



# Circular codes in archaeal genomes

Gabriel Frey, Christian J. Michel\*

*Equipe de Bioinformatique Théorique, LSIT (UMR CNRS-ULP 7005), Université Louis Pasteur de Strasbourg, Pôle API,  
Boulevard Sébastien Brant, 67400 Illkirch, France*

Received 7 November 2002; received in revised form 27 February 2003; accepted 5 March 2003

## Abstract

A new statistical method associating each trinucleotide with a frame is developed for identifying circular codes. Its sensibility allows the detection of several circular codes in the (protein coding) genes of archaeal genomes. Several properties of these circular codes are described, in particular the lengths of the minimal windows to retrieve the construction frames, a new definition of a parameter for measuring some probabilities of words generated by the circular codes, and the types of nucleotides in the trinucleotide sites. Some biological consequences are presented in Discussion.

© 2003 Elsevier Ltd. All rights reserved.

*Keywords:* Circular codes; Archaeal genomes; Genes; Statistical method; Trinucleotide; Frame

## 1. Introduction

The trinucleotide distribution in the (protein coding) genes is not random. Indeed, the synonymous codons (codons coding for the same amino acid) do not occur with the same frequencies in the genes of different species (Grantham et al., 1980; Sharp and Matassi, 1994; Antezana and Kreitman, 1999). This codon usage bias depends on several biological factors, mainly the gene function, the DNA spatial structure (e.g. stacking energies, curvature, superhelicity) and the evolutionary process. The gene function appears to be the most important factor (Ma et al., 2002) compared to the spatial (e.g. Tazi and Bird, 1990) and evolutionary (e.g. Smith et al., 1985) ones. Several processes involved in the gene function may determine the codon usage bias in order to increase translational efficiency and accuracy (Table 1).

The statistical methods leading to these previous results are usually simple and mainly based on the occurrence frequencies of basic motifs in genes, such as the mononucleotides, the dinucleotides, the codons, etc. In this line of research, we have added a new method called frame trinucleotide frequency (FTF) which

analyses the trinucleotide occurrence frequencies in the 3 frames of genes (Arquès and Michel, 1996). By convention, the reading frame established by the codon ATG is frame 0, and frames 1 and 2 are the reading frames shifted by 1 and 2 nucleotides in the 5'–3' direction, respectively. Then, a preferential frame for the 64 trinucleotides can be deduced by assigning to each trinucleotide the frame associated with its highest occurrence frequency. Totally unexpected, by excluding the identical trinucleotides (AAA, CCC, GGG and TTT) and with a few exceptions, this approach identifies the same 3 subsets of 20 trinucleotides per frame in the 2 gene populations  $P$  of eukaryotes EUK and prokaryotes PRO. These 3 sets,  $\mathcal{X}_0(\text{EUK}) = \mathcal{X}_0(\text{PRO})$  (resp.  $\mathcal{X}_1(\text{EUK}) = \mathcal{X}_1(\text{PRO})$  and  $\mathcal{X}_2(\text{EUK}) = \mathcal{X}_2(\text{PRO})$ ) associated with the frame 0 (resp. 1 and 2) have several interesting properties, in particular the property of circular code which is briefly recalled here without mathematical notations. The other properties, such as permutation and complementarity, are detailed in Arquès and Michel (1996, 1997).

*1.1. The 3 sets  $\mathcal{X}_0(P)$ ,  $\mathcal{X}_1(P)$  and  $\mathcal{X}_2(P)$ ,  
 $P = \{\text{EUK}, \text{PRO}\}$ , are (maximal) circular codes*

A circular code is a set of words on an alphabet such as any word written on a circle (the next letter after the last letter of the word being the first letter) has at most

\*Corresponding author. Tel.: +33-3-90-24-44-62; fax: +33-3-90-24-44-55.

*E-mail addresses:* [frey@dpt-info.u-strasbg.fr](mailto:frey@dpt-info.u-strasbg.fr) (G. Frey),  
[michel@dpt-info.u-strasbg.fr](mailto:michel@dpt-info.u-strasbg.fr) (C.J. Michel).

Table 1  
Gene function process involved in the codon usage bias

Type of gene function process	Some references
tRNA pool and codon–anticodon binding strength Codon context	<a href="#">Ikemura (1981)</a> ; <a href="#">Ikemura (1985)</a>
Genome signature Global/local G + C content and mutational biases	<a href="#">Yarus and Folley (1984)</a> ; <a href="#">Shpaer (1986)</a> ; <a href="#">Gutman and Hatfield (1989)</a> ; <a href="#">Berg and Silva (1997)</a> <a href="#">Campbell et al. (1999)</a> <a href="#">Jukes and Bhushan (1986)</a> ; <a href="#">Sharp and Matassi (1994)</a> ; <a href="#">Sueoka (1992)</a>
Gene expression mRNA stability and transcription rate	<a href="#">Grantham et al. (1981)</a> <a href="#">Andersson and Kurland (1990)</a>
Effects of DNA polymerase and repair systems, methylation, CpG islands	<a href="#">Hanai and Wada (1988)</a>
Tertiary structure of the protein coded by the gene	<a href="#">Gu et al. (2002)</a>

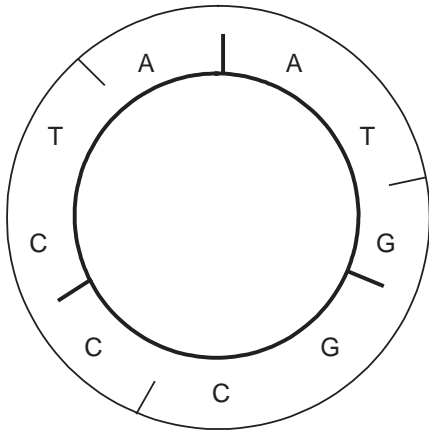


Fig. 1. The set  $\mathcal{X} = \{\text{AAT, ATG, CCT, CTA, GCC, GGC}\}$  is not a circular code as the word  $w = \text{ATGGCCCTA}$ , written on a circle, can be factorized into words of  $\mathcal{X}$  according to 2 different ways: ATG, GCC, CTA (thick line) and AAT, GGC, CCT) (thin line).

one decomposition into words of the circular code. As an example, let the set  $\mathcal{X}$  be composed of the 6 following words:  $\mathcal{X} = \{\text{AAT, ATG, CCT, CTA, GCC, GGC}\}$  and the word  $w$ , be a series of the 9 following letters:  $w = \text{ATGGCCCTA}$ . The word  $w$ , written on a circle, can be factorized into words of  $\mathcal{X}$  according to 2 different ways: ATG, GCC, CTA and AAT, GGC, CCT (Fig. 1). Therefore,  $\mathcal{X}$  is not a circular code. In contrast, if the set  $\mathcal{X}'$  obtained by replacing the last word GGC of  $\mathcal{X}$  by GTC is considered, i.e.  $\mathcal{X}' = \{\text{AAT, ATG, CCT, CTA, GCC, GTC}\}$ , then it does not exist such an ambiguous word with  $\mathcal{X}'$ . Then,  $\mathcal{X}'$  is a circular code. The construction frame of a word generated by a concatenation of the words of the circular code can be retrieved after the reading, anywhere in the word, of a certain number of nucleotides depending on the code. This

series of nucleotides is called the window of the circular code.

For a set of words with 3 letters (trinucleotides) on a 4-letter alphabet, a circular code contains at most 20 words and is maximal. Indeed, a necessary, but not sufficient, condition for a code to be circular is the absence of 2 permuted words in the code (Section 3.3 for details).

### 1.2. The code $\mathcal{X}_0(P)$ is a $\mathcal{C}^3$ code

Recall the definition of the trinucleotide (left circular) permutation: The (left circular) permutation  $\mathcal{P}$  of the trinucleotide  $w_0 = l_0l_1l_2$ ,  $l_0l_1l_2 \in \{\text{A, C, G, T}\}$ , is the permuted trinucleotide  $\mathcal{P}(w_0) = w_1 = l_1l_2l_0$ , e.g.  $\mathcal{P}(\text{AAC}) = \text{ACA}$ , and  $\mathcal{P}(\mathcal{P}(w_0)) = \mathcal{P}(w_1) = w_2$ , e.g.  $\mathcal{P}(\mathcal{P}(\text{AAC})) = \text{CAA}$ .

This definition is naturally extended to the trinucleotide set permutation: The permutation  $\mathcal{P}$  of a set of trinucleotides is the permuted trinucleotide set obtained by the permutation  $\mathcal{P}$  of all its trinucleotides.

Each circular code can be deduced from the permutation of another circular code. The code  $\mathcal{X}_1(P)$  associated with the frame defined as the (left) permutation of the reading frame, i.e. the frame 1, can be obtained by the permutation of the code  $\mathcal{X}_0(P)$ . Similarly, the code  $\mathcal{X}_2(P)$  can be deduced from the permutation of the code  $\mathcal{X}_1(P)$  and the code  $\mathcal{X}_0(P)$ , from the code  $\mathcal{X}_2(P)$ . As the code  $\mathcal{X}_0(P)$  is coding for the reading frame in genes, it is considered as the main code and then called  $\mathcal{C}^3$  code (maximal circular code with 2 permuted maximal circular codes).

*Remark:* The property that  $\mathcal{X}_0(P)$  is a circular code, does not necessarily imply that  $\mathcal{X}_1(P)$  and  $\mathcal{X}_2(P)$  are also circular codes.

Another  $\mathcal{C}^3$  code without the complementary property has been found in the mitochondrial genes ([Arquès and Michel, 1997](#)).

The main biological property directly related to a circular code is the ability to retrieve the reading frame which might be involved in the transcription and the translation apparatus ([Arquès and Michel, 1996, 1997](#)).

### 1.3. The developed method

The method FTF, which was used for determining the  $\mathcal{C}^3$  codes in the gene populations of eukaryotes, prokaryotes and mitochondria, presents exceptions for the assignments of preferential frames to some trinucleotides. These exceptions are usually observed with borderline cases existing when the occurrence frequencies of a trinucleotide are similar in the 3 frames. They have been treated “by hand” ([Arquès and Michel, 1996, 1997](#) for details).

The method proposed here, called frame permuted trinucleotide frequency (FPTF), generalizes and automates the statistical approach FTF and increases the detection sensibility of circular codes in genes. As 2 permuted trinucleotides cannot simultaneously belong to a circular code, the 60 trinucleotides (without AAA, CCC, GGG and TTT) are gathered in 20 classes of 3 trinucleotides invariant by permutation. The method FPTF considers both the preferential frame of a trinucleotide by comparing its occurrence frequencies in the 3 frames and the preferential permuted trinucleotide in a frame by comparing the occurrence frequencies of the 3 permuted trinucleotides in this frame. A statistical function, taking into account these 2 parameters, allows to associate each trinucleotide with a unique frame.

Otherwise, in contrast to the method FTF which has been used with gene populations, i.e. genes from different genomes, the method FPTF is applied here to the genomes, precisely to their complete sets of genes in both strands. So far, no circular code has been searched in the archaeal genomes, members of the third domain of life.

#### 1.4. The archaeal genomes

Archaea have features that are either unique, typically prokaryotic or typically eukaryotic (e.g. [Bernander, 2000](#); [Woese, 2000](#); [Forterre, 2001](#)). They possess a prokaryotic mode of cellular organization, e.g. no nuclear envelope, circular DNA molecules organized similar to those of prokaryotes, etc. On the other hand, they present many eukaryotic similarities in their replication, transcription and translation processes, e.g. introns in tRNA genes, protein synthesis initiation with unformylated methionine, etc.

Most of them are extremophiles, i.e. adapted to live in extreme environment such as high or low temperature, acid or salinity water, etc. Archaea are often viewed as predominant over prokaryotes in all hostile environments. Their variability in their morphology, metabolic pathway and genetic information makes the construction of their phylogenetic tree difficult.

The actual hypothesis divides the archaeal domain into 3 main phyla:

- The Euryarchaeota (EA) containing the methanogens (strict anaerobes and methane producers), the halophiles (strict aerobes living in high-salt environments), some thermoacidophiles (aerobes living in hot and acidic environments) and some hyperthermophiles (living in hot environments).
- The Crenarchaeota (CA) principally composed of hyperthermophiles.

- The Korarchaeota deduced from small subunit ribosomal DNA sequences and for which no organism has been sequenced yet.

Circular codes are searched in the 16 archaeal complete genomes sequenced at the time of writing this article which comprise 12 Euryarchaeota and 4 Crenarchaeota. Their complete sets of genes are analysed in both strands. The method FPTF identifies several circular codes in archaea. Several properties of these codes are described. In particular, a new definition of a parameter allows to measure some probabilities of words generated by the circular codes. Some biological consequences in archaea are presented in Discussion with respect to the 3 2-letter genetic alphabets (purine/pyrimidine, amino/ceto, strong/weak interaction) and the genetic code.

## 2. Method

The methods developed for identifying circular codes in genes consider the permuted trinucleotides. On the genetic alphabet {A,C,G,T}, there are 60 trinucleotides  $w \in \mathcal{T} = \{\text{AAA}, \dots, \text{TTT}\} - \{\text{AAA}, \text{CCC}, \text{GGG}, \text{TTT}\}$  which can be gathered in 20 sets of 3 permuted trinucleotides. In genes, a trinucleotide  $w$  can be read in 3 frames  $p \in \{0, 1, 2\}$  and then noted  $w^p$  with  $p = 0$ : reading frame established by the start trinucleotide ATG, and  $p = 1$  (resp.  $p = 2$ ): reading frame shifted by 1 (resp. 2) nucleotide in the 5'–3' direction. Therefore, there are  $60 \times 3 = 180$  trinucleotides  $w^p$ .

The first method developed, called FTF ([Arquès and Michel, 1996](#)), computes the 180 occurrence frequencies  $o(w^p)$  of the trinucleotides  $w^p$  in genes and assigns a preferential frame  $\text{phasepref}$  to the trinucleotide  $w$ , noted  $w^{\text{phasepref}}$ , such as the occurrence frequency of  $w$  in this frame is the highest one among the 3 possible frames  $p$

$$\text{phase}(w) = k \text{ such as } o(w^k) = \text{MAX}_{p=0}^2 \{o(w^p)\},$$

$$w^{\text{phasepref}} = w^{\text{phase}(w)}. \quad (1)$$

However, as the method FTF is based on the simple parameter  $o(w^p)$ , some trinucleotides are associated with 2 frames  $p$  and  $p'$ , i.e.  $|o(w^p) - o(w^{p'})| \leq \varepsilon$  where  $\varepsilon = 0.2\%$  with the circular code of eukaryotes/prokaryotes ([Arquès and Michel, 1996](#)) and  $\varepsilon = 0.4\%$  with the circular code of mitochondria ([Arquès and Michel, 1997](#)). Furthermore, in these borderline cases, a different preferential frame cannot be assigned to each permuted trinucleotide in order to verify the hypotheses of a potential  $\mathcal{C}^3$  code. These exceptions with the circular codes of eukaryotes/prokaryotes and mitochondria have been treated “by hand” ([Arquès and Michel, 1996, 1997](#)).

In order to have a general and automatic approach for the frame assignment to a trinucleotide, we propose here a quantitative method which considers simultaneously the 3 permuted trinucleotides for assigning a different preferential frame to each permuted trinucleotide. This method called FPTF will identify several circular codes in archaeal genes.

As there are 3 permuted trinucleotides  $w_i$ ,  $i \in \{0, 1, 2\}$ , in 3 different frames  $p \in \{0, 1, 2\}$ , there are 9 possible trinucleotides  $w_i^p$  and therefore  $\binom{9}{3} = 84$  possible sets  $\mathcal{S}^j$ ,  $j \in \{1, \dots, 84\}$ , of 3 trinucleotides. The 84 sets  $\mathcal{S}^j$  are noted  $\mathcal{S}84$  and defined as follows:

$$\begin{aligned} \mathcal{S}84 &= \{ \{w_0^0, w_0^1, w_0^2\}, \dots, \{w_0^0, w_0^2, w_0^1\}, \dots, \{w_0^0, w_1^0, w_1^1\}, \dots, \\ &\quad \{w_0^0, w_1^1, w_1^2\}, \dots, \{w_0^0, w_1^2, w_0^1\}, \dots, \{w_0^0, w_2^0, w_2^1\}, \dots, \\ &\quad \{w_0^1, w_0^2, w_1^0\}, \dots, \{w_0^2, w_0^1, w_1^1\}, \dots, \{w_1^0, w_1^1, w_1^2\}, \dots, \\ &\quad \{w_1^1, w_1^2, w_2^0\}, \dots, \{w_1^2, w_2^0, w_2^1\}, \dots, \{w_2^0, w_2^1, w_2^2\} \} \\ &= \{ \mathcal{S}^1, \dots, \mathcal{S}^{84} \}. \end{aligned}$$

By convention, the trinucleotides  $w_i^p$ ,  $w_{i'}^{p'}$  and  $w_{i''}^{p''}$ ,  $i, i', i'', p, p', p'' \in \{0, 1, 2\}$ , in a set  $\mathcal{S}^j = \{w_i^p, w_{i'}^{p'}, w_{i''}^{p''}\}$  are in the lexicographical order.

Among these 84 sets  $\mathcal{S}^j$ , the 3 sets of left permuted trinucleotides are:

- the trinucleotide  $w_0$  in frame 0, i.e.  $w_0^0$ , the permuted trinucleotide  $w_1$  in frame 1, i.e.  $w_1^1$ , and the permuted trinucleotide  $w_2$  in frame 2, i.e.  $w_2^2$ , leading to the set  $\mathcal{S}^{22} = \{w_0^0, w_1^1, w_2^2\}$ ;
- the trinucleotide  $w_0$  in frame 1, i.e.  $w_0^1$ , the permuted trinucleotide  $w_1$  in frame 2, i.e.  $w_1^2$ , and the permuted trinucleotide  $w_2$  in frame 0, i.e.  $w_2^0$ , leading to the set  $\mathcal{S}^{44} = \{w_0^1, w_1^2, w_2^0\}$ ;
- the trinucleotide  $w_0$  in frame 2, i.e.  $w_0^2$ , the permuted trinucleotide  $w_1$  in frame 0, i.e.  $w_1^0$ , and the permuted trinucleotide  $w_2$  in frame 1, i.e.  $w_2^1$ , leading to the set  $\mathcal{S}^{53} = \{w_0^2, w_1^0, w_2^1\}$ .

These 3 sets  $\mathcal{S}^{22}$ ,  $\mathcal{S}^{44}$  and  $\mathcal{S}^{53}$  associate each trinucleotide with a preferential frame and each frame with a permuted trinucleotide by respecting the definition of the trinucleotide permutation (Section 1). Therefore, in these 3 sets, a relation between a trinucleotide and its frame, allows to deduce (by permutation) the 2 other relations between the permuted trinucleotides and their frames. A statistical analysis will show that these 3 sets occur with the highest value in all the archaeal genomes (see below the Remark concerning the sets  $\mathcal{S}^j$ ).

A statistical function is developed for quantifying these 84 sets  $\mathcal{S}^j$  and determining the preferential set among  $\mathcal{S}^{22}$ ,  $\mathcal{S}^{44}$  and  $\mathcal{S}^{53}$ . Let  $o(w_i^p)$ ,  $i, p \in \{0, 1, 2\}$ , be

the occurrence probability of a trinucleotide  $w_i^p$  in genes. In order to weight the probability  $o(w_i^p)$  of  $w_i^p$  by its probabilities in the 3 frames, the following probability function  $P$  allows to associate a preferential frame with a trinucleotide

$$P(w_i^p) = \frac{o(w_i^p)}{\sum_{p=0}^2 o(w_i^p)}. \quad (2)$$

Similarly, in order to weight the probability  $o(w_i^p)$  of  $w_i^p$  by its permuted trinucleotide probabilities, the following probability function  $Q$  allows to associate a preferential permuted trinucleotide with a frame

$$Q(w_i^p) = \frac{o(w_i^p)}{\sum_{i=0}^2 o(w_i^p)}. \quad (3)$$

*Remark:* In a genome with hundreds of genes, the denominators  $DEN(P)$  and  $DEN(Q)$  are different from 0. Indeed, there is no stop codon  $\tilde{w}_0 \in \{TAA, TAG, TGA\}$  in frame 0 in genes. Therefore,  $o(\tilde{w}_0^0) = 0$ . As  $\tilde{w}_0$  occurs in frames 1 and 2, then  $DEN(P) = \sum_{p=0}^2 o(\tilde{w}_0^p) = \sum_{p=1}^2 o(\tilde{w}_0^p) > 0$ . On the other hand, as  $\mathcal{P}(\tilde{w}_0)$  and  $\mathcal{P}(\mathcal{P}(\tilde{w}_0))$  occur in frame 0, then  $DEN(Q) = \sum_{i=0}^2 o(\tilde{w}_i^0) = \sum_{i=1}^2 o(\tilde{w}_i^0) > 0$ . These 2 inequalities could obviously not be verified with one gene of short length.

The statistical function for selecting the set  $\mathcal{S}$  of 3 trinucleotides  $\{w_{i_0}^{p_0}, w_{i_1}^{p_1}, w_{i_2}^{p_2}\}$ , which has the strongest weight compared to the frame and the permuted trinucleotide, is defined as follows:

$$\begin{aligned} F(\mathcal{S}^{N(i_0, p_0, i_1, p_1, i_2, p_2)}) &= F(w_{i_0}^{p_0}, w_{i_1}^{p_1}, w_{i_2}^{p_2}) \\ &= \frac{1}{3} \sum_{j=0}^2 \frac{1}{2} (P(w_{i_j}^{p_j}) + Q(w_{i_j}^{p_j})) \\ &= \frac{1}{6} \sum_{j=0}^2 (P(w_{i_j}^{p_j}) + Q(w_{i_j}^{p_j})) \end{aligned} \quad (4)$$

with

$$\begin{aligned} N(i_0, p_0, i_1, p_1, i_2, p_2) &= \frac{1}{6} I_0 (I_0^2 - 24I_0 + 191) \\ &\quad - \frac{1}{2} I_1 (I_1 - 17) + I_2 - 45 \end{aligned}$$

and  $I_0 = 3i_0 + p_0 + 1$ ,  $I_1 = 3i_1 + p_1 + 1$ ,  $I_2 = 3i_2 + p_2 + 1$ ,  $i_0, p_0, i_1, p_1, i_2, p_2 \in \{0, 1, 2\}$  verifying  $I_0 < I_1 < I_2$ .

**Property 1.** If the 9 occurrence probabilities  $o(w_i^p)$  in a set  $\mathcal{S}^j$  are identical, i.e.  $o(w_i^p) = o(w_{i'}^{p'}) \forall i, i', p, p' \in \{0, 1, 2\}$ , then  $F(\mathcal{S}^j) = 1/3 \forall j \in \{1, \dots, 84\}$  (proof obvious).

The set  $\mathcal{S}^{\max}$ , having the highest value with the function  $F$  among the 84 sets  $\mathcal{S}^j$ , is

$$\mathcal{S}^{\max} = \mathcal{S}^k \quad \text{such as } F(\mathcal{S}^k) = \text{MAX}_{j=1}^{84} \{F(\mathcal{S}^j)\}. \quad (5)$$

**Remark concerning the sets  $\mathcal{S}^j$ :** Let  $\mathcal{S}_k, k \in \{1, \dots, 20\}$ , be a set of 3 permuted trinucleotides in the lexicographical order among the 20 possible ones, i.e.  $\mathcal{S}_1 = \{AAC, ACA, CAA\}, \dots, \mathcal{S}_{20} = \{GTT, TTG, TGT\}$ . The 9 trinucleotides  $w_i^p$  of a given set  $\mathcal{S}_k$  can be partitioned into

$$\frac{1}{3!} \binom{9}{3} \binom{6}{3} = 280$$

different sets of 3 sets of 3 trinucleotides. Such a partition is noted

$$\begin{aligned} E_k^{(j_0, j_1, j_2)} &= \left\{ \left\{ w_{i_0}^{p_0}, w_{i_1}^{p_1}, w_{i_2}^{p_2} \right\}, \left\{ w_{i_0}^{p'_0}, w_{i_1}^{p'_1}, w_{i_2}^{p'_2} \right\}, \right. \\ &\quad \left. \left\{ w_{i_0}^{p''_0}, w_{i_1}^{p''_1}, w_{i_2}^{p''_2} \right\} \right\} \\ &= \left\{ \mathcal{S}_k^{j_0}, \mathcal{S}_k^{j_1}, \mathcal{S}_k^{j_2} \right\}, \end{aligned}$$

$k \in \{1, \dots, 20\}$ ,  $j_0, j_1, j_2 \in \{1, \dots, 84\}$  and  $j_0 \neq j_1 \neq j_2$ , by associating each set of 3 trinucleotides with a set  $\mathcal{S}_k^j$  such as  $\mathcal{S}_k^{j_l} \cap \mathcal{S}_k^{j_{l'}} = \emptyset$ ,  $l, l' \in \{0, 1, 2\}$  and  $l \neq l'$ . The particular partition studied here is

$$\begin{aligned} E_k^{(22, 44, 53)} &= \left\{ \mathcal{S}_k^{22}, \mathcal{S}_k^{44}, \mathcal{S}_k^{53} \right\} \\ &= \left\{ \left\{ w_0^0, w_1^1, w_2^2 \right\}, \left\{ w_0^1, w_1^2, w_2^0 \right\}, \left\{ w_0^2, w_1^0, w_2^1 \right\} \right\}. \end{aligned}$$

For a given set  $\mathcal{S}_k$ , the partition  $E_k^{(j_0, j_1, j_2)}$  is evaluated by the function  $F$  as follows:

$$F\left(E_k^{(j_0, j_1, j_2)}\right) = \text{MAX}_{l=0}^2 \left\{ F\left(\mathcal{S}_k^{j_l}\right) \right\} \quad (6)$$

and the partition  $E^{(j_0, j_1, j_2)}$  for a genome  $\mathcal{G}$  with 20 sets  $\mathcal{S}_k$ , each corresponding to a trinucleotide permutation class, by the mean  $M$

$$M\left(E^{(j_0, j_1, j_2)}, \mathcal{G}\right) = \frac{1}{20} \sum_{k=1}^{20} F\left(E_k^{(j_0, j_1, j_2)}\right). \quad (7)$$

The computation of  $M\left(E^{(j_0, j_1, j_2)}, \mathcal{G}\right)$  in each genome  $\mathcal{G}$  shows that  $M\left(E^{(22, 44, 53)}, \mathcal{G}\right)$  has the highest value among the 280 partitions  $E^{(j_0, j_1, j_2)}$  in the 16 archaeal genomes (data not shown). This strong statistical result for the partition  $E^{(22, 44, 53)}$  explains the choice of the set  $\mathcal{S}^j = \left\{ w_0^p, w_1^{p'}, w_2^{p''} \right\}$ ,  $p, p', p'' \in \{0, 1, 2\}$  and  $p \neq p' \neq p''$ , for determining the set  $\mathcal{S}^{\text{max}}$ . Finally, the 3 sets  $\mathcal{S}^j$  in the partition  $E^{(22, 44, 53)}$  associate each permuted trinucleotide with a different frame such as the frame of a trinucleotide can be deduced from the frame of a permuted trinucleotide.

Finally, the preferential set  $\mathcal{S}^{\text{pref}}$  among  $\mathcal{S}^{22} = \left\{ w_0^0, w_1^1, w_2^2 \right\}$ ,  $\mathcal{S}^{44} = \left\{ w_0^1, w_1^2, w_2^0 \right\}$  and  $\mathcal{S}^{53} = \left\{ w_0^2, w_1^0, w_2^1 \right\}$  in the partition  $E^{(22, 44, 53)}$ , is identified by comparing the sets  $\mathcal{S}^{22}$ ,  $\mathcal{S}^{44}$  and  $\mathcal{S}^{53}$  with the set  $\mathcal{S}^{\text{max}}$

$$\mathcal{S}^{\text{pref}} = \mathcal{S}^k$$

such as

$$F(\mathcal{S}^{\text{max}}) - F(\mathcal{S}^k) = \text{MIN}_{j=22, 44, 53} \left\{ F(\mathcal{S}^{\text{max}}) - F(\mathcal{S}^j) \right\}. \quad (8)$$

$\mathcal{S}^{\text{pref}}$  is the set with the highest value among the sets  $\mathcal{S}^{22}$ ,  $\mathcal{S}^{44}$  and  $\mathcal{S}^{53}$ . In most cases (see below),  $\mathcal{S}^{\text{max}} = \mathcal{S}^{\text{pref}}$  (i.e.  $F(\mathcal{S}^{\text{max}}) - F(\mathcal{S}^{\text{pref}}) = 0$ ).

Sixteen archaeal genomes  $\mathcal{G}$  are studied:

- *Aeropyrum* (pernix) with 2694 genes containing 1916 kb and noted AP (Crenarchaeota CA),
- *Archeoglobus* (fulgidus) with 2407 genes containing 1989 kb and noted AG (Euryarchaeota EA),
- *Halobacterium* (sp.NCR-1) with 2058 genes containing 1761 kb and noted HB (EA),
- *Methanococcus* (jannashii) with 1709 genes containing 1444 kb and noted MC (EA),
- *Methanopyrus* (kandleri) with 1678 genes containing 1492 kb and noted MP (EA),
- *Methanosarcina acetivorans* with 4440 genes containing 4162 kb and noted MSA (EA),
- *Methanosarcina mazei* with 3371 genes containing 3065 kb and noted MSM (EA),
- *Methanothermobacter* (thermautotrophicus) with 1868 genes containing 1575 kb and noted MT (EA),
- *Pyrobaculum* (aerophilum) with 2605 genes containing 1968 kb and noted PB (CA),
- *Pyrococcus abyssi* with 1762 genes containing 1606 kb and noted PCA (EA),
- *Pyrococcus furiosus* with 2060 genes containing 1740 kb and noted PCF (EA),
- *Pyrococcus horikoshii* with 2058 genes containing 1704 kb and noted PCH (EA),
- *Sulfolobus solfataricus* with 2994 genes containing 2525 kb and noted SLS (CA),
- *Sulfolobus tokodaii* with 2826 genes containing 2276 kb and noted SLT (CA),
- *Thermoplasma acidophilium* with 1478 genes containing 1359 kb and noted TPA (EA),
- *Thermoplasma volcanium* with 1523 genes containing 1353 kb and noted TPV (EA).

In all these genomes, the genes beginning obligatory with a start codon are extracted from both strands.

These large populations allow to have stable frequencies leading to significant statistical results.

### 3. Results

#### 3.1. Identification of 3 subsets of 20 trinucleotides in the 3 frames of the genes in the archaeal genomes

The 180 occurrence frequencies  $o(w_i^p)$  of the 20 sets of 3 permuted trinucleotides  $w_i$  in the 3 frames  $p$  are computed in all the genes of the 16 archaeal genomes (Table 2).

**Remark:** The frequencies of the 3 stop codons TAA, TAG, and TGA in frame 0 are obviously equal to 0 in all the genomes (Table 2).

Table 2

Occurrence frequencies  $o(w_i^p)$  (in % and rounded at 0.1%) of the 20 sets of 3 permuted trinucleotides  $w_i$  in the 3 frames  $p$  of all the genes in the 16 archaeal genomes

	<i>Aeropyrum</i> (AP)			<i>Archaeoglobus</i> (AG)			<i>Halobacterium</i> (HB)			<i>Methanococcus</i> (MC)			<i>Methanopyrus</i> (MP)			<i>Methanosarcina acetivorans</i> (MSA)		
	0	1	2	0	1	2	0	1	2	0	1	2	0	1	2	0	1	2
AAC	1.6	0.5	1.0	2.2	0.5	1.4	2.0	1.0	0.5	1.6	1.0	2.2	1.7	1.1	0.6	2.2	1.6	1.6
ACA	1.1	2.2	0.5	1.0	2.9	0.3	0.4	2.5	0.7	2.0	1.6	0.9	0.4	1.6	0.6	1.6	2.2	0.9
CAA	0.3	1.0	1.7	0.3	2.0	2.6	0.4	0.3	1.7	0.9	3.6	1.4	0.3	0.2	1.9	0.7	2.4	2.3
AAG	2.8	1.0	1.1	4.5	1.8	1.8	1.1	0.9	0.1	3.1	4.6	2.3	3.1	1.8	0.3	2.3	3.6	1.6
AGA	1.1	2.6	2.0	2.2	3.6	2.4	0.1	1.7	0.6	2.8	2.9	6.1	0.6	2.5	1.6	1.3	2.2	3.3
GAA	1.2	0.6	2.0	2.8	2.1	3.5	1.6	0.2	1.8	5.1	2.5	2.0	2.6	0.6	2.9	5.5	1.8	2.5
AAT	0.5	0.3	0.9	1.0	0.7	2.8	0.1	0.3	0.2	3.7	2.3	5.0	0.3	0.3	0.6	2.3	1.4	2.7
ATA	3.2	1.2	0.5	2.3	0.7	0.5	0.1	0.1	0.1	4.6	4.4	1.5	1.0	0.5	0.2	2.1	2.4	0.9
TAA	0.0	1.6	1.1	0.0	2.2	1.2	0.0	0.1	0.1	0.0	5.5	4.8	0.0	0.5	0.8	0.0	2.1	2.0
ACC	1.6	2.0	0.9	1.2	1.4	0.5	3.2	3.2	0.4	0.4	0.5	0.8	1.9	2.4	1.2	1.7	1.7	0.7
CCA	0.9	2.4	1.8	0.8	1.5	0.9	0.4	2.1	2.4	2.3	0.5	0.5	0.5	1.5	1.1	0.6	1.5	1.2
CAC	1.3	0.8	1.7	1.1	0.5	1.3	2.1	0.8	3.6	0.4	0.5	0.4	1.5	0.9	1.8	0.8	0.6	1.2
ACG	1.1	2.4	0.3	1.2	2.9	0.1	3.0	7.6	0.5	0.2	0.7	0.0	2.0	4.7	1.0	0.8	1.9	0.3
CGA	0.3	1.3	3.0	0.1	1.7	3.1	0.5	2.7	10.1	0.0	0.2	0.5	1.1	2.1	6.4	0.3	0.9	2.6
GAC	2.6	0.4	1.3	2.7	0.3	0.9	8.1	0.4	2.5	1.0	0.3	0.4	4.2	1.0	1.8	2.4	0.4	1.2
ACT	0.9	1.1	1.1	0.8	1.7	0.8	0.3	1.3	0.7	1.5	1.2	0.8	0.4	1.3	1.5	1.4	1.2	1.4
CTA	1.8	1.9	1.6	0.5	0.5	1.4	0.1	0.1	1.2	0.9	1.5	0.8	0.9	0.5	0.9	0.5	1.0	1.0
TAC	2.1	1.2	0.8	2.8	0.4	0.6	2.3	0.2	0.2	1.0	0.6	1.0	2.5	1.3	0.6	1.7	0.7	1.5
AGC	2.6	2.1	2.1	1.9	2.3	1.6	1.6	2.3	0.7	0.5	0.8	2.2	0.8	2.7	1.1	1.2	1.3	1.9
GCA	1.4	1.6	1.0	2.1	1.6	1.3	0.8	1.9	1.5	2.3	0.5	0.7	0.6	0.9	1.0	2.6	1.4	1.1
CAG	1.6	1.7	2.8	1.5	1.8	2.4	2.2	0.6	1.2	0.5	2.3	0.4	1.1	0.6	1.2	1.9	2.3	1.9
AGG	4.6	4.1	1.7	3.0	4.7	1.3	0.1	3.1	0.2	1.0	2.2	2.1	1.7	4.8	1.2	1.4	2.2	2.2
GGA	1.2	2.4	3.8	2.5	1.8	4.4	0.5	1.1	4.2	3.6	1.0	2.2	1.8	2.0	6.4	2.7	1.6	3.4
GAG	5.5	1.3	3.3	6.1	1.9	2.9	5.1	0.3	0.7	3.5	2.5	1.2	7.4	1.3	1.5	2.5	1.6	1.0
AGT	0.7	1.0	1.5	0.5	1.5	1.6	0.3	1.4	0.4	1.1	1.2	2.6	0.4	1.6	1.3	1.0	1.0	1.6
GTA	1.5	0.9	0.9	1.1	0.4	0.8	0.2	0.1	1.1	1.5	1.0	0.5	1.5	1.0	1.4	1.8	0.9	0.7
TAG	0.0	3.2	1.8	0.0	1.3	0.5	0.0	0.1	0.1	0.0	3.7	1.6	0.0	1.6	0.4	0.0	1.6	0.5
ATC	1.3	0.8	0.6	2.1	0.9	0.5	3.3	0.7	0.3	1.1	1.0	1.1	3.4	1.1	0.5	2.5	1.3	0.9
TCA	0.7	2.4	0.6	0.9	3.0	0.8	0.2	2.6	0.3	1.5	0.8	0.9	0.4	2.6	0.7	1.2	2.4	1.1
CAT	0.6	0.5	2.2	0.4	0.6	2.8	0.2	0.1	2.5	1.0	1.5	1.2	0.4	0.2	1.9	0.9	0.8	2.1
ATG	2.0	1.3	0.3	2.5	2.1	0.4	1.7	0.4	0.1	2.2	4.7	0.7	1.7	0.7	0.3	2.3	3.0	0.5
TGA	0.0	2.1	1.8	0.0	3.2	3.9	0.0	1.7	0.9	0.0	1.7	5.5	0.0	2.9	1.4	0.0	2.3	4.1
GAT	1.3	0.3	2.0	2.2	0.6	2.9	0.9	0.2	2.2	4.5	1.1	2.2	1.6	0.5	3.3	2.9	0.5	2.3
ATT	0.8	0.3	0.5	2.8	0.8	0.9	0.2	0.2	0.1	4.9	2.6	4.1	0.5	0.4	0.4	2.7	1.6	1.3
TTA	0.5	1.6	0.4	0.6	1.8	1.0	0.1	0.1	0.1	5.3	4.5	1.7	0.5	0.5	0.4	0.9	2.3	1.2
TAT	1.3	0.9	1.9	0.8	0.5	1.0	0.2	0.1	0.1	3.4	2.3	4.1	0.3	0.4	0.7	2.1	0.9	2.4
CCG	1.3	2.5	0.9	1.1	1.7	0.2	2.4	7.3	2.9	0.1	0.2	0.0	2.6	4.4	2.2	1.0	2.0	0.8
CGC	0.7	1.2	2.3	0.2	0.8	1.5	3.3	3.3	7.6	0.0	0.0	0.3	1.6	1.7	3.3	0.5	0.4	0.9
GCC	3.7	2.1	2.2	2.0	0.6	1.2	5.8	2.1	2.0	0.6	0.2	0.5	3.4	1.0	1.4	1.8	0.8	0.8
CCT	1.8	1.9	4.2	0.8	1.3	2.0	0.1	0.8	2.6	0.9	0.5	0.3	0.6	1.0	2.9	1.4	1.1	2.6
CTC	3.9	1.4	2.1	2.5	1.0	1.2	4.0	0.2	1.3	0.3	0.8	0.6	3.5	0.7	1.4	1.8	1.3	1.2
TCC	1.6	2.5	1.1	0.9	1.1	1.3	1.2	2.7	0.5	0.3	0.5	1.4	1.4	3.0	1.1	1.5	1.8	1.4
CGG	0.7	2.2	2.3	0.1	1.8	2.1	2.1	5.2	5.5	0.0	0.2	0.3	2.2	4.1	3.5	0.6	1.6	2.3
GGC	3.3	2.2	3.7	1.8	1.1	2.3	5.0	1.6	4.1	0.4	0.3	1.1	1.9	1.2	3.3	1.5	1.0	1.5
GCG	2.0	2.1	0.6	1.7	1.4	0.2	6.0	5.5	2.5	0.2	0.2	0.0	3.4	2.6	1.9	0.7	1.0	0.4
CGT	0.4	0.8	2.2	0.1	0.4	2.2	0.4	2.1	5.6	0.0	0.2	0.4	1.1	1.4	4.7	0.4	0.3	1.2
GTC	2.3	0.5	1.0	1.5	0.3	0.5	5.0	0.4	1.6	0.5	0.2	0.2	3.8	1.1	1.1	1.6	0.5	0.6
TCG	1.1	2.8	0.3	0.7	3.1	0.1	1.9	8.5	0.5	0.1	0.4	0.0	1.4	6.3	1.0	0.7	2.1	0.3
CTG	2.8	2.2	1.0	2.6	2.1	0.9	3.2	0.7	0.9	0.2	2.2	0.4	3.5	0.6	0.9	2.4	2.3	1.1
TGC	0.6	1.6	1.4	0.9	1.7	2.5	0.4	2.6	0.8	0.4	0.6	2.0	0.9	1.9	0.6	0.7	1.5	2.7
GCT	2.5	1.3	3.3	2.1	1.1	3.2	0.5	0.9	3.8	2.5	0.5	0.6	0.9	0.7	3.3	1.8	0.8	1.9
CTT	1.7	0.6	1.6	2.4	0.7	2.2	0.3	0.1	1.0	0.9	1.3	1.5	0.6	0.4	1.4	3.0	1.3	1.4
TTC	2.1	1.0	0.5	2.8	2.0	0.9	2.9	0.5	0.2	0.9	1.1	1.0	2.6	0.9	0.5	2.0	1.7	2.0
TCT	0.9	1.9	1.7	0.6	1.6	1.9	0.1	1.5	0.5	1.1	1.1	0.7	0.3	1.5	0.9	1.2	1.6	1.9
GGT	1.5	1.1	2.8	1.3	0.7	2.6	0.7	1.0	3.1	1.3	0.4	1.1	2.2	1.3	3.8	1.2	0.7	1.7
GTG	2.9	1.2	0.7	2.2	1.1	0.7	3.8	1.1	0.8	0.6	1.2	0.3	4.3	1.7	1.1	1.4	1.1	0.3
TGG	1.3	3.7	1.4	1.0	2.8	1.7	1.1	3.5	0.6	0.7	1.8	3.2	1.2	4.5	0.7	1.1	2.4	1.4
GTT	2.1	0.3	1.3	3.9	0.4	1.4	0.7	0.1	2.6	4.4	0.9	1.7	1.0	0.4	2.3	2.1	0.8	1.1
TTG	0.7	2.1	0.3	1.1	5.1	0.3	0.7	0.7	0.1	1.9	5.3	0.7	1.1	0.5	0.3	0.8	4.2	0.3
TGT	0.3	1.1	2.1	0.3	0.7	2.2	0.3	1.6	0.5	0.9	0.8	2.8	0.5	1.4	0.8	0.6	0.8	2.3

Table 2 (continued)

	<i>Methanosarcina mazei</i> (MSM)			<i>Methanothermobacter</i> (MT)			<i>Pyrobaculum</i> (PB)			<i>Pyrococcus abyssi</i> (PCA)			<i>Pyrococcus furiosus</i> (PCF)			<i>Pyrococcus horikoshii</i> (PCH)		
	0	1	2	0	1	2	0	1	2	0	1	2	0	1	2	0	1	2
AAC	2.0	1.5	1.5	2.2	1.5	1.4	1.6	0.8	0.9	2.0	0.8	1.8	1.7	1.0	1.9	1.7	1.0	2.0
ACA	1.7	2.1	0.9	2.1	2.9	0.8	1.1	2.0	0.9	0.8	2.3	0.4	1.6	2.2	0.6	1.1	1.7	0.5
CAA	0.5	2.5	2.0	0.1	2.2	1.9	1.0	1.4	1.9	0.7	2.2	2.2	1.0	3.1	2.0	0.8	2.6	1.7
AAG	2.2	3.9	1.7	2.8	0.9	1.9	2.8	2.0	1.2	5.4	2.9	2.2	4.3	4.0	2.5	4.6	3.2	2.4
AGA	1.4	2.2	3.8	1.2	3.3	2.3	1.9	2.3	3.2	1.7	4.0	3.5	2.8	3.2	4.8	2.0	3.4	3.6
GAA	5.7	1.8	2.4	3.2	1.2	2.6	2.0	0.9	1.9	3.6	2.1	3.5	4.9	2.7	2.6	4.1	2.2	3.0
AAT	2.4	1.4	3.0	1.1	1.0	2.1	1.0	1.0	2.6	1.3	0.9	3.4	1.8	1.2	4.4	1.9	1.1	3.7
ATA	2.4	2.5	0.9	4.3	1.2	0.7	2.9	1.6	1.1	4.9	1.5	0.7	4.1	1.8	1.1	4.5	2.2	0.9
TAA	0.0	2.4	2.3	0.0	2.3	1.4	0.0	2.7	2.0	0.0	4.0	2.4	0.0	4.2	2.3	0.0	4.3	2.9
ACC	1.4	1.7	0.7	1.9	2.4	1.3	0.8	1.2	0.4	1.1	1.1	0.8	0.8	0.9	0.8	1.1	0.9	0.9
CCA	0.7	1.4	1.2	1.1	2.3	1.7	0.8	1.8	1.1	1.7	1.3	0.9	2.0	1.1	0.9	1.8	1.2	0.8
CAC	0.8	0.7	1.1	1.1	1.3	1.8	1.1	0.3	1.5	0.9	0.5	0.8	0.8	0.7	0.8	0.7	0.6	0.7
ACG	0.6	1.8	0.3	0.5	1.7	0.3	1.2	3.6	0.3	1.2	2.4	0.1	0.6	1.5	0.1	1.0	1.5	0.1
CGA	0.2	0.7	2.1	0.1	0.4	1.8	0.2	1.0	2.7	0.1	1.4	2.6	0.1	0.6	1.5	0.1	1.1	1.7
GAC	2.4	0.4	1.1	2.9	0.7	1.2	2.8	0.4	1.2	1.9	0.4	0.9	1.6	0.5	0.8	1.2	0.4	0.8
ACT	1.4	1.1	1.4	0.5	1.4	2.2	1.4	1.4	1.0	1.1	1.5	1.4	1.5	1.4	1.5	1.4	1.4	1.5
CTA	0.4	1.0	0.8	0.5	0.2	1.4	1.3	1.4	1.2	1.9	1.4	1.1	1.8	1.3	1.0	1.8	1.7	0.9
TAC	1.5	0.7	1.5	2.3	1.2	0.5	2.6	1.0	0.8	2.4	0.9	0.7	1.8	0.8	1.0	1.9	1.0	1.0
AGC	1.3	1.3	2.0	1.1	1.1	1.4	1.2	1.7	1.5	1.4	1.9	2.0	1.1	1.5	2.3	1.2	1.6	2.3
GCA	2.9	1.4	1.1	2.9	1.4	0.7	1.5	1.4	1.0	1.7	0.8	0.8	2.5	0.7	0.7	1.9	0.6	0.7
CAG	2.0	2.7	2.1	1.8	3.4	2.9	1.1	1.8	1.5	0.9	1.6	1.2	0.8	2.3	1.0	0.8	1.7	0.9
AGG	1.7	2.1	2.5	4.1	4.3	1.8	2.7	3.5	1.4	3.7	4.2	1.8	2.1	3.0	2.2	3.0	3.4	1.9
GGA	2.8	1.6	3.1	2.1	2.4	4.6	1.3	1.7	3.3	3.0	2.3	4.0	3.7	1.7	2.9	3.2	2.1	3.5
GAG	2.4	1.7	1.0	5.0	0.8	2.2	5.0	1.6	2.4	5.3	1.6	2.5	4.0	2.5	1.9	4.2	1.8	2.1
AGT	0.9	0.9	1.7	0.8	1.0	1.1	0.4	1.4	2.5	0.7	1.5	2.5	1.0	1.3	3.1	1.0	1.3	2.5
GTA	1.9	0.8	0.6	1.0	1.0	0.6	2.1	0.6	1.1	1.4	0.7	1.0	1.9	0.7	0.8	1.8	0.9	0.8
TAG	0.0	1.7	0.4	0.0	1.6	0.3	0.0	3.2	1.1	0.0	3.7	1.7	0.0	3.6	1.7	0.0	3.6	1.9
ATC	2.2	1.2	0.9	2.2	0.7	1.5	1.1	0.8	0.6	1.5	1.1	0.6	1.4	1.3	0.6	1.7	1.4	0.8
TCA	1.4	2.3	1.1	2.1	4.1	0.6	0.7	1.6	0.6	0.9	2.2	1.0	1.0	1.9	1.1	1.2	1.9	1.1
CAT	0.9	0.8	2.0	0.8	1.3	4.1	0.4	0.5	1.9	0.6	0.8	2.3	0.7	1.0	2.0	0.8	1.0	2.0
ATG	2.4	3.0	0.5	2.9	3.4	0.6	1.8	1.2	0.3	2.3	2.4	0.4	2.2	3.1	0.5	2.3	3.0	0.5
TGA	0.0	2.2	4.5	0.0	2.8	5.4	0.0	2.1	2.1	0.0	2.3	3.5	0.0	2.0	4.1	0.0	2.1	3.9
GAT	2.9	0.5	2.2	3.1	0.7	2.9	1.5	0.5	1.5	2.8	0.7	3.2	2.8	0.9	2.2	3.1	0.8	2.8
ATT	3.0	1.5	1.3	1.2	0.6	1.3	2.3	1.1	1.3	2.1	1.0	1.7	3.3	1.2	2.2	2.6	1.3	2.1
TTA	0.9	2.6	1.3	0.4	1.3	0.5	2.1	2.4	0.9	1.7	2.9	1.1	2.0	3.1	1.1	2.2	3.4	1.2
TAT	2.0	0.9	2.7	1.0	1.1	1.4	1.7	1.4	1.9	1.4	1.1	1.6	2.2	1.2	2.0	2.0	1.4	2.3
CCG	0.8	1.6	0.7	0.6	1.8	0.7	1.4	2.7	0.6	0.7	1.3	0.1	0.3	0.8	0.1	0.5	1.2	0.1
CGC	0.5	0.4	0.9	0.4	0.2	1.2	0.9	1.9	3.1	0.1	0.4	1.0	0.1	0.2	0.6	0.1	0.3	0.7
GCC	1.6	0.8	0.8	2.7	0.9	0.8	3.9	1.2	2.3	1.9	0.5	1.0	1.4	0.4	0.8	1.6	0.4	0.7
CCT	1.7	0.9	2.5	0.9	1.4	4.1	1.0	1.4	1.8	0.9	1.0	2.3	1.1	0.9	1.7	1.2	1.2	2.1
CTC	1.7	1.4	1.1	3.5	0.3	1.3	2.0	0.9	1.3	2.1	0.9	1.0	1.7	1.3	0.9	1.8	1.1	1.2
TCC	1.3	1.8	1.4	1.3	2.6	0.7	0.8	1.3	0.9	0.7	1.5	1.6	0.6	1.3	1.8	1.0	1.5	1.9
CGG	0.4	1.3	2.1	0.5	1.3	1.9	0.6	2.3	2.5	0.1	1.1	1.2	0.1	0.6	0.7	0.1	1.0	0.9
GGC	1.5	1.0	1.4	1.7	0.9	2.4	3.7	1.8	4.2	1.1	0.9	2.0	0.8	0.7	1.6	0.7	0.8	1.5
GCG	0.6	1.0	0.4	0.7	0.9	0.3	2.9	2.7	0.9	1.0	0.9	0.1	0.6	0.5	0.1	0.7	0.6	0.1
CGT	0.4	0.3	1.0	0.5	0.1	1.4	0.3	1.2	2.9	0.1	0.5	2.2	0.1	0.3	1.2	0.1	0.4	1.6
GTC	1.5	0.5	0.6	2.1	0.6	0.6	1.8	0.3	0.9	1.5	0.4	0.3	1.1	0.4	0.3	1.1	0.4	0.3
TCG	0.5	1.8	0.3	0.3	2.0	0.1	0.8	2.3	0.2	0.5	2.4	0.1	0.3	1.3	0.1	0.5	1.7	0.1
CTG	2.3	2.7	1.1	2.4	2.3	1.0	1.9	1.8	0.7	1.1	1.6	0.5	0.8	2.0	0.5	0.9	1.7	0.5
TGC	0.7	1.4	2.9	0.8	1.1	2.3	0.5	1.5	1.2	0.3	0.9	1.6	0.2	0.8	2.1	0.3	0.8	1.8
GCT	2.0	0.8	1.9	1.1	0.7	1.4	1.6	1.0	2.8	2.0	0.7	2.2	2.1	0.6	1.7	2.2	0.7	1.8
CTT	3.5	1.4	1.2	2.5	0.2	1.6	1.5	0.9	2.1	2.1	0.9	2.4	2.4	1.1	2.1	2.6	1.2	2.4
TTC	1.8	1.8	2.1	2.6	0.9	0.8	1.7	0.9	0.6	2.7	1.8	0.9	1.8	2.1	1.1	2.2	1.9	1.3
TCT	1.3	1.4	2.1	0.6	1.8	1.1	1.0	1.5	1.2	0.7	1.7	1.4	0.9	1.5	1.9	1.0	1.7	1.6
GGT	1.3	0.6	1.6	2.5	0.8	2.8	0.5	0.8	2.1	1.5	0.9	2.2	1.2	0.7	1.5	1.5	0.7	1.9
GTG	1.3	1.2	0.3	2.3	2.1	0.4	3.8	0.3	0.9	1.4	1.0	0.4	1.4	1.2	0.4	1.2	1.1	0.4
TGG	0.9	2.3	1.4	0.8	3.5	1.9	1.5	4.0	1.1	1.2	2.2	2.5	1.2	2.2	2.8	1.2	2.0	2.5
GTT	2.3	0.8	1.0	2.4	0.6	0.8	1.9	0.4	2.2	3.9	0.5	1.9	3.6	0.5	1.6	3.6	0.6	1.7
TTG	0.6	4.6	0.2	0.2	4.2	0.1	1.8	2.9	0.4	1.5	3.7	0.4	1.4	4.8	0.4	1.2	4.1	0.5
TGT	0.6	0.7	2.6	0.5	0.6	2.4	0.4	1.7	1.8	0.3	0.8	1.3	0.4	0.8	2.2	0.3	0.7	1.6

Table 2 (continued)

	<i>Sulfolobus solfataricus</i> (SLS)			<i>Sulfolobus tokodaii</i> (SLT)			<i>Thermoplasma acidophilum</i> (TPA)			<i>Thermoplasma volcanium</i> (TPV)		
	0	1	2	0	1	2	0	1	2	0	1	2
AAC	1.7	1.2	1.7	1.4	1.3	1.8	2.4	0.7	1.4	2.2	1.0	1.6
ACA	1.4	1.5	1.1	1.7	1.1	1.1	1.4	3.0	0.7	1.7	2.3	0.9
CAA	1.6	2.3	1.5	1.6	2.5	1.0	0.3	2.2	2.5	0.8	2.6	2.2
AAG	3.8	3.2	2.3	2.8	4.1	2.4	3.9	1.2	1.6	3.5	2.6	2.0
AGA	2.5	2.7	4.1	2.6	2.2	4.9	1.6	3.6	2.4	1.7	2.7	3.7
GAA	3.9	2.0	2.6	4.7	2.1	2.1	2.8	1.7	3.4	4.0	1.8	2.7
AAT	3.3	1.8	4.6	3.5	2.1	4.9	1.9	1.2	2.2	2.6	1.6	3.5
ATA	5.0	3.8	1.5	5.0	4.1	1.6	5.6	1.8	1.4	5.6	2.9	1.5
TAA	0.0	5.1	4.6	0.0	5.4	4.9	0.0	2.7	1.1	0.0	3.9	2.5
ACC	0.7	0.6	1.3	0.4	0.5	1.4	1.4	1.1	0.8	0.8	1.0	0.9
CCA	1.6	0.9	0.5	1.7	0.5	0.4	1.3	2.2	0.9	1.5	1.2	0.8
CAC	0.5	0.8	0.7	0.4	0.9	0.5	0.8	0.8	1.4	0.7	0.9	0.9
ACG	0.6	1.8	0.2	0.4	1.2	0.1	1.4	2.5	0.3	0.9	2.3	0.3
CGA	0.1	0.7	1.3	0.1	0.4	0.9	0.2	1.5	2.7	0.2	1.0	2.1
GAC	1.3	0.5	0.8	1.0	0.5	0.6	2.0	0.3	1.1	1.8	0.3	0.9
ACT	2.0	1.3	1.6	2.2	1.1	1.6	0.7	1.1	1.3	1.4	1.2	1.5
CTA	1.9	2.5	0.8	1.7	2.9	0.6	0.6	0.6	1.3	1.3	1.8	1.0
TAC	1.7	1.6	1.7	1.3	1.6	1.9	2.6	1.2	0.9	2.2	1.3	1.5
AGC	0.7	0.7	2.1	0.6	0.7	2.4	1.6	1.2	1.3	1.3	1.0	1.8
GCA	1.9	0.5	0.6	2.1	0.4	0.5	2.2	1.9	1.3	2.4	1.4	1.0
CAG	0.5	1.9	0.6	0.5	2.2	0.5	1.9	2.0	2.6	1.3	2.5	1.3
AGG	1.8	2.5	2.3	1.2	1.8	2.5	2.6	3.3	1.5	2.0	2.3	2.0
GGA	2.6	1.5	2.5	2.6	1.2	1.7	2.2	1.7	3.8	2.2	1.5	2.6
GAG	3.0	1.7	1.6	2.3	1.8	1.1	3.2	0.9	1.6	2.3	1.2	1.2
AGT	1.7	1.3	3.0	1.7	1.0	3.2	0.5	0.9	1.2	0.8	1.0	2.0
GTA	2.8	1.7	0.8	3.0	1.7	0.6	1.4	0.6	1.0	2.5	0.9	0.9
TAG	0.0	4.7	2.3	0.0	4.9	2.1	0.0	2.3	0.5	0.0	3.3	1.2
ATC	1.1	1.4	1.3	0.9	1.4	1.3	2.2	2.2	2.1	1.5	1.7	1.9
TCA	1.6	1.3	1.2	1.7	1.0	1.3	1.6	3.6	0.9	1.9	2.0	0.9
CAT	0.8	1.4	1.5	1.0	1.6	1.1	0.9	1.5	4.2	0.9	1.5	2.6
ATG	2.1	3.2	0.5	2.1	3.7	0.6	3.1	3.1	0.5	2.7	3.3	0.5
TGA	0.0	2.1	3.7	0.0	1.7	4.2	0.0	3.2	2.9	0.0	2.2	3.5
GAT	3.4	1.0	2.1	3.7	1.0	1.7	3.8	1.2	3.8	3.7	1.0	2.6
ATT	3.4	2.2	3.4	4.0	2.6	3.7	1.3	1.5	1.7	2.2	1.9	2.3
TTA	4.1	3.7	1.8	4.8	4.3	2.1	0.4	1.2	1.0	1.4	2.8	1.4
TAT	3.1	2.5	3.2	3.6	2.5	4.1	2.0	1.8	1.7	2.6	2.2	3.1
CCG	0.4	0.6	0.1	0.4	0.4	0.0	1.3	1.9	0.3	0.7	1.0	0.2
CGC	0.1	0.3	0.7	0.1	0.2	0.4	0.5	1.3	1.4	0.2	0.7	0.9
GCC	0.7	0.2	0.6	0.5	0.2	0.5	2.1	1.0	1.1	1.3	0.7	1.0
CCT	1.2	0.9	0.8	1.4	0.6	0.6	0.8	0.9	1.9	1.2	0.8	1.5
CTC	0.7	1.0	0.9	0.6	1.1	0.8	1.9	0.6	1.0	1.2	0.8	0.9
TCC	0.8	0.7	1.4	0.5	0.6	1.5	1.6	1.2	1.3	1.0	1.0	1.2
CGG	0.1	0.6	0.7	0.0	0.4	0.7	0.3	1.8	2.3	0.2	1.0	1.5
GGC	0.7	0.5	1.1	0.5	0.4	0.8	2.5	1.0	2.0	1.9	0.8	1.4
GCG	0.7	0.5	0.1	0.4	0.4	0.1	1.4	1.6	0.5	0.8	1.1	0.3
CGT	0.2	0.7	1.2	0.1	0.4	0.7	0.3	0.7	2.1	0.3	0.6	1.5
GTC	0.7	0.5	0.6	0.6	0.4	0.4	1.6	0.3	0.8	1.1	0.3	0.7
TCG	0.5	1.0	0.1	0.3	0.7	0.1	1.2	2.5	0.3	0.9	1.7	0.2
CTG	0.6	2.0	0.2	0.4	2.3	0.2	2.8	1.9	0.5	1.2	2.5	0.4
TGC	0.2	0.7	1.7	0.2	0.6	2.0	0.5	1.6	2.2	0.4	1.2	2.3
GCT	2.3	0.4	0.9	2.6	0.4	0.7	1.3	0.7	2.2	1.9	0.6	1.6
CTT	1.5	1.5	1.9	1.9	1.9	1.6	2.0	0.7	1.5	2.5	1.2	1.6
TTC	1.8	1.4	1.6	1.4	1.6	1.9	3.2	1.5	1.6	2.2	1.3	2.0
TCT	1.5	1.3	1.4	1.9	1.2	1.5	1.0	1.6	2.0	1.7	1.2	1.7
GGT	2.2	0.7	1.3	2.4	0.6	0.9	1.5	0.6	2.0	1.7	0.7	1.3
GTG	1.2	1.4	0.3	0.8	1.5	0.2	2.1	1.1	0.3	1.2	1.2	0.3
TGG	1.1	1.8	2.3	1.0	1.2	2.3	0.9	2.8	2.0	0.8	1.9	2.4
GTT	2.8	1.0	1.9	3.0	1.1	1.6	2.1	0.5	1.3	2.4	0.7	1.3
TTG	1.6	3.0	0.6	1.0	3.5	0.6	0.7	3.0	0.1	1.2	3.6	0.3
TGT	0.4	0.9	2.0	0.5	0.8	2.5	0.1	1.1	1.9	0.2	1.0	2.4



Formula (8) using these frequencies  $o(w_i^p)$  identifies 20 preferential sets  $\mathcal{S}^{\text{pref}}$  of permuted trinucleotides in each archaeal genome (Table 3). The preferential sets  $\mathcal{S}^{\text{pref}}$  have 218 out of  $20 \times 16 = 320$ , i.e. 68%, the highest value (first rank  $Rk$ ) with the function  $F$  among the 84 sets  $\mathcal{S}^j$  in the 16 archaeal genomes and have 265 out of 320, i.e. 83%, the first 3 values (Table 3).

The 2 strong statistical results, highest value for the partition  $E^{(22,44,53)}$  (see above) and more than 2 out of 3 the highest value for the set  $\mathcal{S}^{\text{pref}}$ , are a consequence of the distribution of the 60 trinucleotides to the 3 frames which is revealed by the sensibility of the method developed.

These preferential sets  $\mathcal{S}^{\text{pref}}$  lead to 3 subsets  $\mathcal{X}_0(\mathcal{G})$ ,  $\mathcal{X}_1(\mathcal{G})$  and  $\mathcal{X}_2(\mathcal{G})$  of 20 trinucleotides associated with the frames 0, 1 and 2, respectively, in the 16 archaeal genomes  $\mathcal{G}$ . Table 3 indicates directly the subset  $\mathcal{X}_0(\mathcal{G})$ , e.g.

$$\begin{aligned} \mathcal{X}_0(\text{AP}) = \{ & \text{AAC, AAG, ATA, ACC,} \\ & \text{GAC, TAC, AGC, GAG,} \\ & \text{GTA, ATC, ATG, ATT,} \\ & \text{GCC, CTC, GCG, GTC,} \\ & \text{CTG, TTC, GTG, GTT} \} \end{aligned}$$

in the genome  $\mathcal{G} = \text{AP}$ . The 2 subsets  $\mathcal{X}_1(\mathcal{G})$  and  $\mathcal{X}_2(\mathcal{G})$  are deduced from the subset  $\mathcal{X}_0(\mathcal{G})$  by the permutation property (Section 3.2).

The same 3 subsets of 20 trinucleotides are identified with the 2 archaeal genomes  $\mathcal{G} = \{\text{AP, AG}\}$  (Table 3):  $\mathcal{X}_i(\text{AP}) = \mathcal{X}_i(\text{AG})$ ,  $i \in \{0, 1, 2\}$ . Therefore, 15 subsets  $\mathcal{X}_0(\mathcal{G})$  are identified in the archaeal genomes.

### 3.2. Permutation property

$\mathcal{P}(\mathcal{X}_0(\mathcal{G})) = \mathcal{X}_1(\mathcal{G})$  and  $\mathcal{P}(\mathcal{X}_1(\mathcal{G})) = \mathcal{X}_2(\mathcal{G})$  in the 16 archaeal genomes  $\mathcal{G}$  (by construction of the method).  $\mathcal{X}_0(\mathcal{G})$  generates  $\mathcal{X}_1(\mathcal{G})$  by 1 permutation and  $\mathcal{X}_2(\mathcal{G})$  by 2 permutations.

### 3.3. Circular code property

#### 3.3.1. Definition of a circular code

Recall of a few notations:  $\mathcal{A}$  being a finite alphabet,  $\mathcal{A}^*$  denotes the words on  $\mathcal{A}$  of finite length including the empty word of length 0 and  $\mathcal{A}^+$ , the words on  $\mathcal{A}$  of finite length  $\geq 1$ . Let  $w_1 w_2$  be the concatenation of the 2 words  $w_1$  and  $w_2$ .

A subset  $\mathcal{X}$  of  $\mathcal{A}^+$  is a circular code if  $\forall n, m \geq 1$  and  $x_1, x_2, \dots, x_n \in \mathcal{X}$ ,  $y_1, y_2, \dots, y_m \in \mathcal{X}$  and  $r \in \mathcal{A}^*$ ,  $s \in \mathcal{A}^+$ , the equalities  $s x_2 x_3 \dots x_n r = y_1 y_2 \dots y_m$  and  $x_1 = r s$  imply  $n = m$ ,  $r = 1$  and  $x_i = y_i$ ,  $1 \leq i \leq n$  (Béal, 1993; Berstel and Perrin, 1985) (Fig. 2). In other terms, every word on  $\mathcal{A}$  “written on a circle” has at most one decomposition (factorization) on  $\mathcal{X}$ . Therefore, the construction frame of any word generated by a circular code  $\mathcal{X}$  (precisely,

generated by any concatenation of words of a circular code) can be retrieved as there is a unique decomposition on  $\mathcal{X}$ .

In the following,  $\mathcal{X}$  will be a set of words of length 3 as a gene is a concatenation of trinucleotides.

By excluding all the trinucleotides  $w = \text{lll}$  and by gathering the remaining trinucleotides in sets of 3 trinucleotides so that, in each set, the 3 trinucleotides are deduced from each other by permutation, a potential circular code, has at most one trinucleotide per set. This rule applied to the alphabet  $\mathcal{A} = \{\text{A, C, G, T}\}$  leads to 20 sets of 3 permuted trinucleotides and therefore, to  $3^{20} \sim 3.5$  billions potential circular codes.

A complete study of circular codes with trinucleotides on the reduced alphabet  $\mathcal{A} = \{\text{R, Y}\}$  is presented in Arquès and Michel (1997). It allows to introduce simply the different concepts, e.g. the concept of decomposition and some properties of circular codes. It is done “by hand” as there are only  $3^2 = 9$  potential circular codes (2 sets of 3 permuted trinucleotides). Obviously, such a study cannot be realized on the alphabet  $\mathcal{A} = \{\text{A, C, G, T}\}$  as there are 12 964 440 circular codes among  $3^{20} \sim 3.5$  billions potential ones (Arquès and Michel, 1997). The search of circular codes on the alphabets at 3 (and more) letters needs the introduction of a few classical definitions and results in coding theory, in particular the concept of the flower automaton (Béal, 1993; Berstel and Perrin, 1985).

#### 3.3.2. Identification of circular codes in the archaeal genomes

The subsets  $\mathcal{X}_0(\mathcal{G})$ ,  $\mathcal{X}_1(\mathcal{G})$  and  $\mathcal{X}_2(\mathcal{G})$  are maximal (i.e. 20 trinucleotides) circular codes in the 16 archaeal genomes  $\mathcal{G}$ .

A probabilistic model, based on the nucleotide frequencies with a hypothesis of absence of correlation between successive bases on a DNA strand, has been proposed by Koch and Lehmann (1997) for constructing and identifying circular codes. Recently, we have given a mathematical proof that their model can only generate a subset of circular codes which are obtained by the flower automaton method, in particular their model cannot generate the circular code of eukaryotes/prokaryotes (Lacan and Michel, 2001). So far, the flower automaton method is the most simple and complete one for identifying circular codes.

The flower automaton  $\mathcal{F}(\mathcal{X})$  associated with a subset  $\mathcal{X}$  of  $\mathcal{A}^+$  has a particular state (labeled 1 in Fig. 3) and cycles issued from this state 1 and labeled by words of  $\mathcal{X}$  (Fig. 3).

Therefore, to prove that “ $\mathcal{X}(\mathcal{G})$  is a circular code” is equivalent to prove that  $\mathcal{F}(\mathcal{X}(\mathcal{G}))$  does not contain 2 cycles labeled with the same word.

An example with the set  $\mathcal{X}_0(\text{AP}) = \mathcal{X}_0(\text{AG})$  in the archaeal genomes AP and AG gives the main steps (without details) for proving that a set of words is a

Table 3

Identification of 20 preferential sets  $\mathcal{S}^{\text{pref}}$  of permuted trinucleotides in each archaeal genome by giving the trinucleotide  $w^0$  in frame 0 (reading frame) in each set  $\mathcal{S}_k$ . The preferential sets  $\mathcal{S}^{\text{pref}}$  have 218 out of  $20 \times 16 = 320$ , i.e. 68%, the highest value (first rank  $Rk$ ) with the function  $F$  among the 84 sets  $\mathcal{S}^j$  in the 16 archaeal genomes and have 265 out of 320, i.e. 83%, the first 3 values. The subsets  $\mathcal{X}_0(\mathcal{G})$  in frame 0 in the 16 archaeal genomes  $\mathcal{G}$  are directly deduced, e.g.  $\mathcal{X}_0(\text{AP}) = \{\text{AAC, AAG, ATA, ACC, GAC, TAC, AGC, GAG, GTA, ATC, ATG, ATT, GCC, CTC, GCG, GTC, CTG, TTC, GTG, GTT}\}$  in the genome  $\mathcal{G} = \text{AP}$ .

$\mathcal{S}_k$	AP		AG		HB		MC		MP		MSA		MSM		MT	
	$w^0$	$Rk$	$w^0$	$Rk$	$w^0$	$Rk$	$w^0$	$Rk$	$w^0$	$Rk$	$w^0$	$Rk$	$w^0$	$Rk$	$w^0$	$Rk$
(AAC,ACA,CAA)	AAC	1	AAC	1	AAC	1	ACA	1	AAC	1	AAC	1	AAC	6	AAC	1
(AAG,AGA,GAA)	AAG	1	AAG	1	AAG	2	GAA	1	AAG	1	GAA	1	GAA	1	AAG	2
(AAT,ATA,TAA)	ATA	1	ATA	1	AAT	33	ATA	1	ATA	3	ATA	2	ATA	1	ATA	1
(ACC,CCA,CAC)	ACC	2	ACC	2	ACC	4	CCA	1	ACC	7	ACC	2	ACC	7	ACC	4
(ACG,CGA,GAC)	GAC	1	GAC	1	GAC	1	GAC	1	GAC	1	GAC	1	GAC	1	GAC	1
(ACT,CTA,TAC)	TAC	8	TAC	1	TAC	1	ACT	1	TAC	6	TAC	12	ACT	6	TAC	3
(AGC,GCA,CAG)	AGC	2	AGC	19	CAG	2	GCA	1	CAG	2	GCA	1	GCA	1	GCA	2
(AGG,GGA,GAG)	GAG	1	GAG	1	GAG	1	GGA	7	GAG	1	GAG	1	GAG	3	GAG	1
(AGT,GTA,TAG)	GTA	1	GTA	1	TAG	7	GTA	1	GTA	3	GTA	1	GTA	1	GTA	1
(ATC,TCA,CAT)	ATC	1	ATC	1	ATC	1	TCA	2	ATC	1	ATC	1	ATC	1	ATC	1
(ATG,TGA,GAT)	ATG	1	ATG	5	ATG	1	GAT	1	ATG	1	GAT	1	GAT	1	GAT	1
(ATT,TTA,TAT)	ATT	1	ATT	1	TAT	3	ATT	3	ATT	3	ATT	1	ATT	1	ATT	1
(CCG,CGC,GCC)	GCC	1	GCC	1	GCC	1	GCC	2	GCC	1	GCC	1	GCC	1	GCC	1
(CCT,CTC,TCC)	CTC	1	CTC	3	CTC	1	CCT	1	CTC	1	CTC	1	CTC	1	CTC	1
(CGG,GGC,GCG)	GCG	10	GCG	4	GGC	6	GCG	3	GCG	1	GGC	2	GGC	1	GGC	7
(CGT,GTC,TCG)	GTC	1	GTC	1	GTC	1	GTC	1	GTC	1	GTC	1	GTC	1	GTC	1
(CTG,TGC,GCT)	CTG	3	CTG	5	CTG	1	GCT	1	CTG	1	GCT	3	GCT	2	GCT	6
(CTT,TTC,TCT)	TTC	1	TTC	9	TTC	1	TTC	13	TTC	1	CTT	4	CTT	2	TTC	2
(GGT,GTG,TGG)	GTG	1	GTG	1	GTG	1	GGT	1	GTG	1	GTG	1	GTG	1	GTG	1
(GTT,TTG,TGT)	GTT	1	GTT	1	TTG	1	GTT	1	TTG	1	GTT	1	GTT	1	GTT	1

$\mathcal{S}_k$	PB		PCA		PCF		PCH		SLS		SLT		TPA		TPV	
	$w^0$	$Rk$	$w^0$	$Rk$	$w^0$	$Rk$	$w^0$	$Rk$	$w^0$	$Rk$	$w^0$	$Rk$	$w^0$	$Rk$	$w^0$	$Rk$
(AAC,ACA,CAA)	AAC	1	AAC	1	ACA	7	ACA	8	ACA	3	ACA	1	AAC	1	AAC	3
(AAG,AGA,GAA)	GAA	8	AAG	1	GAA	1	GAA	10	GAA	2	GAA	1	AAG	1	GAA	3
(AAT,ATA,TAA)	ATA	1	ATA	1	ATA	1	ATA	1	ATA	1	ATA	2	ATA	1	ATA	1
(ACC,CCA,CAC)	ACC	6	ACC	25	CCA	21	CCA	18	CCA	1	CCA	1	ACC	1	CCA	10
(ACG,CGA,GAC)	GAC	1	GAC	1	GAC	1	GAC	1	GAC	1	GAC	1	GAC	1	GAC	1
(ACT,CTA,TAC)	TAC	1	TAC	7	TAC	22	CTA	36	ACT	1	ACT	1	TAC	2	ACT	16
(AGC,GCA,CAG)	GCA	3	GCA	1	GCA	1	GCA	1	GCA	1	GCA	1	AGC	6	GCA	1
(AGG,GGA,GAG)	GAG	1	GAG	1	GAG	2	GAG	1	GAG	1	GGA	2	GAG	1	GAG	1
(AGT,GTA,TAG)	GTA	1	GTA	1	GTA	1	GTA	1	GTA	1	GTA	1	GTA	1	GTA	1
(ATC,TCA,CAT)	ATC	1	ATC	1	ATC	1	ATC	1	TCA	6	TCA	2	ATC	1	ATC	7
(ATG,TGA,GAT)	ATG	5	GAT	2	GAT	1	GAT	1	GAT	1	GAT	1	ATG	6	GAT	1
(ATT,TTA,TAT)	ATT	1	ATT	1	ATT	1	ATT	1	ATT	9	ATT	5	TAT	19	ATT	2
(CCG,CGC,GCC)	GCC	1	GCC	1	GCC	1	GCC	1	GCC	1	GCC	1	GCC	1	GCC	1
(CCT,CTC,TCC)	CTC	1	CTC	1	CTC	4	CTC	2	CCT	1	CCT	1	CTC	2	CTC	2
(CGG,GGC,GCG)	GCG	7	GCG	2	GCG	1	GCG	1	GCG	2	GGC	8	GGC	2	GGC	1
(CGT,GTC,TCG)	GTC	1	GTC	1	GTC	1	GTC	1	GTC	1	GTC	1	GTC	1	GTC	1
(CTG,TGC,GCT)	CTG	1	GCT	1	GCT	1	GCT	1	GCT	1	GCT	1	CTG	5	GCT	1
(CTT,TTC,TCT)	TTC	1	TTC	1	CTT	1	TTC	6	TTC	7	TCT	1	TTC	7	CTT	18
(GGT,GTG,TGG)	GTG	1	GTG	2	GTG	10	GTG	9	GGT	1	GGT	1	GTG	1	GGT	3
(GTT,TTG,TGT)	GTT	4	GTT	1	GTT	1	GTT	1	GTT	1	GTT	1	GTT	1	GTT	1

circular code. The flower automaton  $\mathcal{F}(\mathcal{X}_0(\text{AP}))$  of  $\mathcal{X}_0(\text{AP})$  is constructed by reading the 20 words of  $\mathcal{X}_0(\text{AP})$  in the way given in Fig. 3. As all words of  $\mathcal{X}_0(\text{AP})$  are trinucleotides, i.e. of length equal to 3, the states of  $\mathcal{F}(\mathcal{X}_0(\text{AP}))$  can be associated with frames: state 1 with frame 0, states 2–5 with frame 1 and states 6–15 with frame 2. The search of 2 cycles labeled with the

same word in  $\mathcal{F}(\mathcal{X}_0(\text{AP}))$  can be realized by comparing all the words starting from the frames 0 (state 1) and 1 (states 2–5) (Fig. 4). The absence of 2 cycles is verified when all the compared words are “non-extensible” (Fig. 4).

The implementation in Java of the flower automaton algorithm (not described here) demonstrates that the 15

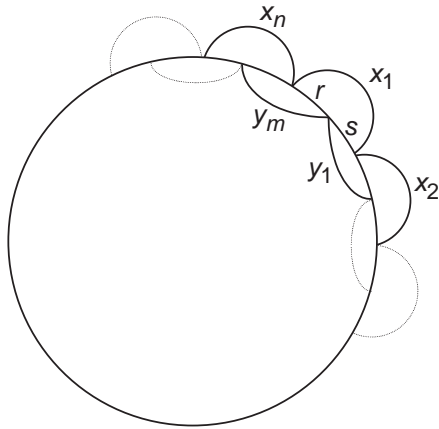


Fig. 2. A representation of the definition of a circular code.

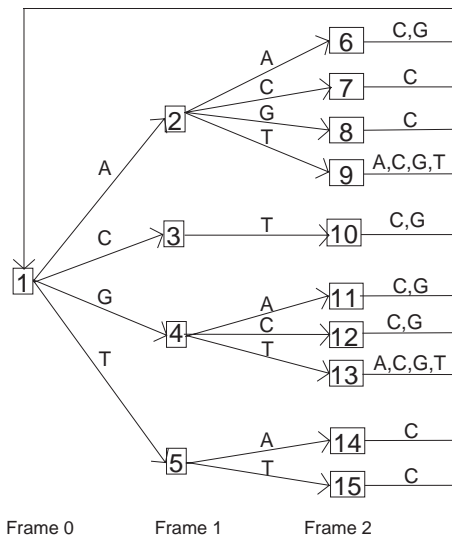


Fig. 3. Flower automaton  $\mathcal{F}(X_0(AP))$  of the archaeal  $\mathcal{C}^3$  code  $X_0(AP)$ .

sets  $X_0(\mathcal{G})$  are (maximal) circular code with 2 permuted (maximal) circular codes  $X_1(\mathcal{G})$  and  $X_2(\mathcal{G})$ . Therefore, 15  $\mathcal{C}^3$  codes  $X_0(\mathcal{G})$  in reading frame are identified in the archaeal genomes.

*Remark:* Nine  $\mathcal{C}^3$  codes  $X_0(\mathcal{G})$  in the 9 genomes  $\mathcal{G} = \{\{AP, AG\}, MP, MSA, MSM, MT, PB, PCA, PCH, SLT\}$  are identified with a strong statistical significance. Indeed, their  $20 \times 9 = 180$  preferential sets  $\mathcal{S}_k^{\text{pref}}$  of permuted trinucleotides have values  $F(\mathcal{S}_k^{\text{pref}}) \gg 1/3$  (data not shown), i.e. different from the random case (Property 1 in the Method).

Five  $\mathcal{C}^3$  codes  $X_0(\mathcal{G})$  have one among 20 set  $\mathcal{S}_k$ , named  $\tilde{\mathcal{S}}_k$ , which is not preferential as its value  $F(\tilde{\mathcal{S}}_k) \approx 1/3$  is close to the random case (Property 1 in Method), precisely  $F(\tilde{\mathcal{S}}_k) < 0.37$  (data not shown). In these sets  $\tilde{\mathcal{S}}_k$ , the occurrence probabilities of the 3 permuted trinucleotides are similar in the 3 frames. These sets are:  $\tilde{\mathcal{S}}_4$  in  $\mathcal{G} = TPV$ ,  $\tilde{\mathcal{S}}_6$  in  $\mathcal{G} = PCF$ ,  $\tilde{\mathcal{S}}_{12}$

in  $\mathcal{G} = TPA$  and  $\mathcal{S}_{18}$  in  $\mathcal{G} = \{MC, SLS\}$  (data not shown).

The  $\mathcal{C}^3$  code  $X_0(HB)$  has 2 sets  $\tilde{\mathcal{S}}_k$ :  $\tilde{\mathcal{S}}_3 = \{AAT, ATA, TAA\}$  and  $\tilde{\mathcal{S}}_{12} = \{ATT, TTA, TAT\}$  (Table 2).  $\tilde{\mathcal{S}}_3$  and  $\tilde{\mathcal{S}}_{12}$  are the only sets among 20 only made of A and T. As the genome HB has the highest GC content (68%; Ng et al., 2000) compared to the other genomes, its frequencies of A and T are low (data not shown).

These sets  $\tilde{\mathcal{S}}_k$  are rare as there are in total 7 sets  $\tilde{\mathcal{S}}_k$  among  $15 \times 20 = 300$  sets  $\tilde{\mathcal{S}}_k^{\text{pref}}$ , i.e. about 2%. In these sets  $\tilde{\mathcal{S}}_k$ , the preferential set  $\mathcal{S}_k^{\text{pref}}$  is assigned to one of the 3 sets  $\mathcal{S}^{22}$ ,  $\mathcal{S}^{44}$  and  $\mathcal{S}^{53}$  leading to a  $\mathcal{C}^3$  code  $X_0(\mathcal{G})$ . Note that a solution has always been found with these 7 sets  $\tilde{\mathcal{S}}_k$  as 15  $\mathcal{C}^3$  codes  $X_0(\mathcal{G})$  have been identified.

### 3.3.3. Identical circular codes on both DNA strands in the archaeal genomes

In each archaeal genome, the same  $\mathcal{C}^3$  code  $X_0(\mathcal{G})$  is identified in genes of the direct and complementary DNA strands as the occurrence frequencies of the 64 trinucleotides in the 3 frames of genes are almost identical in both strands (data not shown).

### 3.4. Construction frame determination property in the archaeal circular codes

The property of the determination of the construction frame (decomposition) of a word generated by any concatenation of words of a circular code (more simply, generated by a circular code), is a consequence of the circular code property. The principle is presented briefly with the previous example of the circular code  $X_0(AP) = X_0(AG)$  in the archaeal genomes AP and AG (see Béal, 1993; Berstel and Perrin, 1985 for the details). The unicity of the decomposition of a word into words of a circular code is not obvious.

For example, the word  $w' = TAGCGACCT$  of length 9 can be decomposed into trinucleotides of  $X_0(AP)$  in 2 ways:

- $\{12, 4\}T\{9, 13\}A1G4C12G1A2C7C1T5$  with  $l_1TA, GCG, ACC, Tl_2l_3 \in X_0(AP)$
- $\{9, 13\}T1A2G8C1G4A11C1C3T10$  with  $l_1l_2T, AGC, GAC, CTl_3 \in X_0(AP)$

where the states of the flower automaton  $\mathcal{F}(X_0(AP))$  of  $X_0(AP)$  separate the letters of  $w'$  (Fig. 3). Therefore, the word  $w'$  has 2 construction frames: frame 1 (initial states  $\{2,4\}$  of  $\mathcal{F}(X_0(AP))$ ) or frame 2 (initial states  $\{9,13\}$  of  $\mathcal{F}(X_0(AP))$ ) (Fig. 3). With the example chosen, the right factor GCGACCT of  $w'$  is one of the 4 longest words of length 7 identified in Fig. 4. It begins at the states 1 and 2 in  $\mathcal{F}(X_0(AP))$  (Fig. 3). The left factor TA of  $w'$  can be concatenated with this right factor by ending the letter A of TA at the states 1 and 2 in  $\mathcal{F}(X_0(AP))$  (Fig. 3).

The unicity of such a decomposition is proved by using the properties of the flower automaton  $\mathcal{F}(X_0(\mathcal{G}))$

Words starting from the states 1 and 2 in the flower automaton  $\mathcal{F}(\mathcal{X}_0(\text{AP}))$

A  
AC  
ACC  
**ACCT**  
AG  
AGC  
**AGCT**  
C  
G  
GC  
GCC  
**GCCT**  
GCG  
GCGA  
GCGAC  
GCGACC  
**GCGACCT**  
GCGAG  
GCGAGC  
**GCGAGCT**  
**GCGC**  
GCGT  
GCGTA  
GCGTAC  
**GCGTACT**  
GCGTT  
GCGTTC  
**GCGTTCT**

T  
TA  
TAC  
**TACT**  
TT  
TTC  
**TTCT**

Words starting from the states 1 and 3 in the flower automaton  $\mathcal{F}(\mathcal{X}_0(\text{AP}))$

**T**

Words starting from the states 1 and 4 in the flower automaton  $\mathcal{F}(\mathcal{X}_0(\text{AP}))$

A  
AC  
ACC  
**ACCT**  
AG  
AGC  
**AGCT**  
C  
T  
TA  
TAC  
**TACT**  
TT  
TTC  
**TTCT**

Words starting from the states 1 and 5 in the flower automaton  $\mathcal{F}(\mathcal{X}_0(\text{AP}))$

A  
AC  
ACC  
**ACCT**  
**T**

Fig. 4. Absence of 2 cycles labeled with the same word in the flower automaton  $\mathcal{F}(\mathcal{X}_0(\text{AP}))$  proving that the set  $\mathcal{X}_0(\text{AP})$  in the archaeal genome AP is a circular code. The words in bold are “non-extensible”.

of the circular code  $\mathcal{X}(\mathcal{G})$ . If any letter of a word generated by  $\mathcal{X}(\mathcal{G})$  can be associated with a unique state of  $\mathcal{F}(\mathcal{X}(\mathcal{G}))$ , then a frame can be assigned to this letter as each state of  $\mathcal{F}(\mathcal{X}(\mathcal{G}))$  is related to a frame. In the case of a flower automaton  $\mathcal{F}$ , there exists a number

$n(\mathcal{X}(\mathcal{G}))$  so that for a word of length  $\geq n(\mathcal{X}(\mathcal{G}))$ , all paths associated with this word have the same terminal state, i.e. a unique state for a letter of this word is identified. Thus, the construction frame of the word can be determined and the word can be decomposed into

Table 4

Lengths of the minimal windows to retrieve the construction frames 0 (reading frame), 1 and 2 with the 15 archaeal circular codes  $\mathcal{X}_0(\mathcal{G})$ ,  $\mathcal{X}_1(\mathcal{G})$  and  $\mathcal{X}_2(\mathcal{G})$ , respectively

$\mathcal{G}$	AP,AG	HB	MC	MP	MSA	MSM	MT	PB	PCA	PCF	PCH	SLS	SLT	TPA	TPV
$\mathcal{X}_0(\mathcal{G})$	10	9	7	10	10	9	10	7	10	7	10	7	5	7	9
$\mathcal{X}_1(\mathcal{G})$	10	10	10	8	10	11	10	8	8	6	8	10	7	7	11
$\mathcal{X}_2(\mathcal{G})$	7	11	7	11	10	7	10	10	10	6	9	10	6	7	10

words of the circular code. In other words, the length of the minimal window to retrieve the construction frame is  $n(\mathcal{X}(\mathcal{G}))$  letters.

With the previous example,  $n(\mathcal{X}_0(\text{AP})) = 10$  (first line and column in Table 4), i.e. the length of the minimal window to retrieve the construction frame with  $\mathcal{X}_0(\text{AP})$  is 10 letters. The word  $w'$  of length 9 was chosen in the set of the longest words with 2 construction frames. The concatenation of only one letter  $l$  to  $w'$  leads to the word  $w = w'l$  of length 10 with a unique construction frame as 10 is the length of the minimal window. For example with  $l = \{A, T\}$ , the construction frame of  $w$  is unique and equal to 1 as there is no edge labeled  $\{A, T\}$  leaving the state 10 in  $\mathcal{F}(\mathcal{X}_0(\text{AP}))$  (Fig. 3). Then, the unique decomposition of  $w$  into trinucleotides of  $\mathcal{X}_0(\text{AP})$  is  $l_1\text{TA}, \text{GCG}, \text{ACC}, \text{T}\{A, T\}l_2$ .

The program developed in Java testing all possible paths in a flower automaton (not described here) allows to compute the lengths  $n(\mathcal{X}(\mathcal{G}))$  of the minimal windows of the archaeal circular codes  $\mathcal{X}_0(\mathcal{G})$ ,  $\mathcal{X}_1(\mathcal{G})$  and  $\mathcal{X}_2(\mathcal{G})$  to retrieve the construction frames 0, 1 and 2, respectively (Table 4). These lengths vary between 5 for  $\mathcal{X}_0(\text{SLT})$  and 11 for  $\mathcal{X}_1(\text{MSM})$ ,  $\mathcal{X}_1(\text{TPV})$ ,  $\mathcal{X}_2(\text{HB})$  and  $\mathcal{X}_2(\text{MK})$  (Table 4).

### 3.5. Probabilities of words generated by the archaeal circular codes

A new definition of a simple parameter is introduced for measuring some probabilities of words generated by the circular codes. Let  $N(\mathcal{L})$  be the number of words of a given length  $\mathcal{L}$ , i.e.  $N(\mathcal{L}) = 4^\mathcal{L}$  on the alphabet  $\{A, C, G, T\}$ . Let  $\mathcal{F}_p(\mathcal{L}, \mathcal{X}(\mathcal{G}))$  be the set of words of a given length  $\mathcal{L}$  obtained from all the states of the frame  $p \in \{0, 1, 2\}$ , of the flower automaton  $\mathcal{F}$  of the circular code  $\mathcal{X}(\mathcal{G})$ . By using the classical Poincaré formula, the number  $N(\mathcal{L}, \mathcal{X}(\mathcal{G}))$  of words of a given length  $\mathcal{L}$  generated by the circular code  $\mathcal{X}(\mathcal{G})$  is equal to

$$N(\mathcal{L}, \mathcal{X}(\mathcal{G})) = \text{CARD} \left( \bigcup_{p=0}^2 \mathcal{F}_p(\mathcal{L}, \mathcal{X}(\mathcal{G})) \right) = \sum_{p=0}^2 (-1)^p S_p \quad (9)$$

with  $S_p = \sum_{0 \leq i_0 < \dots < i_p \leq 2} \text{CARD}(\mathcal{F}_{i_0}(\mathcal{L}, \mathcal{X}(\mathcal{G})) \cap \dots \cap \mathcal{F}_{i_p}(\mathcal{L}, \mathcal{X}(\mathcal{G})))$ .

Then, its probability  $\mathcal{P}(\mathcal{L}, \mathcal{X}(\mathcal{G}))$  is equal to

$$\mathcal{P}(\mathcal{L}, \mathcal{X}(\mathcal{G})) = \frac{N(\mathcal{L}, \mathcal{X}(\mathcal{G}))}{N(\mathcal{L})} \quad (10)$$

The words with a unique construction frame generated by the circular code  $\mathcal{X}(\mathcal{G})$  are the words obtained from all the states of a unique frame, i.e. the words with an empty intersection in the 3 sets  $\mathcal{F}_p(\mathcal{L}, \mathcal{X}(\mathcal{G}))$ . Therefore, the number  $M(\mathcal{L}, \mathcal{X}(\mathcal{G}))$  of words of a given length  $\mathcal{L}$  with a unique construction frame generated by the circular code  $\mathcal{X}(\mathcal{G})$ , is equal to

$$M(\mathcal{L}, \mathcal{X}(\mathcal{G})) = N(\mathcal{L}, \mathcal{X}(\mathcal{G})) - \text{CARD} \left( \bigcup_{0 \leq i_0 < i_1 \leq 2} (\mathcal{F}_{i_0}(\mathcal{L}, \mathcal{X}(\mathcal{G})) \cap \mathcal{F}_{i_1}(\mathcal{L}, \mathcal{X}(\mathcal{G}))) \right) = \sum_{p=0}^2 (-1)^p (p+1) S_p \quad (11)$$

and its probability  $\mathcal{Q}(\mathcal{L}, \mathcal{X}(\mathcal{G}))$ , to

$$\mathcal{Q}(\mathcal{L}, \mathcal{X}(\mathcal{G})) = \frac{M(\mathcal{L}, \mathcal{X}(\mathcal{G}))}{N(\mathcal{L})} \quad (12)$$

With the previous example of the  $\mathcal{C}^3$  code  $\mathcal{X}_0(\text{AP}) = \mathcal{X}_0(\text{AG})$  in the archaeal genomes AP and AG, Table 5 shows the probabilities  $\mathcal{P}(\mathcal{L}, \mathcal{X}_0(\text{AP}))$  and  $\mathcal{Q}(\mathcal{L}, \mathcal{X}_0(\text{AP}))$  by varying the word length  $\mathcal{L} \in \{3, 11\}$ . With  $\mathcal{L} = 3$ , only 1 word, GGG, cannot be generated by the flower automaton  $\mathcal{F}(\mathcal{X}_0(\text{AP}))$  (Fig. 3). Therefore,  $\mathcal{P}(3, \mathcal{X}_0(\text{AP})) = (4^3 - 1)/4^3 \approx 0.984$ . With  $\mathcal{L} = 9$ , 4 words has 2 construction frames (Fig. 3): TAGCGACT, TAGCGAGCT, TAGCGTACT and TAGCGTTCT with the frames 1 (initial states  $\{2, 4\}$  in  $\mathcal{F}(\mathcal{X}_0(\text{AP}))$ ) and 2 (initial states  $\{9, 13\}$  in  $\mathcal{F}(\mathcal{X}_0(\text{AP}))$ ). Note that the suffixes of length 7 of these 4 words are the 4 words with the greatest length in Fig. 4. No word has 2 construction frames with  $\mathcal{L} = 10$  which is the length of the minimal window to retrieve the construction frame 0 with  $\mathcal{X}_0(\text{AP})$  (Table 4).

The 2 probabilities  $\mathcal{P}(\mathcal{L}, \mathcal{X}_0(\mathcal{G}))$  and  $\mathcal{Q}(\mathcal{L}, \mathcal{X}_0(\mathcal{G}))$  with the 15 archaeal  $\mathcal{C}^3$  codes  $\mathcal{X}_0(\mathcal{G})$  (in reading frame) by varying the word length  $\mathcal{L} \in \{3, 11\}$  are given in

Table 5

Probability  $\mathcal{P}(\mathcal{L}, \mathcal{X}_0(\text{AP}))$  (%) of words generated by the archaeal  $\mathcal{C}^3$  code  $\mathcal{X}_0(\text{AP})$  (in reading frame) and probability  $\mathcal{Q}(\mathcal{L}, \mathcal{X}_0(\text{AP}))$  (%) of words with a unique construction frame also generated by the same code  $\mathcal{X}_0(\text{AP})$  by varying the word length between 3 and 11

Length $\mathcal{L}$	Number of words not generated	Number of words in 1 frame	Number of words in 2 frames	Number of words in 3 frames	Sum	Probability $\mathcal{P}(\mathcal{L}, \mathcal{X}_0(\text{AP}))$ (%)	Probability $\mathcal{Q}(\mathcal{L}, \mathcal{X}_0(\text{AP}))$ (%)
3	1	32	29	2	64	98.4	50.0
4	43	176	37	0	256	83.2	68.8
5	360	628	36	0	1024	64.8	61.3
6	2205	1862	29	0	4096	46.2	45.5
7	11396	4976	12	0	16384	30.4	30.4
8	51544	13984	8	0	65536	21.4	21.3
9	223748	38392	4	0	262144	14.6	14.6
10	948576	100000	0	0	1048576	9.5	9.5
11	3914304	280000	0	0	4194304	6.7	6.7

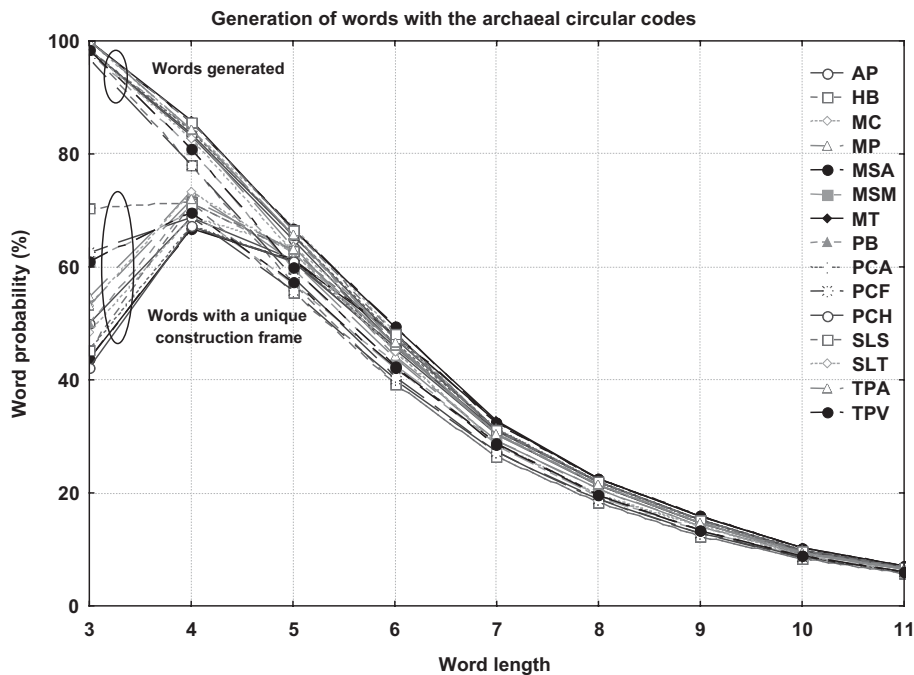


Fig. 5. Probability  $\mathcal{P}(\mathcal{L}, \mathcal{X}_0(\mathcal{G}))$  (%) of words generated by the 15 archaeal  $\mathcal{C}^3$  codes  $\mathcal{X}_0(\mathcal{G})$  (in reading frame) and probability  $\mathcal{Q}(\mathcal{L}, \mathcal{X}_0(\mathcal{G}))$  (%) of words with a unique construction frame also generated by the same codes  $\mathcal{X}_0(\mathcal{G})$  by varying the word length between 3 and 11.

Fig. 5. For all the codes, the probabilities  $\mathcal{P}(\mathcal{L}, \mathcal{X}_0(\mathcal{G}))$  decrease from a mean value of  $98.9 \pm 0.8\%$ , where 0.8 represents the average deviation, at  $\mathcal{L} = 3$ , to a mean value of  $6.5 \pm 0.4\%$  at  $\mathcal{L} = 11$ . Five codes in the archaeal genomes  $\mathcal{G} = \{\text{MC}, \text{MSA}, \text{MT}, \text{SLS}, \text{TPA}\}$  generate all the 64 trinucleotides (100% at  $\mathcal{L} = 3$ ). For all the codes, the probabilities  $\mathcal{Q}(\mathcal{L}, \mathcal{X}_0(\mathcal{G}))$  increase from a mean value of  $52.3 \pm 6.6\%$  at  $\mathcal{L} = 3$ , to a maximum with a mean value of  $69.7 \pm 1.9\%$  at  $\mathcal{L} = 4$ , then decrease and begin to merge with the probabilities  $\mathcal{P}(\mathcal{L}, \mathcal{X}_0(\mathcal{G}))$  at  $\mathcal{L} = 5$ . Therefore, almost all the words of lengths  $\geq 5$  generated by the 15 archaeal  $\mathcal{C}^3$  codes allow to retrieve the reading frame in genes. The number

of words with a unique construction frame also differs between the different codes. For example at  $\mathcal{L} = 5$ ,  $\mathcal{Q}(5, \mathcal{X}_0(\text{MC})) = 64\%$  while  $\mathcal{Q}(5, \mathcal{X}_0(\text{HB})) = 55\%$ , i.e. the code  $\mathcal{X}_0(\text{MC})$  generates 84 more words of length 5 with a unique construction frame compared to the code  $\mathcal{X}_0(\text{HB})$ . The code  $\mathcal{X}_0(\text{MC})$  may lead to a greater gene diversity. Otherwise, an analytical model of gene evolution based on an independent mixing of the trinucleotides of the  $\mathcal{C}^3$  code of eukaryotes/prokaryotes allows to simulate the actual genes of eukaryotes/prokaryotes (Arquès et al., 1998). Therefore, the probabilities  $\mathcal{P}(\mathcal{L}, \mathcal{X}_0(\mathcal{G}))$  and  $\mathcal{Q}(\mathcal{L}, \mathcal{X}_0(\mathcal{G}))$  may represent evolutionary parameters of the word diversity of the



Table 6 (continued)

		AG, AP	HB	MC	MP	MSA	MSM	MT	PB	PCA	PCF	PCH	SLS	SLT	TPA	TPV
YYR	CCA	1	1	0	1	1	1	1	1	1	0	0	0	0	1	1
	CCG	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	CTA	2	2	1	2	2	1	2	2	2	0	0	1	1	2	1
	CTG	0	0	1	0	1	1	1	0	1	1	1	1	1	0	1
	TCA	1	1	0	1	1	1	1	1	1	1	1	0	0	1	1
	TCG	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	TTA	1	2	1	1	1	1	1	1	1	1	1	1	1	1	1
	TTG	1	0	1	0	1	1	1	1	1	1	1	1	1	1	1
	MEAN	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
YYY	CCC	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—
	CCT	2	2	0	2	2	2	2	2	2	2	2	0	0	2	2
	CTC	0	0	1	0	0	0	0	0	0	0	0	1	1	0	0
	CTT	2	2	1	2	0	0	2	2	2	0	2	1	1	2	0
	TCC	1	1	2	1	1	1	1	1	1	1	1	2	2	1	1
	TCT	1	1	0	1	2	2	1	1	1	2	1	0	0	1	2
	TTC	0	0	2	0	1	1	0	0	0	1	0	2	2	0	1
	TTT	—	—	—	—	—	—	—	—	—	—	—	—	—	—	—
	MEAN	0 1 2	0 1 2	0 1 2	0 1 2	0 1 2	0 1 2	0 1 2	0 1 2	0 1 2	0 1 2	0 1 2	0 1 2	0 1 2	0 1 2	0 1 2

different  $\mathcal{C}^3$  codes, and in particular of the trinucleotide diversity.

3.6. Types of nucleotides in the trinucleotide sites of the archaeal circular codes

The 4 types of nucleotides in the 15 archaeal  $\mathcal{C}^3$  codes  $\mathcal{X}_0(\mathcal{G})$  (in reading frame) occur in:

- the first trinucleotide site except T with the 3 genomes  $\mathcal{G} = \{\text{MSM, PCF, TPV}\}$ ,
- the second trinucleotide site except G with the 5 genomes  $\mathcal{G} = \{\text{MP, PB, PCA, PCF, PCH}\}$ ,
- the third trinucleotide site except A with the genome  $\mathcal{G} = \text{HB}$  and G with the genome  $\mathcal{G} = \text{SLT}$ .

Therefore, only one type of nucleotide T (resp. G) may not occur in the first (resp. second) site, but 2 types of nucleotides A and G may not occur in the third site. In other words, G may not occur in 2 sites, the second and the third ones, and C always occurs in the 3 sites. The code  $\mathcal{X}_0(\text{PCF})$  has no trinucleotides T and G in the first and second sites, respectively. It codes for the smallest number of amino acids among all the archaeal codes (see Discussion).

Similar rules can be deduced obviously with the 15 archaeal codes  $\mathcal{X}_1(\mathcal{G})$  and  $\mathcal{X}_2(\mathcal{G})$  by permutation.

4. Discussion

4.1. Common codons in the archaeal circular codes

The common and rare codons in the 15 archaeal  $\mathcal{C}^3$  codes  $\mathcal{X}_0(\mathcal{G})$  (in reading frame) are (from Table 3):

- 20 codons are absent (occurrence number  $Nb = 0$  of codons in the 15 codes)
- 11 codons are very rare ( $1 \leq Nb \leq 3$ )
- 9 codons are rare ( $4 \leq Nb \leq 7$ )
- 10 codons are common ( $8 \leq Nb \leq 11$ ): AAC, ACC, GAA, GAT, GCA, GCG, GCT, GTG, TAC, TTC
- 7 codons are very common ( $12 \leq Nb \leq 14$ ): ATA, ATC, ATT, CTC, GAG, GTA, GTT
- 3 codons are always present ( $Nb = 15$ ): GAC, GCC, GTC which code for Asp, Ala and Val, respectively.

4.2. Consequences on the 2-letter alphabets of the archaeal genomes

The 3 circular codes  $\mathcal{X}_0(\mathcal{G})$ ,  $\mathcal{X}_1(\mathcal{G})$  and  $\mathcal{X}_2(\mathcal{G})$  in an archaeal genome  $\mathcal{G}$  have 20 trinucleotides in the frames 0, 1 and 2, respectively. Therefore, a preferential frame for the 8 2-letter trinucleotides on a given 2-letter alphabet can be deduced by considering for each 2-letter trinucleotide, the average frame associated with the 8 4-letter trinucleotides specified on the 2-letter trinucleotide and belonging to the code  $\mathcal{X}_0(\mathcal{G})$  (frame 0),  $\mathcal{X}_1(\mathcal{G})$  (frame 1) and  $\mathcal{X}_2(\mathcal{G})$  (frame 2).

On the alphabet  $\{\text{R, Y}\}$  ( $\text{R} = \text{purine} = \{\text{A, G}\}$ ,  $\text{Y} = \text{pyrimidine} = \{\text{C, T}\}$ ), the trinucleotide RYY occurs in frame 0 with all the 15 archaeal codes (Table 6). With the 2 genomes  $\mathcal{G} = \{\text{MSM, TPV}\}$ , all the 8 4-letter trinucleotides specified on RYY are in frame 0. The pattern RYY is also deduced from the circular code of mitochondria (Arquès and Michel, 1997). Its permuted trinucleotides YYR and YRY occur obviously in frames 1 and 2, respectively (Table 6). The trinucleotide RYR also occurs in frame 0 with 13 archaeal codes and its permuted trinucleotides YRR and RRY in frames 1



Table 7

Amino acids coded by the 15 archaeal circular codes  $\mathcal{X}_0(\mathcal{G})$ ,  $\mathcal{X}_1(\mathcal{G})$  and  $\mathcal{X}_2(\mathcal{G})$  according to the universal genetic code. The numbers 0 (in bold), 1 and 2 are associated with the frames 0 (reading frame), 1 and 2 of the codes  $\mathcal{X}_0(\mathcal{G})$ ,  $\mathcal{X}_1(\mathcal{G})$  and  $\mathcal{X}_2(\mathcal{G})$ , respectively

		AG, AP	HB	MC	MP	MSA	MSM	MT	PB	PCA	PCF	PCH	SLS	SLT	TPA	TPV
Lysine, Lys, K	AAG	<b>0</b>	<b>0</b>	1	<b>0</b>	1	1	<b>0</b>	1	<b>0</b>	1	1	1	1	<b>0</b>	1
Methionine, Met, M	ATG	<b>0</b>	<b>0</b>	1	<b>0</b>	1	1	1	<b>0</b>	1	1	1	1	1	<b>0</b>	1
Phenylalanine, Phe, F	TTC	<b>0</b>	<b>0</b>	2	<b>0</b>	1	1	<b>0</b>	<b>0</b>	<b>0</b>	1	<b>0</b>	2	2	<b>0</b>	1
Tryptophan, Trp, W	TGG	1	1	2	1	1	1	1	1	1	1	1	2	2	1	2
Asparagine, Asn, N	AAC	<b>0</b>	<b>0</b>	2	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	2	2	2	2	<b>0</b>	<b>0</b>
	AAT	2	<b>0</b>	2	2	2	2	2	2	2	2	2	2	2	2	2
Aspartic acid, Asp, D	GAC	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>
	GAT	2	2	<b>0</b>	2	<b>0</b>	<b>0</b>	<b>0</b>	2	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	2	<b>0</b>
Cysteine, Cys, C	TGC	1	1	2	1	2	2	2	1	2	2	2	2	2	1	2
	TGT	2	1	2	1	2	2	2	2	2	2	2	2	2	2	2
Glutamic acid, Glu, E	GAA	2	2	<b>0</b>	2	<b>0</b>	<b>0</b>	2	<b>0</b>	2	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	2	<b>0</b>
	GAG	<b>0</b>	<b>0</b>	1	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	1	<b>0</b>	<b>0</b>
Glutamine, Gln, Q	CAA	2	2	1	2	2	2	2	2	2	1	1	1	1	2	2
	CAG	2	<b>0</b>	1	<b>0</b>	1	1	1	1	1	1	1	1	1	2	1
Histidine, His, H	CAC	2	2	1	2	2	2	2	2	2	1	1	1	1	2	2
	CAT	2	2	1	2	2	2	2	2	2	2	2	1	1	2	2
Tyrosine, Tyr, Y	TAC	<b>0</b>	<b>0</b>	2	<b>0</b>	<b>0</b>	2	<b>0</b>	<b>0</b>	<b>0</b>	1	1	2	2	<b>0</b>	2
	TAT	2	<b>0</b>	2	2	2	2	2	2	2	2	2	2	2	2	2
Glycine, Gly, G	GGA	2	2	<b>0</b>	2	2	2	2	2	2	2	2	2	<b>0</b>	2	2
	GGC	2	<b>0</b>	2	2	<b>0</b>	<b>0</b>	<b>0</b>	2	2	2	2	2	<b>0</b>	<b>0</b>	<b>0</b>
	GGT	2	2	<b>0</b>	2	2	2	2	2	2	2	2	<b>0</b>	<b>0</b>	2	<b>0</b>
Isoleucine, Ile, I	ATA	<b>0</b>	1	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>
	ATC	<b>0</b>	<b>0</b>	2	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	2	2	<b>0</b>	<b>0</b>
	ATT	<b>0</b>	1	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>
Proline, Pro, P	CCA	1	1	<b>0</b>	1	1	1	1	1	1	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	1	1
	CCG	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	CCT	2	2	<b>0</b>	2	2	2	2	2	2	2	2	<b>0</b>	<b>0</b>	2	2
Alanine, Ala, A	GCA	1	2	<b>0</b>	2	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	1	<b>0</b>
	GCC	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>
	GCG	<b>0</b>	1	<b>0</b>	<b>0</b>	1	1	1	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	1	1	1
	GCT	2	2	<b>0</b>	2	<b>0</b>	<b>0</b>	<b>0</b>	2	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	2	<b>0</b>
Threonine, Thr, T	ACA	1	1	<b>0</b>	1	1	1	1	1	1	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	1	1
	ACC	<b>0</b>	<b>0</b>	2	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	2	2	2	2	<b>0</b>	<b>0</b>
	ACG	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	ACT	1	1	<b>0</b>	1	1	<b>0</b>	1	1	1	2	2	<b>0</b>	<b>0</b>	1	<b>0</b>
Valine, Val, V	GTA	<b>0</b>	2	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>
	GTC	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>
	GTG	<b>0</b>	<b>0</b>	1	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	1	1	<b>0</b>	1
	GTT	<b>0</b>	2	<b>0</b>	2	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>
Arginine, Arg, R	AGA	1	1	2	1	2	2	1	2	1	2	2	2	2	1	2
	AGG	1	1	2	1	1	1	1	1	1	1	1	1	2	1	1
	CGA	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
	CGC	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
	CGG	1	2	1	1	2	2	2	1	1	1	1	1	2	2	2
	CGT	2	2	2	2	2	2	2	2	2	2	2	2	2	2	2
Leucine, Leu, L	CTA	2	2	1	2	2	1	2	2	2	<b>0</b>	<b>0</b>	1	1	2	1
	CTC	<b>0</b>	<b>0</b>	1	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	<b>0</b>	1	1	<b>0</b>	<b>0</b>
	CTG	<b>0</b>	<b>0</b>	1	<b>0</b>	1	1	1	<b>0</b>	1	1	1	1	1	<b>0</b>	1
	CTT	2	2	1	2	<b>0</b>	<b>0</b>	2	2	2	<b>0</b>	2	1	1	2	<b>0</b>
	TTA	1	2	1	1	1	1	1	1	1	1	1	1	1	1	1
	TTG	1	<b>0</b>	1	<b>0</b>	1	1	1	1	1	1	1	1	1	1	1
Serine, Ser, S	AGC	<b>0</b>	1	2	1	2	2	2	2	2	2	2	2	2	<b>0</b>	2
	AGT	2	1	2	2	2	2	2	2	2	2	2	2	2	2	2
	TCA	1	1	<b>0</b>	1	1	1	1	1	1	1	1	<b>0</b>	<b>0</b>	1	1
	TCC	1	1	2	1	1	1	1	1	1	1	1	2	2	1	1
	TCG	1	1	1	1	1	1	1	1	1	1	1	1	1	1	1
	TCT	1	1	<b>0</b>	1	2	2	1	1	1	2	1	<b>0</b>	<b>0</b>	1	2
Stop codon, ochre	TAA	1	2	1	1	1	1	1	1	1	1	1	1	1	1	1
Stop codon, amber	TAG	1	<b>0</b>	1	1	1	1	1	1	1	1	1	1	1	1	1
Stop codon, opal	TGA	1	1	2	1	2	2	2	1	2	2	2	2	2	1	2

Table 8

Classification of the archaeal genomes according to the number of amino acids coded by the 15  $\mathcal{C}^3$  codes  $\mathcal{X}_0(\mathcal{G})$  (in reading frame) (deduced from Table 7). The top part gives the 6 amino acids always coded by these codes

	PCF	PCH	MC, SLS, SLT	MSM, TPV	MSA	PB	PCA	MT	AG, AP	MP	HB	TPA
Aspartic acid, Asp, D	+	+	+	+	+	+	+	+	+	+	+	+
Glutamic acid, Glu, E	+	+	+	+	+	+	+	+	+	+	+	+
Isoleucine, Ile, I	+	+	+	+	+	+	+	+	+	+	+	+
Alanine, Ala, A	+	+	+	+	+	+	+	+	+	+	+	+
Threonine, Thr, T	+	+	+	+	+	+	+	+	+	+	+	+
Valine, Val, V	+	+	+	+	+	+	+	+	+	+	+	+
Lysine, Lys, K	–	–	–	–	–	–	+	+	+	+	+	+
Methionine, Met, M	–	–	–	–	–	+	–	–	+	+	+	+
Phenylalanine, Phe, F	–	+	–	–	–	+	+	+	+	+	+	+
Asparagine, Asn, N	–	–	–	+	+	+	+	+	+	+	+	+
Glutamine, Gln, Q	–	–	–	–	–	–	–	–	–	+	+	–
Tyrosine, Tyr, Y	–	–	–	–	+	+	+	+	+	+	+	+
Glycine, Gly, G	–	–	+	+	+	–	–	+	–	–	+	+
Proline, Pro, P	+	+	+	–	–	–	–	–	–	–	–	–
Leucine, Leu, L	+	+	–	+	+	+	+	+	+	+	+	+
Serine, Ser, S	–	–	+	–	–	–	–	–	+	–	–	+
Total number of AA	8	9	9	9	10	11	11	12	13	13	14	14

and 2, respectively (Table 6). There is no preferential frame for RRR and YYY. Therefore, the archaeal genes follow the pattern  $\text{RYN} = \{\text{RYR}, \text{RYY}\}$  ( $\text{N} = \{\text{R}, \text{Y}\}$ ) contrary to the classical one RNY deduced from the circular code of eukaryotes/prokaryotes (Arquès and Michel, 1996; Eigen and Schuster, 1978). Otherwise, the pattern RYN is a maximal circular code on the alphabet  $\{\text{R}, \text{Y}\}$  but not a self-complementary one.

This analysis on the 2 other 2-letter alphabets  $\{\text{K}, \text{M}\}$  ( $\text{K} = \text{ceto} = \{\text{G}, \text{T}\}$ ,  $\text{M} = \text{amino} = \{\text{A}, \text{C}\}$ ) and  $\{\text{S}, \text{W}\}$  ( $\text{S} = \text{strong interaction} = \{\text{C}, \text{G}\}$ ,  $\text{W} = \text{weak interaction} = \{\text{A}, \text{T}\}$ ) does not reveal a preferential 2-letter pattern as the 4-letter trinucleotides in frame 0 are spread on several 2-letter trinucleotides (data not shown). Therefore, in archaeal genes, the 2-letter alphabet close to the alphabet  $\{\text{A}, \text{C}, \text{G}, \text{T}\}$  is  $\{\text{R}, \text{Y}\}$ .

#### 4.3. Consequences on the genetic code of the archaeal genomes

Four amino acids (AA) are never coded by the 15 archaeal  $\mathcal{C}^3$  codes  $\mathcal{X}_0(\mathcal{G})$  (in reading frame): Arg, Cys, His and Trp (Table 7). These 4 AA have a complex chemical structure in terms of their numbers of atoms or cycles. Indeed, Trp is the single AA with 2 cycles. Arg and His are the 2 most complex positively charged (basic) polar AA: Arg has the greatest number of atoms in the side chain (without cycle) and His has a cycle which leads, under certain conditions, to a dual acid/base property allowing catalytic reactions in the active sites of enzymes. Cys can form disulfide linkages by reaction with another Cys.

Six amino acids are always coded by the 15 archaeal codes: Ala, Asp, Glu, Ile, Thr and Val (Table 7). Ala and Asp represent the complete group of negatively

charged (acidic) polar AA. These 6 AA are equally represented in the 2 classes of aminoacyl-tRNA synthetases with a class I containing Glu, Ile and Val, and a class II holding Ala, Asp and Thr (reviewed in Schimmel et al., 1993; Hartman, 1995; Saks and Sampson, 1995). Leu close to Ile is coded by 12 archaeal codes.

These archaeal codes code for a number of AA varying from 8 AA with  $\mathcal{G} = \text{PCF}$  to 14 AA with  $\mathcal{G} = \{\text{HB}, \text{TPA}\}$  (Table 8). Circular codes coding a small number of AA may be more “primitive”. Pro, which is based on a ring structure that includes the central carbon atom, is only coded by the primitive archaeal codes (number of AA  $\leq 9$  with  $\mathcal{G} = \{\text{PCF}, \text{PCH}, \text{MC}, \text{SLS}, \text{SLT}\}$ ). Two AA, Gly and Ser, are coded both by primitive and evolved circular codes. Lys is the single basic AA encoded by archaeal codes, and furthermore, only by the evolved ones (number of AA  $\geq 11$  with  $\mathcal{G} = \{\text{PCA}, \text{MT}, \text{AG}, \text{AP}, \text{MP}, \text{HB}, \text{TPA}\}$ ). Met, Phe, Asn, Gln and Tyr are mainly coded by the evolved archaeal codes.

#### Acknowledgements

We thank the referee for his advice.

#### References

- Andersson, S.G.E., Kurland, C.G., 1990. Codon preferences in free living microorganisms. *Microbiol. Rev.* 54, 198–210.
- Antezana, M.A., Kreitman, M., 1999. The nonrandom location of synonymous codons suggests that reading frame-independent forces have patterned codon preferences. *J. Mol. Evol.* 49, 36–43.

- Arquès, D.G., Michel, C.J., 1996. A complementary circular code in the protein coding genes. *J. theor. Biol.* 182, 45–58.
- Arquès, D.G., Michel, C.J., 1997. A circular code in the protein coding genes of mitochondria. *J. theor. Biol.* 189, 273–290.
- Arquès, D.G., Fallot, J.-P., Michel, C.J., 1998. An evolutionary analytical model of a complementary circular code simulating the protein coding genes, the 5' and 3' regions. *Bull. Math. Biol.* 60, 163–194.
- Béal, M.-P., 1993. *Codage Symbolique*. Masson, Paris.
- Berg, O.G., Silva, P.J.N., 1997. Codon bias in *Escherichia coli*: the influence of codon context on mutation and selection. *Nucleic Acids Res.* 25, 1397–1404.
- Bernander, R., 2000. Chromosome replication, nucleotid segregation and cell division in archaea. *Trends Microbiol.* 8, 278–283.
- Berstel, J., Perrin, D., 1985. *Theory of Codes*. Academic Press, New York.
- Campbell, A., Mrázek, J., Karlin, S., 1999. Genomic signature comparisons among prokaryote, plasmid, and mitochondrial DNA. *Proc. Natl Acad. Sci. USA* 96, 9184–9189.
- Eigen, M., Schuster, P., 1978. The hypercycle. A principle of natural self-organization. Part C: the realistic hypercycle. *Naturwissenschaften* 65, 341–369.
- Forterre, P., 2001. Genomics and early cellular evolution. The origin of the DNA world. *C. R. Acad. Sci. III* 324, 1067–1076.
- Grantham, R., Gautier, C., Gouy, M., Mercier, R., Pave, A., 1980. Codon catalog usage and the genome hypothesis. *Nucleic Acids Res.* 8, r49–r62.
- Grantham, R., Gautier, C., Gouy, M., Jacobzone, M., Mercier, R., 1981. Codon catalog usage is a genome strategy modulated for gene expressivity. *Nucleic Acids Res.* 9, r43–r74.
- Gu, W.J., Ma, J.M., Zhou, T., Sun, X., Lu, Z.H., 2002. Studies on the codon usage bias of genes coding for proteins with different tertiary structure. *Acta Biophys. Sin.* 18, 81–86.
- Gutman, G.A., Hatfield, G.W., 1989. Nonrandom utilization of codons pairs in *Escherichia coli*. *Proc. Natl Acad. Sci. USA* 86, 3699–3703.
- Hanai, R., Wada, A., 1988. The effects of guanine and cytosine variation on dinucleotide frequency and amino acid composition in the human genome. *J. Mol. Evol.* 27, 321–325.
- Hartman, H., 1995. Speculations on the origin of the genetic code. *J. Mol. Evol.* 40, 541–544.
- Ikemura, T., 1981. Correlation between the abundance of *Escherichia coli* transfer RNAs and the occurrence of the respective codons in its protein genes: a proposal for a synonymous codon choice that is optimal for the *E. coli* translational system. *J. Mol. Biol.* 151, 389–409.
- Ikemura, T., 1985. Codon usage and tRNA content in unicellular and multicellular organisms. *Mol. Biol. Evol.* 2, 12–34.
- Jukes, T.H., Bhushan, V., 1986. Silent nucleotide substitutions and G+C content of some mitochondrial and bacterial genes. *J. Mol. Evol.* 24, 39–44.
- Koch, A.J., Lehmann, J., 1997. About a symmetry of the genetic code. *J. theor. Biol.* 189, 171–174.
- Lacan, J., Michel, C.J., 2001. Analysis of a circular code model. *J. theor. Biol.* 213, 159–170.
- Ma, J., Zhou, T., Gu, W., Sun, X., Lu, Z., 2002. Cluster analysis of the codon use frequency of MHC genes from different species. *BioSystems* 65, 199–207.
- Ng, W.V., et al., 2000. Genome sequence of *Halobacterium* species NRC-1. *Proc. Natl Acad. Sci. USA* 97, 12176–12181.
- Saks, M.E., Sampson, J.R., 1995. Evolution of tRNA recognition systems and tRNA gene sequences. *J. Mol. Evol.* 40, 509–518.
- Schimmel, P., Giegé, R., Moras, D., Yokoyama, S., 1993. An operational RNA code for amino acids and possible relationship to genetic code. *Proc. Natl Acad. Sci. USA* 90, 8763–8768.
- Sharp, P.M., Matassi, G., 1994. Codon usage and genome evolution. *Curr. Opin. Genet. Dev.* 4, 851–860.
- Shpaer, E.G., 1986. Constraints on codon context in *Escherichia coli* genes. Their possible role in modulating the efficiency of translation. *J. Mol. Biol.* 188, 555–564.
- Smith, T.F., Ralph, W.W., Goodman, M., Czelusniak, J., 1985. Codon usage in the vertebrate haemoglobins and its implications. *Mol. Biol. Evol.* 2, 390–398.
- Sueoka, N., 1992. Directional mutation pressure, selection constraints, and genetic equilibria. *J. Mol. Evol.* 34, 95–114.
- Tazi, J., Bird, A., 1990. Alternative chromatin structure at CpG islands. *Cell* 60, 909–920.
- Woese, C.R., 2000. Interpreting the universal phylogenetic tree. *Proc. Natl Acad. Sci. USA* 97, 8392–8396.
- Yarus, M., Folley, L.S., 1984. Sense codons are found in specific contexts. *J. Mol. Biol.* 182, 529–540.