



Brief communication

Identification of all trinucleotide circular codes

Christian J. Michel^{a,*}, Giuseppe Pirillo^{b,c}^a Equipe de Bioinformatique Théorique, FDBT, LSIT (UMR CNRS-ULP 7005), Université de Strasbourg, Pôle API, Boulevard Sébastien Brant, 67400 Illkirch, France^b Consiglio Nazionale delle Ricerche, Istituto di Analisi dei Sistemi ed Informatica "Antonio Ruberti", Unità di Firenze, Dipartimento di Matematica "U.Dini", viale Morgagni 67/A, 50134 Firenze, Italia^c Université de Marne-la-Vallée, 5 boulevard Descartes, 77454 Marne-la-Vallée Cedex 2, France

ARTICLE INFO

Article history:

Received 31 January 2010

Received in revised form 18 March 2010

Accepted 21 March 2010

Keywords:

Circular code

Trinucleotide

Growth function

Necklace

Gene

Genetic code

ABSTRACT

A new trinucleotide proposition is proved here and allows all the trinucleotide circular codes on the genetic alphabet to be identified (their numbers and their sets of words). This new class of genetic motifs, i.e. circular codes (or synchronizing genetic motifs), may be involved in the structure and the origin of the genetic code, and in reading frames of genes.

© 2010 Elsevier Ltd. All rights reserved.

1. Introduction

We continue our study of the properties of trinucleotide circular codes. For 50 years, codes, comma-free codes and circular codes have been mathematical objects studied in biology, mainly to understand the structure and the origin of the genetic code as well as the reading frame (construction) of genes, see the pioneer works (Crick et al., 1957; Golomb et al., 1958a,b). In order to have an intuitive meaning of these notions, codes are written on a straight line while comma-free codes and circular codes are written on a circle, but in both cases, unique decipherability is required.

The genetic code based on 64 trinucleotides is a code in the sense of language theory, more precisely a uniform code (Berstel and Perrin, 1985), but not a circular code (Lassez, 1976) (see Remark 2 below). Before the discovery of the genetic code, Crick et al. (1957) proposed a maximal comma-free code of 20 trinucleotides for coding the 20 amino acids. In 1996, a maximal circular code X_0 of 20 trinucleotides was identified statistically on a large gene population of eukaryotes and also on a large gene population of prokaryotes (Arquès and Michel, 1996):

$$X_0 = \{AAC, AAT, ACC, ATC, ATT, CAG, CTC, CTG, GAA, GAC, GAG, GAT, GCC, GGC, GGT, GTA, GTC, GTT, TAC, TTC\}$$

* Corresponding author.

E-mail addresses: michel@dpt-info.u-strasbg.fr (C.J. Michel), pirillo@math.unifi.it (G. Pirillo).

This code X_0 has remarkable properties. For example, X_0 is self-complementary: 10 trinucleotides are complementary to the 10 other trinucleotides, e.g. AAC is complementary to GTT, AAT to ATT, etc. The two sets of 20 trinucleotides, called X_1 and X_2 , obtained by a simple shift operation of X_0 , one and two letters, respectively, are also maximal circular codes (Arquès and Michel, 1996). This surprising result, still mysterious, was cited/discussed in research works in mathematics/computer science and mainly in theoretical biology, e.g. (Koch and Lehman, 1997; Béal and Senellart, 1998; Bassino, 1999; Štambuk, 1999; Jolivet and Rothen, 2001; Nikolaou and Almirantis, 2003; May et al., 2004; Lassez et al., 2007; Pirillo, 2003; José et al., 2009). Its main biological consequence would be that genes have (or had) two codes: the classical genetic code to code the amino acids and a circular code to retrieve the reading frames of genes. Therefore, the computational study of trinucleotide circular codes is particularly important in biology.

The determinations of very small classes of trinucleotide circular codes, precisely the 99,320 self-complementary trinucleotide circular codes (Pirillo and Pirillo, 2005) and about 559 millions trinucleotide comma-free codes (Michel et al., 2008a), were obtained by using the classical flower automaton algorithm (Berstel and Perrin, 1985). We recently identified a relation between these two classes of trinucleotide codes by constructing a hierarchy of codes that are closed by the comma-free codes and the circular codes (Michel et al., 2008b). The whole class of all the trinucleotide circular codes is identified in this paper (their numbers and their sets of words). This problem has a computational complexity with

an order of magnitude significantly higher than the two previous cases (more than 200 times). Indeed, about 116 billion trinucleotide circular codes are identified. The proof of a new trinucleotide proposition (**Proposition 3**), which appears obvious afterwards, allows the computational problem associated with the general case to be solved. Thus, this short **Proposition 3** which can easily be programmed, allows circular codes (synchronizing genetic motifs) on the genetic alphabet to be identified.

2. Definitions

Let \mathcal{A} denote a finite alphabet, \mathcal{A}^* , the set of all words over \mathcal{A} and \mathcal{A}^+ , the set of all words over \mathcal{A} except the empty word ε . Given a subset X of \mathcal{A}^* , X^n is the set of the words over \mathcal{A} which is the product of n words from X , i.e. $X^n = \{x_1x_2 \cdots x_n | x_i \in X\}$.

There is a correspondence between the genetic and language-theoretic concepts. The letters (or nucleotides or bases) define the genetic alphabet $\mathcal{A}_4 = \{A, C, G, T\}$. The set of non-empty words (resp. words) over \mathcal{A}_4 is denoted by \mathcal{A}_4^+ (resp. \mathcal{A}_4^*). The set of the 16 words of length two (or dinucleotides or dileters) is denoted by \mathcal{A}_4^2 . The set of the 64 words of length three (or trinucleotides or trileters) is denoted by \mathcal{A}_4^3 . The total order over the alphabet \mathcal{A}_4 is $A < C < G < T$. Consequently, \mathcal{A}_4^+ is lexicographically ordered: given two words $u, v \in \mathcal{A}_4^+$, u is smaller than v in lexicographical order, written $u < v$, if and only if either u is a proper prefix of v or there exist $x, y \in \mathcal{A}_4$, $x < y$, and $r, s, t \in \mathcal{A}_4^*$ such that $u = rxs$ and $v = ryt$.

Definition 1. Code: A set X of \mathcal{A}^+ is a code over \mathcal{A} if for each $x_1, \dots, x_n, x'_1, \dots, x'_m \in X$, $n, m \geq 1$, the condition $x_1 \cdots x_n = x'_1 \cdots x'_m$ implies $n = m$ and $x_i = x'_i$ for $i = 1, \dots, n$.

Remark 1. The set \mathcal{A}_4^3 itself is a code. More precisely, it is a uniform code (Berstel and Perrin, 1985).

Notation 1. Consequently, any non-empty subset of \mathcal{A}_4^3 is a code called trinucleotide code in this paper.

Definition 2. Trinucleotide circular code: A trinucleotide code $X \in \mathcal{A}_4^3$ is circular if for each $x_1, \dots, x_n, x'_1, \dots, x'_m \in X$, $n, m \geq 1$, $p \in \mathcal{A}_4^*$, $s \in \mathcal{A}_4^*$, the conditions $sx_2 \cdots x_n p = x'_1 \cdots x'_m$ and $x_1 = ps$ imply $n = m$, $p = \varepsilon$ and $x_i = x'_i$ for $i = 1, \dots, n$.

Remark 2. \mathcal{A}_4^3 is obviously not a trinucleotide circular code.

Definition 3. Maximal trinucleotide circular code: A trinucleotide circular code $X \in \mathcal{A}_4^3$ is maximal if for each $x \in \mathcal{A}_4^3$, $X \cup \{x\}$ is not a trinucleotide circular code.

Remark 3. Any trinucleotide circular code with 20 words is maximal. Therefore, the lengths of trinucleotide circular codes vary between 1 and 20.

3. Propositions

Proposition 1. The number of trinucleotide circular codes of length 1 is equal to 60.

Proof. Obvious. \square

Proposition 2. The number of trinucleotide circular codes of length 20 is equal to 12,964,440.

Proof. This number was obtained in 1996 by using the flower automaton algorithm (Table 2(d) in Arquès and Michel, 1996). \square

In order to compute the growth function of trinucleotide circular codes for all lengths $l = 1, \dots, 20$, we extend the necklace definition (Pirillo, 2003; Michel et al., 2008b). $l_1, l_2, \dots, l_{n-1}, l_n, \dots$ are

letters in \mathcal{A}_4 , $d_1, d_2, \dots, d_{n-1}, d_n, \dots$ are dileters in \mathcal{A}_4^2 and n is an integer satisfying $n \geq 2$.

Definition 4. Letter Dileter Continued Closed Necklaces (LDCCN): We say that the ordered sequence $l_1, d_1, l_2, d_2, \dots, d_{n-1}, l_n, d_n, l_{n+1}$ is an $(n+1)$ LDCCN for a subset $X \subset \mathcal{A}_4^3$ if $l_1 d_1, l_2 d_2, \dots, l_n d_n \in X$ and $d_1 l_2, d_2 l_3, \dots, d_{n-1} l_n, d_n l_{n+1} \in X$ and $l_1 = l_{n+1}$.

Notation 2. An $(n+1)$ LDCCN $l_1, d_1, l_2, d_2, \dots, d_{n-1}, l_n, d_n, l_{n+1}$ is denoted by $[l_1, d_1, l_2, d_2, \dots, d_{n-1}, l_n, d_n]$. Accordingly: a 2LDCCN, i.e. $[l_1, d_1]$, has the form l_1, d_1, l_1 ; a 3LDCCN, i.e. $[l_1, d_1, l_2, d_2]$, has the form l_1, d_1, l_2, d_2, l_1 ; a 4LDCCN, i.e. $[l_1, d_1, l_2, d_2, l_3, d_3]$, has the form $l_1, d_1, l_2, d_2, l_3, d_3, l_1$; a 5LDCCN, i.e. $[l_1, d_1, l_2, d_2, l_3, d_3, l_4, d_4]$, has the form $l_1, d_1, l_2, d_2, l_3, d_3, l_4, d_4, l_1$.

Proposition 3. Let X be a trinucleotide circular code. The following conditions are equivalent.

- (i) X is a trinucleotide circular code.
- (ii) X has no nLDCCN for any integer $n \in \{2, 3, 4, 5\}$.

Proof. (i) \Rightarrow (ii). By way of contradiction, suppose that X has some nLDCCN for some integer $n \in \{2, 3, 4, 5\}$.

If it is a 2LDCCN then $l_1, d_1, l_1, d_1, l_1, d_1, l_1, d_1, l_1$ is a 5LDCCN for X .

If it is a 3LDCCN then $l_1, d_1, l_2, d_2, l_1, d_1, l_2, d_2, l_1$ is a 5LDCCN for X .

If it is a 4LDCCN then $l_1, d_1, l_2, d_2, l_3, d_3, l_1, d_1, l_2$ is a 5LDCCN for X .

If it is a 5LDCCN then $l_1, d_1, l_2, d_2, l_3, d_3, l_4, d_4, l_1$ is a 5LDCCN for X .

In each of these four cases, by **Proposition 1**, X is not a trinucleotide circular code. Contradiction.

(ii) \Rightarrow (i). By way of contradiction, suppose that X is not a trinucleotide circular code. By **Proposition 1**, X has a 5LDCCN, say $l_1, d_1, l_2, d_2, l_3, d_3, l_4, d_4, l_5$. As \mathcal{A}_4 has four letters, then $l_i = l_j$ for some i, j , $1 \leq i < j \leq 5$.

If $j - i = 4$ then $l_1 = l_5$ and $[l_1, d_1, l_2, d_2, l_3, d_3, l_4, d_4]$ is a 5LDCCN for X .

If $j - i = 3$ then $[l_i, d_i, l_{i+1}, d_{i+1}, l_{i+2}, d_{i+2}]$ is a 4LDCCN for X .

If $j - i = 2$ then $[l_i, d_i, l_{i+1}, d_{i+1}]$ is a 3LDCCN for X .

If $j - i = 1$ then $[l_i, d_i]$ is a 2LDCCN for X .

In each of these four cases, by **Proposition 1**, there is a contradiction with (ii). \square

Necklace algorithm (principle): This new **Proposition 3** is used to compute all the trinucleotide circular codes (growth function for all lengths $l = 1, \dots, 20$). The principle of this necklace algorithm is simple. If the algorithm identifies a necklace i LDCCN for a given $i \in \{2, 3, 4, 5\}$ in a code, then it is not circular and the algorithm stops avoiding to analyse the next necklaces j LDCCN for $j > i$ and $j \in \{2, 3, 4, 5\}$.

4. Results

Table 1 shows the number $\text{Nb}(l)$ of trinucleotide circular codes of length l . The growth function has a minimum number $\text{NbMin} = 60$ at $l = 1$ and a maximum number $\text{NbMax} = 23,403,485,556$ at $l = 13$. **Fig. 1** associated with **Table 1** gives the graphical distribution of trinucleotide circular codes. The distribution is asymmetric with respect to NbMax at $l = 13$. The numbers of codes of $l = 13$ and $l = 14$ are close. There are $\text{NbPot}(l) = \binom{20}{l} \times 3^l$ potential

trinucleotide circular codes of length $l \in \{1, 20\}$. Therefore, the probability $\text{Pr}(l)$ of a trinucleotide circular code of length l is equal to $\text{Pr}(l) = \text{Nb}(l)/\text{NbPot}(l)$. **Table 1** and **Fig. 1** also show this proba-

Table 1
Growth function of trinucleotide circular codes. The 1st, 2nd and 3rd rows give their lengths l , their occurrence numbers $Nb(l)$ and their probabilities $Pr(l)$, respectively.

l	1	2	3	4	5	6	7	8	9	10
$Nb(l)$	60	1,704	30,432	382,164	3,568,212	25,507,512	141,639,780	614,568,102	2,086,742,208	5,542,646,244
$Pr(l)$	1	9.96×10^{-1}	9.89×10^{-1}	9.74×10^{-1}	9.47×10^{-1}	9.03×10^{-1}	8.35×10^{-1}	7.44×10^{-1}	6.31×10^{-1}	5.08×10^{-1}
l	11	12	13	14	15	16	17	18	19	20
$Nb(l)$	11,503,061,124	18,615,667,124	23,403,485,556	22,700,634,924	16,787,523,072	9,279,022,320	3,708,717,048	1,012,099,740	168,726,792	12,964,440
$Pr(l)$	3.87×10^{-1}	2.78×10^{-1}	1.89×10^{-1}	1.22×10^{-1}	7.55×10^{-2}	4.45×10^{-2}	2.52×10^{-2}	1.37×10^{-2}	7.26×10^{-3}	3.72×10^{-3}

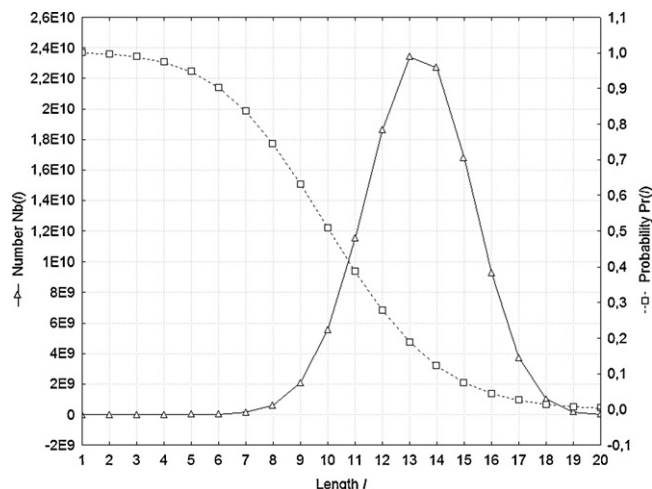


Fig. 1. Growth function of trinucleotide circular codes. The x-axis represents their length l , the primary y-axis, their number $Nb(l)$ and the secondary y-axis, their probability $Pr(l)$.

bility distribution which decreases from 1 at $l = 1$ to 3.72×10^{-3} at $l = 20$.

5. Conclusion

After a first computer result in 1996 giving the number (and its sets of words) of trinucleotide circular codes of length $l = 20$ (maximal codes) (Arquès and Michel, 1996), the trinucleotide circular codes for all lengths $l = 1, \dots, 20$ are identified here (numbers and sets of words). In particular, the 12,964,440 maximal trinucleotide circular codes obtained with the flower automaton algorithm are retrieved with this new necklace algorithm. The necklace programming allows the determination of these synchronizing genetic motifs which may be involved, in particular, in the structure and the origin of the genetic code, and in the reading frames of genes.

Acknowledgement

The second author thanks the Dipartimento di matematica “U. Dini” for giving him a friendly hospitality.

References

- Arquès, D.G., Michel, C.J., 1996. A complementary circular code in the protein coding genes. *J. Theor. Biol.* 182, 45–58.
- Bassino, F., 1999. Generating function of circular codes. *Adv. Appl. Math.* 22, 1–24.
- Béal, M.-P., Senellart, J., 1998. On the bound of the synchronization delay of a local automaton. *Theor. Comput. Sci.* 205, 297–306.
- Berstel, J., Perrin, D., 1985. *Theory of Codes*. Academic Press, London.
- Crick, F.H.C., Griffith, J.S., Orgel, L.E., 1957. Codes without commas. *Proc. Natl. Acad. Sci.* 43, 416–421.
- Golomb, S.W., Gordon, B., Welch, L.R., 1958a. Comma-free codes. *Can. J. Math.* 10, 202–209.
- Golomb, S.W., Welch, L.R., Delbrück, M., 1958b. Construction and properties of comma-free codes. *Biol. Med. Dan. Vid. Selsk.*, 23.
- Jolivet, R., Rothen, F., 2001. Peculiar symmetry of DNA sequences and evidence suggesting its evolutionary origin in a primeval genetic code. In: Ehrenfreund, P., Angerer, O., Battrick, B. (Eds.), *Proceedings of the First European Workshop in Exo-/astro-biology*. ESA SP-496 Noordwijk, 173–176.
- José, M.V., Govezensky, T., García, J.A., Bobadilla, J.R., 2009. On the evolution of the standard genetic code: vestiges of critical scale invariance from the RNA world in current prokaryote genomes. *PLoS ONE* 4 (2), e4340.
- Koch, A.J., Lehman, J., 1997. About a symmetry of the genetic code. *J. Theor. Biol.* 189, 171–174.
- Lassez, J.-L., 1976. Circular codes and synchronization. *Int. J. Comput. Syst. Sci.* 5, 201–208.
- Lassez, J.-L., Rossi, R.A., Bernal, A.E., 2007. Crick’s hypothesis revisited: the existence of a universal coding frame. *IEEE AINAW’07*.

- May, E.E., Vouk, M.A., Bitzer, D.L., Rosnick, D.I., 2004. An error-correcting framework for genetic sequence analysis. *J. Franklin Inst.* 341, 89–109.
- Michel, C.J., Pirillo, G., Pirillo, M.A., 2008a. Varieties of comma-free codes. *Comput. Math. Appl.* 55, 989–996.
- Michel, C.J., Pirillo, G., Pirillo, M.A., 2008b. A relation between trinucleotide comma-free codes and trinucleotide circular codes. *Theor. Comput. Sci.* 401, 17–25.
- Nikolaou, C., Almirantis, Y., 2003. Mutually symmetric and complementary triplets: difference in their use distinguish systematically between coding and non-coding genomic sequences. *J. Theor. Biol.* 223, 477–487.
- Pirillo, G., 2003. A characterization for a set of trinucleotides to be a circular code. In: Pellegrini, C., Cerrai, P., Freguglia, P., Benci, V., Israel, G. (Eds.), *Determinism, Holism, and Complexity*. Kluwer.
- Pirillo, G., Pirillo, M.A., 2005. Growth function of self-complementary circular codes. *Biol. Forum* 98, 97–110.
- Štambuk, N., 1999. On circular coding properties of gene and protein sequences. *Croat. Chem. Acta* 72, 999–1008.