



HAL
open science

Evolutionary conservation and functional implications of circular code motifs in eukaryotic genomes

Gopal Dila, Christian Michel, Olivier Poch, Raymond Ripp, Julie D Thompson

► **To cite this version:**

Gopal Dila, Christian Michel, Olivier Poch, Raymond Ripp, Julie D Thompson. Evolutionary conservation and functional implications of circular code motifs in eukaryotic genomes. *BioSystems*, 2019, 175, pp.57-74. 10.1016/j.biosystems.2018.10.014 . hal-02381474

HAL Id: hal-02381474

<https://hal.science/hal-02381474>

Submitted on 4 Dec 2020

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

Evolutionary conservation and functional implications of circular code motifs in eukaryotic genomes

Gopal Dila, Christian J. Michel, Olivier Poch, Raymond Ripp, Julie D. Thompson

CSTB, ICube

CNRS, University of Strasbourg

300 Boulevard Sébastien Brant

67400 Illkirch, France

Email: d.gopal@outlook.com, c.michel@unistra.fr, olivier.poch@unistra.fr, raymond.ripp@unistra.fr,
thompson@unistra.fr

Keywords: circular code motifs; genome evolution; genetic code; gene expression.

Abstract

A set X of 20 trinucleotides has been found to have the highest average occurrence in the reading frame, compared to the two shifted frames, of genes of bacteria, archaea, eukaryotes, plasmids and viruses (Michel, 2017, 2015; Arquès and Michel, 1996). This set X has an interesting mathematical property, since X is a maximal C^3 self-complementary trinucleotide circular code (Arquès and Michel, 1996). Furthermore, any motif obtained from this circular code X has the capacity to retrieve, maintain and synchronize the reading frame in genes. In a recent study of the X motifs in the complete genome of the yeast, *Saccharomyces cerevisiae*, it was shown that they are significantly enriched in the reading frame of the genes (protein-coding regions) of the genome (Michel *et al.*, 2017). It was suggested that these X motifs may be evolutionary relics of a primitive code originally used for gene translation. The aim of this paper is to address two questions: are X motifs conserved during evolution? and do they continue to play a functional role in the processes of genome decoding and protein production? In a large scale analysis involving complete genomes from four mammals and nine different yeast species, we highlight specific evolutionary pressures on the X motifs in the genes of all the genomes, and identify important new properties of X motif conservation at the level of the encoded amino acids. We then compare the occurrence of X motifs with existing experimental data concerning protein expression and protein production, and report a significant correlation between the number of X motifs in a gene and increased protein abundance. In a general way, this work suggests that motifs from circular codes, i.e. motifs having the property of reading frame retrieval, may represent functional elements located within the coding regions of extant genomes.

1. Introduction

The same set X of trinucleotides (also known as codons) was identified in average in genes (reading frame) of bacteria, archaea, eukaryotes, plasmids and viruses (Michel 2017, 2015; Arquès and Michel, 1996). It contains the 20 following trinucleotides

$$X = \{AAC, AAT, ACC, ATC, ATT, CAG, CTC, CTG, GAA, GAC, \\ GAG, GAT, GCC, GGC, GGT, GTA, GTC, GTT, TAC, TTC\} \quad (1)$$

and codes the 12 following amino acids (three and one letter notation)

$$\mathcal{X} = \{Ala, Asn, Asp, Gln, Glu, Gly, Ile, Leu, Phe, Thr, Tyr, Val\} \\ = \{A, N, D, Q, E, G, I, L, F, T, Y, V\}. \quad (2)$$

This set X has several strong mathematical properties. In particular, it is self-complementary, i.e. 10 trinucleotides of X are complementary to the other 10 trinucleotides of X , e.g. $AAC \in X$ is complementary to $GTT \in X$, and it is a circular code. A circular code is defined as a set of words such that any motif obtained from this set, allows to retrieve, maintain and synchronize the original (construction) frame. Thus, the circular code X may represent a self-correcting property of the genetic code. Indeed, it has been proposed recently that the circular code X may participate in the regulation of gene transcription (El Houmami and Seligmann, 2017). Other correction properties of the genetic code

have also been proposed with roles in gene translation, including compensation of tRNA misloading (Seligmann, 2011), prevention of protein misfolding (Seligmann and Warthi, 2017), or termination of translation after ribosomal frameshifting (Seligmann and Pollock, 2004).

Motifs from the circular code X (denoted (1) above) having this frame retrieval property are called X motifs. Since 1996, the theory of circular codes in genes has mainly been developed by analysing the classes, the numbers and the mathematical properties of circular codes using probability-statistics, combinatorics and graph theory (reviews in Michel, 2008, and Fimmel and Strüngmann, 2018). More recently, the circular code theory was applied to the complete genome sequence of a living organism, namely the eukaryote *Saccharomyces cerevisiae* (Michel *et al.*, 2017). It was shown that X motifs from the circular code X (1) were significantly enriched in the genes (protein-coding regions) of the genome. The authors hypothesized that the X motifs may be evolutionary relics of a primitive code originally used for translation.

In this article, we describe a large-scale study of the X motifs in two independent sets of complete genomes. The first set is composed of four mammal genomes, representing highly evolved species and closely related genomes. The second set is built from nine yeast genomes, representing the simplest eukaryotes with more divergent genome sequences. Each set includes a well-studied and annotated 'reference' genome: the human genome for the first set and the *Saccharomyces cerevisiae* genome for the second set. We first highlight specific evolutionary pressures on the X motifs in the genes of both sets of genomes, and identify important new properties of X motif conservation at the level of the encoded amino acids. Thus, the 20 trinucleotides of the circular code X (1) are grouped according to the amino acids they encode, leading to a new hypothesis for the evolution of the genetic code where each amino acid was coded by the most constrained circular codes, namely strong comma-free and comma-free codes.

Then, we investigate the potential functional role of X motifs in the regulation of gene expression. To achieve this, we compare the occurrence of X motifs with existing experimental data concerning protein expression and production, and report a significant correlation between the number of X motifs in a gene and increased protein abundance. Taken together, the results represent compelling evidence suggesting that X motifs may indeed contribute to the complex mechanisms of protein synthesis in extant genomes.

2. Method

After recalling a few basic definitions of circular codes, we define three classes of motifs: the X motifs constructed from the circular code X (1) identified in genes, as well as non- X motifs and random motifs used to evaluate the significance of X motifs. The statistical analyses of these motifs are based on very simple statistics, namely frequencies and mean frequencies, leading to clear biological results.

2.1. Definitions of circular code

We recall a few definitions without detailed explanation (i.e. without examples and figures) that are necessary for understanding the main properties of the X motifs obtained from the trinucleotide circular code X identified in genes (Michel, 2017, 2015; Arquès and Michel, 1996).

Notation 1. Let us denote the nucleotide 4-letter alphabet $B = \{A, C, G, T\}$ where A stands for adenine, C stands for cytosine, G stands for guanine and T stands for thymine. The trinucleotide set over B is denoted by $B^3 = \{AAA, \dots, TTT\}$. The set of non-empty words (words, respectively) over B is denoted by B^+ (B^* , respectively).

Definition 1. A set $S \subseteq B^+$ is a *code* if, for each $x_1, \dots, x_n, y_1, \dots, y_m \in S, n, m \geq 1$, the condition $x_1 \cdots x_n = y_1 \cdots y_m$ implies $n = m$ and $x_i = y_i$ for $i = 1, \dots, n$.

Definition 2. Any non-empty subset of the code B^3 is a code and called *trinucleotide code*.

Definition 3. The genetic code is a trinucleotide code. It defines a surjective map $\mathcal{g}: \tilde{B}^3 \rightarrow P$ where $\tilde{B}^3 = B^3 \setminus \{TAA, TAG, TGA\}$ and P is the set of the 20 peptide components (amino acids). We also use the following notation: a sequence s of trinucleotides, i.e. a gene, codes a sequence noted $\mathcal{g}(s)$ of amino acids, i.e. a protein.

Example 1. $\mathcal{g}(GGA) = Gly$, $\mathcal{g}^{-1}(Gly) = \{GGA, GGC, GGG, GGT\}$ and $\mathcal{g}(GACATCCTG) = DIL$ where D, I and L are amino acids.

Definition 4. A trinucleotide code $X \subseteq B^3$ is *circular* if, for each $x_1, \dots, x_n, y_1, \dots, y_m \in X, n, m \geq 1, r \in B^*, s \in B^+$, the conditions $sx_2 \cdots x_n r = y_1 \cdots y_m$ and $x_1 = rs$ imply $n = m, r = \varepsilon$ (empty word) and $x_i = y_i$ for $i = 1, \dots, n$.

We briefly recall the proof used to determine whether a code is circular or not, with the most recent and powerful approach which relates an oriented (directed) graph to a trinucleotide code.

Definition 5. (Fimmel *et al.*, 2016). Let $X \subseteq B^3$ be a trinucleotide code. The directed graph $\mathcal{G}(X) = (V(X), E(X))$ associated with X has a finite set of vertices $V(X)$ and a finite set of oriented edges $E(X)$ (ordered pairs $[v, w]$ where $v, w \in X$) defined as follows:

$$\begin{cases} V(X) = \{N_1, N_3, N_1N_2, N_2N_3: N_1N_2N_3 \in X\} \\ E(X) = \{[N_1, N_2N_3], [N_1N_2, N_3]: N_1N_2N_3 \in X\} \end{cases}$$

The theorem below gives a relation between a trinucleotide code which is circular and its associated graph.

Theorem 1. (Fimmel *et al.*, 2016). Let $X \subseteq B^3$ be a trinucleotide code. The following statements are equivalent:

- (i) The code X is circular.
- (ii) The graph $\mathcal{G}(X)$ is acyclic.

We also recall the results that characterize the comma-free codes and the strong comma-free codes by the longest paths in their associated graphs.

Theorem 2. (Fimmel *et al.*, 2016). Let $X \subseteq B^3$ be a trinucleotide circular code. The following statements are equivalent:

- (i) X is comma-free.

(ii) The longest path in $\mathcal{G}(X)$ is of length at most 2.

Theorem 3. (Fimmel *et al.*, 2017). Let $X \subseteq B^3$ be a trinucleotide circular code. The following statements are equivalent:

(i) X is strong comma-free.

(ii) The longest path in $\mathcal{G}(X)$ is of length at most 1.

Thus, the reading frame is retrieved after the reading of 2 nucleotides with motifs from a strong comma-free code, of 3 nucleotides (1 trinucleotide) with motifs from a comma-free code and of at most 13 nucleotides (4 trinucleotides + 1 nucleotide) with motifs from circular codes.

The trinucleotide set X coding the reading frame in genes of bacteria, archaea, eukaryotes, plasmids and viruses is a maximal (20 trinucleotides) C^3 (the 2 shifted codes by permutation of X are also circular) self-complementary (10 trinucleotides of X are complementary to the 10 other trinucleotides of X) trinucleotide circular code (Michel 2017, 2015; Arquès and Michel, 1996).

2.2. Definition of X motifs, non- X motifs and random motifs

Definition 6. As in Michel *et al.* (2017), a X motif $m(X)$ constructed from the circular code X (1), is a word with cardinality $4 \leq c \leq 20$ trinucleotides and length $l \geq c \geq 4$ trinucleotides. Here, we consider only the X motifs $m(X)$ found in reading frame of genes.

Indeed, the X motifs $m(X)$ have a cardinality $c \leq 20$ trinucleotides as the circular code X (1) has 20 trinucleotides. The minimal length $l = 4$ trinucleotides was chosen based on the requirement for 13 nucleotides in order to retrieve the reading frame. The class of motifs of X with cardinality $c < 4$ are excluded here because they are mostly associated with the “pure” trinucleotide repeats often found in non-coding regions of the genome (Michel *et al.*, 2017; El Soufi and Michel, 2017).

The fundamental property of a X motif $m(X)$ is the ability to retrieve, synchronize and maintain the reading frame. Indeed, a window of 13 nucleotides located anywhere in a sequence generated from the circular code X (1) is sufficient to retrieve the reading (correct, construction) frame of the sequence. It is important to stress again that this window for retrieving the reading frame in a sequence can be located anywhere in the sequence, i.e. no other frame signal, including start and stop trinucleotides, is required to identify the reading frame.

Example 2. For the convenience of the reader, we give an example of a X motif $m(X) = m_1$ from the circular code X (1) in a sequence $s = \dots AAAGGTGCCGAAGCCCTGGAGGAAAAG \dots$. In the sequence s , there is a X motif $m_1 = GGTGCCGAAGCCCTGGAGGAA$ of cardinal $c = 5$ trinucleotides $\{CTG, GAA, GAG, GCC, GGT\}$ and length $l = 7$ trinucleotides. Note that m_1 cannot be extended to the left or to the right in s due to the presence of the periodic trinucleotide AAA (left) and the trinucleotide AAG (right) which both do not belong to the circular code X . Then, the reading frame of the sequence s can easily be deduced from the X motif m_1 : $\dots, AAA, GGT, GCC, GAA, GCC, CTG, GAG, GAA, AAG, \dots$

Definition 7. For simplification reasons, a non- X motif $m(\bar{X})$ is any word of any cardinality and length constructed from the nucleotide 4-letter alphabet $B = \{A, C, G, T\}$ except the X motifs $m(X)$ defined in

Definition 6. As for the X motifs $m(X)$, we only consider the non- X motifs $m(\bar{X})$ found in reading frame of genes.

Note that using this simplified notation, the class of motifs of X of cardinality $c < 4$ trinucleotides and length $l < 4$ trinucleotides belong to the non- X motifs $m(\bar{X})$.

In order to evaluate the statistical significance of X motifs in genes, we define, as in Michel *et al.* (2017), random motifs.

Definition 8. A R random motif $m(R)$ constructed from a random code R , is a word with cardinality $4 \leq c \leq 20$ trinucleotides and length $l \geq c \geq 4$ trinucleotides. A random code R is generated according to the properties of X , except its circularity property:

- (i) R has a cardinality equal to 20 trinucleotides;
- (ii) The total number of each nucleotide A, C, G and T in R is equal to 15 (note that $20 \times 3 = 15 \times 4$);
- (iii) R has no stop trinucleotides $\{TAA, TAG, TGA\}$ and no periodic trinucleotides $\{AAA, CCC, GGG, TTT\}$;
- (iv) R is not a circular code. Its associated graph $\mathcal{G}(R)$ is cyclic ($\mathcal{G}(R)$ being not shown).

As for the X motifs $m(X)$ and the non- X motifs $m(\bar{X})$, we only consider the R random motifs $m(R)$ found in reading frame of genes.

In order to obtain a statistically significant distribution of random codes, a set of 100 (different) random codes R are generated according to Definition 8. Examples of such random codes are given in Appendix in Michel *et al.* (2017).

Definition 9. We say that a letter $N_i \in B$ belongs to a motif m of length $l(m)$ if $1 \leq i \leq l(m)$.

Notation 2. $\mathcal{S} = \{X, \bar{X}, R\}$ denotes the three trinucleotide codes associated with the studied motifs $m(X)$, $m(\bar{X})$ and $m(R)$.

2.3. Multiple alignment of genes

In the following sections, we briefly recall the multiple alignment of genes and the notations used. A reference gene sequence $s_1 = \mathbb{R}$ (by convention here the reference sequence is the first sequence in the alignment) is aligned with its orthologous corresponding $n - 1$ genes s_2, \dots, s_n where $s_2, \dots, s_n \in B^+$. The genes s_1, s_2, \dots, s_n have respective lengths $|s_1|, |s_2|, \dots, |s_n|$. Note that orthologous genes originate from a common DNA ancestral sequence and diverged after a speciation event.

A gene multiple alignment s_1, s_2, \dots, s_n , $n \geq 2$, is a mapping z on the alphabet $(B \cup \{\varepsilon\})^n \setminus (\{\varepsilon\})^n$ whose projection on the 1st component is s_1 , up to the projection on the n th component is s_n . Thus, a gene multiple alignment z of letter length l is noted

$$z = \begin{pmatrix} \bar{N}_{11} & \cdots & \bar{N}_{l1} \\ \bar{N}_{12} & \cdots & \bar{N}_{l2} \\ \vdots & \vdots & \vdots \\ \bar{N}_{1n} & \cdots & \bar{N}_{ln} \end{pmatrix}$$

with the reference sequence $\mathbb{R} = s_1 = \bar{N}_{11}, \dots, \bar{N}_{l1}$, up to the sequence $s_n = \bar{N}_{1n}, \dots, \bar{N}_{ln}$ such that the nucleotide $\bar{N}_{ji} \in B \cup \{\varepsilon\}$ for $i = 1, \dots, n$ and $j = 1, \dots, l$ and, where ε being classically associated with the gap symbol "-" or ".". An aligned tuple $(\bar{N}_{j1}, \dots, \bar{N}_{jn})$ at the j th position such that $\bar{N}_{j1}, \bar{N}_{jn} \in B$ with

$\bar{N}_{j1} \neq \bar{N}_{ji}$ and $i \geq 2$ denotes the substitution of the j th nucleotide \bar{N}_{j1} of \mathbb{R} by the j th nucleotide \bar{N}_{ji} of s_i . An aligned tuple $(\bar{N}_{j1}, \dots, \bar{N}_{ji}, \dots, \bar{N}_{jn})$ such that $\bar{N}_{j1} \in B$ and $\bar{N}_{ji} \in \{\varepsilon\}$ with $i \geq 2$ denotes the deletion of the j th nucleotide \bar{N}_{j1} of \mathbb{R} . An aligned tuple $(\bar{N}_{j1}, \dots, \bar{N}_{ji}, \dots, \bar{N}_{jn})$ such that $\bar{N}_{j1} \in \{\varepsilon\}$ and $\bar{N}_{ji} \in B$ with $i \geq 2$ denotes the insertion of the j th nucleotide \bar{N}_{ji} of s_i .

The X motifs $m(X)$, the non- X motifs $m(\bar{X})$ and the R random motifs $m(R)$ located in the gene multiple alignment belong to the alphabet $B \cup \{\varepsilon\}$, i.e. they may contain gaps. Those which are located in the gene sequence s_i , $i = 1, \dots, n$, are noted $m(X, s_i)$, $m(\bar{X}, s_i)$ and $m(R, s_i)$, respectively, in particular, $m(X, \mathbb{R})$, $m(\bar{X}, \mathbb{R})$ and $m(R, \mathbb{R})$, respectively, in the reference gene \mathbb{R} . [Figure 1](#) shows an example of a part of a gene multiple alignment.

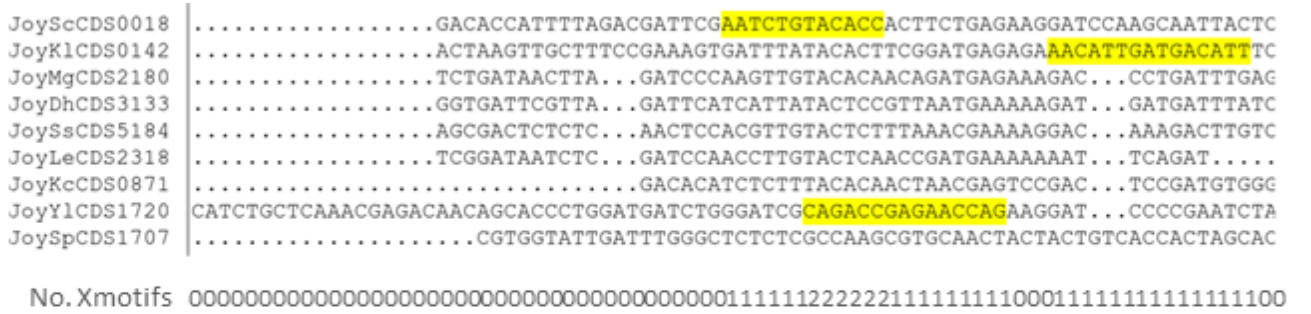


Figure 1. Screenshot of a yeast gene multiple alignment. The X motifs $m(X)$ in the reading frame of genes (Definition 6: cardinality $4 \leq c \leq 20$ trinucleotides and length $l \geq c \geq 4$ trinucleotides) are coloured in yellow. The reference (first) gene $s_1 = \mathbb{C}$ (*Saccharomyces Cerevisiae Sc*) contains one X motif $m(X, \mathbb{C})$ in reading frame and two non- X motifs $m(\bar{X}, \mathbb{C})$ in reading frame (Definition 7) without colour. Two X motifs are identified in the reading frame of two other yeast genes: $m(X, \mathbb{L})$ (*Kluyveromyces lactis Kl*, $s_2 = \mathbb{L}$) and $m(X, s_8)$ (*Yarrowia lipolytica Yl*). The number of X motifs (last row) is used in the calculation of the positional conservation parameter (Section 2.4).

2.4. Positional conservation parameter of X motifs and random motifs

Here, we consider whether the position of X motifs is preserved within the genes from different organisms. To do this, for each column of a gene alignment, the number of organisms with a X motif at this position was calculated. For example, in [Figure 1](#), the number of organisms with a X motif is equal to 0, 1, or 2. This number was normalized by the number of organisms having a nucleotide at that position in the alignment and not a gap.

Formally, we define a simple statistical parameter for analysing the positional conservation of motifs in the reference genes in the multiple alignments.

Definition 10. The positional conservation score $Ppc(m)$ of all motifs $m = (m(S, \mathbb{R}))$, $S \in \{X, R\}$ for studying X motifs and R random motifs, of letter lengths $l(m)$, m on the alphabet $B \cup \{\varepsilon\}$ (with gaps), in the reference genes $s_1 = \mathbb{R}$ in all the gene multiple alignments s_1, s_2, \dots, s_n is equal to

$$Ppc(m) = Ppc(m(S, \mathbb{R})) = \frac{1}{\sum_{m \in \mathbb{R}} l(m)} \sum_{m \in \mathbb{R}} \sum_{j=1}^{l(m)} \frac{1}{Nb_j} \sum_{i=1}^n \delta_{i,j}$$

where

$$\delta_{i,j} = \begin{cases} 1 & \text{if } N_{ji} \in B \text{ and } N_{ji} \in m(\mathcal{S}, s_i) \\ 0 & \text{otherwise} \end{cases},$$

Nb_j is the number of nucleotides without gaps at position j in the multiple alignment of n genes, $2 \leq Nb_j \leq n$ for $j = 1, \dots, l(m)$. The condition $N_{ji} \in B$ and $N_{ji} \in m(\mathcal{S}, s_i)$ signifies that the letter N_{ji} at the j th position in the gene s_i is a nucleotide and not a gap, and belongs to a motif m .

Remark 1. The positional score $Ppc(m) \in]0,1]$. The positional conservation Ppc of the motif m is the lowest in the alignment when $Ppc(m) \approx 0$, i.e. when the motif m in the reference genome is aligned with zero motifs in the other genomes. The positional conservation Ppc of the motif m is the highest in the alignment when $Ppc(m) = 1$ corresponding to the case where all genes without gaps have X motifs in the same position as the reference genes.

2.5. Pairwise alignment parameters of X motifs and non- X motifs

A pairwise alignment is a multiple alignment z with $n = 2$ sequences of letter length l such that their nucleotides $N \in B \cup \{\epsilon\}$ (with gaps). Several classical pairwise alignment parameters are used to estimate the conservation of a pairwise alignment, including (i) the percentage of alignment positions that contain identical nucleotides, and (ii) the ratio of synonymous to non-synonymous substitutions. These parameters are briefly recalled in the following definitions and are illustrated with examples from [Figure 2](#), and [Table 1](#) and [Table 2](#).

```
uc001acd.3_hg38      | CTACATCCCGGGCACGGACATCCTGGACCTGGAGAACCAGCGAGAAAACCTGGAGCAGCCATTCTGAGTGTGTTCA
uc001acd.3_tupBel1  | CTACATCCCTGGGACGGACATCCCAGGCCTGGACAGTCAGCGAGAGAACCCTGGAGCAGCCATTCTGAGTGTGTTCA
uc001acd.3_mm10     | CTACATCCCTGGGACGGACATCCCAGGCCTGGACAGTCAGCGAGAAAACCTGGAAACAGCCATTCTGAGTGTATTCA
uc001acd.3_canFam3  | CTACATCCCTGGGACGGACATCCCAGGCCTGGAGAGCCCGCGAGAAAACCTGGAAACAGCCATTCTGAGTGTGTTCA
```

Figure 2. Screenshot of a mammal gene multiple alignment. The X motifs $m(X)$ in the reading frame of genes (Definition 6: cardinality $4 \leq c \leq 20$ trinucleotides and length $l \geq c \geq 4$ trinucleotides) are coloured in yellow. The reference (first) gene $s_1 = \mathbb{H}$ (*Homo sapiens hg38*) contains two X motifs $m(X, \mathbb{H})$ in reading frame and three non- X motifs $m(\bar{X}, \mathbb{H})$ in reading frame (Definition 7) without colour. The 2nd non- X motif $m(\bar{X}, \mathbb{H})$ is composed of one trinucleotide $CGA \notin X$ (1). Three X motifs are identified in the reading frame of three other mammal genes: $m(X, s_2)$ (*Tupaia Belangeri tupBel1*), $m(X, \mathbb{M})$ (*Mus musculus mm10*, $s_3 = \mathbb{M}$) and $m(X, s_4)$ (*Canis Lupus Familiaris canFam3*).

Reference gene \mathbb{H}	1st X motif $m(X, \mathbb{H})$	2nd X motif $m(X, \mathbb{H})$
Protein $\varphi(s)$ of \mathbb{H}	<u>D</u> <u>I</u> <u>L</u> <u>D</u> <u>L</u> <u>E</u> <u>N</u> <u>Q</u>	<u>E</u> <u>N</u> <u>L</u> <u>E</u> <u>Q</u>
Gene s of \mathbb{H}	<u>GAC</u> <u>ATC</u> <u>CTG</u> <u>GAC</u> <u>CTG</u> <u>GAG</u> <u>AAC</u> <u>CAG</u>	<u>GAA</u> <u>AAC</u> <u>CTG</u> <u>GAG</u> <u>CAG</u>
Gene s' of \mathbb{M}	<u>GAC</u> <u>ATC</u> <u>CCG</u> <u>GGC</u> <u>CCA</u> <u>GAA</u> <u>CAT</u> <u>CAC</u>	<u>GAA</u> <u>AAC</u> <u>CTG</u> <u>GAA</u> <u>CAG</u>
Protein $\varphi(s')$ of \mathbb{M}	<u>D</u> <u>I</u> <u>P</u> <u>G</u> <u>P</u> <u>E</u> <u>H</u> <u>H</u>	<u>E</u> <u>N</u> <u>L</u> <u>E</u> <u>Q</u>

Table 1. From [Figure 2](#), the alignment of the two X motifs $m(X, \mathbb{H})$ in reading frame of total length 39 nucleotides of the reference gene $s_1 = \mathbb{H}$ (*Homo sapiens hg38*) and the gene \mathbb{M} (*Mus musculus mm10*). There are 30 identical nucleotide pairs and 9 different nucleotide pairs (underlined). The protein alignment associated with the gene alignment is given by applying the universal genetic code map φ (Definition 3) to each trinucleotide.

Reference gene \mathbb{H}	1st non- X motif $m(\bar{X}, \mathbb{H})$	2nd $m(\bar{X}, \mathbb{H})$	3rd non- X motif $m(\bar{X}, \mathbb{H})$
Protein $\varphi(s)$ of \mathbb{H}	<u>Y</u> <u>I</u> <u>P</u> <u>G</u> <u>T</u>	<u>R</u>	<u>P</u> <u>F</u> <u>L</u> <u>S</u> <u>V</u> <u>F</u>
Gene s of \mathbb{H}	TAC ATC CCG GGC ACG	CGA	CCA TTC CTG AGT GTG TTC
Gene s' of \mathbb{M}	TAC ATC CCT GGG ACG	<u>CCA</u>	CCA TTC CTG AGT <u>GTA</u> TTC
Protein $\varphi(s')$ of \mathbb{M}	<u>Y</u> <u>I</u> <u>P</u> <u>G</u> <u>T</u>	<u>P</u>	<u>P</u> <u>F</u> <u>L</u> <u>S</u> <u>V</u> <u>F</u>

Table 2. From [Figure 2](#), the alignment of the three non- X motifs $m(\bar{X}, \mathbb{H})$ in reading frame of total length 36 nucleotides of the reference gene $s_1 = \mathbb{H}$ (*Homo sapiens hg38*) and the gene \mathbb{M} (*Mus musculus mm10*). There are 32 identical nucleotide pairs and 4 different nucleotide pairs (underlined). The protein alignment associated with the gene alignment is given by applying the universal genetic code map φ (Definition 3) to each trinucleotide.

Definition 11. The percentage $Pid(m)$ of identical nucleotides of all motifs $m = (m(\mathcal{S}, \mathbb{R}))$, $\mathcal{S} \in \{X, \bar{X}\}$ for studying X motifs and non- X motifs, of letter lengths $l(m)$, m on the alphabet $B \cup \{\varepsilon\}$ (with gaps), in the reference genes $s_1 = \mathbb{R}$ in all the gene pairwise alignments s_1 and s_2 is equal to

$$Pid(m) = Pid(m(\mathcal{S}, \mathbb{R})) = \frac{1}{\sum_{m \in \mathbb{R}} l(m)} \sum_{m \in \mathbb{R}} \sum_{i=1}^{l(m)} \delta_i$$

where the operator δ_i , $1 \leq i \leq l(m)$, associated with a pair of letters N is defined by

$$\delta_i = \begin{cases} 1 & \text{if } N_{i1} \in B \text{ and } N_{i1} = N_{i2}. \\ 0 & \text{otherwise} \end{cases}$$

Example 3. From [Table 1](#), $Pid(m(X, \mathbb{H})) = 30/39 = 76.92\%$.

Definition 12. Let $f_i(c)$, $g_i(c)$ respectively, be the fraction of synonymous, non-synonymous respectively, potential substitutions at the i th site, $i = 1, 2, 3$, of a given codon $c = N_1N_2N_3$. Then, the numbers $Ns(c)$, $Nns(c)$ respectively, of synonymous, non-synonymous respectively, sites for a given codon c , are defined according to Nei and Gojobori (1986) by $Ns(c) = \sum_{i=1}^3 f_i(c)$ and $Nns(c) = \sum_{i=1}^3 g_i(c) = \sum_{i=1}^3 (1 - f_i(c)) = 3 - Ns(c)$.

Example 4. In the case of the codon $c = CTG$ coding the amino acid $\varphi(c) = Leu$, $f_1(Leu) = \frac{1}{3}$ as only the 1st site substitution $CTG \rightarrow TTG$ is synonymous out of ATG , GTG and TTG , $f_2(Leu) = 0$ as there is no 2nd site synonymous substitution out of CAG , CCG and CGG , and $f_3(Leu) = \frac{3}{3} = 1$ as all the 3rd site substitutions are synonymous out of CTA , CTC and CTT . Then, $Ns(Leu) = \frac{1}{3} + 0 + 1 = \frac{4}{3}$ and $Nns(Leu) = 3 - \frac{4}{3} = \frac{5}{3}$.

The definitions of $Ns(c)$ and $Nns(c)$ for a given codon are naturally extended to a motif m .

Definition 13. The potential numbers $Ns(m)$, $Nns(m)$ respectively, of synonymous, non-synonymous respectively, sites for a motif $m = (m(\mathcal{S}, \mathbb{R}))$, $\mathcal{S} \in \{X, \bar{X}\}$ for studying X motifs and non- X motifs, of letter length l , m on the alphabet B (without gaps), are equal to $Ns(m) = \sum_{c \in m} Ns(c)$ where $Ns(c)$ is defined in Definition 12, and $Nns(m) = l - Ns(m)$. Then, $Ns(m)$ and $Nns(m)$ are computed for all motifs m in the reference sequence s_1 of the gene pairwise alignments.

Example 5. From [Table 1](#), the potential numbers $Ns(m)$ and $Nns(m)$ of synonymous and non-synonymous sites for the two X motifs are $Ns(m(X, \mathbb{H})) = \frac{1}{3} + \frac{2}{3} + \frac{4}{3} + \frac{1}{3} + \frac{4}{3} + \frac{1}{3} + \frac{1}{3} + \frac{1}{3} + \frac{1}{3} + \frac{4}{3} + \frac{1}{3} + \frac{1}{3} = \frac{23}{3} \approx 7.67$ and $Nns(m(X, \mathbb{H})) = 39 - \frac{23}{3} = \frac{94}{3} \approx 31.33$.

Definition 14. Let $Os(m)$, $Ons(m)$ respectively, be the observed numbers of synonymous, non-synonymous respectively, substitutions of a reference motif $m = (m(\mathcal{S}, \mathbb{R}))$, $\mathcal{S} \in \{X, \bar{X}\}$ in the gene $s_1 = \mathbb{R}$ for studying X motifs and non- X motifs, in the motif m' in gene s_2 of letter lengths $l(m) = l(m')$, m, m' on the alphabet B (without gaps), in all gene pairwise alignments s_1 and s_2 .

Remark 2. $Os(m) + Ons(m) = l(m) - \sum_{i=1}^{l(m)} \delta_i$ where δ_i is defined in Definition 11.

Example 6. From [Table 1](#), $Os(m(X, \mathbb{H})) = 4$ (four synonymous substitutions: the 3rd site of $CTG (L)$, $GAG (E) \rightarrow GAA (E)$, the 3rd site of $AAC (N)$ and $GAG (E) \rightarrow GAA (E)$) and $Ons(m(X, \mathbb{H})) = 5$ (five non-synonymous substitutions: $CTG (L) \rightarrow CCG (P)$, $GAC (D) \rightarrow GGC (G)$, $CTG (L) \rightarrow CCA (P)$, $AAC (N) \rightarrow CAT (H)$ and $CAG (Q) \rightarrow CAC (H)$).

Definition 15. The percentages $Ps(m)$, $Pns(m)$ respectively, of synonymous, non-synonymous respectively, substitutions of a reference motif $m = (m(\mathcal{S}, \mathbb{R}))$, $\mathcal{S} \in \{X, \bar{X}\}$ in the gene $s_1 = \mathbb{R}$ for studying X motifs and non- X motifs, in the motif m' in the gene s_2 of letter lengths $l(m) = l(m')$, m, m' on the alphabet B (without gaps), in all gene pairwise alignments s_1 and s_2 , are equal to $Ps(m) = Os(m)/Ns(m)$ and $Pns(m) = Ons(m)/Nns(m)$ where $Os(m)$ and $Ons(m)$ are defined in Definition 14, and $Ns(m)$ and $Nns(m)$ in Definition 13.

Example 7. From [Table 1](#), $Ps(m(X, \mathbb{H})) = \frac{4}{23} = \frac{12}{23} \approx 0.52$ and $Pns(m(X, \mathbb{H})) = \frac{5}{94} = \frac{15}{94} \approx 0.16$.

Example 8 summarizes the parameters for the three non- X motifs of [Table 2](#).

Example 8. From [Table 2](#), $Pid(m(\bar{X}, \mathbb{H})) = 32/36 = 88.89\%$, $Ns(m(\bar{X}, \mathbb{H})) = \frac{1}{3} + \frac{2}{3} + \frac{3}{3} + \frac{3}{3} + \frac{3}{3} + \frac{4}{3} + \frac{3}{3} + \frac{1}{3} + \frac{4}{3} + \frac{1}{3} + \frac{3}{3} + \frac{1}{3} = \frac{29}{3} \approx 9.67$, $Nns(m(\bar{X}, \mathbb{H})) = 3 \times 12 - \frac{29}{3} = \frac{79}{3} \approx 26.33$, $Os(m(\bar{X}, \mathbb{H})) = 3$ (three synonymous substitutions: $CCG (P) \rightarrow CCT (P)$, $GGC (G) \rightarrow GGG (G)$ and $GTG (V) \rightarrow GTA (V)$ and $Ons(m(\bar{X}, \mathbb{H})) = 1$ (one non-synonymous substitution: $CGA (R) \rightarrow CCA (P)$), $Ps(m(\bar{X}, \mathbb{H})) = \frac{3}{29} = \frac{9}{29} \approx 0.31$ and $Pns(m(\bar{X}, \mathbb{H})) = \frac{1}{79} = \frac{3}{79} \approx 0.04$.

2.6. Codon substitution matrix of X motifs and random motifs

We define a codon (trinucleotide) substitution matrix $\mathbf{A}(m(\mathcal{S}, \mathbb{R}))$, $\mathcal{S} \in \{X, R\}$ for studying X motifs and R random motifs, in the reference genes $s_1 = \mathbb{R}$ in all the gene multiple alignments s_1, s_2, \dots, s_n . The codon substitution matrix $\mathbf{A}(m(\mathcal{S}, \mathbb{R})) = [a_{ij}]_{1 \leq i \leq 64, 1 \leq j \leq 64}$ of size 64×64 (square matrix) where the 64 rows and the 64 columns are associated with the 64 codons B^3 , has element $a_{ij} = Nb(\mathcal{S}[j] \rightarrow B^3[i])$ in row i and column j referring to the number of substitutions of codon $\mathcal{S}[j]$ (j th codon of \mathcal{S}) of the motifs m (in the reference genes \mathbb{R}) by the aligned codon $B^3[i]$ (i th codon of B^3) of the $n - 1$ genes s_2, \dots, s_n .

We define the normalized matrix $\mathbf{B}(m(\mathcal{S}, \mathbb{R})) = [b_{ij}]_{1 \leq i \leq 64, 1 \leq j \leq 64}$ with element $b_{ij} = a_{ij}/a_{.j}$ where $a_{.j} = \sum_{k=1}^{64} a_{kj}$ ($a_{.j} \neq 0$), for $1 \leq i \leq 64$ and $1 \leq j \leq 64$, such that it is stochastic in column. The column normalization, rather than a full matrix normalization, allows the codons to be compared whatever the codon usage.

Remark 3. The elements a_{ii} of \mathbf{A} and b_{ii} of \mathbf{B} can be indexed either by numbers or by the codons B^3 .

Remark 4. The diagonal elements a_{ii} of \mathbf{A} and b_{ii} of \mathbf{B} can be different from 0.

Remark 5. For a given code \mathcal{S} , i.e. the circular code X or a given random code R , with $Card(\mathcal{S}) = 20$ codons \mathcal{S} (Section 2.2), the matrices \mathbf{A} and \mathbf{B} have 20 non-empty codon columns and $64 - 20 = 44$ empty codon columns. However, the 20 codon columns of \mathbf{A} and \mathbf{B} vary with each (different) random code R , which obviously differ from the 20 codon columns of the circular code X .

Example 9. An example of construction of the matrices \mathbf{A} and \mathbf{B} is given from the alignment in [Table 3](#). The first codon column leads to the submatrix \mathbf{A} given in [Table 4](#). The procedure is iterated for each codon column and leads to the matrix \mathbf{A} given in [Table 5](#). The normalized matrix \mathbf{B} is given in [Table 6](#).

s_1	GAG	GAC	ATC	CTG	GAC	CTG	AAC	CAG
s_2	GAC	GAC	ATC	CCA	GGC	CTG	AGT	CAG
s_3	GAA	GAC	ATC	CCG	GGC	CCA	CAT	CAC
s_4	GAG	GAC	ATC	CGG	GGC	CTG	AGC	CCG

Table 3. Example of a multiple alignment of four genes where $s_1 = \mathbb{R}$ is the reference gene.

\mathbf{A}	GAG
GAA	1 GAG \rightarrow GAA
GAC	1 GAG \rightarrow GAC
GAG	1 GAG \rightarrow GAG

Table 4. Codon substitution submatrix \mathbf{A} of the first codon column from example of [Table 3](#).

\mathbf{A}	AAC	ATC	CAG	CTG	GAC	GAG
AGC	1					
AGT	1					
ATC		3				
CAC			1			
CAG			1			
CAT	1					
CCA				2		
CCG			1	1		
CGG				1		
CTG				2		
GAA						1
GAC					3	1
GAG						1
GGC					3	

Table 5. Codon substitution matrix \mathbf{A} from example of [Table 3](#). The remaining codon rows and columns equal to 0 are not shown.

B	AAC	ATC	CAG	CTG	GAC	GAG
AGC	1/3					
AGT	1/3					
ATC		1				
CAC			1/3			
CAG			1/3			
CAT	1/3					
CCA				1/3		
CCG			1/3	1/6		
CGG				1/6		
CTG				1/3		
GAA						1/3
GAC					1/2	1/3
GAG						1/3
GGC					1/2	
Sum	1	1	1	1	1	1

Table 6. Normalized matrix **B** from example of [Table 3](#). The remaining codon rows and columns equal to 0 are not shown.

2.7. Amino acid conservation parameter of X motifs and random motifs

We define a simple statistical parameter for analysing the conservation of X motifs and R random motifs for the 12 amino acids \mathcal{X} (2) coded by the circular code X , in the reference genes $s_1 = \mathbb{R}$ in all the gene multiple alignments s_1, s_2, \dots, s_n .

Definition 16. The percentage $Paac(m(\mathcal{S}, \mathbb{R}), p)$ of conservation of X codons per amino acid p (peptide component) coded by all the motifs $m = (m(\mathcal{S}, \mathbb{R}))$, $\mathcal{S} \in \{X, R\}$ for studying X motifs and R random motifs, in the reference genes $s_1 = \mathbb{R}$ in all the gene multiple alignments s_1, s_2, \dots, s_n , is equal to

$$Paac(m(\mathcal{S}, \mathbb{R}), p) = \frac{1}{Card(\mathcal{G}^{-1}(p) \cap \mathcal{S})} \sum_{\substack{i,j \in \mathcal{G}^{-1}(p) \\ i,j \in \mathcal{S}}} b_{ij}(m(\mathcal{S}, \mathbb{R}))$$

where $p \in \mathcal{X} = \{A, D, E, F, G, I, L, N, Q, T, V, Y\}$ (2), $b_{ij}(m(\mathcal{S}, \mathbb{R}))$ is the element of the normalized matrix **B** of \mathcal{S} defined in Section 2.6 and the inverse genetic code map \mathcal{G}^{-1} defined in Definition 3.

Definition 17. The mean percentage $\bar{Paac}(m(\mathcal{S}, \mathbb{R}), \mathcal{X})$ of conservation of X codons in the 12 amino acids \mathcal{X} (2) coded by all the motifs $m = (m(\mathcal{S}, \mathbb{R}))$, $\mathcal{S} \in \{X, R\}$, in the reference genes $s_1 = \mathbb{R}$ in all the gene multiple alignments s_1, s_2, \dots, s_n , is equal to

$$\bar{Paac}(m(\mathcal{S}, \mathbb{R}), \mathcal{X}) = \frac{1}{Card(\mathcal{X})} \sum_{p \in \mathcal{X}} Paac(m(\mathcal{S}, \mathbb{R}), p)$$

where $p \in \mathcal{X} = \{A, D, E, F, G, I, L, N, Q, T, V, Y\}$ (2) and $Paac(m(\mathcal{S}, \mathbb{R}), p)$ defined in Definition 16.

Remark 6. The mean percentage \bar{Paac} give the same statistical weight for each amino acid.

Definition 18. To achieve a strong statistical significance, we use the information from the 100 random codes R (R_1, \dots, R_{100}), and not only one random code, the percentage $Paac(\bar{m}(R, \mathbb{R}), p)$ of conservation of X codons per amino acid p coded by the R mean random motifs $\bar{m}(R, \mathbb{R})$ in the reference genes $s_1 = \mathbb{R}$ in all the gene multiple alignments s_1, s_2, \dots, s_n , is equal to

$$Paac(\bar{m}(R, \mathbb{R}), p) = \frac{1}{\sum_{k=1}^{100} \delta_k} \sum_{k=1}^{100} Paac(m(R_k, \mathbb{R}), p)$$

where $p \in \mathcal{X} = \{A, D, E, F, G, I, L, N, Q, T, V, Y\}$ (2), $Paac(m(R_k, \mathbb{R}), p)$ defined in Definition 16 and $\delta_k = 1$ if $\mathcal{G}^{-1}(p) \cap R_k \neq \emptyset$ (i.e. the random code R_k can code the amino acid p) and $\delta_k = 0$ otherwise.

Remark 7. For the R mean random motifs $\bar{m}(R, \mathbb{R})$, we only analyse in the mean matrix $\bar{\mathbf{B}}$ the trinucleotides coding the 12 amino acids \mathcal{X} of the circular code X .

Definition 19. The mean percentage $\bar{Paac}(\bar{m}(R, \mathbb{R}), \mathcal{X})$ of conservation of X codons in the 12 amino acids \mathcal{X} (2) coded by the R mean random motifs $\bar{m}(R, \mathbb{R})$ in the reference genes $s_1 = \mathbb{R}$ in all the gene multiple alignments s_1, s_2, \dots, s_n , is equal to

$$\bar{Paac}(\bar{m}(R, \mathbb{R}), \mathcal{X}) = \frac{1}{Card(\mathcal{X})} \sum_{p \in \mathcal{X}} Paac(\bar{m}(R, \mathbb{R}), p)$$

where $p \in \mathcal{X} = \{A, D, E, F, G, I, L, N, Q, T, V, Y\}$ (2) and $Paac(\bar{m}(R, \mathbb{R}), p)$ defined in Definition 18.

2.8. Data

In order to increase the significance of the results, two classes of independent alignments are investigated. The first class of alignments is based on four mammal genomes, which represent highly evolved species and closely related genomes. The second class of alignments is built from nine yeast genomes, which represent the simplest eukaryotes and are more divergent. The human genome for the first class and the *Saccharomyces cerevisiae* genome for the second class are taken as reference genomes as they are well documented model organisms.

2.8.1. Mammal gene alignments

From the mammalian gene multiple alignments available on the UCSC site (https://bds.mpi-cbg.de/hillerlab/144VertebrateAlignment_CESAR/, Sharma and Hiller, 2017), we have used genes from four well annotated genomes. [Table 7](#) shows some summary statistics of the four selected mammal genomes.

Genome name	Identification	Number of genes	Nucleotide length of genes
<i>Canis lupus familiaris</i>	<i>canFam3</i>	21,137	34,379,490
<i>Homo sapiens</i>	<i>hg38</i> (H)	22,352	36,808,167
<i>Mus musculus</i>	<i>mm10</i> (M)	20,178	33,519,381
<i>Tupaia belangeri</i>	<i>tupBel1</i>	18,485	23,387,559

Table 7. Genomes of four mammals.

H. sapiens (*hg38*, H) is taken as the reference genome and is present in each of the 22,352 gene alignments. One or two corresponding genes from the three other species may be missing, in which case the corresponding genes are replaced by gaps in the alignment.

2.8.2. Yeast gene alignments

For the yeast multiple alignments, the protein sequences of nine different yeasts and the localization of the corresponding nucleic acid sequence on the chromosomes (Table 8) are obtained from the NCBI Genbank (<https://www.ncbi.nlm.nih.gov/genbank/>).

Genome name	Identification	Number of genes	Nucleotide length of genes
<i>Debaryomyces hansenii</i>	<i>Dh</i>	6288	7,506,066
<i>Kluyveromyces lactis</i>	<i>Kl</i> (L)	5085	7,729,998
<i>Kuraishia capsulata</i>	<i>Kc</i>	5989	6,911,424
<i>Lodderomyces elongisporus</i>	<i>Le</i>	5799	7,110,237
<i>Meyerozyma guilliermondii</i>	<i>Mg</i>	5920	6,633,972
<i>Saccharomyces cerevisiae</i>	<i>Sc</i> (C)	6008	8,246,529
<i>Scheffersomyces stipitis</i>	<i>Ss</i>	5818	6,991,422
<i>Schizosaccharomyces pombe</i>	<i>Sp</i>	4980	5,614,506
<i>Yarrowia lipolytica</i>	<i>Yl</i>	6472	6,762,072

Table 8. Genomes of nine yeasts.

S. cerevisiae (Sc, C) is taken as the reference genome. A BLAST (Altschul *et al.*, 1997) database of all protein sequences of these nine organisms is created. For all protein sequences of *S. cerevisiae*, a BLAST search in this database is performed. Then, the protein alignments containing from 2 to 9 sequences are obtained using ClustalW (Thompson *et al.*, 1994). The corresponding nucleic sequence alignments are created by localizing each amino acid on the genome. The BLAST searches, alignments and some data analyses were performed using our in-house software platform Gscope (R. Ripp, unpublished, details in Section 2.9).

2.9. Software development

In a nucleic sequence alignment, the *X* and *R* random motifs are localized in the genes using a program developed in the Java language (El Soufi and Michel, 2017). The program takes optional parameters that define the minimum cardinality *c* (in trinucleotides) and the length *l* (in trinucleotides) of the *X* and *R* motifs searched. The *X* and *R* motifs verify Definition 6 (cardinality $c \geq 4$ trinucleotides and with any length $l \geq c \geq 4$ trinucleotides). Although the *X* and *R* motifs are contiguous in the gene sequences, gaps may be inserted during the alignment process.

Gscope is an integrated platform allowing the analysis of all kinds of genomic data. It is written in Tcl/Tk and runs under all operating systems. It is specifically designed to perform high throughput analyses. Gscope includes the tools necessary to create the basic data, analysis tools and visualization interfaces. It also allows the creation of SQL relational databases and the querying and display of the available information through a web based interface (Wscope).

3. Results

The results presented below are based on basic frequency statistics and their biological significance is clear. In order to evaluate the statistical significance of the different results presented below, we chose an approach that involved comparing the results obtained for the X motifs with those obtained for R random motifs generated by 100 (different) random codes R . This approach avoids the problems associated with defining statistical hypotheses about the nucleotide composition, the length and the random model of the different regions of the genome. The main disadvantage of our approach is the additional computational resources required to obtain the results for the 100 random codes.

This section is divided into two main parts. In the first part, we estimate the evolutionary conservation of X motifs in two large-scale sets of genes from mammal and yeast species ($\sim 20,000$ and ~ 6000 genes, respectively). In the second part, we evaluate the potential functional activity of the X motifs, by correlating them with existing experimental data.

3.1. Evolutionary conservation of X motifs in mammal and yeast genes

3.1.1. Enrichment of X motifs in mammal and yeast genes

We first investigated the occurrence number and codon length of X motifs in the two sets of genes. [Figure 3](#) and [Figure 4](#) show a very strong enrichment of X motifs in both mammal and yeast genes compared to the R random motifs from the 100 (different) random codes R . The number of X motifs in mammal genes is equal to 173,390, compared to a mean number of 60,330 R motifs. This difference is significant according to a one-sided Student's t -test with value $p \approx 10^{-82}$. The number of X motifs in yeast genes is equal to 35,833, compared to a mean number of 15,853 R motifs. Again, this difference is significant according to a one-sided Student's t -test with value $p \approx 10^{-75}$. This result is an additional and strong confirmation of the enrichment of X motifs in genes previously observed in the yeast *S. cerevisiae* (Michel *et al.*, 2017).

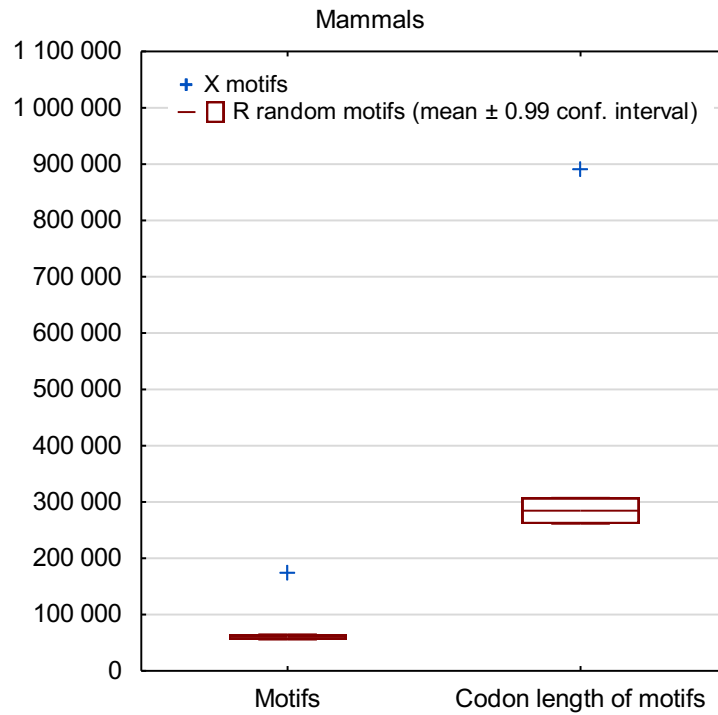


Figure 3. Comparison of the number of *X* and *R* random motifs and their codon lengths in mammalian genes. The number of *X* motifs is represented with a blue cross. The distribution of the *R* random motifs from the 100 random codes *R* is indicated by boxplots representing the mean and ± 0.99 confidence interval.

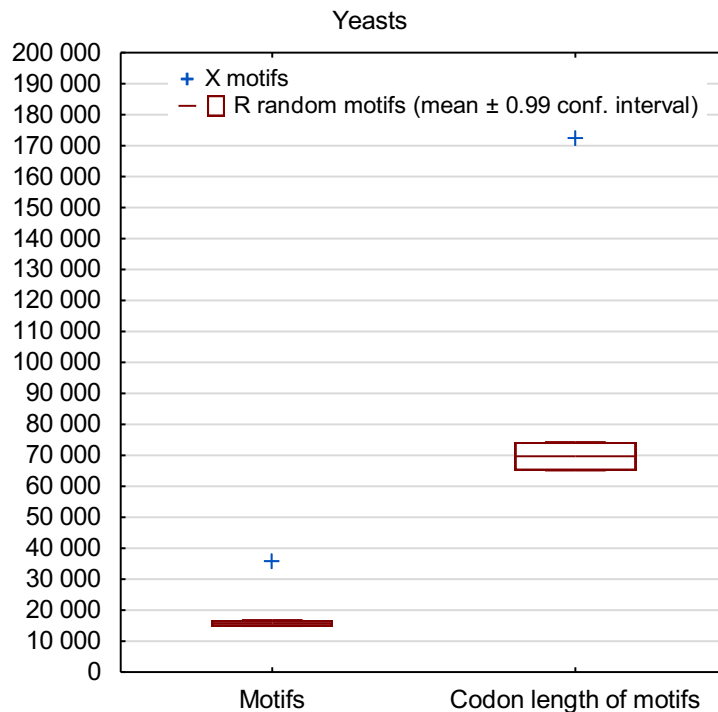


Figure 4. Comparison of the number of *X* and *R* random motifs and their codon lengths in the yeast genes. The number of *X* motifs is represented with a blue cross. The distribution of the *R* random motifs from the 100 random codes *R* is indicated by boxplots representing the mean and ± 0.99 confidence interval.

3.1.2. Positional conservation of X motifs in mammal and yeast genes

We then used the multiple alignments corresponding to the 22,352 mammal genes and the 6008 yeast genes, to calculate the positional conservation scores $Ppc(m)$ (Definition 10; $m = (m(\mathcal{S}, \mathbb{R}))$, $\mathcal{S} \in \{X, R\}$) for the X motifs and the R random motifs from the 100 (different) random codes R , shown in [Figure 5](#) and [Figure 6](#), respectively. The positional conservation score Ppc measures the number of motifs that are found in the same columns in a given multiple alignment. For both mammals and yeasts, the number of X motifs with the highest positional conservation score $Ppc = 1$ was higher than the number of R motifs. In contrast, the number of X motifs with the lowest positional conservation score $Ppc < 0.25$ was much lower than the number of R motifs. A one sample Wilcoxon signed rank indicated that the X motifs and the R motifs have significantly different medians with two-sided values $p = 0.031$ for the mammals and $p = 0.016$ for the yeasts.

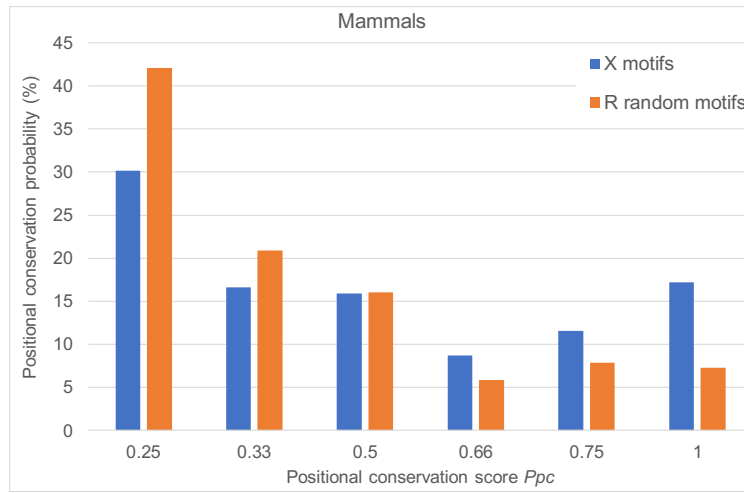


Figure 5. Positional conservation probability (%) of X motifs and R random motifs in the mammal gene multiple alignments with respect to the human reference genes \mathbb{H} as a function of the score $Ppc(m)$ (Definition 10; $m = (m(\mathcal{S}, \mathbb{H}))$, $\mathcal{S} \in \{X, R\}$) varying from 0 (no conservation in the alignment) to 1 (highest conservation in the alignment).

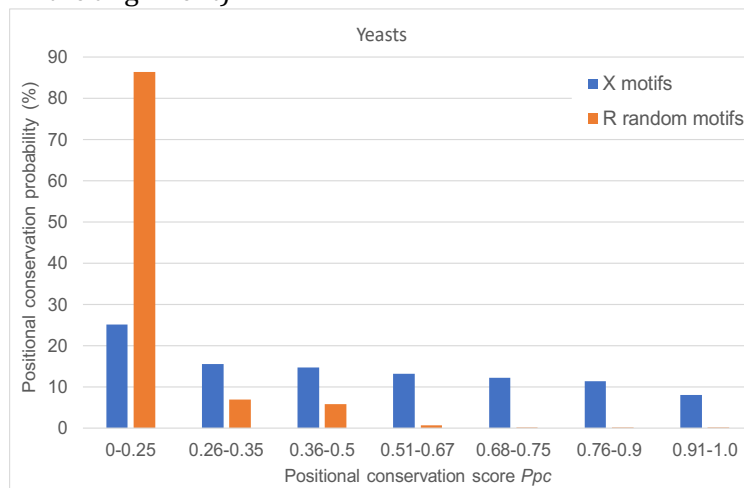


Figure 6. Positional conservation probability (%) of X motifs and R random motifs in the yeast gene multiple alignments with respect to the *S. cerevisiae* reference genes \mathbb{C} as a function of the score $Ppc(m)$ (Definition 10; $m = (m(\mathcal{S}, \mathbb{C}))$, $\mathcal{S} \in \{X, R\}$) varying from 0 (no conservation in the alignment) to 1 (highest conservation in the alignment).

We conclude that X motifs are more likely to be preserved in the same position in the orthologous genes of mammals and yeasts than R random motifs.

3.1.3. Sequence conservation of X motifs in mammal and yeast genes

In the previous section, we showed that X motifs tend to occur at the same positions in orthologous genes from different organisms. To investigate the level of sequence conservation within the X motifs that are found at the same position, we computed several classical pairwise alignment parameters that were defined in Section 2.5.

We first calculated the percentage $Pid(m)$ (Definition 11; $m = (m(\mathcal{S}, \mathbb{R}))$, $\mathcal{S} \in \{X, \bar{X}\}$), of identical nucleotides in the aligned X motifs and compared it to the percentage of identical nucleotides in the aligned non- X motifs (alignment columns with no X motifs). For this initial analysis, we selected two organisms from each of the mammal and yeast gene sets. For the mammals, 14,681 gene pairwise alignments containing both human \mathbb{H} and mouse \mathbb{M} genes were used, and for the yeasts, 1088 gene pairwise alignments containing both *S. cerevisiae* \mathbb{C} and *K. lactis* \mathbb{L} genes were used. The Pid observed in X motifs was 87.44% for \mathbb{H} - \mathbb{M} alignments, and 59.88% for \mathbb{C} - \mathbb{L} alignments. In comparison, the Pid observed in non- X motifs was 77.56% for \mathbb{H} - \mathbb{M} alignments, and 53.94% for \mathbb{C} - \mathbb{L} alignments. For \mathbb{H} - \mathbb{M} alignments ($n = 14,681$), a χ^2 test shows a strongly significant difference between the Pid values of X motifs (87.44%) and non- X motifs (77.56%) with one-sided value $p \approx 10^{-110}$. For \mathbb{C} - \mathbb{L} alignments ($n = 1088$), a χ^2 test shows a significant difference between the Pid values of X motifs (59.88%) and non- X motifs (53.94%) with a one-sided value $p \approx 0.005$. Thus, the sequences of X motifs are generally more conserved, i.e. evolve more slowly, than the remainder of the gene alignments.

This increased conservation of X motifs indicates that their sequences are maintained during natural selection, and may reflect differences in the strengths of positive selection or purifying selection. To understand the relative contributions of these different modes of selection better, we calculated the ratio $Pns(m)/Ps(m)$ (Definition 15; $m = (m(\mathcal{S}, \mathbb{R}))$, $\mathcal{S} \in \{X, \bar{X}\}$), of non-synonymous to synonymous substitutions for \mathbb{H} - \mathbb{M} and \mathbb{C} - \mathbb{L} alignments (Table 9 and Table 10). This ratio is commonly used to infer purifying ($Pns/Ps < 1$) or positive ($Pns/Ps > 1$) selection in genes. It is important to note that a non-synonymous substitution implies a change in the amino acid in the translated protein, while a synonymous substitution only changes the codon: the original codon is replaced by another codon coding for the same amino acid.

\mathbb{H} - \mathbb{M} alignment	Nns	Ns	Ons	Os	Pns	Ps	Pns/Ps
X motifs	1,611,224	480,358	99,670	184,643	0.06	0.38	0.16
Non- X motifs	19,772,931	8,225,136	1,524,889	2,572,797	0.08	0.31	0.25

Table 9. Comparison of non-synonymous and synonymous substitutions for X motifs and non- X motifs in pairs of aligned genes in human \mathbb{H} and mouse \mathbb{M} . $Nns(m)$, $Ns(m)$ respectively, are the potential numbers of non-synonymous, synonymous respectively, sites for the motifs $m = (m(\mathcal{S}, \mathbb{H}))$, $\mathcal{S} \in \{X, \bar{X}\}$, (Definition 13). $Ons(m)$, $Os(m)$ respectively, are the observed numbers of non-synonymous, synonymous respectively, substitutions of the motifs m (Definition 14). $Pns(m)$, $Ps(m)$ respectively, are

the percentages of non-synonymous, synonymous respectively, substitutions of the motifs m (Definition 15).

ℂ-ℒ alignment	Nns	Ns	Ons	Os	Pns	Ps	Pns/Ps
X motifs	369,426	93,981	103,766	80,081	0.28	0.85	0.33
Non- X motifs	5,310,908	1,973,266	1,580,781	1,362,399	0.30	0.69	0.43

Table 10. Comparison of non-synonymous and synonymous substitutions for X motifs and non- X motifs in pairs of aligned genes in *S. cerevisiae* ℂ and *K. lactis* ℒ. $Nns(m)$, $Ns(m)$ respectively, are the potential numbers of non-synonymous, synonymous respectively, sites for the motifs $m = (m(\mathcal{S}, \mathcal{C}))$, $\mathcal{S} \in \{X, \bar{X}\}$, (Definition 13). $Ons(m)$, $Os(m)$ respectively, are the observed numbers of non-synonymous, synonymous respectively, substitutions of the motifs m (Definition 14). $Pns(m)$, $Ps(m)$ respectively, are the percentages of non-synonymous, synonymous respectively, substitutions of the motifs m (Definition 15).

When we compare the rates of non-synonymous and synonymous substitutions, Pns and Ps respectively, for ℍ-ℳ than ℂ-ℒ, the values of Pns and Ps are obviously lower in both X motifs and non- X motifs, due to the smaller phylogenetic distance between ℍ and ℳ. In other words, ℍ and ℳ are more closely related than ℂ and ℒ, so we would expect less substitutions, both synonymous and non-synonymous. Also, the values of Pns are lower than Ps for X motifs and non- X motifs in both sets of genes. Again, this is expected since non-synonymous substitutions have a larger effect on the translated protein and occur less often than synonymous substitutions.

Importantly, we observe significantly lower ratios Pns/Ps of non-synonymous to synonymous substitutions in X motifs than in non- X motifs, suggesting more evolutionary constraints on X motifs. Furthermore, while the proportion Pns of non-synonymous substitutions is lower between X motifs and non- X motifs, the proportion Ps of synonymous substitutions is higher in X motifs than in non- X motifs. This result motivated the studies presented in the next section, which were designed to analyse in more detail the specific selective constraints in X motifs.

3.1.4. Synonymous substitutions of trinucleotides in X motifs

Given a multiple alignment of orthologous gene sequences, our goal is to determine whether trinucleotides in X motifs are conserved beyond what would be expected by chance if they were evolving only under the selective pressure on the amino acid they encode. Therefore, we chose to consider only those positions in the alignment with a conserved amino acid, i.e. involving only synonymous substitutions.

We first calculated the codon substitution matrices $\mathbf{A}(m)$ (defined in Section 2.6) for the X motifs in all the mammal and yeast gene multiple alignments. These two matrices $\mathbf{A}(m)$ are shown in Appendix (Table 19 and Table 20). We then normalized the columns of $\mathbf{A}(m)$ to produce the two normalized codon substitution matrices $\mathbf{B}(m)$ (defined in Section 2.6) for the X motifs, and extracted the rows and columns of the matrices $\mathbf{B}(m)$ that correspond to the synonymous substitutions of the X codons, as shown in Appendix (Table 21 and Table 22). We also calculated the equivalent submatrices $\mathbf{B}(m)$ for the R random motifs. Finally, we calculated the percentages $Paac(m(X, \mathbb{R}), p)$ (Definition 16) of conservation

of X codons per amino acid $p \in \mathcal{X}$ in [Table 21](#) and [Table 22](#), for the two mammal and yeast gene multiple alignments, as summarized in [Table 11](#) and [Table 12](#). In addition, we provide in these [Table 11](#) and [Table 12](#), the mean percentages $\bar{P}aac(m(X, \mathbb{R}), \mathcal{X})$ (Definition 17) of conservation of the 12 amino acids \mathcal{X} (2) for mammals and yeasts. For comparison, the values of $Paac(\bar{m}(R, \mathbb{R}), p)$ (Definition 18) and $\bar{P}aac(\bar{m}(R, \mathbb{R}), \mathcal{X})$ (Definition 19) calculated for the R mean random motifs from the 100 random codes R are reported. As these $Paac$ and $\bar{P}aac$ values for R motifs are mean values, their distributions with a ± 0.99 confidence interval are shown in [Figure 7](#) and [Figure 8](#).

	Mean	A	D	E	F	G	I	L	N	Q	T	V	Y
$m(X, \mathbb{H})$	78.1	66.1	89.4	90.3	78.9	77.3	84.9	78.3	85.6	80.7	63.5	65.7	76.1
$\bar{m}(R, \mathbb{H})$	67.6	60.2	73.0	76.8	77.7	68.1	68.0	67.1	70.5	71.3	58.1	62.6	74.6

Table 11. For the mammal gene multiple alignments with respect to the human reference genes $s_1 = \mathbb{H}$, mean percentage $\bar{P}aac(m(X, \mathbb{H}), \mathcal{X})$ (Definition 17) of conservation of X codons for the 12 amino acids \mathcal{X} (2) coded by the X motifs and percentages $Paac(m(X, \mathbb{H}), p)$ (Definition 16) of conservation per amino acid $p \in \mathcal{X}$ (first row). Mean percentage $\bar{P}aac(\bar{m}(R, \mathbb{H}), \mathcal{X})$ (Definition 19) of conservation of X codons for the 12 amino acids \mathcal{X} (2) coded by the R mean random motifs (from the 100 random codes R) and percentages $Paac(\bar{m}(R, \mathbb{H}), p)$ (Definition 18) of conservation per amino acid $p \in \mathcal{X}$ (second row).

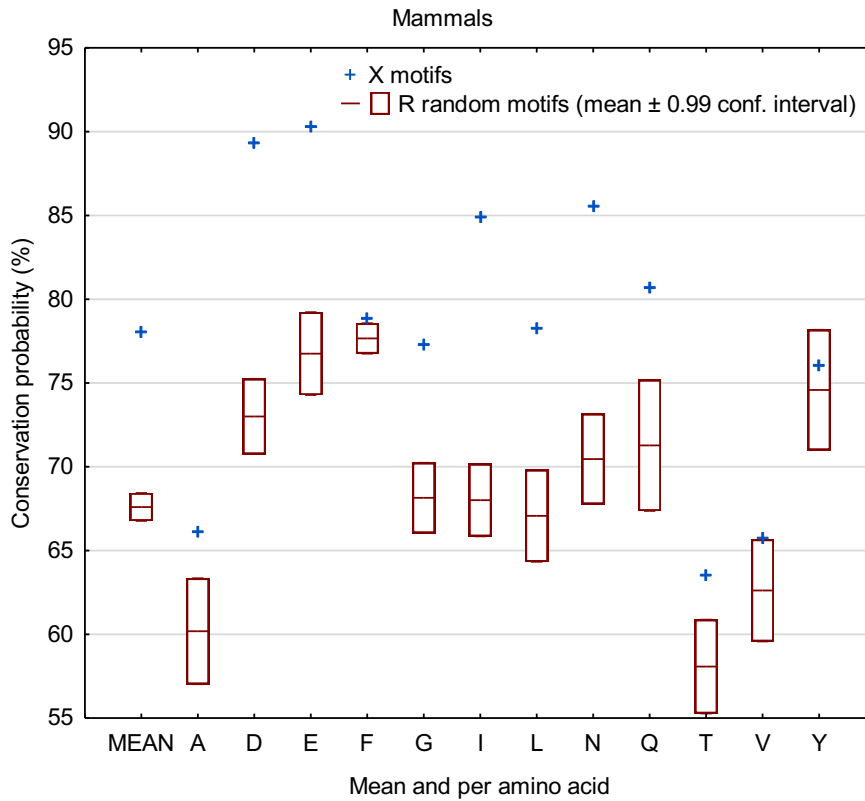


Figure 7 (associated with [Table 11](#)). For the mammal gene multiple alignments with respect to the human reference genes $s_1 = \mathbb{H}$, mean percentage $\bar{P}aac(m(X, \mathbb{H}), \mathcal{X})$ (Definition 17) of conservation of X codons for the 12 amino acids \mathcal{X} (2) coded the X motifs and percentages $Paac(m(X, \mathbb{H}), p)$ (Definition 16) of conservation per amino acid $p \in \mathcal{X}$ (blue cross). The distribution of the R mean random motifs (from the 100 random codes R) is indicated by boxplots representing the mean percentage $\bar{P}aac(\bar{m}(R, \mathbb{H}), \mathcal{X})$ (Definition 19) and percentages $Paac(\bar{m}(R, \mathbb{H}), p)$ (Definition 18) with a ± 0.99 confidence interval.

A new and strong property is identified with the X motifs of the circular code X . The average percentage ($\bar{P}aac$) conservation of X codons is significantly higher in X motifs than the conservation observed in the R mean random motifs in the mammal gene alignments (one-sided Student's t -test with $p \approx 10^{-55}$) ([Table 11](#) and [Figure 7](#)). Furthermore, this is true for 11 out of 12 amino acids (percentage $Paac$). For the amino acid Y , the conservation of X codons in X motifs is higher than in R motifs although the difference is not significant at 0.99. This new property can be formalized simply by the following inequalities:

$$\begin{cases} \bar{P}aac(m(X, \mathbb{H}), X) > \bar{P}aac(\bar{m}(R, \mathbb{H}), X) \\ Paac(m(X, \mathbb{H}), p) > Paac(\bar{m}(R, \mathbb{H}), p) \quad \forall p \in X \end{cases}$$

	Mean	A	D	E	F	G	I	L	N	Q	T	V	Y
$m(X, \mathbb{C})$	29.3	16.1	45.3	41.8	29.9	40.4	39.3	10.6	33.5	13.6	14.9	33.6	33.0
$\bar{m}(R, \mathbb{C})$	22.3	19.1	26.9	25.2	29.8	27.1	20.3	21.8	22.3	20.9	15.6	18.2	33.0

Table 12. For the yeast gene multiple alignments with respect to the *S. cerevisiae* reference genes $s_1 = \mathbb{C}$, mean percentage $\bar{P}aac(m(X, \mathbb{C}), X)$ (Definition 17) of conservation of X codons for the 12 amino acids X (2) coded by the X motifs and percentages $Paac(m(X, \mathbb{C}), p)$ (Definition 16) of conservation per amino acid $p \in X$ (first row). Mean percentage $\bar{P}aac(\bar{m}(R, \mathbb{C}), X)$ (Definition 19) of conservation of X codons for the 12 amino acids X (2) coded by the R mean random motifs (from the 100 random codes R) and percentages $Paac(\bar{m}(R, \mathbb{C}), p)$ (Definition 18) of conservation per amino acid $p \in X$ (second row).

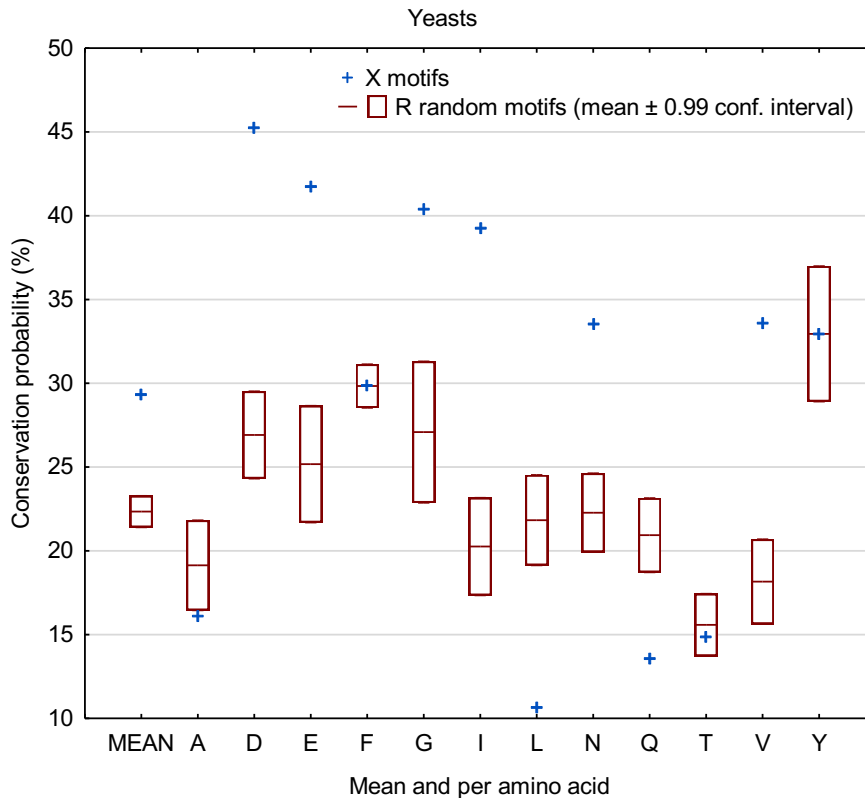


Figure 8 (associated with [Table 12](#)). For the yeast gene multiple alignments with respect to the *S. cerevisiae* reference genes $s_1 = \mathbb{C}$, mean percentage $\bar{P}aac(m(X, \mathbb{C}), X)$ (Definition 17) of conservation of X codons for the 12 amino acids X (2) coded by the X motifs and percentages $Paac(m(X, \mathbb{C}), p)$ (Definition 16) of conservation per amino acid $p \in X$ (blue cross). The distribution of the R mean random motifs (from the 100 random codes R) is indicated by boxplots representing the mean

percentage $\bar{P}aac(\bar{m}(R, \mathbb{C}), \mathcal{X})$ (Definition 19) and percentages $Paac(\bar{m}(R, \mathbb{C}), p)$ (Definition 18) with a ± 0.99 confidence interval.

The average percentage ($\bar{P}aac$) conservation of X codons is significantly higher in X motifs than in the R mean random motifs in the yeast gene alignments (one-sided Student's t -test with $p \approx 10^{-35}$) (Table 12 and Figure 8). For 6 out of 12 amino acids, the conservation (percentage $Paac$) of X codons in X motifs is higher than in R motifs. In contrast, for the amino acids A, L, Q and T , the conservation is lower than in the R motifs. For the amino acids F and L , the conservation is similar to the R motifs. This property can also be summarized by the following inequalities:

$$\begin{cases} \bar{P}aac(m(X, \mathbb{C}), \mathcal{X}) > \bar{P}aac(\bar{m}(R, \mathbb{C}), \mathcal{X}) \\ Paac(m(X, \mathbb{C}), p) > Paac(\bar{m}(R, \mathbb{C}), p) \quad \forall p \in \mathcal{X} \setminus \{A, L, Q, T\} \end{cases}$$

The conservation values of the motifs observed in the yeast alignments is lower than in the human alignments. This is not surprising since it is known that the yeasts diverged much earlier (more synonymous and non-synonymous substitutions) than the mammals included in this study. This evolutionary diversity in yeasts may also explain the exception with the four amino acids observed with this simple statistical parameter $Paac$. It should also be stressed that the identified conservation property of X codons in X motifs with respect to the amino acids is independent of the codon usage, the GC content, the nucleotide composition, the length of genes, etc.

3.1.5. A new hypothesis of evolution of the genetic code: union of circular codes associated with each amino acid

The statistical analyses performed in the previous sections show that the 20 trinucleotides of the circular code X (1) are strongly linked to the amino acids they encode, thus, leading to a partition of the 20 trinucleotides of X into 12 trinucleotide classes, each trinucleotide class being associated with an amino acid $p \in \mathcal{X}$ (2). This property leads us to propose that the extant genetic code may result from a union of circular codes: the subcodes of the circular code X (1) associated with each amino acid (Table 13 and Figure 9). Remember that a subcode of a circular code, is also circular. Interestingly, classes of circular codes with the strongest constraints of reading frame retrieval, i.e. the strong comma-free and comma-free codes (Theorem 3 and Theorem 2), can code an amino acid. Using these theorems, we determine the circular class of each trinucleotide code involved in Table 13 (an initial approach developed in Michel, 2014, Section 3.4, Table 6).

The evolution of the genetic code may thus have started from the circular codes with the strongest constraints, i.e. the strong comma-free codes and the comma-free codes with motifs retrieving the reading frame after the reading of 2 and 3 nucleotides, i.e. a nucleotide length of a codon or anticodon. It is tempting to suggest that these circular codes may have emerged independently in different "primitive soups". However, such strongly constrained coding systems may not have been viable in the long term. By relaxing the constraints, they may have evolved to circular codes having flexible motifs for retrieving the reading frame after the reading of at most 13 nucleotides, and to non-circular codes

without the ability to retrieve the reading frame. Among the 12 amino acids \mathcal{X} (2) coded by the circular code X (1), 10 amino acids are coded by strong comma-free codes and 2 amino acids E_X and L_X of X , by comma-free codes (Table 13 and Figure 9). In the extant genetic code, only 3 amino acids D , N and Q are still coded by strong comma-free codes, 6 amino acids A , E , I , T , V and Y , by comma-free codes, 1 amino acid L , by a circular code, and 2 amino acids F and G , by simple codes (not circular). The union of circular codes allows to extend the amino acid coding. For example, the union of the strong comma-free code $Q_X = \{CAG\}$ of X and the strong comma-free code $\{CAA\}$ leads to the strong comma-free code $Q = \{CAA, CAG\}$ of the genetic code, etc. Obviously, the union of 2 comma-free codes does not imply that the resulting code is comma-free, see for example the case of the amino acid L (Table 13). The 8 remaining amino acids could have been generated by mutations in circular codes. The extant genetic code is a code from a mathematical point of view (Definition 3), however it is not circular, i.e. it does not have the ability to retrieve the reading frame in genes after its circularity property loss.

AA	Circular code X	Class	Union	Class	Genetic code	Class
Asn	$N_X = \{AAC, AAT\}$	SCF			$N = \{AAC, AAT\}$	SCF
Asp	$D_X = \{GAC, GAT\}$	SCF			$D = \{GAC, GAT\}$	SCF
Gln	$Q_X = \{CAG\}$	SCF	$\{CAA\}$	SCF	$Q = \{CAA, CAG\}$	SCF
Glu	$E_X = \{GAA, GAG\}$	CF			$E = \{GAA, GAG\}$	CF
Phe	$F_X = \{TTC\}$	SCF	$\{TTT\}$	NC	$F = \{TTC, TTT\}$	NC
Tyr	$Y_X = \{TAC\}$	SCF	$\{TAT\}$	CF	$Y = \{TAC, TAT\}$	CF
Ile	$I_X = \{ATC, ATT\}$	SCF	$\{ATA\}$	CF	$I = \{ATA, ATC, ATT\}$	CF
Ala	$A_X = \{GCC\}$	SCF	$\{GCA, GCG, GCT\}$	CF	$A = \{GCA, GCC, GCG, GCT\}$	CF
Gly	$G_X = \{GGC, GGT\}$	SCF	$\{GGA, GGG\}$	NC	$G = \{GGA, GGC, GGG, GGT\}$	NC
Thr	$T_X = \{ACC\}$	SCF	$\{ACA, ACG, ACT\}$	CF	$T = \{ACA, ACC, ACG, ACT\}$	CF
Val	$V_X = \{GTA, GTC, GTT\}$	SCF	$\{GTG\}$	CF	$V = \{GTA, GTC, GTG, GTT\}$	CF
Leu	$L_X = \{CTC, CTG\}$	CF	$\{CTA, CTT, TTA, TTG\}$	CF	$L = \{CTA, CTC, CTG, CTT, TTA, TTG\}$	C

Table 13. Classes of codes (non-circular NC, circular C, comma-free CF, strong comma-free SCF) of the 12 amino acids \mathcal{X} (2) with respect to the circular code X (1) and the universal genetic code.

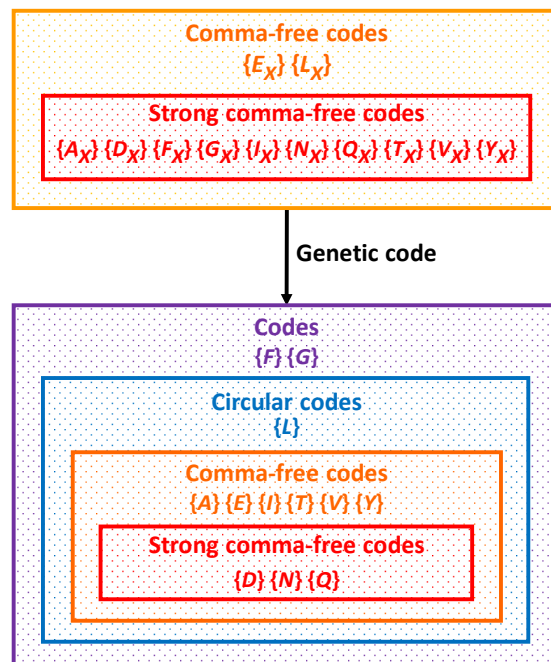


Figure 9 (associated with Table 13). Evolution of the genetic code by union of circular codes associated with each amino acid from the circular code X (1).

3.2. Functionality of *X* motifs in extant genomes

In the previous section, we identified specific evolutionary constraints, suggesting that the *X* motifs in the genomes included in this study have evolved under purifying selection. Indeed, the nucleotides in the *X* motifs display a considerable excess of synonymous substitutions compared to the non-*X* motifs. Furthermore, the average conservation of codons in *X* motifs is significantly higher than expected if the substitution process was random. These results suggest a possible functional role of *X* motifs, presumably as elements of the complex genome decoding system. In order to investigate the potential effects of *X* motifs on the translation of protein-coding genes, we compared the frequency of *X* motifs in the genes of the four mammalian and nine yeast species with existing experimental data on protein expression and protein production. We will show that the experimental data can generally be explained by circular code motifs, i.e. motifs having the property of reading frame retrieval.

3.2.1. Dicodons associated with reduced protein production are not located in *X* motifs

Recently, experimental studies in *S. cerevisiae* (Gamble *et al.*, 2016) were performed to investigate the effects of different codons on translation efficiency. The authors measured the expression levels of more than 35,000 synthetic protein variants in which three adjacent codons of the coding sequence were randomized. No individual codons had consistent effects on gene expression. However, 17 pairs of adjacent codons (called in the following dicodons) were identified that, when they were present in-frame in the coding sequence, reduced the expression level of the genes. This list is recalled in [Table 14](#). In this list, we identified the codons belonging to the circular code *X*.

Dicodon	Class	Dicodon	Class
<i>AGGCGA</i>	<i>NN</i>	<i>CGAGCG</i>	<i>NN</i>
<i>AGGCGG</i>	<i>NN</i>	<i>CTCCCG</i>	<i>XN</i>
<i>ATACGA</i>	<i>NN</i>	<i>CTGATA</i>	<i>XN</i>
<i>ATACGG</i>	<i>NN</i>	<i>CTGCCG</i>	<i>XN</i>
<i>CGAATA</i>	<i>NN</i>	<i>CTGCCA</i>	<i>XN</i>
<i>CGACCG</i>	<i>NN</i>	<i>GTACCG</i>	<i>XN</i>
<i>CGACGA</i>	<i>NN</i>	<i>GTACGA</i>	<i>XN</i>
<i>CGACGG</i>	<i>NN</i>	<i>GTGCCA</i>	<i>NN</i>
<i>CGACTG</i>	<i>NX</i>		

Table 14. List of the 17 dicodons that reduced the expression level of the genes (Gamble *et al.*, 2016) (1st and 3rd columns). Class of the dicodons according to its codons belonging to the circular code *X* (symbol *X*) or not (symbol *N*) (2nd and 4th columns).

Surprisingly, none of these 17 dicodons are composed of two *X* codons meaning that they cannot be located in a *X* motif.

3.2.2. Correlation of *X* motifs with dicodons associated with low and high protein production

Following the work of Gamble *et al.* (2016), Diambre (2017) performed a statistical analysis of dicodon usage frequencies over two sets of proteins: a low protein abundance (PA) set and a high PA set, from

nine diverse organisms including three prokaryotes, one plant, one yeast (*S. cerevisiae*), and two multicellular eukaryotes and two mammals. The working hypothesis was that sequences encoding abundant proteins should be optimized, in the sense of translation efficiency. He found an important bias of dicodon usage depending on PA and determined which dicodons were statistically associated with low or high abundance. These usage preferences cannot be explained by the frequency usage of the single codons. The statistical analysis of coding sequences of nine organisms reveals that in many cases dicodon preferences are shared between related organisms.

Dicodon	Class	Dicodon	Class
AAAATA	NN	CAGAAA	XN
AATGCA	XN	GAAAGT	XN
AATTGG	XN	GAACTA	XN
AGTAAG	NN	GCATTT	NN
AGTGTT	NN	TATAAA	NN
ATAGGT	NX	TATCCG	NN
ATTAAA	XN	TTTCAG	NX
CAAAGT	NN	TTTTTT	NN

Table 15. List of the 16 dicodons with low protein abundance (Diambra, 2017) (1st and 3rd columns). Class of the dicodons according to its codons belong to the circular code *X* (symbol *X*) or not (symbol *N*) (2nd and 4th columns).

In addition to the 17 previous dicodons (identified in *S. cerevisiae*), this study identified 16 new dicodons (Table 15) associated with low protein abundance in a number of different organisms. Again, these 16 dicodons cannot be located in a *X* motif. Thus, there are 33 low abundance dicodons that support the circular code theory.

Furthermore, the study revealed 40 dicodons shared between different organisms and preferentially used by high abundance proteins (Table 16). Importantly, 27 of these 40 dicodons (67.5%) are potentially in *X* motifs.

Dicodon	Class	Dicodon	Class	Dicodon	Class	Dicodon	Class
AACAAC	XX	ACCTTC	XX	GACACC	XX	GTCACC	XX
AACAAG	XN	ATCAAC	XX	GACTAC	XX	GTCATC	XX
AACACC	XX	ATCAAG	XN	GATGCT	XN	GTTGCC	XX
AAGTCC	NN	ATCACC	XX	GCCAAC	XX	TACAAC	XX
ACCAAC	XX	ATCATC	XX	GCCAAG	XN	TACAAG	XN
ACCAAG	XN	ATTGCC	XX	GCCACC	XX	TCCACC	NX
ACCACC	XX	CCACCA	NN	GCCATC	XX	TTCAAC	XX
ACCATC	XX	CGTCGT	NN	GCCGCC	XX	TTCAAG	XN
ACCATT	XX	GACAAC	XX	GGTGTC	XX	TTCACC	XX
ACCGCC	XX	GACAAG	XN	GTCAAG	XN	TTCATC	XX

Table 16. List of the 40 dicodons with high protein abundance (Diambra, 2017) (1st, 3rd, 5th and 7th columns). Class of the dicodons according to its codons belong to the circular code *X* (symbol *X*) or not (symbol *N*) (2nd, 4th, 6th and 8th columns).

3.2.3. Classification of genes as low or high abundance according to the circular code theory

The experimental and statistical results of Gamble *et al.* (2016) and Diambre (2017) ([Table 14](#), [Table 15](#) and [Table 16](#)) can be summarized in the following [Table 17](#).

	XX	{NN,NX,XN}	Total
Low abundance protein	0	33	33
High abundance protein	27	13	40
Total	27	46	73

Table 17. Contingency table of low/high abundance protein and presence/absence of dicodons XX (deduced from [Table 14](#), [Table 15](#) and [Table 16](#)).

A χ^2 test shows a strongly significant relation between the presence/absence of dicodons XX and protein abundancy with a one sided value $p \approx 10^{-9}$ ([Table 17](#)). The following probabilities can be easily deduced from [Table 17](#):

$$P(\text{Low abundance protein} \mid XX) = 0/33 = 0\%,$$

$$P(\text{High abundance protein} \mid XX) = 27/40 = 67.5\%.$$

Thus, the presence-absence of XX dicodons in a gene is an important and new factor in the classification of genes as low or high abundance.

3.2.4. Presence of X motifs in wild type genes and genes optimized to increase expression

The SGDB database (Wu *et al.*, 2007) contains gene expression data for genes that have been experimentally re-engineered to increase gene expression. Generally, this is achieved by replacing codons in the wild type gene with optimal codons for the expression system (i.e. replace rare codons with the most frequently used codons in the organism). We only considered the re-engineered genes that did not involve non-synonymous changes. Thus, we analysed 42 re-engineered genes that had increased expression and 4 re-engineered genes had no significant increase in expression. We searched for X motifs and R random motifs (from the 100 random codes) in the wild type genes and the genes optimized for gene expression. Then, we calculated the mean number and the mean nucleotide length of X and R motifs per sequence ([Table 18](#)).

		Mean number of X motifs	Mean number of R motifs	Mean length of X motifs	Mean length of R motifs
42 genes with increased expression	Wild type	5.4	3.6	86.1	53.7
	Optimized gene	11.2	3.7	188.6	58.2
3 genes with no increased expression	Wild type	5.3	2.6	80.0	35.8
	Optimized gene	5.0	3.8	80.0	55.6

Table 18. Mean number and mean nucleotide length of X and R random motifs per wild type gene and per optimized gene from the SGDB database (Wu *et al.*, 2007).

For the re-engineered genes that did not present an increase expression, we observe a non-significant difference in the mean number ($5.0 - 5.3 = -0.3$, one tailed Wilcoxon test with value $p = 0.50$) and no

difference in the mean length ($80.0 - 80.0 = 0$) of X motifs between the optimized genes and the wild type genes. These differences are also not significant for R motifs ($3.8 - 2.6 = 1.2$ and $55.6 - 35.8 = 19.8$, respectively, data not shown). In contrast, for the re-engineered genes that resulted in increased expression, the optimized genes have significantly more X motifs ($11.2 - 5.4 = 5.8$, one tailed Wilcoxon test with value $p \approx 10^{-6}$) and the X motifs covered a larger proportion of the genes ($188.6 - 86.1 = 102.5$, one tailed Wilcoxon test with value $p \approx 10^{-6}$) for most genes (Figure 10). These differences are not observed with the R motifs (one tailed Wilcoxon test, p values equal to 0.24 and 0.12, respectively). Thus, this important result suggest a potential new strategy for the efficient gene optimization.

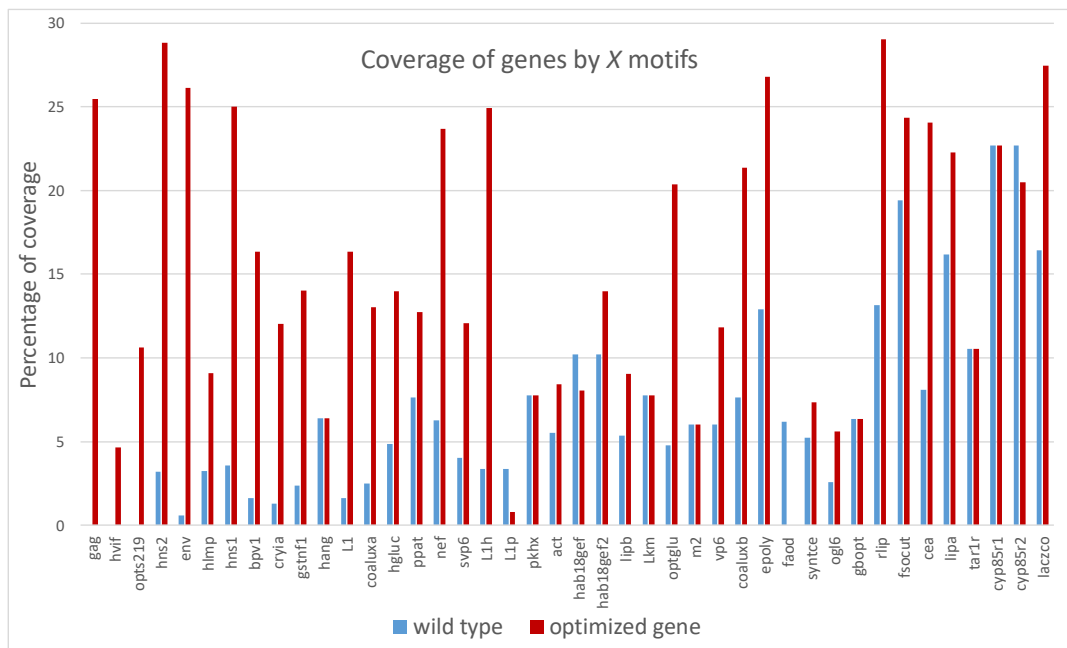


Figure 10. Percentage coverage (total length of X motifs divided by the total length of genes) of 42 wild type and optimized genes by X motifs.

4. Conclusion

The work described in this paper addressed two questions: are X motifs conserved during evolution? and do they continue to play a functional role in the processes of genome decoding and protein synthesis?

We performed a large scale study involving the complete genomes of four mammals and nine yeast species. The organisms chosen represent a large phylogenetic distribution, and a wide variety of gene structures, ranging from the simple, single exon genes of *S. cerevisiae* to the highly complex intron/exon structure of human genes. To avoid any bias towards a specific sequence alignment algorithm or evolutionary model, the multiple alignments of the gene sequences were obtained by two different methods. First, high quality mammal gene alignments were obtained from a previous independent study (Sharma and Hiller, 2017) of genome annotation methods. Second, multiple alignments of the yeast genes were constructed using a simple protein alignment method (ClustalW, Thompson *et al.*, 1994).

Furthermore, well characterized, well annotated genomes (Human and *S. cerevisiae*) were chosen to ensure high quality gene models.

In a preliminary analysis, we identified a strongly significant enrichment of *X* motifs (number and length) in both mammal and yeast genes, confirming our previous findings in *S. cerevisiae* (Michel *et al.*, 2017). We then calculated a number of different measures of evolutionary conservation, and showed that the *X* motifs are more conserved than the rest of the gene sequences, with a lower ratio Pns/Ps of non-synonymous to synonymous substitutions, indicative of purifying selection. These results were found to hold in both the mammal and yeast gene alignments. We then performed a more in-depth investigation of the synonymous substitutions in *X* motifs. At this stage, we modelled the evolutionary processes at the codon level, since it is important to account precisely for the protein-coding constraints on each nucleotide site. We demonstrated that the sequence conservation observed in the *X* motifs is the result of two types of selective pressure. The first type is the pressure to maintain the amino acids of the proteins encoded by the genes. The second type of selective pressure applies only to *X* motifs and highlights a new conservation property of *X* motifs per amino acid, which led us to propose a novel hypothesis for the evolution of the genetic code as a union of circular codes associated with each amino acid.

The increased conservation of *X* motifs and the specific evolutionary constraints suggest that *X* motifs may represent an additional, overlapping function within the protein-coding regions of genomes. Indeed, the genetic code establishes the rules to translate the 64 possible codons into the 20 amino acids and a stop signal. It is well known that the genetic code is degenerate and that different synonymous codons encoding the same amino acid are not used with the same frequency in different species. The genetic code also contains information that influences the rate and efficiency of translation, although the mechanisms of codon-mediated regulation are still not clear (Brule and Greyhack, 2017). Many recent studies have been performed to try to explain the different codon usages observed and their effects on translation. In particular, it has been shown that the efficiency of translating a particular codon is influenced by the nature of the immediately adjacent flanking codons (Gamble *et al.*, 2016; Diambra, 2017; Chevance and Hughes, 2017). These studies are mostly based on statistical and/or experimental analyses of gene sequences without being related to the results to a theoretical model. Here, we have investigated the pertinence of the circular code theory to explain the observations.

For example, in two related studies (Gamble *et al.*, 2016; Diambra, 2017), a total of 33 dicodons were found to be associated with low protein abundance, and 40 dicodons associated with high abundance proteins. We identified a significant correlation between the protein abundance level and dicodons belonging to the circular code *X*. To further investigate this link between the presence of *X* motifs in a gene and the expression level, we compared a set of re-engineered genes (that had been experimentally optimized for increased expression, using synonymous substitutions to replace rare codons with more frequent ones) with the original wild type genes. Again, we identified a significant correlation between the number and length of *X* motifs and protein expression levels. These results, taken together, suggest

that increasing the proportion of X motifs in a gene may represent an important new strategy for their efficient optimization.

The molecular mechanisms underlying the functional correlations observed here remain to be elucidated. However, it has been observed previously that short X motifs have also been conserved in many transfer RNAs (tRNAs) and ribosomal RNAs (rRNAs) (Michel, 2012, 2013; El Soufi and Michel, 2014, 2015). In particular, the universally conserved nucleotides A1492, A1493 G530 in the ribosome decoding center are located in short X motifs. Given the self-complementary property of the circular code X , it is possible that there is some kind of interaction between the X motifs in the protein-coding genes and the ribosomal X motifs in order to maintain the correct reading frame during translation.

The results presented here indicate that the circular code motifs may explain how the choice of different synonymous codons within genes and between species, known as codon usage bias, impacts nucleic acid stability, protein levels, structure, function and evolution. Further investigation, and *in vitro* or *in vivo* experimental validation, will be required to refine our hypothesis for the evolution of the genetic code and the proposed functional role in the translation of genes.

REFERENCES

- Altschul S.F., Madden T.L., Schäffer A.A., Zhang J., Zhang Z., Miller W., Lipman D.J. (1997). Gapped BLAST and PSI-BLAST: a new generation of protein database search programs. *Nucleic Acids Research* 125, 3389-3402.
- Arquès D.G., Michel C.J. (1996). A complementary circular code in the protein coding genes. *Journal of Theoretical Biology* 182, 45-58.
- Brule C.E., Grayhack E.J. (2017). Synonymous codons: choose wisely for expression. *Trends in Genetics* 33, 283-297.
- Chevance F.F.V., Hughes K.T. (2017). Case for the genetic code as a triplet of triplets. *Proceedings of the National Academy of Sciences U.S.A.*, 1614896114.
- Diambra L.A. (2017). Differential bicodon usage in lowly and highly abundant proteins. *PeerJ* 5:e3081.
- El Houmami N., Seligmann H. (2017). Evolution of nucleotide punctuation marks: from structural to linear signals. *Frontiers in Genetics* 8, 36.
- El Soufi K., Michel C.J. (2014). Circular code motifs in the ribosome decoding center. *Computational Biology and Chemistry* 52, 9-17.
- El Soufi K., Michel C.J. (2015). Circular code motifs near the ribosome decoding center. *Computational Biology and Chemistry* 59, 158-176.
- El Soufi K., Michel C.J. (2017). Unitary circular code motifs in genomes of eukaryotes. *Biosystems* 153, 45-62.
- Fimmel E., Michel C.J., Strüngmann L. (2016). n -Nucleotide circular codes in graph theory. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences* 374, 20150058.
- Fimmel E., Michel C.J., Strüngmann L. (2017). Strong comma-free codes in genetic information. *Bulletin of Mathematical Biology* 79, 1796-1819.
- Fimmel E., Strüngmann L. (2018). Mathematical fundamentals for the noise immunity of the genetic code. *Biosystems* 164, 186-198.
- Gamble C.E., Brule C.E., Dean K.M., Fields S., Grayhack E.J. (2016). Adjacent codons act in concert to modulate translation efficiency in yeast. *Cell* 166, 679-690.
- Michel C.J. (2008). A 2006 review of circular codes in genes. *Computer and Mathematics with Applications* 55, 984-988.
- Michel C.J. (2012). Circular code motifs in transfer and 16S ribosomal RNAs: A possible translation code in genes. *Computational Biology and Chemistry* 37, 24-37.
- Michel C.J. (2013). Circular code motifs in transfer RNAs. *Computational Biology and Chemistry* 45, 17-29.
- Michel C.J. (2014). A genetic scale of reading frame coding. *Journal of Theoretical Biology* 355, 83-94.
- Michel C.J. (2015). The maximal C^3 self-complementary trinucleotide circular code X in genes of bacteria, eukaryotes, plasmids and viruses. *Journal of Theoretical Biology* 380, 156-177.

- Michel C.J. (2017). The maximal C^3 self-complementary trinucleotide circular code X in genes of bacteria, archaea, eukaryotes, plasmids and viruses. *Life* 7, 20, 1-16.
- Michel C.J., Nguefack Ngoune V., Poch O., Ripp R., Thompson J.D. (2017). Enrichment of circular code motifs in the genes of the yeast *Saccharomyces cerevisiae*. *Life* 7, 52, 1-20.
- Nei M., Gojobori T. (1986). Simple methods for estimating the numbers of synonymous and nonsynonymous nucleotide substitutions. *Molecular Biology Evolution* 3, 418-26.
- Seligmann H., Pollock DD. (2004). The ambush hypothesis: hidden stop codons prevent off-frame gene reading. *DNA and Cell Biology* 23, 701-705.
- Seligmann H. (2011). Error compensation of tRNA misacylation by codon-anticodon mismatch prevents translational amino acid misinsertion. *Computational Biology and Chemistry* 35, 81-95.
- Seligmann H., Warthi G. (2017). Genetic Code Optimization for Cotranslational Protein Folding: Codon Directional Asymmetry Correlates with Antiparallel Betasheets, tRNA Synthetase Classes. *Computational and Structural Biotechnology Journal* 15, 412-424.
- Sharma V., Hiller M. (2017). Increased alignment sensitivity improves the usage of genome alignments for comparative gene annotation. *Nucleic Acids Research* 45, 8369-8377.
- Thompson J.D., Higgins D.G., Gibson T.J. (1994). CLUSTAL W: improving the sensitivity of progressive multiple sequence alignment through sequence weighting, position-specific gap penalties and weight matrix choice. *Nucleic Acids Research*.22, 4673-4680.
- Wu G., Zheng Y., Qureshi I., Zin H.T., Beck T., Bulka B., Freeland S.J. (2007). SGDB: a database of synthetic genes re-designed for optimizing protein over-expression. *Nucleic Acids Research* 35, D76-D79.

<i>A</i>		<i>D</i>		<i>E</i>		<i>F</i>		<i>G</i>			<i>I</i>		
<i>GCA GCC GCG GCT</i>		<i>GAC GAT</i>		<i>GAA GAG</i>		<i>TTC TTT</i>		<i>GGA GGC GGG GGT</i>			<i>ATA ATC ATT</i>		
<i>GCA</i>		<i>GAC</i> 72.9	22.0	<i>GAA</i> 67.3	13.0	<i>TTC</i> 78.9		<i>GGA</i>			<i>ATA</i>		
<i>GCC</i>	66.1	<i>GAT</i> 16.7	67.1	<i>GAG</i> 22.7	77.7	<i>TTT</i>		<i>GGC</i>	69.6	21.9	<i>ATC</i>	73.1	22.4
<i>GCG</i>		Sum 89.6	89.1	Sum 90.0	90.7	Sum 78.9		<i>GGG</i>			<i>ATT</i>	12.5	61.8
<i>GCT</i>		Mean 89.4		Mean 90.3		Mean 78.9		<i>GGT</i>	10.7	52.4	Sum 85.6	84.2	
Sum	66.1							Sum	80.3	74.3	Mean 84.9		
Mean	66.1							Mean	77.3				

<i>L</i>					<i>N</i>		<i>Q</i>		<i>T</i>			<i>V</i>			<i>Y</i>	
<i>CTA CTC CTG CTT TTA TTG</i>					<i>AAC AAT</i>		<i>CAA CAG</i>		<i>ACA ACC ACG ACT</i>			<i>GTA GTC GTG GTT</i>			<i>TAC TAT</i>	
<i>CTA</i>					<i>AAC</i> 71.1	22.2	<i>CAA</i>		<i>ACA</i>			<i>GTA</i> 42.7	1.9	2.1	<i>TAC</i> 76.1	
<i>CTC</i>	68.1	4.6			<i>AAT</i> 15.5	62.3	<i>CAG</i>	80.7	<i>ACC</i>	63.5		<i>GTC</i> 6.2	64.2	14.8	<i>TAT</i>	
<i>CTG</i>	10.9	73.0			Sum 86.6	84.5	Sum 80.7		<i>ACG</i>			<i>GTG</i>			Sum 76.1	
<i>CTT</i>					Mean 85.6		Mean 80.7		<i>ACT</i>			<i>GTT</i> 3.9	7.5	53.9	Mean 76.1	
<i>TTA</i>									Sum 63.5			Sum 52.8	73.6	70.9		
<i>TTG</i>									Mean 63.5			Mean 65.7				
Sum	79.0	77.6														
Mean		78.3														

Table 21. For the mammal gene multiple alignments with respect to the human reference genes $s_1 = \mathbb{H}$, codon substitution submatrices of $\mathbf{B}(m(X, \mathbb{H}))$ (in %) of X motifs $m(X, \mathbb{H})$ (Section 2.6) for the 12 amino acids $p \in \mathcal{X} = \{A, D, E, F, G, I, L, N, Q, T, V, Y\}$ coded by the circular code X (1).

A		D		E		F		G			I	
GCA GCC GCG GCT		GAC GAT		GAA GAG		TTC TTT		GGA GGC GGG GGT			ATA ATC ATT	
GCA		GAC 19.9	17.9	GAA 26.3	23.6	TTC 29.9		GGA			ATA	
GCC	16.1	GAT 25.6	27.1	GAG 17.3	16.4	TTT		GGC	10.1	9.1	ATC	18.5 16.3
GCG		Sum 45.4	45.1	Sum 43.6	39.9	Sum 29.9		GGG			ATT	20.7 23.1
GCT		Paac 45.3		Mean 41.8		Mean 29.9		GGT	21.9	39.7	Sum	39.2 39.4
Sum	16.1							Sum	32.0	48.8	Mean	39.3
Paac	16.1							Mean	40.4			

L					N		Q		T			V			Y	
CTA CTC CTG CTT TTA TTG					AAC AAT		CAA CAG		ACA ACC ACG ACT			GTA GTC GTG GTT			TAC TAT	
CTA					AAC 20.6	16.2	CAA		ACA			GTA 5.2	5.0	5.0	TAC 33.0	
CTC	5.6	6.2			AAT 15.0	15.3	CAG 13.6		ACC 14.9			GTC 8.4	13.3	11.9	TAT	
CTG	4.2	5.2			Sum 35.6	31.5	Sum 13.6		ACG			GTG			Sum 33.0	
CTT					Mean 33.5		Mean 13.6		ACT			GTT 13.5	18.8	19.7	Mean 33.0	
TTA									Sum 14.9			Sum 27.1	37.2	36.6		
TTG									Mean 14.9			Mean 33.6				
Sum	9.9	11.4														
Mean		10.6														

Table 22. For the yeast gene multiple alignments with respect to the *S. cerevisiae* reference genes $s_1 = \mathbb{C}$, codon substitution submatrices of $\mathbf{B}(m(X, \mathbb{C}))$ (in %) of X motifs $m(X, \mathbb{C})$ (Section 2.6) for the 12 amino acids $p \in \mathcal{X} = \{A, D, E, F, G, I, L, N, Q, T, V, Y\}$ coded by the circular code X (1).