# Genomic Abelian Finite Groups

Robersy Sanchez [1] and Jesús Barreto[2]

[1]Department of Biology. Pennsylvania State University, University Park, PA 16802.
E-mail: rus547@psu.edu
ORCID: https://orcid.org/0000-0002-5246-1453

[2]Universidad Central "Marta Abreu" de Las Villas. Santa Clara. Cuba.
E-mail: barretouclv@gmail.com

[1] Corresponding author:
rus547@psu.edu

## Abstract

Experimental studies reveal that genome architecture splits into natural domains suggesting a well-structured genomic architecture, where, for each species, genome populations are integrated by individual mutational variants. Herein, we show that the architecture of population genomes from the same or closed related species can be quantitatively represented in terms of the direct sum of homocyclic abelian groups defined on the genetic code, where populations from the same species lead to the same canonical decomposition into $p$-groups. This finding unveils a new ground for the application of the abelian group theory to genomics and epigenomics, opening new horizons for the study of the biological processes (at genomic scale) and provides new lens for genomic medicine.

**Keywords**: Genomics, Genetic code, Abelian groups, genome algebra

# 1   Introduction

28

29   The analysis of the genome architecture is one of biggest challenges for the current and future

30   genomics. Current bioinformatic tools make possible faster genome annotation process than some

31   years ago [2]. Current experimental genomic studies suggest that genome architectures must obey

32   specific mathematical biophysics rules [3–6].

33   Experimental results points to an injective relationship: *DNA sequence → 3D chromatin*

34   *architecture* [3,4,6], and failures of DNA repair mechanisms in preserving the integrity of the DNA

35   sequences lead to dysfunctional genomic rearrangements which frequently are reported in several

36   diseases [5]. Hence, some hierarchical logic is inherent to the genetic information system that makes

37   it feasible for mathematical studies. In particular, there exist mathematical biology reasons to analyze

38   the genetic information system as a communication system [7–10].

39   Under the assumption that current forms of life evolved from simple primordial cells with very

40   simple genomic structure and robust coding apparatus, the genetic code is a fundamental link to the

41   primeval form of live, which played an essential role on the primordial architecture. The genetic

42   code, the set of biochemical rules used by living cells to translate information encoded within genetic

43   material into proteins, sets the basis for our understanding of the mathematical logic inherent to the

44   genetic information system [9,11]. The genetic code is the cornerstone of live on earth. Not a single

45   form of live could evolve or exist, as we currently know it, without the genetic code.

46   Several genetic code algebraic structures has been introduced to study effect of the quantitative

47   relationship between the coding apparatus and the mutational process on protein-coding regions [12–

48   16]. Formally the genetic code only is limited to translated coding regions where the number of RNA

49   bases is a multiple of 3. However, as suggested in reference [17], the difficulties in prebiotic synthesis

50   of the nucleosides components of RNA (nucleo-base + sugar) and suggested that some of the original

51   bases may not have been the present purines or pyrimidines [18]. Piccirilli et al. [19] demonstrated

52   that the alphabet can in principle be larger. Switzer et al. [20] have shown an enzymatic incorporation

53   of new functionalized bases into RNA and DNA. This expanded the genetic alphabet from 4 to 5 or

54    more letters, which permits new base pairs, and provides RNA molecules with the potential to greatly

55    increase their catalytic power.

56         It is important to notice that even in the current (*friendly*) environmental conditions not a single

57    cell can survive without a DNA repair enzymatic machinery and that such an enzymatic machinery

58    did not existed at all in the primaeval forms of live. Here, we are confronting the *chicken and egg*

59    problem. To date, the best solution (to our knowledge) is the admission of alternative base-pairs in

60    the primordial DNA alphabet which, as suggested in the studies on the prebiotic chemistry, could

61    contribute to the thermal and general physicochemical stability of the primordial DNA molecules.

62         Several algebraic structures have been proposed including an additional letter into the DNA

63    alphabet: A, C, G, T. The new letter (D) stands for current insertion deletion/mutations or for

64    alternative wobble base pairing, which would be a relict fingerprint from primordial enzymes derived

65    from a more degenerated ancestral genetic code [17,21,22]. Supporting evidence for the existence of

66    a more degenerated ancestral genetic code built up on a larger alphabet is found in the tRNA anticodon

67    region permitting wobble base pairing by including, e.g., bases such as: inosine (in eukariotes),

68    agmatidine (in archaea), and lysidine (in bacteria), which has been proposed as evolutionary solutions

69    to the need for lower the high translational noise connected to the reading of the AUA and AUG

70    codons [23,24]. Additionally, various alternative base pairs like methylated cytosine and adenine are

71    still present in the current genomes playing an important role in the epigenetic adaptation of

72    organismal populations to the continuous environmental changes [10,25].

73         Cytosine DNA methylation results from the addition of methyl groups to cytosine C5 residues,

74    and the configuration of methylation within a genome provides trans-generational epigenetic

75    information. These epigenetic modifications can influence the transcriptional activity of the

76    corresponding genes, or maintain genome integrity by repressing transposable elements and affecting

77    long-term gene silencing mechanisms [26,27].

78         In this scenario, we shall show that all possible DNA molecules and, consequently, genomes

79    can be described by way of finite abelian groups which can be split into the direct sum of homocyclic

80    2-groups and 5-groups defined on the genetic code. A homocyclic group is a direct sum of cyclic

81     groups of the same order. Any finite abelian group can be decomposed into a direct sum of

82     homocyclic $p$-groups [28], i.e., a group in which the order of every element is a power of a primer number

83     $p$.

84          The genetic code algebraic structures under scrutiny in the mentioned references covered rings

85     and vector spaces with a common feature, the corresponding additive group is an abelian group of

86     prime-power order. Next, to help a better comprehension of the current work, a brief introductory

87     summary on these groups is provided. Results presented here generalizes the application of the

88     genetic code algebras (reported in several publications) to the whole genome.

**1.1    Reported genetic code abelian groups relevant for the current study**

89

90     Herein, we assume that readers are familiar with the definition of abelian group, which otherwise can

91     be found in textbooks and elsewhere including Wikipedia. Nevertheless, all the abelian groups

92     discussed here are isomorphic to the well-known abelian groups of integer module $n$, which are easily

93     apprehended by a college-average educated mind. For example, the abelian group defined on the set

94     {0, 1, 2, 3, 4}, which corresponds to the group of integer modulo 5 ( $\mathbb{Z}_5$ ), where $(2 + 1) \bmod 5 = 3$,

95     $(1 + 3) \bmod 5 = 4$, $(2 + 3) \bmod 5 = 0$, etc. The subjacent biophysical and biochemical reasonings to

96     define the algebraic operations on the set of DNA bases and on the codon set were given in references

97     [12,14,17].

*1.1.1   The $\mathbb{Z}_{64}-algebras$ of the genetic code ($C_g$)*

98

99     The $\mathbb{Z}_{64}-algebras$ of the genetic code ($C_g$) and gene sequences were stated several years ago. In the

100    $\mathbb{Z}_{64}-algebra$ $C_g$ the sum operation, defined on the codon set, is a manner to consecutively obtain all

101    codons from the codon AAC (UUG) in such a way that the genetic code will represent a non-

102    dimensional code scale of amino acids interaction energy in proteins.

103         A description of the genetic code abelian finite group ($C_g$, +) can be found in [12]. Group

104    $\left( C_g, + \right)$ is isomorphic to the group on the set $\mathbb{Z}_{2^6}$ (the sum of integer modulo 64), which formally

105    will be expressed as $\left( C_g, + \right) \cong \left( \mathbb{Z}_{2^6}, + \right)$. This group on the set $\mathbb{Z}_{2^6}$ (the sum of integer modulo 64).

106     The mapping of the set of codons $X_1 X_2 X_3 \in C_g$ into the set $\mathbb{Z}_{2^6}$ is straightforward after consider the

107     bijection $A \leftrightarrow 0, C \leftrightarrow 1, G \leftrightarrow 2, U \leftrightarrow 3$ and the function $g(x) = 4x_1 + 16x_2 + x_3$. For example:

$$
\begin{array}{llll}
AGC & \leftrightarrow & 33 \\
+\,UGU & \leftrightarrow & +47 \\
\hline
ACA & \leftrightarrow & 16 \bmod 64
\end{array}
\qquad
\begin{array}{llll}
AGC & \leftrightarrow & 33 \\
+\,ACU & \leftrightarrow & +18 \\
\hline
AUU & \leftrightarrow & 51 \bmod 64
\end{array}
\qquad
\begin{array}{llll}
GGC & \leftrightarrow & 41 \\
CUA & \leftrightarrow & +52 \\
\hline
UCC & \leftrightarrow & 29 \bmod 64
\end{array}
$$

108

109     The $Z_{64}$-algebra $C_g$, however, is limited to protein-coding regions, while it is well known that,

110     in eukaryotes, only a small fraction of the genome −about 3%− called open reading frame (ORF)

111     encodes for proteins [18]. Since non-coding DNA sequences can have a base pairs number not

112     multiple of three, complete chromosomes and genomes cannot be described by means of group

113     $(C_g, +)$. In addition, natural genomic variations that includes insertions and deletion mutations (indel

114     mutations) across individuals from the same population and close-related populations from different

115     species cannot be represented with group $(C_g, +)$.

116     *1.1.2    The $\left(\mathbb{Z}_2^6, +\right)$ group of the genetic code ($C_g$)*

117     Group $(C_g, +)$ is the additive group of a module over a ring, which however, do not conform to a

118     vector space. To build a genetic code vector space, a Galois field ($GF(4)$) structure in the ordered base

119     set $B = \{G, U, A, C\}$ was introduced in reference [14]. In particular, an isomorphism with the Galois

120     field is defined by means of its binary representation $\mathbb{Z}_2 \times \mathbb{Z}_2 = \{(0, 0),\ (0, 1), (1, 0), (1, 1)\}$, i.e. a

121     unique $GF(4)$ up to isomorphism exists, such that a bijection $f : \{G, U, A, C\} \leftrightarrow \mathbb{Z}_2 \times \mathbb{Z}_2$ from the

122     DNA base set $B = \{G, U, A, C\}$ to the set of binary duplets $(\alpha_1, \alpha_2)$ is stated., where $\alpha_i \in \mathbb{Z}_2 = \{0,1\}$

123     , for $i \in \{1, 2\}$. For example, the bijection $f$ is defined as:

124     $$f(G) = (0, 0), f(U) = (0, 1), f(A) = (1, 0), f(C) = (1, 1).$$

125     The additive group of bases is the Klein four-group, which is defined by the group presentation:

126     $V = \{U, A \mid U + U = A + A = C + C = G,\ A + U = C\}$, i.e., $(B, +) \cong \left(\mathbb{Z}_2^2, +\right)$. Next, the abelian group

127     on the set of codons $B^3$ was defined as the direct third power $B^3 = B \times B \times B$ of the group $(B, +)$, i.e.

128    $(B^3,+) = (B,+) \times (B,+) \times (B,+)$, which is isomorphic to the group: $(\mathbb{Z}_2^6,+) = (\mathbb{Z}_2^2,+) \times (\mathbb{Z}_2^2,+) \times (\mathbb{Z}_2^2,+)$

129    , i.e., $(B^3,+) \cong (\mathbb{Z}_2^6,+)$. The sum operation on the set $(B^3,+)$ follows from the sum operation by

130    coordinates.

131       As pointed out before by Crick, the first two bases of codons determine the physicochemical

132    properties of aminoacids [29]. The four encoded amino acids of every class are either the same or

133    show very similar physicochemical properties. This genetic code regularity is captured by the quotient

134    group $B^3/G_{GGA}$, where $G_{GGA}$ is a subgroup of $B^3$ integrated by the elements {GGG,GGA} (the

135    elements of the quotient group $B^3/G_{GGA}$ are given in Table 5 from [14]). The quotient group

136    $B^3/G_{GGA}$ is isomorphic to group $(\mathbb{Z}_2^5,+) = (\mathbb{Z}_2^2,+) \times (\mathbb{Z}_2^2,+) \times (\mathbb{Z}_2,+)$. Each element of this group

137    represents an equivalence class of codons. Two triplets $X_1X_2X_3$ and $Y_1Y_2Y_3$ are equivalent if, and

138    only if, the difference $X_1X_2X_3 + Y_1Y_2Y_3 \in G_{DDA}$. In biological terms, substitution mutations

139    involving codons from the same class will not alter (or at least no substantially alter in most of the

140    cases) the physicochemical properties of the encoded protein domains, since in the worst scenario

141    involves aminoacids with very close physicochemical properties, with the exception of codon for

142    aminoacid tryptophan.

143    *1.1.3   The $\mathbb{Z}_{125}$ group of the extended genetic code ($C_e$)*

144    The extension of the *genetic code group* $(C_g,+)$ follows straightforward from the extension of the

145    codon set, which is easily accomplished extending the source alphabet of the standard genetic code:

146    {A, C, G, U} and, consequently, extending the base triplet set (extended triplet) as $X_1X_2X_3$, $X_i \in$ {D,

147    A, C, G, U} [22]. The new algebraic structure $(C_e,+)$ is isomorphic to the abelian group defined on

148    the set $\mathbb{Z}_{5^3}$ (the sum of integer modulo 125), formally, $(C_e,+) \cong (\mathbb{Z}_{5^3},+)$. The mapping of the set of

149    codons $X_1X_2X_3 \in C_e$ into the set $\mathbb{Z}_{5^3}$ is straightforward after consider the bijection

150 $D \leftrightarrow 0, A \leftrightarrow 1, C \leftrightarrow 2, G \leftrightarrow 3, U \leftrightarrow 4$ and the function $g(x) = 5x_1 + 25x_2 + x_3$ (see Table 1). For

151 example:

| AGC ↔ 82 | AGC ↔ 82 | GGC ↔ 92 |
|---|---|---|
| + UGU ↔ +99 | + DCU ↔ +54 | + CUD ↔ +110 |
| ACA ↔ 56 mod 125 | CDA ↔ 11 mod 125 | DGC ↔ 77 mod 125 |

152

153 **Table 1**. Ordered set of extended triplets corresponding to the elements from $\mathbb{Z}_{5^3}$

| a | D I | D III | A I | A III | | C I | C III | | G I | G III | | U I | U III | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| D | 0 | DDD | 25 | DAD | | 50 | DCD | | 75 | DGD | | 100 | DUD | | D |
| | 1 | DDA | 26 | DAA | | 51 | DCA | | 76 | DGA | | 101 | DUA | | A |
| | 2 | DDC | 27 | DAC | | 52 | DCC | | 77 | DGC | | 102 | DUC | | C |
| | 3 | DDG | 28 | DAG | | 53 | DCG | | 78 | DGG | | 103 | DUG | | G |
| | 4 | DDU | 29 | DAU | | 54 | DCU | | 79 | DGU | | 104 | DUU | | U |
| A | 5 | ADD | 30 | AAD | | 55 | ACD | | 80 | AGD | | 105 | AUD | | D |
| | 6 | ADA | 31 | AAA | K | 56 | ACA | T | 81 | AGA | R | 106 | AUA | I | A |
| | 7 | ADC | 32 | AAC | N | 57 | ACC | T | 82 | AGC | S | 107 | AUC | I | C |
| | 8 | ADG | 33 | AAG | K | 58 | ACG | T | 83 | AGG | R | 108 | AUG | M | G |
| | 9 | ADU | 34 | AAU | N | 59 | ACU | T | 84 | AGU | S | 109 | AUU | I | U |
| C | 10 | CDD | 35 | CAD | | 60 | CCD | | 85 | CGD | | 110 | CUD | | D |
| | 11 | CDA | 36 | CAA | Q | 61 | CCA | P | 86 | CGA | R | 111 | CUA | L | A |
| | 12 | CDC | 37 | CAC | H | 62 | CCC | P | 87 | CGC | R | 112 | CUC | L | C |
| | 13 | CDG | 38 | CAG | Q | 63 | CCG | P | 88 | CGG | R | 113 | CUG | L | G |
| | 14 | CDU | 39 | CAU | H | 64 | CCU | P | 89 | CGU | R | 114 | CUU | L | U |
| G | 15 | GDD | 40 | GAD | | 65 | GCD | | 90 | GGD | | 115 | GUD | | D |
| | 16 | GDA | 41 | GAA | E | 66 | GCA | A | 91 | GGA | G | 116 | GUA | V | A |
| | 17 | GDC | 42 | GAC | D | 67 | GCC | A | 92 | GGC | G | 117 | GUC | V | C |
| | 18 | GDG | 43 | GAG | E | 68 | GCG | A | 93 | GGG | G | 118 | GUG | V | G |
| | 19 | GDU | 44 | GAU | D | 69 | GCU | A | 94 | GGU | G | 119 | GUU | V | U |
| U | 20 | UDD | 45 | UAD | | 70 | UCD | | 95 | UGD | | 120 | UUD | | D |
| | 21 | UDA | 46 | UAA | Stop | 71 | UCA | S | 96 | UGA | Stop | 121 | UUA | L | A |
| | 22 | UDC | 47 | UAC | Y | 72 | UCC | S | 97 | UGC | C | 122 | UUC | F | C |
| | 23 | UDG | 48 | UAG | Stop | 73 | UCG | S | 98 | UGG | W | 123 | UUG | L | G |
| | 24 | UDU | 49 | UAU | Y | 74 | UCU | S | 99 | UGU | C | 124 | UUU | F | U |

154 ᵃ Bijection between the base-triplets set and the elements from sets $\mathbb{Z}_{5^3}$ as given in [22].

155

156 *1.1.4 The $\left(\mathbb{Z}_5^3, +\right)$ group of the extended genetic code ($C_e$)*

157 The Galois field $GF(5)$ of the DNA set of bases $\mathfrak{B} = \{D, A, C, G, U\}$ was introduced in reference [17].

158 This structure led to the definition of a $\mathbb{Z}_5$−vector space $\mathfrak{B}^3$ over the set $\mathfrak{B}^3 = \mathfrak{B} \times \mathfrak{B} \times \mathfrak{B}$

159 isomorphic to the set $\mathbb{Z}_5^3 = \mathbb{Z}_5 \times \mathbb{Z}_5 \times \mathbb{Z}_5$ [17,30]. But here, we are interested only in the abelian

160 groups $\left(\mathfrak{B}, +\right)$ and $\left(\mathfrak{B}^3, +\right)$. After the bijection $D \leftrightarrow 0, A \leftrightarrow 1, C \leftrightarrow 2, G \leftrightarrow 3, U \leftrightarrow 4$, the sum

161 operation of two DNA bases follows from the sum operation on the Galois field $GF(5)$ (i.e., on $\mathbb{Z}_5$,

162    the sum of integers modulo 5). For example, $C + U \leftrightarrow (2 + 4) \mod 5 = 1 \leftrightarrow A$. The sum operation

163    on the set $\mathfrak{B}^3$ follows from the sum operation by coordinates.

164         It is worthy to notice that there 24 way to define each one of the above mentioned algebraic

165    structures [30,31]. Nevertheless, for each defined genetic code group, there is only one (genetic code

166    abelian group) up to isomorphism, which lead to their representation as an abelian group, where the

167    sum operation corresponds to the sum of integer modulo $n \in \{2, 2^6, 5, 5^3\}$.

## 2    The General Theoretical Model

169    Herein, it will be showed that, in a general scenario, the whole genome population from any species

170    or close related species, can be algebraically represented as a direct sum of abelian cyclic groups or

171    more specifically abelian $p$-groups. Basically, we propose the representation of multiple sequence

172    alignments (MSA) of length $N$ as the direct sum:

173
$$G = \left( \mathbb{Z}_{p_1} \right)^{n_1} \oplus \left( \mathbb{Z}_{p_2} \right)^{n_2} \oplus \cdots \oplus \left( \mathbb{Z}_{p_k} \right)^{n_k} \qquad [1]$$

174    Where $p_i \in \{2, 5, 2^6, 5^3\}$ and $N = n_1 + n_2 + \ldots + n_k$. Here, we assume the usual definition of direct sum

175    of groups [32]. Let $B_i$ $(i \in I = \{1, \ldots, n\})$ be a family of subgroups of $G$, subject to the following two

176    conditions:

177         1)   $\sum B_i = G$. That is, $B_i$ together generates $G$.

178         2)   For every $i \in I$:   $B_i \cap \sum B_j = 0$.

179    Then, it is said that $G$ is the direct sum of its subgroups $B_i$, which formally is expressed by the

180    expression: $G = \underset{i}{\oplus} B_i$ or $G = B_1 \oplus \ldots \oplus B_n$.

181         In superior organisms, genomic DNA sequences are integrated by intergenic regions and gene

182    regions. The former are the larger regions, while the later includes the protein-coding regions as

183    subsets. The MSA of DNA and protein-coding sequences reveals allocations of the nucleotide bases

184    and aminoacids into stretched of *strings*. The alignment of these stretched would indicate the presence

185     of substitution, *indel* mutations. As a result, the alignment of a whole chromosome DNA sequences

186     from several individuals from the same or close-related species can be split into well-defined

187     subregions or domains, and each one of them can be represented as homocyclic abelian groups, i.e.,

188     a cyclic group of *prime-power* order (Fig. 1). As a result, each DNA sequence is represented as a *N*-

189     dimensional vector with numerical coordinates representing bases and codons.

190

191 **Fig.1**. A typical DNA multiple sequence alignment (MSA) including segments of protein-coding
192 regions. A MSA would include the presence of substitution, insertion and deletion mutations (*indel*
193 mutations). The aligned sequences can be grouped into blocks, which can be algebraically represented
194 by abelian groups.
195

196     An intuitive mathematical representation of MSA is implicit in Fig.1, with following

197     observations:

198       a) Every DNA sequence from the MSA and every subsequence on it can be represented as a

199        vector with element coordinates defined in some abelian group. For example,

200        $\left(C_g, +\right) \cong \left(\mathbb{Z}_{64}, +\right)$, the first five codons from the first DNA sequence from Fig. 1,

201        $\{\text{ATA, CCC, ATG, GCC, AAC}\} \in \left(C_g, +\right)$, can be represented by the vector of integers:

202        $\{48, 21, 50, 25, 1\}$ where each coordinate is an element from group $\left(\mathbb{Z}_{64}, +\right)$ (see Table 1

203        from reference [12] and the introduction section).

204       b)    Any MSA can be algebraically represented as a symbolic composition of abelian

205        groups each one of them is isomorphic to an abelian group of integers module *n*. Such a

206        composition can be algebraically represented as a direct sum of homocyclic abelian groups.

207        For example, the multiple sequence alignment from Fig. 1 can be represented by the direct

208        sum of abelian groups:

209
$$G = \left(\mathbb{Z}_{2^6}\right)^5 \oplus \left(\mathbb{Z}_5\right)^8 \oplus \left(\mathbb{Z}_{5^3}\right)^5 \oplus \left(\mathbb{Z}_5\right)^7 \oplus \left(\mathbb{Z}_{5^3}\right)^4 \qquad [2]$$

210        In more specific scenario, the multiple sequence alignment from Fig. 1 can be represented by

211     the direct sum of abelian 2-groups and 5-groups:

212
$$G = (\mathbb{Z}_2^6)^5 \oplus (\mathbb{Z}_5)^8 \oplus (\mathbb{Z}_{5^3})^5 \oplus (\mathbb{Z}_5)^7 \oplus (\mathbb{Z}_2^6)^4 \qquad [3]$$

213 Or strictly as the direct sum of abelian 5-groups:

214
$$G = (\mathbb{Z}_5^3)^5 \oplus (\mathbb{Z}_5)^8 \oplus (\mathbb{Z}_{5^3})^5 \oplus (\mathbb{Z}_5)^7 \oplus (\mathbb{Z}_5^3)^4 \qquad [4]$$

215 Although the above *direct sums* of abelian groups provides a useful compact representation of

216 MSA, for application purposes to genomics, we would also consider to use the concept of direct

217 product (*cartesian sum or complete direct sums*) [32]. Next, let $S$ be a set of abelian cyclic groups

218 identified in the MSA $M$ of length $N$ (i.e., every DNA sequence from $M$ has $N$ bases). Let $\ell_i$ the

219 number of bases or triples of bases covered on $M$ by group $S_i \in S$ where $\sum_i \ell_i = N$. Hence, each

220 DNA sequence on the $M$ can be represented by a cartesian product $(b_1, \ldots, b_n)$ where $b_i \in S_i$

221 $(i = 1, \ldots, n)$ and $n = |S|$. Let $G_i$ be a group defined on the set of all elements $(0, \ldots, 0, b_i, 0, \ldots 0)$

222 where $b_i \in S_i$ stands on the $i^{th}$ place and 0 everywhere else. It is clear that $S_i \cong G_i$. In this context,

223 the set of all vectors $(b_1, \ldots, b_n)$ with equality and addition of vectors defined coordinate-wise

224 becomes a group ($G$) named direct product (cartesian sum) of groups $S_i$ ($G_i$), i.e.:

225
$$G = \otimes_i S_i = \oplus_i G_i \qquad [5]$$

226 An illustration of the cartesian sum application was given above in observation a).

227 **3 Results**

228 Results essentially comprise an application of the fundamental theorem of abelian finite groups

229 [28,32]. By this theorem every finite abelian group $G$ is isomorphic to a direct sum of cyclic groups

230 of prime-power order of the form:

231
$$G = \mathbb{Z}_{p_1^{\alpha_1}} \oplus \mathbb{Z}_{p_2^{\alpha_2}} \oplus \cdots \oplus \mathbb{Z}_{p_n^{\alpha_n}} \qquad [6]$$

232 Or (in short) $G = \oplus_{i=1}^n \mathbb{Z}_{p_i^{\alpha_i}}$, where the $p_i$'s are primes (not necessarily distinct), $\alpha_i \in \mathbb{N}$ and $\mathbb{Z}_{p_i^{\alpha_i}}$ is

233 the group of integer module $p_i^{\alpha_i}$. The abelian group representation of the MSA from Fig. 1 given by

234 expressions [1] and [2] correspond to the cases where the finite abelian group $G$ is a direct sum of

235  *prime-power order*, while expression [3] reflects the fact that any finite abelian group can be

236  decomposed into a direct sum of homocyclic *p*-groups [28,32], in this $p = 5$.

237  As is showed in Fig 1, this abelian group is a heterocyclic group that split into a direct sum of

238  homocyclic *prime-power order*, each one of them split into the direct sum of cyclic *p*-groups with

239  same order. For example, in expression [4] we have the subgroup: $\left(\mathbb{Z}_5^3\right)^4 = \oplus_{i=1}^{12} \mathbb{Z}_5$, which is a direct

240  sum of 12 homocyclic 5-groups $\left(\mathbb{Z}_5, +\right) \cong \left(\mathcal{B}, +\right)$. The case of $\mathbb{Z}_{2^6}$ representation of the genetic code

241  (as given in [12]) is less evident. It follows from the fact that the genetic code table is integrated by

242  16 subsets of codons with form $K = \left\{XY\text{A}, XY\text{C}, XY\text{G}, XY\text{U}\right\}$, where $X \in B$ and $Y \in B$ are fixed,

243  the sum operation on each set $K$ is defined by coordinates as in the set of bases $\left(B, \otimes\right)$, and codon

244  $XY$A is taken as identity element. For example, $K=$ {CGA, CGC, CGG, CGU} with codon CGA as

245  identity element. In other words, $\left(K, +\right) \cong \left(B, \otimes\right) \cong \left(\mathbb{Z}_2^2, +\right)$, which corresponds to the Klein four

246  group as defined on $\mathbb{Z}_2^2$.

247  Notice that for each fixed length *N* we can build manifold heterocyclic groups $S_i$, and each one

248  of them can have different decomposition into *p*-groups. So, each group $S_i$ could be characterized by

249  means of their corresponding canonical decomposition into *p*-groups. This last detail is exemplified

250  in Fig. 2, where an exon region from the enzyme *phospholipase B domain containing-2* (PLBD2)

251  simultaneously encodes information for several aminoacids and carries the footprint to be targeted by

252  the transcription factor REST. Four possible group representations for this exon subregion are

253  suggested in the top of the figure (panel **a**). These types of protein-coding regions are called *duons*,

254  since their base-triplets encode information not only for aminoacids but also for transcription

255  enhancers [33–35].

256
257  **Fig. 2**. The DNA sequence motifs targeted by transcription factors usually integrate genomic building
258  block across several species. **a**, DNA sequence alignment of the protein-coding sequences from
259  phospholipase B domain containing-2 (PLBD2) carrying the footprint sequence motif recognized
260  (targeted) by the Silencing Transcription factor (REST), also known as Neuron-Restrictive Silencer
261  Factor (NRSF) REST (NRSF). **b**. Sequence logo of the footprint motif recognized REST (NRSF) on
262  the exons (derived from TF2DNA dataset [36]). **c**, Translation of the codon sequences using the one-
263  letter symbol of the aminoacids.

12

264     The group representation is particularly interesting for the analysis of DNA sequence motifs,

265     which typically are highly conserved across the species. As suggested in Fig. 2, there are some

266     subregions of DNA or protein sequences where there are few or not gaps introduced and mostly

267     substitution mutations are found. Such subregions conform blocks that can cover complete DNA

268     sequence motifs targeted by DNA biding proteins like transcription factors, which are identifiable by

269     bioinformatic algorithm like BLAST [37]. Herein, the case of double coding called our attention,

270     where the DNA sequence simultaneously encode information transcription factor targeted sequence

271     motif and the codon sequence encoding for aminoacids. Notice that the abelian group

272     $\left(C_g,+\right)\cong\left(\mathbb{Z}_{64},+\right)$ defined on the standard genetic code is enough to quantitatively describe these

273     motifs (Fig. 2). However, a further application of group theory together with additional knowledge

274     on the biological function this motif can lead to a more specific decomposition into abelian groups.

275     No matter how complex a genomic region might be, it has an abelian group representation.

276     As shown in Fig. 3, two different protein-coding (gene) models from two different genome

277     populations can lead to the same direct sum of abelian $p$-groups and the same final aminoacids

278     sequence (protein). The respective exon regions have different lengths and gaps ("-", representing

279     base D in the extended genetic code) were added to exons 1 and 2 (from panel **a**) to preserve the

280     reading frame in the group representation (after transcription and splicing gaps are removed). Both

281     gene models, from panel **a** and **b**, however, lead to the same direct sum of abelian $5$-groups:

282     $$\left(\mathbb{Z}_5\right)^{n+7}\oplus\left(\mathbb{Z}_5^3\right)^3\oplus\left(\mathbb{Z}_5\right)^{3+m+m+2}\oplus\left(\mathbb{Z}_5^3\right)^3\oplus\left(\mathbb{Z}_5\right)^{n+8}.$$

283

284     **Fig. 3**. Two different protein-coding (gene) models can lead to the same abelian group representation
285     and the same protein sequence. A dummy intron was drawn carrying the typical sequence motif
286     targeted by the spliceosome the donor ($GUR$) and acceptor ($Y^m AG$) sites, where $R\in\left\{A,G\right\}$ (purines)
287     and $Y\in\left\{C,U\right\}$, $X$ stands for any base, and $n$ and $m$ indicate the number of bases present in the
288     corresponding sub-sequences (pyrimidines). **a**, A gene model based on a *dummy* consensus sequence
289     where gaps representing base D from the extended genetic code were added to preserve the coding
290     frame, which naturally is restored by splicing soon after Transcription. **b**, A gene model where both
291     exons, 1 and 2, carries a complete set of three codons (base-triplets). Both models, from panels **a** and
292     **b**, leads to the same canonical direct sum of abelian 5-groups.
293

294    An example considering changes on the gene-body reading frames as those introduced by

295    alternative splicing is shown in Fig. 4. Gene-bodies with annotated alternative splicing can easily be

296    represented by any of the groups $(\mathbb{Z}_{5^3})^n$ or $(\mathbb{Z}_5^3)^n$ (Fig.4a). The splicing scenario can include enhancer

297    regions as well (Fig.4b).

298

299    **Fig. 4**. The abelian group representation of a given genome only depend on our current knowledge
300    on its annotation. **a,** the alternative splicing specified for an annotated gene model does not alter the
301    abelian group representation and only would add information for the decomposition of the existing
302    cyclic groups into subgroups. **b**, a more complex gene model including detailed information on the
303    promoter regions. A GC box (G5MG4CU) motif is located upstream of a TATA box (TATAWAW)
304    motif in the promoter region.  The GC box is commonly the binding site for Zinc finger proteins,
305    particularly, Sp1 transcription factors. A putative GC box was included in exon 2, which is an atypical
306    scenario, but it can be found, e.g., in the second exon from the gene encoding for sphingosine kinase
307    1 (SPHK1), transcript variant 2 (NM_182965, CCDS11744.1).

308

309    As commented in the introduction, cytosine DNA methylation is implicitly included in

310    extended base-triple group representation. Typically, methylation analysis of methylomes is

311    addressed to identify methylation changes induced by, for example, environmental changes,

312    lifestyles, age, or diseases. So, in this case the letter D stands for methylated cytosine, since only

313    epigenetic changes are evaluated. A concrete example with two genes from patients with pediatric

314    acute lymphoblastic leukemia (PALL) is presented in Fig. 5.

315

316    **Fig. 5**. Vector representation of differentially methylated exons regions from genes EGEL7 and
317    P2RY1 from patients with pediatric acute lymphoblastic leukemia (PALL). **a**. Segment of exon-6
318    from gene EGFL7. **b**. Segment of exon-1 from gene P2RY1. Methylated cytosines are highlighted in
319    yellow background. In PALL patients, gene EGEL7 mostly hypomethylated and gene P2RY1 mostly
320    hypermethylated in respect to healthy individuals (WT). The encoded aminoacid sequence is given
321    using the one letter symbols. Both genes, EGEL7 and P2RY1, were identified in the top ranked list
322    of differentially methylated genes integrating clusters of hubs in the protein-protein interaction
323    networks from PALL reported in reference [38].

324

325    It is obvious that the MSA from a whole genome population derives from the MSA of every

326    genomic region, from the same or closed related species. At this point, it is worthy to recall that there

327    is not, for example, just one human genome or just one from any other species, but populations of

328    human genomes and genomes populations from other species. Since every genomic region can be

329    represented by the direct sum of abelian cyclic groups of prime-power order, then the whole genome

330    population from individuals from the same or closed related species can be represented as an abelian

331    group, which will be, in turns, the direct sum of abelian cyclic groups of prime-power order. Hence,

332    results lead us to the representation of whole genomes populations of individuals from the same

333    species or close related species (as suggested in Fig.1) by means of direct sum of their group

334    representation into abelian cyclic groups. A general illustration of this modelling would be, for

335    example:

336
$$S = \ldots \oplus (\mathbb{Z}_{5^3})^{n_1} \oplus \overbrace{(\mathbb{Z}_{2^6})^{m_1}}^{motif} \oplus (\mathbb{Z}_{5^3})^{n_2} \oplus \ldots \oplus \overbrace{(\mathbb{Z}_{2}^{2})^{m_2}}^{domain} \oplus \ldots \oplus \overbrace{(\mathbb{Z}_{5^3})^{n_p}}^{domain} \oplus \overbrace{(\mathbb{Z}_{2^6})^{m_p}}^{motif} \ldots [7]$$

337    That is, the fundamental theorem of abelian finite groups has an equivalent in genomics.

338    **Theorem 1**. The genomic architecture from a genome population can be quantitatively represented

339    as an abelian group isomorphic to a direct sum of cyclic groups of prime-power order.

340        The proof of this theorem is self-evident across the discussion and examples presented here.

341    Basically, the group representations of the genetic code lead to the group representations of local

342    genomic domains in terms of cyclic groups of prime-power order, for example, $\left(\mathfrak{B}^3, +\right) \cong \left(\mathbb{Z}_5^3, +\right)$

343    or $\left(C_e, +\right) \cong \left(\mathbb{Z}_{5^3}, +\right)$, till covering the whole genome. As for any finite abelian group, the abelian

344    group representation of genome populations can be expressed in terms a direct sum of abelian cyclic

345    groups of prime-power order. Any new discovering on the annotation of given genome population

346    will only split an abelian group, already defined on some genomic domain/region, into the direct sum

347    of abelian subgroups ∎.

348    ## 4   Discussions

349    Under the assumption that the current forms of life are the result of an evolutionary process started

350    from very simple primordial cells, the current non-coding DNA must be the relict footprint of multiple

351    recombination of ancient DNA domains in all the permissible forms, which in ancient times were

352    rules by an ancient genetic code. In consequence, on this scenario, the group representations of the

353    genetic code are logically extended from relatively small local DNA domains to the whole genome.

354     Examples shown in Fig. 1 to 4 indicates whatever would be the genomic architecture for given

355     species, the observed variations in the individual populations and in populations from closed related

356     species, it can be quantitatively described as the direct sum of abelian cyclic groups. The

357     discovering/annotation of new genomic features will only lead to the decomposition of previous

358     known abelian cyclic groups representing some genomic subregion into direct sums of subgroups. In

359     such algebraic representation DNA sequence motifs for which only substitution mutations happened

360     are specifically represented by the abelian group $\left(C_g,+\right)\cong\left(\mathbb{Z}_{64},+\right)$, in protein coding regions, and

361     by any or combination of groups $\left(B,+\right)\cong\left(\mathbb{Z}_2^2,+\right)$, $\left(B^2,+\right)\cong\left(\mathbb{Z}_2^4,+\right)$ or $B^3/G_{\mathrm{GGA}}\cong\left(\mathbb{Z}_2^5,+\right)$ in

362     non-protein coding regions.

363     Results indicate that the genome architecture of whole populations can be quantitively studied

364     in the framework of abelian group theory. Two sets of MSA, $S_1$ and $S_2$, could split into different

365     cyclic groups and, however, these sequences can be isomorphic between them because have the same

366     canonical decomposition. Particularly, the genetic code abelian group $\left(\mathfrak{B}^3,+\right)\cong\left(\mathbb{Z}_5^3,+\right)$ is enough

367     for an algebraic representation of the genome population from the same species or close related

368     species. However, such a decomposition is biologically poor and, as suggested in Figs. 4 to 5, masks

369     relevant biological features from the genome architecture. A further decomposition into the direct

370     sum of abelian groups will only depends on our knowledge on the genome annotation for specified

371     species.

372     As suggested in Figs. 3 and 4, base D from the extended genetic code (represented as gaps in

373     the MSA) results useful preserving the information on the natural reading frame in the abelian group

374     representation. It is worthy to notices that, for the transcriptional and splicing enzymatic machinery,

375     the information on the reading frame preservation is already in the sequence. Molecular machines

376     perform precise logical operations [39], which in this case result in a sort of molecular *enthymeme*

377     (logical) operation where the conclusion is omitted obeying the principle of cellular economy. In

378     other words, in the algebraic representation of gene and genome populations base D carries real

379     biological information.

380    From several examples provided here, it is clear that there exists a language for the genome

381    architecture that can be represented in terms of sums of finite abelian groups. The future developments

382    of genome annotation from several species can certainly lead to the discovery of logical rules of a

383    such language determining the possible viable variations in the populations. As suggested in Fig. 5,

384    the identification of quotient groups (at larger scale) can permit the stratification of large genome

385    population into equivalence classes corresponding to individual subpopulations, each one of them

386    carrying particular viable variations of species genome architecture.

387    As indicated in reference [12], natural genomic rearrangement like DNA recombination and

388    translocation at structural and functional domain can be represented as group automorphisms and

389    endomorphisms. Biologically, such description corresponds to the fact that the new genetic

390    information is recreated, simply, by way of reorganization of the genetic material in the chromosomes

391    of living organisms [5,40]. The analysis and discussion on the application of the endomorphism ring

392    theory to describe the dynamics of genome population is a promising subject for further studies.

393    Particularly promising is the application of the genomic abelian groups on epigenomic studies,

394    which results when base D stands for the methylated cytosine. As suggested in Fig.5, a precise

395    decomposition of methylation motif into the direct sum of abelian finite group can leads to their

396    classification into unambiguous equivalence classes. This open the doors for the application of based

397    machine-learning bioinformatic approaches to study the methylation changes induced on individual

398    populations by, e.g., environmental changes, aging process and diseases, which is of particular

399    interest in genomic medicine [41].

400    Results presented here would have considerable positive impact on current molecular

401    evolutionary biology, which heavily relies on evolutionary null hypotheses about the past. As

402    suggested in reference [30], the genomic abelian groups open new horizons for the study of the

403    molecular evolutionary stochastic processes (at genomic scale) and with relevant biomedical

404    applications, founded on a deterministic ground, which only depends on the physicochemical

405    properties of DNA bases and aminoacids. In this case, the only molecular evolutionary hypothesis

406    needed about the past is a fact, the existence of the genetic code.

## 5   Conclusions

407

408   Results to date indicate that the genetic code and, ultimately, the physicochemical properties of DNA

409   bases on which the genetic code algebraic structure are defined, has a deterministic effect or at least

410   partially rules on the current genome architectures, in such a way that the abelian group

411   representations of the genetic code are logically extended to the whole genome. In consequence, the

412   fundamental theorem of abelian finite groups can be applied to the whole genome. This result opens

413   new horizons for further genomics studies with the application of the abelian group theory, which

414   currently is well developed and well documented [32,42].

415   Results suggest that the architecture of current population genomes is quite far from

416   randomness and obeys deterministic rules. Although the random nature of the mutational process,

417   only a small fraction of mutations is fixed in genomic populations. In particular, fixation events are

418   ruled by the genetic code architecture, which as shown by Sanchez (2018), it can be simulated as an

419   optimization process by using genetic algorithms [30]. This points to the study of the dynamics of

420   genome populations as a stochastic deterministic process. Genome stochasticity derives from the

421   stochasticity of mutational process and from the stochasticity of biochemical reactions, which gives

422   rise to a rich population diversity and phenotypic plasticity that help to prevent population extinction.

423   The deterministic part derives from its architecture, which can be represented in terms of a canonical

424   direct sum of homocyclic abelian groups derived from the genetic code, hold for all the individuals

425   from the same population/species.

## References

426

427   1.    Sanchez R, Barreto J, Morgado E, Grau R. Abelian Finite Group of DNA Genomic Sequences.
428         arXiv   Quant   Methods.   2005;   1-6.   https://arxiv.org/abs/q-bio/0512042.   Available:
429         https://www.researchgate.net/publication/2183122_Abelian_Finite_Group_of_DNA_Genom
430         ic_Sequences
431   2.    Whalen S, Truty RM, Pollard KS. Enhancer-promoter interactions are encoded by complex
432         genomic signatures on looping chromatin. Nat Genet. 2016;48: 488–496. doi:10.1038/ng.3539
433   3.    Nuebler J, Fudenberg G, Imakaev M, Abdennur N, Mirny LA. Chromatin organization by an
434         interplay of loop extrusion and compartment segregation. Proc Natl Acad Sci U S A.

435         2018;115: E6697–E6706. doi:10.1073/pnas.1717730115

436    4.    Rowley MJ, Corces VG. Organizational principles of 3D genome architecture. Nat Rev Genet.
437         Springer US; 2018;19: 789–800. doi:10.1038/s41576-018-0060-8

438    5.    Piazza A, Heyer WD. Homologous Recombination and the Formation of Complex Genomic
439         Rearrangements.    Trends    Cell    Biol.    Elsevier    Ltd;    2019;29:    135–149.
440         doi:10.1016/j.tcb.2018.10.006

441    6.    Zheng H, Xie W. The role of 3D genome organization in development and cell differentiation.
442         Nat Rev Mol Cell Biol. Springer US; 2019;20: 535–550. doi:10.1038/s41580-019-0132-4

443    7.    Schneider TD. Evolution of biological information. Nucleic Acids Res. 2000;28: 2794–9.
444         Available:
445         http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=102656&tool=pmcentrez&rende
446         rtype=abstract

447    8.    Yockey HP. Origin of life on earth and Shannon's theory of communication. Comput Chem.
448         2000;24: 105–123. doi:10.1016/S0097-8485(99)00050-9

449    9.    Sanchez R, Grau R. A genetic code Boolean structure. II. The genetic information system as
450         a  Boolean  information  system.  Bull  Math  Biol.  2005/07/07.  2005;67:  1017–1029.
451         doi:10.1016/j.bulm.2004.12.004

452    10.   Sanchez R, Mackenzie SA. Information thermodynamics of cytosine DNA methylation.
453         Bardoni  B,  editor.  PLoS  One.  Public  Library  of  Science;  2016;11:  e0150427.
454         doi:10.1371/journal.pone.0150427

455    11.   Sánchez R, Morgado E, Grau R. A genetic code Boolean structure. I. The meaning of Boolean
456         deductions. Bull Math Biol. 2005;67: 1–14. doi:10.1016/j.bulm.2004.05.005

457    12.   Sanchez R, Morgado E, Grau R. Gene algebra from a genetic code algebraic structure. J Math
458         Biol. 2005/07/14. 2005;51: 431–457. doi:10.1007/s00285-005-0332-8

459    13.   Sanchez R, Morgado E, Grau R, Sánchez R. A genetic code Boolean structure. I. The meaning
460         of  Boolean  deductions.  Bull  Math  Biol.  2005/02/05.  2005;67:  1–14.
461         doi:10.1016/j.bulm.2004.05.005

462    14.   Sanchez R, Grau R, Morgado E. A novel Lie algebra of the genetic code over the Galois field
463         of four DNA bases. Math Biosci. 2006;202: 156–174. doi:10.1016/j.mbs.2006.03.017

464    15.   José M V., Zamudio GS, Morgado ER. A unified model of the standard genetic code. R Soc
465         Open Sci. 2017;4: 1–13. doi:10.1098/rsos.160908

466    16.   José M V., Morgado ER, Govezensky T. Genetic Hotels for the Standard Genetic Code:
467         Evolutionary Analysis Based upon Novel Three-Dimensional Algebraic Models. Bull Math
468         Biol. 2011;73: 1443–1476. doi:10.1007/s11538-010-9571-y

469    17.   Sánchez R, Grau R. An algebraic hypothesis about the primeval genetic code architecture.
470         Math    Biosci.    2009/07/18.    2009;221:    60–76.    doi:S0025-5564(09)00114-X    [pii]

471         10.1016/j.mbs.2009.07.001

472   18.   Orgel LE. Prebiotic chemistry and the origin of the RNA world. Crit Rev Biochem Mol Biol.

473         2004;39: 99–123. doi:10.1080/10409230490460765

474   19.   Piccirilli JA, Benner SA, Krauch T, Moroney SE, Benner SA. Enzymatic incorporation of a

475         new base pair into DNA and RNA extends the genetic alphabet. Nature. 1990;343: 33–37.

476         doi:10.1038/343033a0

477   20.   Switzer C, Moronev SE, Benner SA. Enzymatic Incorporation of a New Base Pair into DNA

478         and RNA. J Am Chem Soc. 1989;111: 8322–8323. doi:10.1021/ja00203a067

479   21.   Sanchez R, Grau R, Morgado E. A Novel DNA Sequence Vector Space over an extended

480         Genetic Code Galois Field. MATCH Commun Math Comput Chem. 2006;56: 5–20.

481         Available: http://match.pmf.kg.ac.rs/electronic_versions/Match56/n1/match56n1_5-20.pdf

482   22.   Sanchez R. Evolutionary Analysis of DNA-Protein-Coding Regions Based on a Genetic Code

483         Cube     Metric.     Curr     Top     Med     Chem.     2014;14:     407–417.

484         doi:10.2174/1568026613666131204110022

485   23.   Di Giulio M. LUCA as well as the ancestors of archaea, bacteria and eukaryotes were

486         progenotes: Inference from the distribution and diversity of the reading mechanism of the

487         AUA and AUG codons in the domains of life. BioSystems. Elsevier B.V.; 2020;198: 104239.

488         doi:10.1016/j.biosystems.2020.104239

489   24.   Di Giulio M. Errors of the ancestral translation, LUCA, and nature of its direct descendants.

490         BioSystems. Elsevier B.V.; 2021;206: 104433. doi:10.1016/j.biosystems.2021.104433

491   25.   Smith ZD, Meissner A. DNA methylation: Roles in mammalian development. Nat Rev Genet.

492         Nature Publishing Group; 2013;14: 204–220. doi:10.1038/nrg3354

493   26.   Severin PMD, Zou X, Gaub HE, Schulten K. Cytosine methylation alters DNA mechanical

494         properties. Nucleic Acids Res. 2011;39: 8740–51. doi:10.1093/nar/gkr578

495   27.   Sriraman A, Debnath TK, Xhemalce B, Miller KM. Making it or breaking it: DNA

496         methylation     and     genome     integrity.     Essays     Biochem.     2020;     687–703.

497         doi:10.1042/ebc20200009

498   28.   Fuchs L. Abelian groups. Publishing House of the Hungarian Academy of Sciences.

499         Publishing House of the Hungarian Academy of Sciences; 1958.

500   29.   Crick FHC. The Origin of the Genetic Code. J Mol Biol. 1968;38: 367–379.

501   30.   Sanchez R. Symmetric Group of the Genetic-Code Cubes. Effect of the Genetic-Code

502         Architecture on the Evolutionary Process. MATCH Commun Math Comput Chem. 2018;79:

503         527–560.                                                         Available:

504         http://match.pmf.kg.ac.rs/electronic_versions/Match79/n3/match79n3_527-560.pdf

505   31.   José M V, Morgado ER, Sánchez R, Govezensky T. The 24 Possible Algebraic

506         Representations of the Standard Genetic Code in Six or in Three Dimensions. Adv Stud Biol.

507      2012;4: 119–152. Available: http://www.m-hikari.com/asb/asb2012/asb1-4-2012/joseASB1-
508      4-2012-1.pdf

509  32. Fuchs L. Infinite Abelian Groups, Volume 1. 1st Editio. Academic Press; 1970.

510  33. Stergachis AB, Haugen E, Shafer A, Fu W, Vernot B, Reynolds A, et al. Exonic transcription
511      factor binding directs codon choice and affects protein evolution. Science (80- ). 2013;342:
512      1367–72. doi:10.1126/science.1243490

513  34. Reyna-Llorens I, Burgess SJ, Reeves G, Singh P, Stevenson SR, Williams BP, et al. Ancient
514      duons may underpin spatial patterning of gene expression in C 4 leaves. Proc Natl Acad Sci
515      U S A. 2018;115: 1931–1936. doi:10.1073/pnas.1720576115

516  35. Yadav VK, Smith KS, Flinders C, Mumenthaler SM, De S. Significance of duon mutations in
517      cancer genomes. Sci Rep. Nature Publishing Group; 2016;6: 1–9. doi:10.1038/srep27437

518  36. Pujato M, Kieken F, Skiles AA, Tapinos N, Fiser A. Prediction of DNA binding motifs from
519      3D models of transcription factors; identifying TLX3 regulated genes. Nucleic Acids Res.
520      2014;42: 13500–13512. doi:10.1093/nar/gku1228

521  37. Yang M, Derbyshire MK, Yamashita RA, Marchler-Bauer A. NCBI's Conserved Domain
522      Database and Tools for Protein Domain Analysis. Curr Protoc Bioinforma. 2020;69: 1–25.
523      doi:10.1002/cpbi.90

524  38. Sanchez R, Mackenzie SA. Integrative Network Analysis of Differentially Methylated and
525      Expressed Genes for Biomarker Identification in Leukemia. Sci Rep. 2020;10: 2123.
526      doi:10.1038/s41598-020-58123-2

527  39. Schneider TD. Theory of molecular machines. II. Energy dissipation from molecular
528      machines.    J    Theor    Biol.    1991;148:    125–137.    Available:
529      http://www.ncbi.nlm.nih.gov/pubmed/2016881

530  40. Yu M, Ren B. The three-dimensional organization of mammalian genomes. Annu Rev Cell
531      Dev Biol. 2017;33: 265–289. doi:10.1146/annurev-cellbio-100616-060531

532  41. Salameh Y, Bejaoui Y, El Hajj N. DNA Methylation Biomarkers in Aging and Age-Related
533      Diseases. Front Genet. 2020;11: 1–11. doi:10.3389/fgene.2020.00171

534  42. Fuchs L. Infinite Abelian Groups, Volume 2. Academic Press; 1973.

535

Figure 1

Figure 3

**a**

5'UTR — Exon 1 — Intron 1 — Exon 2 — Intron 2 — Exon 3

$\dots$ ACGC AUGCACC--GUR$X^j$Y$^k$AG-CUGCCCUAC--GUR$X^n$Y$^m$AG-UGAGAACU

$(\mathbb{Z}_5)^{n+4}$ $(\mathbb{Z}_{5^3})^3$ or $(\mathbb{Z}_5^3)^3$ $(\mathbb{Z}_5)^{3+n+m+2}$ $(\mathbb{Z}_{5^3})^4$ or $(\mathbb{Z}_5^3)^4$ $(\mathbb{Z}_5)^{3+n+m+2}$ $(\mathbb{Z}_{5^3})^3$ or $(\mathbb{Z}_5^3)^3$

**Transcription & Alternative Splicing**

| AUG | CAC | CCU | GCC | CUA | CUG | AGA | ACU |
|-----|-----|-----|-----|-----|-----|-----|-----|
| M | H | P | A | L | L | R | T |

| AUG | CAC | CUG | AGA | ACU |
|-----|-----|-----|-----|-----|
| M | H | L | R | T |

**b**

**Gene body**

**Promoter region**

GC Box — TATA Box — 5'UTR — Exon 1 — Intron — Exon 2 — GC Box — 3'UTR

$\dots$G$^5$MG$^4$CU$X^i$TATAWAW$X^j$G$X^k$C ATGCACCC-GUR$X^n$Y$^m$AG--TGCCG$^3$CG$^4$CTGAGA$\dots$

$\mathbb{Z}_2$ $(\mathbb{Z}_5)^2$

$(\mathbb{Z}_5)^5$ $(\mathbb{Z}_5)^4$ $(\mathbb{Z}_5)^i$ $(\mathbb{Z}_5)^7$ $\mathbb{Z}_5$ $(\mathbb{Z}_5)^{k+2}$ $(\mathbb{Z}_{5^3})^3$ or $(\mathbb{Z}_5^3)^3$ $(\mathbb{Z}_5)^{n+m+5}$ $(\mathbb{Z}_{5^3})^3$ or $(\mathbb{Z}_5^3)^3$ $(\mathbb{Z}_5)^{p+2}$

**Transcription & Splicing**

| AUG | CAC | CCU | GCC | GGG | CGG | GGC | UGA |
|-----|-----|-----|-----|-----|-----|-----|-----|
| M | H | P | A | G | R | G | Stop |

Figure 4

Figure 5

Figure 2