exclusion principle (see [20–22]) enables one to utilize the results concerning labeled graphs for enumerating the colorations of the respective unlabeled graphs with a given group of automorphisms (or symmetry group) [20–22], and even with a given monoid (semigroup) of endomorphisms [22]. From the general combinatorial point of view, the same procedure works equally well for s -ary sequences, too.

From among other sequences, we shall pick herein only few [9–12; 23]. In particular, [9–10] investigate the *square-free words*; in these, subwords BbB , wherein B is an arbitrary block and b is the first letter of it, are forbidden. Overlapping words and special circular codes have been considered in [11] and [12], respectively. Finally, uncancelable sequences of the elements of a finite regular monoid R that exclude subsequences of type aba , wherein a and b are inverses in R , are of use in an algebraic treatment of genomic sequences [23], proposed by the present author.

Now we shall turn to posing some problems that follow from the whole text above.

4.3 Open problems

The following problems will represent only a very small part of the problems that could be posed in such a case.

Problem 1. To enumerate s -ary circular sequences of length l ($1 \leq l < s^q$) that are included in all De Bruijn s^q -sets with fixed positive integers s and q ; $s \geq 2$ and $q \geq 1$.

Problem 2. To enumerate subsequences of length l ($1 \leq l < s^q + q - 1$) of all linear s -ary De Bruijn sequences of length $s^q + q - 1$ with fixed positive integers s and q ; $s \geq 2$ and $q \geq 1$.

Problem 3. To enumerate s -ary Kautz sequences of length l ($s \geq 2; l \geq 2$).

Problem 4. To enumerate s -ary ($s \geq 2$) subsequences without subwords of type aba .

Some other sequences, whose consideration is omitted herein, are planned to be considered in our next publications.

Acknowledgments

I am grateful to Dr. Valery Kirzhner, Profs. Edward Trifonov and Alexander Bolshoy for stimulating the idea of preparing this paper. Also, I sincerely thank my anonymous referee for his/her expert work and linguistical help that considerably improved the text.

References

- [1] De Bruijn N. G., A Combinatorial Problem, *Nederl. Akad. Wetensch. Proc.*, 1946, v. 49, p. 758–764. *Indagationes Math.*, 1946, v. 8, p. 461–467.
- [2] Hall M., Jr., *Combinatorial Theory*, John Wiley and Sons, Inc., New York, 1967, p. 91–99.
- [3] Lempel A., m -Ary Closed Sequences, *J. Comb. Theory*, 1971, v. 10, p. 253–258.
- [4] Kautz W. H., in *Design of Optimal Interconnection Networks for Multiprocessors, Architecture and Design of Digital Computers. Nato Advanced Summer Institute*, 1969, p. 249–272.
- [5] Fiol M. A., Alegre I. and Yerra J. L. A., Line Digraph Iterations and the (d, k) Problem for Directed Digraphs, *Proc. 10th Int. Symp. Comp. Arch.*, 1983, p. 174–177.
- [6] Villar J. L., Kautz s -Ary Closed Sequences, in *Combinatorics '88. Proceedings of the International Conference on Incidence Geometries and Combinatorial Structures, Ravello, Italy, May 23–28, 1988* (A. Barlotti, G. Lunardon, F. Mazzocca, N. Melone, D. Olanda, A. Pasini and G. Tallini, eds.), v. 2, Mediterranean Press, 1991, p. 459–469.
- [7] Balaban A. T. and Harary F., Chemical Graphs IV (Aromaticity VI): Dihedral Groups and Monocyclic Aromatic Compounds. *Rev. Roumaine Chim.*, 1967, v. 12, p. 1511–1515.
- [8] Lloyd E. K., Necklace Enumeration with Adjacency Restrictions, in *Combinatorics. The Proceedings of the British Combinatorial Conference held in the University College of Wales, Aberystwyth, 2–6 July 1973* (T. P. McDonough and V. C. Mavron, eds.), Cambridge University Press, Cambridge, 1974, p. 97–102.
- [9] Fife E. D., Irreducible Binary Sequences, in *Combinatorics on Words—Progress and Perspectives* (L. J. Cumming, ed.), Academic Press, Toronto, New York, p. 91–100.
- [10] Shelton R. O., On the Structure and Extendibility of Square-Free Words, *ibid.*, p. 101–118.
- [11] Almeida J., Overlapping of Words in Rational Languages, *ibid.*, p. 119–131.
- [12] Berstel J. and Perrin D., Codes Circulaires, *ibid.*, p. 133–165.
- [13] Harary F., *Graph Theory*, Addison-Wesley Publishing Co., Inc., Reading, Mass., 1969.

- [14] Harary F. and Palmer E. M., *Graphical Enumeration*, Academic Press, New York and London, 1973.
- [15] Tutte W. T., *Graph Theory*, Addison-Wesley, 1984.
- [16] Cvetković D.M., Doob M. and Sachs H., *Spectra of Graphs: Theory and Application*, Academic Press, Berlin, 1980.
- [17] Rosenfeld V. R., Some Spectral Properties of the Arc-Graph, *Commun. Math. Comput. Chem. (MATCH)*, 2001, no. 43, p. 41–48.
- [18] De Bruijn N. G. and van Aardenne-Ehrenfest T., Circuits and Trees in Oriented Graphs, *Simon Stevin*, 1951, v. 28, p. 203–217.
- [19] Hutschenreuther H., Einfacher Beweis des Matrix-Gerüst-Satzes der Netzwerktheorie, *Wiss. Z. TH Ilmenau*, 1967, v. 13, s. 403–404.
- [20] Kerber A., *Algebraic Combinatorics via Finite Group Actions*, Wissenschaftsverlag, Mannheim, Wien, Zürich, 1991.
- [21] Kerber A., *Applied Finite Group Actions*, Springer Verlag, Berlin, Heidelberg, New York, 1999.
- [22] Rosenfeld V. R., Yet Another Generalization of Pólya's Theorem: Enumerating Equivalence Classes of Objects with a Prescribed Monoid of Endomorphisms, *Commun. Math. Comput. Chem. (MATCH)*, 2001, no. 43, p. 111–130.
- [23] Rosenfeld V. R., An Algebraic Model of Closed Loops in Proteins, *Commun. Math. Comput. Chem. (MATCH)*, submitted.
- [24] Trifonov E. N., Making Sense of the Human Genome, in *Structure and Methods*, v. 1, *Human Genome Initiative and DNA Recombination*, Adenine Press, N. Y., 1990, p. 69–76.
- [25] Trifonov E. N., Informational Structure of Genetic Sequences and Nature of Gene Splicing, in *Advances in Biomolecular Simulations*, (E. N. Lavery, J.-L. Rivail and J. Smith, eds.), AIP Conference Proceedings 239, N. Y., 1991, p. 329–338.
- [26] Popov O., Segal D. M. and Trifonov E. N., Linguistic Complexity of Protein Sequences As Compared to Texts of Human Languages, *Biosystems*, v. 38, p. 65–74.
- [27] Bolshoy A., Shapiro K., Trifonov E. N. and Ioshiknes I., Enhancement of the Nucleosomal Pattern in Sequence of Lower Complexity, *Nucleic Acids Research*, 1997, v. 25, no. 16, p. 3248–3254.
- [28] Gabrielian A. and Bolshoy A., Sequence Complexity and DNA Curvature, *Comput. Chem*, 1999, v. 23, p. 263–274.
- [29] Troyanskaya O. G., Arbell O., Koren Y., Landau G. M. and Bolshoy A., Sequence Complexity Profiles of Prokariotic Genomic Sequences: A fast Algorithm for Calculating Linguistic Complexity, *Bioinformatics*, 2002, submitted.