

1

Genomic Abelian Finite Groups

2

3 Robersey Sanchez ¹ and Jesús Barreto²

4

5 ¹Department of Biology. Pennsylvania State University, University Park, PA 16802.

6 E-mail: rus547@psu.edu

7 ORCID: <https://orcid.org/0000-0002-5246-1453>

8

9 ²Universidad Central "Marta Abreu" de Las Villas. Santa Clara. Cuba.

10 E-mail: barretouclv@gmail.com

11

12 ¹ Corresponding author:

13 rus547@psu.edu

14

15 Abstract

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

24

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

26

27

28 **1 Introduction**

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 [1]. Current experimental genomic studies suggest that genome architectures must obey
32 specific mathematical biophysics rules [2–5].

33 Experimental results points to an injective relationship: *DNA sequence* \rightarrow *3D chromatin*
34 *architecture* [2,3,5], 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 [4]. 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
38 analyze the genetic information system as a communication system [6–9].

39 Under the assumption that current forms of life evolved from simple primordial cells with
40 very simple genomic structure and robust coding apparatus, the genetic code is a fundamental link
41 to the primeval form of live, which played an essential role on the primordial architecture. The
42 genetic code, the set of biochemical rules used by living cells to translate information encoded within
43 genetic material into proteins, sets the basis for our understanding of the mathematical logic inherent
44 to the genetic information system [8,10]. The genetic code is the cornerstone of live on earth. Not a
45 single 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
47 quantitative relationship between the coding apparatus and the mutational process on protein-coding
48 regions [11–15]. Formally the genetic code only is limited to translated coding regions where the
49 number of RNA bases is a multiple of 3. However, as suggested in reference [16], the difficulties
50 in prebiotic synthesis of the nucleosides components of RNA (nucleo-base + sugar) and suggested
51 that some of the original bases may not have been the present purines or pyrimidines [17]. Piccirilli
52 et al. [18] demonstrated that the alphabet can in principle be larger. Switzer et al. [19] have shown
53 an enzymatic incorporation of new functionalized bases into RNA and DNA. This expanded the

54 genetic alphabet from 4 to 5 or more letters, which permits new base pairs, and provides RNA
55 molecules with the potential to greatly increase their catalytic power.

56 It is important to notice that even in the current (*friendly*) environmental conditions not a
57 single cell can survive without a DNA repair enzymatic machinery and that such an enzymatic
58 machinery did not exist at all in the primaeval forms of life. Here, we are confronting the *chicken*
59 *and egg* problem. To date, the best solution (to our knowledge) is the admission of alternative base-
60 pairs in the primordial DNA alphabet which, as suggested in the studies on the prebiotic chemistry,
61 could contribute to the thermal and general physicochemical stability of the primordial DNA
62 molecules.

63 Several algebraic structures have been proposed including an additional letter into the DNA
64 alphabet: A, C, G, T. The new letter (D) stands for current insertion deletion/mutations or for
65 alternative wobble base pairing, which would be a relict fingerprint from primordial enzymes
66 derived from a more degenerated ancestral genetic code [16,20,21]. Supporting evidence for the
67 existence of a more degenerated ancestral genetic code built up on a larger alphabet is found in the
68 tRNA anticodon region permitting wobble base pairing by including, e.g., bases such as: inosine (in
69 eukariotes), agmatidine (in archaea), and lysidine (in bacteria), which has been proposed as
70 evolutionary solutions to the need for lower the high translational noise connected to the reading of
71 the AUA and AUG codons [22,23]. Additionally, various alternative base pairs like methylated
72 cytosine and adenine are still present in the current genomes playing an important role in the
73 epigenetic adaptation of organismal populations to the continuous environmental changes [9,24].

74 Cytosine DNA methylation results from the addition of methyl groups to cytosine C5
75 residues, and the configuration of methylation within a genome provides trans-generational
76 epigenetic information. These epigenetic modifications can influence the transcriptional activity of
77 the corresponding genes, or maintain genome integrity by repressing transposable elements and
78 affecting long-term gene silencing mechanisms [25,26].

79 In this scenario, we shall show that all possible DNA molecules and, consequently, genomes
80 can be described by way of finite abelian groups which can be split into the direct sum of homocyclic

81 2-groups and 5-groups defined on the genetic code. A homocyclic group is a direct sum of cyclic
82 groups of the same order. Any finite abelian group can be decomposed into a direct sum of
83 homocyclic p -groups [27], i.e., a group in which the order of every element is a power of a primer number
84 p .

85 The genetic code algebraic structures under scrutiny in the mentioned references covered
86 rings and vector spaces with a common feature, the corresponding additive group is an abelian group
87 of prime-power order. Next, to help a better comprehension of the current work, a brief introductory
88 summary on these groups is provided. Results presented here generalizes the application of the
89 genetic code algebras (reported in several publications) to the whole genome.

90 **1.1 Reported genetic code abelian groups relevant for the current study**

91 Herein, we assume that readers are familiar with the definition of abelian group, which otherwise
92 can be found in textbooks and elsewhere including Wikipedia. Nevertheless, all the abelian groups
93 discussed here are isomorphic to the well-known abelian groups of integer module n , which are
94 easily apprehended by a college-average educated mind. For example, the abelian group defined on
95 the set $\{0, 1, 2, 3, 4\}$, which corresponds to the group of integer modulo 5 (\mathbb{Z}_5), where $(2 + 1) \bmod$
96 $5 = 3$, $(1 + 3) \bmod 5 = 4$, $(2 + 3) \bmod 5 = 0$, etc. The subjacent biophysical and biochemical
97 reasonings to define the algebraic operations on the set of DNA bases and on the codon set were
98 given in references [11,13,16].

99 *1.1.1 The \mathbb{Z}_{64} -algebras of the genetic code (C_g)*

100 The \mathbb{Z}_{64} -algebras of the genetic code (C_g) and gene sequences were stated several years ago. In
101 the \mathbb{Z}_{64} -algebra C_g the sum operation, defined on the codon set, is a manner to consecutively
102 obtain all codons from the codon AAC (UUG) in such a way that the genetic code will represent a
103 non-dimensional code scale of amino acids interaction energy in proteins.

104 A description of the genetic code abelian finite group ($C_g, +$) can be found in [11]. Group
105 $(C_g, +)$ is isomorphic to the group on the set \mathbb{Z}_{2^6} (the sum of integer modulo 64), which formally

106 will be expressed as $(C_g, +) \cong (\mathbb{Z}_{2^6}, +)$. The mapping of the set of codons $X_1X_2X_3 \in C_g$ into the set
 107 \mathbb{Z}_{2^6} is straightforward after consider the bijection $A \leftrightarrow 0, C \leftrightarrow 1, G \leftrightarrow 2, U \leftrightarrow 3$ and the function
 108 $g(x) = 4x_1 + 16x_2 + x_3$. For example:

$$\begin{array}{r} \text{AGC} \leftrightarrow 33 \\ + \text{UGU} \leftrightarrow +47 \\ \hline \text{ACA} \leftrightarrow 16 \pmod{64} \end{array} \quad \begin{array}{r} \text{AGC} \leftrightarrow 33 \\ + \text{ACU} \leftrightarrow +18 \\ \hline \text{AUU} \leftrightarrow 51 \pmod{64} \end{array} \quad \begin{array}{r} \text{GGC} \leftrightarrow 41 \\ \text{CUA} \leftrightarrow +52 \\ \hline \text{UCC} \leftrightarrow 29 \pmod{64} \end{array}$$

109

110 The \mathbb{Z}_{64} -algebra C_g , however, is limited to protein-coding regions, while it is well known that,
 111 in eukaryotes, only a small fraction of the genome –about 3%– called open reading frame (ORF)
 112 encodes for proteins [18]. Since non-coding DNA sequences can have a base pairs number not
 113 multiple of three, complete chromosomes and genomes cannot be described by means of group
 114 $(C_g, +)$. In addition, natural genomic variations that includes insertions and deletion mutations
 115 (indel mutations) across individuals from the same population and close-related populations from
 116 different species cannot be represented with group $(C_g, +)$.

117 *1.1.2 The $(\mathbb{Z}_2^6, +)$ group of the genetic code (C_g)*

118 Group $(C_g, +)$ is the additive group of a module over a ring, which however, do not conform to a
 119 vector space. To build a genetic code vector space, a Galois field ($GF(4)$) structure in the ordered
 120 base set $B = \{G, U, A, C\}$ was introduced in reference [13]. In particular, an isomorphism with the
 121 Galois field is defined by means of its binary representation $\mathbb{Z}_2 \times \mathbb{Z}_2 = \{(0, 0), (0, 1), (1, 0), (1, 1)\}$,
 122 i.e. a 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
 123 the DNA base set $B = \{G, U, A, C\}$ to the set of binary duplets (α_1, α_2) is stated., where
 124 $\alpha_i \in \mathbb{Z}_2 = \{0, 1\}$, for $i \in \{1, 2\}$. For example, the bijection f is defined as:

125
$$f(G) = (0, 0), f(U) = (0, 1), f(A) = (1, 0), f(C) = (1, 1).$$

126 The additive group of bases is the Klein four-group, which is defined by the group
 127 presentation: $V = \{U, A \mid U + U = A + A = C + C = G, A + U = C\}$, i.e., $(B, +) \cong (\mathbb{Z}_2^2, +)$. Next, the

128 abelian group on the set of codons B^3 was defined as the direct third power $B^3 = B \times B \times B$ of the
129 group $(B, +)$, i.e. $(B^3, +) = (B, +) \times (B, +) \times (B, +)$, which is isomorphic to the group:
130 $(\mathbb{Z}_2^6, +) = (\mathbb{Z}_2^2, +) \times (\mathbb{Z}_2^2, +) \times (\mathbb{Z}_2^2, +)$, i.e., $(B^3, +) \cong (\mathbb{Z}_2^6, +)$. The sum operation on the set $(B^3, +)$
131 follows from the sum operation by coordinates.

132 As pointed out before by Crick, the first two bases of codons determine the physicochemical
133 properties of aminoacids [28]. The four encoded amino acids of every class are either the same or
134 show very similar physicochemical properties. This genetic code regularity is captured by the
135 quotient group B^3/G_{GGA} , where G_{GGA} is a subgroup of B^3 integrated by the elements
136 $\{GGG, GGA\}$ (the elements of the quotient group B^3/G_{GGA} are given in Table 5 from [13]). The
137 quotient group 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
138 of this group represents an equivalence class of codons. Two triplets $X_1X_2X_3$ and $Y_1Y_2Y_3$ are
139 equivalent if, and only if, the difference $X_1X_2X_3 + Y_1Y_2Y_3 \in G_{DDA}$. In biological terms, substitution
140 mutations involving codons from the same class will not alter (or at least no substantially alter in
141 most of the cases) the physicochemical properties of the encoded protein domains, since in the worst
142 scenario involves aminoacids with very close physicochemical properties, with the exception of
143 codon for aminoacid tryptophan.

144 1.1.3 The \mathbb{Z}_{125} group of the extended genetic code (C_e)

145 The extension of the *genetic code group* $(C_e, +)$ follows straightforward from the extension of the
146 codon set, which is easily accomplished extending the source alphabet of the standard genetic code:
147 $\{A, C, G, U\}$ and, consequently, extending the base triplet set (extended triplet) as $X_1X_2X_3, X_i \in \{D,$
148 $A, C, G, U\}$ [21]. The new algebraic structure $(C_e, +)$ is isomorphic to the abelian group defined on
149 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
150 of codons $X_1X_2X_3 \in C_e$ into the set \mathbb{Z}_{5^3} is straightforward after consider the bijection

151 $D \leftrightarrow 0, A \leftrightarrow 1, C \leftrightarrow 2, G \leftrightarrow 3, U \leftrightarrow 4$ and the function $g(x) = 5x_1 + 25x_2 + x_3$ (see Table 1). For

152 example:

$$\begin{array}{rcl} \text{AGC} & \leftrightarrow & 82 \\ + \text{UGU} & \leftrightarrow & +99 \\ \hline \text{ACA} & \leftrightarrow & 56 \text{ mod } 125 \end{array} \qquad \begin{array}{rcl} \text{AGC} & \leftrightarrow & 82 \\ + \text{DCU} & \leftrightarrow & +54 \\ \hline \text{CDA} & \leftrightarrow & 11 \text{ mod } 125 \end{array} \qquad \begin{array}{rcl} \text{GGC} & \leftrightarrow & 92 \\ + \text{CUD} & \leftrightarrow & +110 \\ \hline \text{DGC} & \leftrightarrow & 77 \text{ mod } 125 \end{array}$$

153

154 **Table 1.** Ordered set of extended triplets corresponding to the elements from \mathbb{Z}_5^3

a	D		A		C		G		U						
	I	III	I	III	I	III	I	III	I	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	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

155 ^a Bijection between the base-triplets set and the elements from sets \mathbb{Z}_5^3 as given in [21].

156

157 *1.1.4 The $(\mathbb{Z}_5^3, +)$ group of the extended genetic code (C_e)*

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

159 [16]. This structure led to the definition of a \mathbb{Z}_5 -vector space \mathfrak{B}^3 over the set $\mathfrak{B}^3 = \mathfrak{B} \times \mathfrak{B} \times \mathfrak{B}$

160 isomorphic to the set $\mathbb{Z}_5^3 = \mathbb{Z}_5 \times \mathbb{Z}_5 \times \mathbb{Z}_5$ [16,29]. But here, we are interested only in the abelian

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

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

163 , the sum of integers modulo 5). For example, $C + U \leftrightarrow (2 + 4) \bmod 5 = 1 \leftrightarrow A$. The sum operation
164 on the set \mathfrak{B}^3 follows from the sum operation by coordinates.

165 It is worthy to notice that there 24 way to define each one of the above mentioned algebraic
166 structures [29,30]. Nevertheless, for each defined genetic code group, there is only one (genetic code
167 abelian group) up to isomorphism, which lead to their representation as an abelian group, where the
168 sum operation corresponds to the sum of integer modulo $n \in \{2, 2^6, 5, 5^3\}$.

169 2 The General Theoretical Model

170 Herein, it will be showed that, in a general scenario, the whole genome population from any species
171 or close related species, can be algebraically represented as a direct sum of abelian cyclic groups or
172 more specifically abelian p -groups. Basically, we propose the representation of multiple sequence
173 alignments (MSA) of length N as the direct sum:

$$174 \quad G = \left(\mathbb{Z}_{p_1}\right)^{n_1} \oplus \left(\mathbb{Z}_{p_2}\right)^{n_2} \oplus \dots \oplus \left(\mathbb{Z}_{p_k}\right)^{n_k} \quad (1)$$

175 Where $p_i \in \{2, 5, 2^6, 5^3\}$ and $N = n_1 + n_2 + \dots + n_k$. Here, we assume the usual definition of direct sum
176 of groups [31]. Let B_i ($i \in I = \{1, \dots, n\}$) be a family of subgroups of G , subject to the following
177 two conditions:

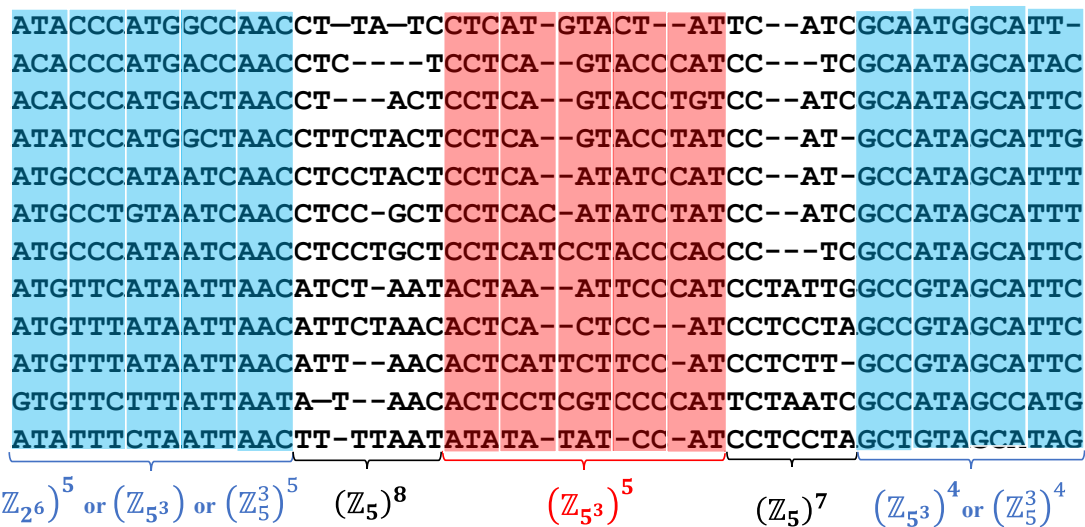
178 1) $\sum B_i = G$. That is, B_i together generates G .

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

180 Then, it is said that G is the direct sum of its subgroups B_i , which formally is expressed by the
181 expression: $G = \bigoplus_i B_i$ or $G = B_1 \oplus \dots \oplus B_n$.

182 In superior organisms, genomic DNA sequences are integrated by intergenic regions and gene
183 regions. The former are the larger regions, while the later includes the protein-coding regions as
184 subsets. The MSA of DNA and protein-coding sequences reveals allocations of the nucleotide bases
185 and aminoacids into stretched of *strings*. The alignment of these stretched would indicate the

186 presence of substitution, *indel* mutations. As a result, the alignment of a whole chromosome DNA
 187 sequences from several individuals from the same or close-related species can be split into well-
 188 defined subregions or domains, and each one of them can be represented as homocyclic abelian
 189 groups, i.e., a cyclic group of *prime-power* order (Fig. 1). As a result, each DNA sequence is
 190 represented as a N -dimensional vector with numerical coordinates representing bases and codons.
 191



192 **Fig.1.** A typical DNA multiple sequence alignment (MSA) including segments of protein-coding
 193 regions. A MSA would include the presence of substitution, insertion and deletion mutations (*indel*
 194 mutations). The aligned sequences can be grouped into blocks, which can be algebraically
 195 represented by abelian groups.
 196
 197

198 An intuitive mathematical representation of MSA is implicit in Fig.1, with following
 199 observations:

- 200 a) Every DNA sequence from the MSA and every subsequence on it can be represented as a
 201 vector with element coordinates defined in some abelian group. For example,
 202 $(C_g, +) \cong (\mathbb{Z}_{64}, +)$, the first five codons from the first DNA sequence from Fig. 1,
 203 $\{ATA, CCC, ATG, GCC, AAC\} \in (C_g, +)$, can be represented by the vector of integers:
 204 $\{48, 21, 50, 25, 1\}$ where each coordinate is an element from group $(\mathbb{Z}_{64}, +)$ (see Table 1
 205 from reference [11] and the introduction section).

206 b) Any MSA can be algebraically represented as a symbolic composition of abelian
 207 groups each one of them is isomorphic to an abelian group of integers module n . Such a
 208 composition can be algebraically represented as a direct sum of homocyclic abelian groups.
 209 For example, the multiple sequence alignment from Fig. 1 can be represented by the direct
 210 sum of abelian groups:

$$211 \quad G = (\mathbb{Z}_{2^6})^5 \oplus (\mathbb{Z}_5)^8 \oplus (\mathbb{Z}_{5^3})^5 \oplus (\mathbb{Z}_5)^7 \oplus (\mathbb{Z}_{5^3})^4 \quad (2)$$

212 In more specific scenario, the multiple sequence alignment from Fig. 1 can be represented by
 213 the direct sum of abelian 2-groups and 5-groups:

$$214 \quad G = (\mathbb{Z}_2^6)^5 \oplus (\mathbb{Z}_5)^8 \oplus (\mathbb{Z}_{5^3})^5 \oplus (\mathbb{Z}_5)^7 \oplus (\mathbb{Z}_2^6)^4 \quad (3)$$

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

$$216 \quad G = (\mathbb{Z}_5^3)^5 \oplus (\mathbb{Z}_5)^8 \oplus (\mathbb{Z}_{5^3})^5 \oplus (\mathbb{Z}_5)^7 \oplus (\mathbb{Z}_5^3)^4 \quad (4)$$

217 Although the above *direct sums* of abelian groups provides a useful compact representation
 218 of MSA, for application purposes to genomics, we would also consider to use the concept of direct
 219 product (*cartesian sum or complete direct sums*) [31]. Next, let S be a set of abelian cyclic groups
 220 identified in the MSA M of length N (i.e., every DNA sequence from M has N bases). Let ℓ_i the
 221 number of bases or triples of bases covered on M by group $S_i \in S$ where $\sum_i \ell_i = N$. Hence, each
 222 DNA sequence on the M can be represented by a cartesian product (b_1, \dots, b_n) where $b_i \in S_i$
 223 ($i = 1, \dots, n$) and $n = |S|$. Let \mathcal{G}_i be a group defined on the set of all elements $(0, \dots, 0, b_i, 0, \dots, 0)$
 224 where $b_i \in S_i$ stands on the i^{th} place and 0 everywhere else. It is clear that $S_i \cong \mathcal{G}_i$. In this context,
 225 the set of all vectors (b_1, \dots, b_n) with equality and addition of vectors defined coordinate-wise
 226 becomes a group (\mathcal{G}) named direct product (cartesian sum) of groups $S_i(\mathcal{G}_i)$, i.e.:

$$227 \quad G = \otimes_i S_i = \oplus_i \mathcal{G}_i \quad (5)$$

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

229 3 Results

230 Results essentially comprise an application of the fundamental theorem of abelian finite groups
231 [27,31]. By this theorem every finite abelian group G is isomorphic to a direct sum of cyclic groups
232 of prime-power order of the form:

$$233 \quad G = \mathbb{Z}_{p_1^{\alpha_1}} \oplus \mathbb{Z}_{p_2^{\alpha_2}} \oplus \cdots \oplus \mathbb{Z}_{p_n^{\alpha_n}} \quad (6)$$

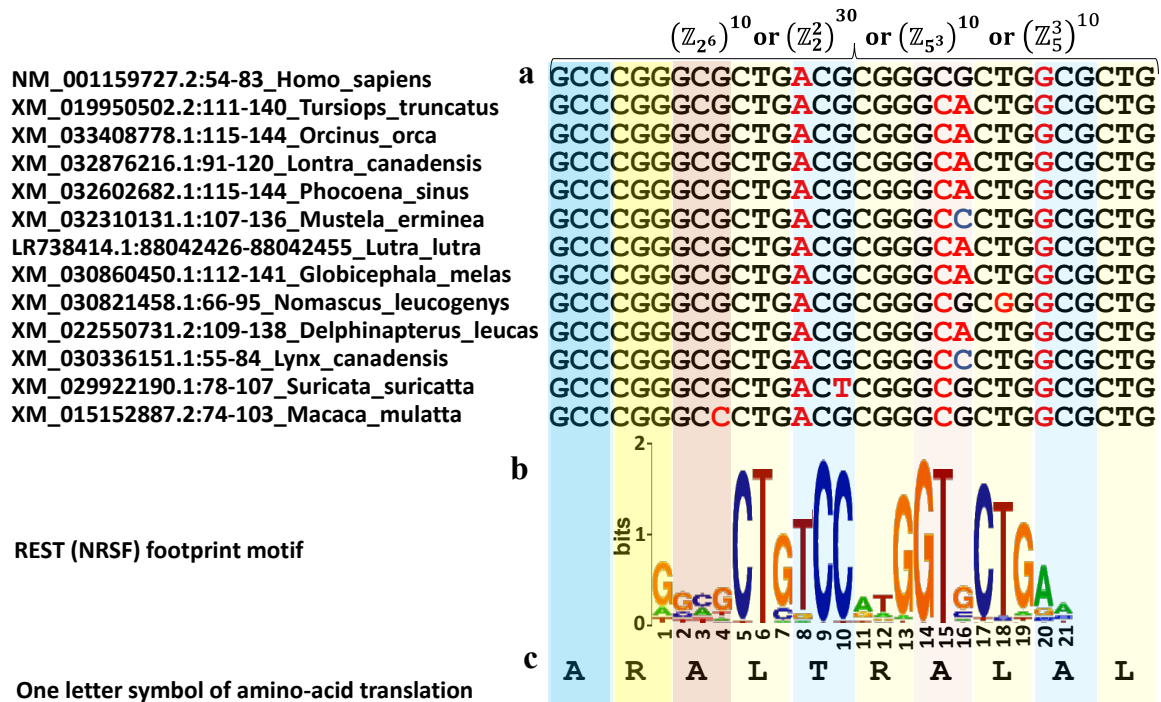
234 Or (in short) $G = \bigoplus_{i=1}^n \mathbb{Z}_{p_i^{\alpha_i}}$, where the p_i 's are primes (not necessarily distinct), $\alpha_i \in \mathbb{N}$ and $\mathbb{Z}_{p_i^{\alpha_i}}$
235 is the group of integer module $p_i^{\alpha_i}$. The abelian group representation of the MSA from Fig. 1 given
236 by expressions (1) and (2) correspond to the cases where the finite abelian group G is a direct sum
237 of *prime-power order*, while expression (3) reflects the fact that any finite abelian group can be
238 decomposed into a direct sum of homocyclic p -groups [27,31], in this $p = 5$.

239 As is showed in Fig 1, this abelian group is a heterocyclic group that split into a direct sum
240 of homocyclic *prime-power order*, each one of them split into the direct sum of cyclic p -groups with
241 same order. For example, in expression [4] we have the subgroup: $(\mathbb{Z}_5^3)^4 = \bigoplus_{i=1}^{12} \mathbb{Z}_5$, which is a
242 direct sum of 12 homocyclic 5-groups $(\mathbb{Z}_5, +) \cong (\mathfrak{B}, +)$. The case of \mathbb{Z}_{2^6} representation of the genetic
243 code (as given in [11]) is less evident. It follows from the fact that the genetic code table is integrated
244 by 16 subsets of codons with form $K = \{XYA, XYC, XYG, XYU\}$, where $X \in B$ and $Y \in B$ are
245 fixed, the sum operation on each set K is defined by coordinates as in the set of bases (B, \otimes) , and
246 codon XYA is taken as identity element. For example, $K = \{CGA, CGC, CGG, CGU\}$ with codon
247 CGA as identity element. In other words, $(K, +) \cong (B, \otimes) \cong (\mathbb{Z}_2^2, +)$, which corresponds to the
248 Klein four group as defined on \mathbb{Z}_2^2 .

249 Notice that for each fixed length N we can build manifold heterocyclic groups S_i , and each
250 one of them can have different decomposition into p -groups. So, each group S_i could be
251 characterized by means of their corresponding canonical decomposition into p -groups. This last

252 detail is exemplified in Fig. 2, where an exon region from the enzyme *phospholipase B domain*
 253 *containing-2* (PLBD2) simultaneously encodes information for several aminoacids and carries the
 254 footprint to be targeted by the transcription factor REST. Four possible group representations for
 255 this exon subregion are suggested in the top of the figure (panel a). These types of protein-coding
 256 regions are called *duons*, since their base-triplets encode information not only for aminoacids but
 257 also for transcription enhancers [32–34].

258



259 **Fig. 2.** The DNA sequence motifs targeted by transcription factors usually integrate genomic
 260 building block across several species. **a**, DNA sequence alignment of the protein-coding sequences
 261 from phospholipase B domain containing-2 (PLBD2) carrying the footprint sequence motif
 262 recognized (targeted) by the Silencing Transcription factor (REST), also known as Neuron-
 263 Restrictive Silencer Factor (NRSF) REST (NRSF). **b**, Sequence logo of the footprint motif
 264 recognized REST (NRSF) on the exons (derived from TF2DNA dataset [35]). **c**, Translation of the
 265 codon sequences using the one-letter symbol of the aminoacids.

266

267

268

269

270

271

272

273

The group representation is particularly interesting for the analysis of DNA sequence motifs,
 which typically are highly conserved across the species. As suggested in Fig. 2, there are some
 subregions of DNA or protein sequences where there are few or not gaps introduced and mostly
 substitution mutations are found. Such subregions conform blocks that can cover complete DNA
 sequence motifs targeted by DNA binding proteins like transcription factors, which are identifiable
 by bioinformatic algorithm like BLAST [36]. Herein, the case of double coding called our attention,

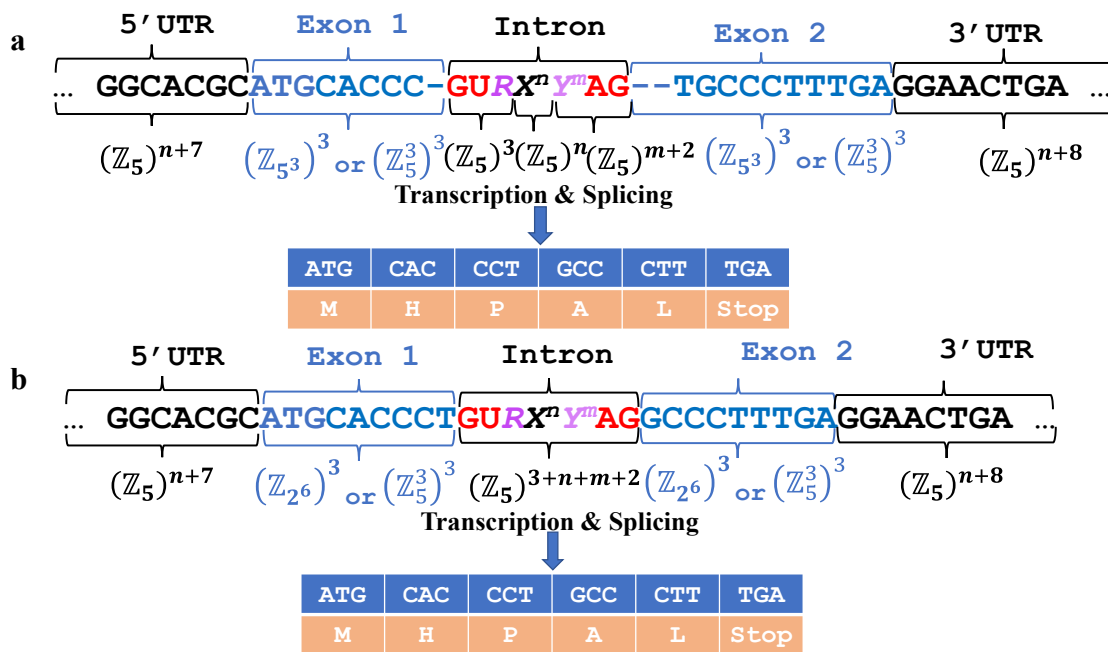
274 where the DNA sequence simultaneously encode information transcription factor targeted sequence
275 motif and the codon sequence encoding for aminoacids. Notice that the abelian group
276 $(C_g, +) \cong (\mathbb{Z}_{64}, +)$ defined on the standard genetic code is enough to quantitatively describe these
277 motifs (Fig. 2). However, a further application of group theory together with additional knowledge
278 on the biological function reveal a more specific decomposition of the motif into abelian groups.

279 No matter how complex a genomic region might be, it has an abelian group representation.
280 As shown in Fig. 3, two different protein-coding (gene) models from two different genome
281 populations can lead to the same direct sum of abelian p -groups and the same final aminoacids
282 sequence (protein). The respective exon regions have different lengths and gaps (“-”, representing
283 base D in the extended genetic code) were added to exons 1 and 2 (from panel **a**) to preserve the
284 reading frame in the group representation (after transcription and splicing gaps are removed). Both
285 gene models, from panel **a** and **b**, however, lead to the same direct sum of abelian 5-groups:

$$286 \quad (\mathbb{Z}_5)^{n+7} \oplus (\mathbb{Z}_5^3)^3 \oplus (\mathbb{Z}_5)^{3+m+m+2} \oplus (\mathbb{Z}_5^3)^3 \oplus (\mathbb{Z}_5)^{n+8}.$$

287 An example considering changes on the gene-body reading frames as those introduced by
288 alternative splicing is shown in Fig. 4. Gene-bodies with annotated alternative splicing can easily be
289 represented by any of the groups $(\mathbb{Z}_{5^3})^n$ or $(\mathbb{Z}_5^3)^n$ (Fig.4a). The splicing scenario can include
290 enhancer regions as well (Fig.4b). As commented in the introduction, cytosine DNA methylation is
291 implicitly included in extended base-triple group representation. Typically, methylation analysis of
292 methylomes is addressed to identify methylation changes induced by, for example, environmental
293 changes, lifestyles, age, or diseases. So, in this case the letter D stands for methylated cytosine, since
294 only epigenetic changes are evaluated. A concrete example with two genes from patients with
295 pediatric acute lymphoblastic leukemia (PALL) is presented in Fig. 5.

296

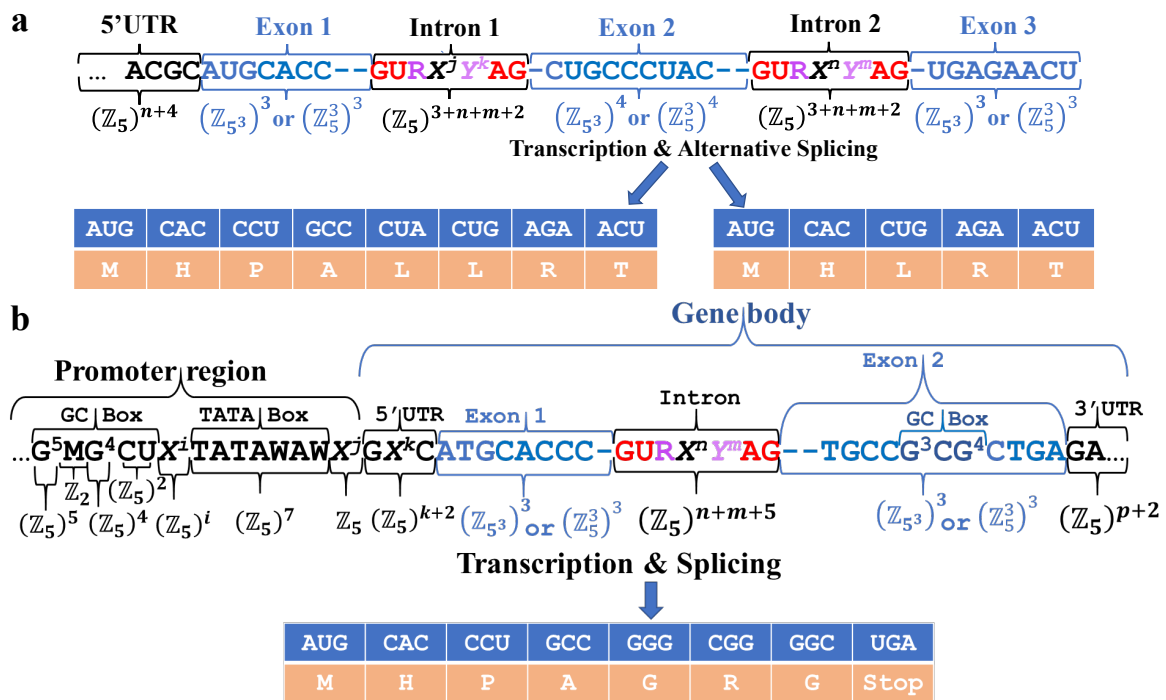


297
298
299
300
301
302
303
304
305
306
307

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

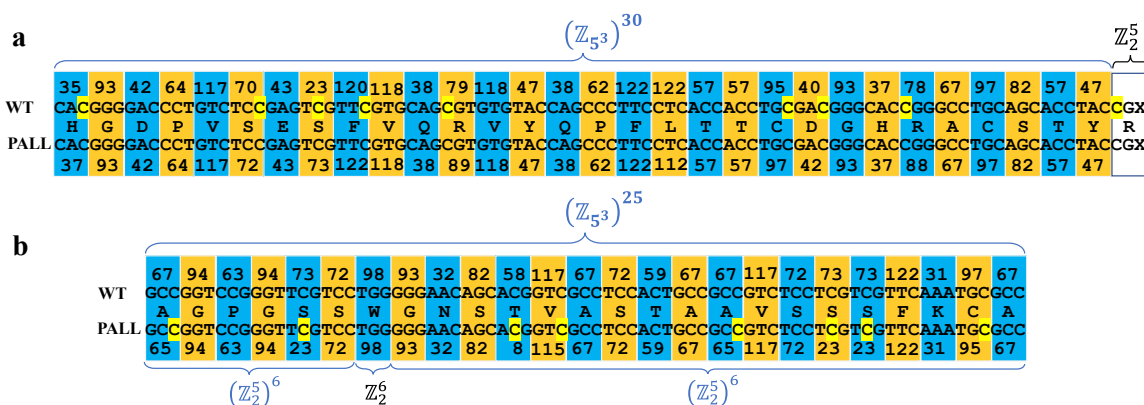
308
309
310
311
312
313
314
315
316
317

It is obvious that the MSA from a whole genome population derives from the MSA of every genomic region, from the same or closed related species. At this point, it is worthy to recall that there is not, for example, just one human genome or just one from any other species, but populations of human genomes and genomes populations from other species. Since every genomic region can be represented by the direct sum of abelian cyclic groups of prime-power order, then the whole genome population from individuals from the same or closed related species can be represented as an abelian group, which will be, in turns, the direct sum of abelian cyclic groups of prime-power order.



318

319 **Fig. 4.** The abelian group representation of a given genome only depend on our current knowledge
 320 on its annotation. **a**, the alternative splicing specified for an annotated gene model does not alter the
 321 abelian group representation and only would add information for the decomposition of the existing
 322 cyclic groups into subgroups. **b**, a more complex gene model including detailed information on the
 323 promoter regions. A GC box (G5MG4CU) motif is located upstream of a TATA box (TATAWAW)
 324 motif in the promoter region. The GC box is commonly the binding site for Zinc finger proteins,
 325 particularly, Sp1 transcription factors. A putative GC box was included in exon 2, which is an
 326 atypical scenario, but it can be found, e.g., in the second exon from the gene encoding for
 327 sphingosine kinase 1 (SPHK1), transcript variant 2 (NM_182965, CCDS11744.1).
 328



329

330 **Fig. 5.** Vector representation of differentially methylated exons regions from genes EGEL7 and
 331 P2RY1 from patients with pediatric acute lymphoblastic leukemia (PALL). **a**. Segment of exon-6
 332 from gene EGFL7. **b**. Segment of exon-1 from gene P2RY1. Methylated cytosines are highlighted
 333 in yellow background. In PALL patients, gene EGEL7 mostly hypomethylated and gene P2RY1
 334 mostly hypermethylated in respect to healthy individuals (WT). The encoded aminoacid sequence
 335 is given using the one letter symbols. Both genes, EGEL7 and P2RY1, were identified in the top
 336 ranked list of differentially methylated genes integrating clusters of hubs in the protein-protein
 337 interaction networks from PALL reported in reference [37].

338 Hence, results lead us to the representation of whole genomes populations of individuals
 339 from the same species or close related species (as suggested in Fig.1) by means of direct sum of
 340 their group representation into abelian cyclic groups. A general illustration of this modelling
 341 would be, for example:

$$342 \quad S = \dots \oplus (\mathbb{Z}_{5^3})^{n_1} \oplus \overbrace{(\mathbb{Z}_{2^6})^{m_1}}^{\text{motif}} \oplus (\mathbb{Z}_{5^3})^{n_2} \oplus \dots \oplus \overbrace{(\mathbb{Z}_2^2)^{m_2}}^{\text{domain}} \oplus \dots \oplus \overbrace{(\mathbb{Z}_{5^3})^{n_p}}^{\text{domain}} \oplus \overbrace{(\mathbb{Z}_{2^6})^{m_p}}^{\text{motif}} \dots (7)$$

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

344 **Theorem 1.** The genomic architecture from a genome population can be quantitatively represented
 345 as an abelian group isomorphic to a direct sum of cyclic groups of prime-power order.

346 The proof of this theorem is self-evident across the discussion and examples presented here.
 347 Basically, the group representations of the genetic code lead to the group representations of local
 348 genomic domains in terms of cyclic groups of prime-power order, for example, $(\mathfrak{B}^3, +) \cong (\mathbb{Z}_5^3, +)$
 349 or $(C_e, +) \cong (\mathbb{Z}_{5^3}, +)$, till covering the whole genome. As for any finite abelian group, the abelian
 350 group representation of genome populations can be expressed in terms a direct sum of abelian cyclic
 351 groups of prime-power order. Any new discovering on the annotation of given genome population
 352 will only split an abelian group, already defined on some genomic domain/region, into the direct
 353 sum of abelian subgroups ■.

354 4 Discussions

355 Under the assumption that the current forms of life are the result of an evolutionary process started
 356 from very simple primordial cells, the current non-coding DNA must be the relict footprint of
 357 multiple recombination of ancient DNA domains in all the permissible forms, which in ancient times
 358 were rules by an ancient genetic code. In consequence, on this scenario, the group representations
 359 of the genetic code are logically extended from relatively small local DNA domains to the whole
 360 genome.

361 Examples shown in Fig. 1 to 4 indicates whatever would be the genomic architecture for given
 362 species, the observed variations in the individual populations and in populations from closed related

363 species, it can be quantitatively described as the direct sum of abelian cyclic groups. The
364 discovering/annotation of new genomic features will only lead to the decomposition of previous
365 known abelian cyclic groups representing some genomic subregion into direct sums of subgroups.
366 In such algebraic representation DNA sequence motifs for which only substitution mutations
367 happened are specifically represented by the abelian group $(C_g, +) \cong (\mathbb{Z}_{64}, +)$, in protein coding
368 regions, and by any or combination of groups $(B, +) \cong (\mathbb{Z}_2^2, +)$, $(B^2, +) \cong (\mathbb{Z}_2^4, +)$ or
369 $B^3/G_{\text{GGA}} \cong (\mathbb{Z}_2^5, +)$ in non-protein coding regions.

370 Results indicate that the genome architecture of whole populations can be quantitatively studied
371 in the framework of abelian group theory. Two sets of MSA, S_1 and S_2 , could split into different
372 cyclic groups and, however, these sequences can be isomorphic between them because have the
373 same canonical decomposition. Particularly, the genetic code abelian group $(\mathfrak{B}^3, +) \cong (\mathbb{Z}_5^3, +)$ is
374 enough for an algebraic representation of the genome population from the same species or close
375 related species. However, such a decomposition is biologically poor and, as suggested in Figs. 4 to
376 5, masks relevant biological features from the genome architecture. A further decomposition into
377 the direct sum of abelian groups will only depends on our knowledge on the genome annotation for
378 specified species.

379 As suggested in Figs. 3 and 4, base D from the extended genetic code (represented as gaps in
380 the MSA) results useful preserving the information on the natural reading frame in the abelian group
381 representation. It is worthy to notices that, for the transcriptional and splicing enzymatic machinery,
382 the information on the reading frame preservation is already in the sequence. Molecular machines
383 perform precise logical operations [38], which in this case result in a sort of molecular *enthymeme*
384 (logical) operation where the conclusion is omitted obeying the principle of cellular economy. In
385 other words, in the algebraic representation of gene and genome populations base D carries real
386 biological information.

387 From several examples provided here, it is clear that there exists a language for the genome
388 architecture that can be represented in terms of sums of finite abelian groups. The future

389 developments of genome annotation from several species can certainly lead to the discovery of
390 logical rules of a such language determining the possible viable variations in the populations. As
391 suggested in reference [13], the identification of quotient groups (at larger scale) can permit the
392 stratification of large genome population into equivalence classes corresponding to individual
393 subpopulations, each one of them carrying particular viable variations of species genome
394 architecture. An illustration on a very simplified example is given in Fig. 3, where the extended base
395 CCD (CC-, potentially encoding for aminoacid proline, P) and codon CCT (CCU, encoding for P)
396 belong to the same equivalent class from the extended genetic code shown in Table 1. In other words,
397 the fact that DNA sequence motifs, domains, genes, chromosome and whole genomes can be
398 algebraically split into sets of equivalence classes gives birth to a new level in the current biological
399 taxonomy, which we would call *genomic algebraical taxonomy of species*.

400 As indicated in reference [11], natural genomic rearrangement like DNA recombination and
401 translocation at structural and functional domain can be represented as group automorphisms and
402 endomorphisms. Biologically, such description corresponds to the fact that the new genetic
403 information is recreated, simply, by way of reorganization of the genetic material in the
404 chromosomes of living organisms [4,39]. The analysis and discussion on the application of the
405 endomorphism ring theory to describe the dynamics of genome population is a promising subject
406 for further studies.

407 Particularly promising is the application of the genomic abelian groups on epigenomic
408 studies, which results when base D stands for the methylated cytosine. As suggested in Fig.5, a
409 precise decomposition of methylation motif into the direct sum of abelian finite group can leads to
410 their classification into unambiguous equivalence classes. This open the doors for the application
411 of based machine-learning bioinformatic approaches to study the methylation changes induced on
412 individual populations by, e.g., environmental changes, aging process and diseases, which is of
413 particular interest in genomic medicine [40].

414 Results presented here would have considerable positive impact on current molecular
415 evolutionary biology, which heavily relies on evolutionary null hypotheses about the past. As

416 suggested in reference [29], the genomic abelian groups open new horizons for the study of the
417 molecular evolutionary stochastic processes (at genomic scale) and with relevant biomedical
418 applications, founded on a deterministic ground, which only depends on the physicochemical
419 properties of DNA bases and aminoacids. In this case, the only molecular evolutionary hypothesis
420 needed about the past is a fact, the existence of the genetic code.

421 **5 Conclusions**

422 Results to date indicate that the genetic code and, ultimately, the physicochemical properties of DNA
423 bases on which the genetic code algebraic structure are defined, has a deterministic effect or at least
424 partially rules on the current genome architectures, in such a way that the abelian group
425 representations of the genetic code are logically extended to the whole genome. In consequence, the
426 fundamental theorem of abelian finite groups can be applied to the whole genome. This result opens
427 new horizons for further genomics studies with the application of the abelian group theory, which
428 currently is well developed and well documented [31,41].

429 Results suggest that the architecture of current population genomes is quite far from
430 randomness and obeys deterministic rules. Although the random nature of the mutational process,
431 only a small fraction of mutations is fixed in genomic populations. In particular, fixation events are
432 ruled by the genetic code architecture, which as shown by Sanchez (2018), it can be simulated as an
433 optimization process by using genetic algorithms [29]. This points to the study of the dynamics of
434 genome populations as a stochastic deterministic process. Genome stochasticity derives from the
435 stochasticity of mutational process and from the stochasticity of biochemical reactions, which gives
436 rise to a rich population diversity and phenotypic plasticity that help to prevent population
437 extinction. The deterministic part derives from its architecture, which can be represented in terms of
438 a canonical direct sum of homocyclic abelian groups derived from the genetic code, hold for all the
439 individuals from the same population/species.

440 **References**

- 441 1. Whalen S, Truty RM, Pollard KS. Enhancer-promoter interactions are encoded by complex
442 genomic signatures on looping chromatin. *Nat Genet.* 2016;48: 488–496.

- 443 doi:10.1038/ng.3539
- 444 2. Nuebler J, Fudenberg G, Imakaev M, Abdennur N, Mirny LA. Chromatin organization by
445 an interplay of loop extrusion and compartmental segregation. *Proc Natl Acad Sci U S A*.
446 2018;115: E6697–E6706. doi:10.1073/pnas.1717730115
- 447 3. Rowley MJ, Corces VG. Organizational principles of 3D genome architecture. *Nat Rev*
448 *Genet*. Springer US; 2018;19: 789–800. doi:10.1038/s41576-018-0060-8
- 449 4. Piazza A, Heyer WD. Homologous Recombination and the Formation of Complex Genomic
450 Rearrangements. *Trends Cell Biol*. Elsevier Ltd; 2019;29: 135–149.
451 doi:10.1016/j.tcb.2018.10.006
- 452 5. Zheng H, Xie W. The role of 3D genome organization in development and cell
453 differentiation. *Nat Rev Mol Cell Biol*. Springer US; 2019;20: 535–550.
454 doi:10.1038/s41580-019-0132-4
- 455 6. Schneider TD. Evolution of biological information. *Nucleic Acids Res*. 2000;28: 2794–9.
456 Available:
457 [http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=102656&tool=pmcentrez&rend](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=102656&tool=pmcentrez&rendertype=abstract)
458 [ertype=abstract](http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=102656&tool=pmcentrez&rendertype=abstract)
- 459 7. Yockey HP. Origin of life on earth and Shannon’s theory of communication. *Comput Chem*.
460 2000;24: 105–123. doi:10.1016/S0097-8485(99)00050-9
- 461 8. Sanchez R, Grau R. A genetic code Boolean structure. II. The genetic information system as
462 a Boolean information system. *Bull Math Biol*. 2005/07/07. 2005;67: 1017–1029.
463 doi:10.1016/j.bulm.2004.12.004
- 464 9. Sanchez R, Mackenzie SA. Information thermodynamics of cytosine DNA methylation.
465 Bardoni B, editor. *PLoS One*. Public Library of Science; 2016;11: e0150427.
466 doi:10.1371/journal.pone.0150427
- 467 10. Sánchez R, Morgado E, Grau R. A genetic code Boolean structure. I. The meaning of
468 Boolean deductions. *Bull Math Biol*. 2005;67: 1–14. doi:10.1016/j.bulm.2004.05.005
- 469 11. Sanchez R, Morgado E, Grau R. Gene algebra from a genetic code algebraic structure. *J Math*
470 *Biol*. 2005/07/14. 2005;51: 431–457. doi:10.1007/s00285-005-0332-8
- 471 12. Sanchez R, Morgado E, Grau R, Sánchez R. A genetic code Boolean structure. I. The
472 meaning of Boolean deductions. *Bull Math Biol*. 2005/02/05. 2005;67: 1–14.
473 doi:10.1016/j.bulm.2004.05.005
- 474 13. Sanchez R, Grau R, Morgado E. A novel Lie algebra of the genetic code over the Galois field
475 of four DNA bases. *Math Biosci*. 2006;202: 156–174. doi:10.1016/j.mbs.2006.03.017
- 476 14. José M V., Zamudio GS, Morgado ER. A unified model of the standard genetic code. *R Soc*
477 *Open Sci*. 2017;4: 1–13. doi:10.1098/rsos.160908
- 478 15. José M V., Morgado ER, Govezensky T. Genetic Hotels for the Standard Genetic Code:

- 479 Evolutionary Analysis Based upon Novel Three-Dimensional Algebraic Models. *Bull Math*
480 *Biol.* 2011;73: 1443–1476. doi:10.1007/s11538-010-9571-y
- 481 16. Sánchez R, Grau R. An algebraic hypothesis about the primeval genetic code architecture.
482 *Math Biosci.* 2009/07/18. 2009;221: 60–76. doi:S0025-5564(09)00114-X [pii]
483 10.1016/j.mbs.2009.07.001
- 484 17. Orgel LE. Prebiotic chemistry and the origin of the RNA world. *Crit Rev Biochem Mol Biol.*
485 2004;39: 99–123. doi:10.1080/10409230490460765
- 486 18. Piccirilli JA, Benner SA, Krauch T, Moroney SE, Benner SA. Enzymatic incorporation of a
487 new base pair into DNA and RNA extends the genetic alphabet. *Nature.* 1990;343: 33–37.
488 doi:10.1038/343033a0
- 489 19. Switzer C, Moronev SE, Benner SA. Enzymatic Incorporation of a New Base Pair into DNA
490 and RNA. *J Am Chem Soc.* 1989;111: 8322–8323. doi:10.1021/ja00203a067
- 491 20. Sanchez R, Grau R, Morgado E. A Novel DNA Sequence Vector Space over an extended
492 Genetic Code Galois Field. *MATCH Commun Math Comput Chem.* 2006;56: 5–20.
493 Available: http://match.pmf.kg.ac.rs/electronic_versions/Match56/n1/match56n1_5-20.pdf
- 494 21. Sanchez R. Evolutionary Analysis of DNA-Protein-Coding Regions Based on a Genetic
495 Code Cube Metric. *Curr Top Med Chem.* 2014;14: 407–417.
496 doi:10.2174/1568026613666131204110022
- 497 22. Di Giulio M. LUCA as well as the ancestors of archaea, bacteria and eukaryotes were
498 progenotes: Inference from the distribution and diversity of the reading mechanism of the
499 AUA and AUG codons in the domains of life. *BioSystems. Elsevier B.V.*; 2020;198: 104239.
500 doi:10.1016/j.biosystems.2020.104239
- 501 23. Di Giulio M. Errors of the ancestral translation, LUCA, and nature of its direct descendants.
502 *BioSystems. Elsevier B.V.*; 2021;206: 104433. doi:10.1016/j.biosystems.2021.104433
- 503 24. Smith ZD, Meissner A. DNA methylation: Roles in mammalian development. *Nat Rev*
504 *Genet.* Nature Publishing Group; 2013;14: 204–220. doi:10.1038/nrg3354
- 505 25. Severin PMD, Zou X, Gaub HE, Schulten K. Cytosine methylation alters DNA mechanical
506 properties. *Nucleic Acids Res.* 2011;39: 8740–51. doi:10.1093/nar/gkr578
- 507 26. Sriraman A, Debnath TK, Xhemalce B, Miller KM. Making it or breaking it: DNA
508 methylation and genome integrity. *Essays Biochem.* 2020; 687–703.
509 doi:10.1042/ebc20200009
- 510 27. Fuchs L. Abelian groups. Publishing House of the Hungarian Academy of Sciences.
511 Publishing House of the Hungarian Academy of Sciences; 1958.
- 512 28. Crick FHC. The Origin of the Genetic Code. *J Mol Biol.* 1968;38: 367–379.
- 513 29. Sanchez R. Symmetric Group of the Genetic-Code Cubes. Effect of the Genetic-Code
514 Architecture on the Evolutionary Process. *MATCH Commun Math Comput Chem.* 2018;79:

- 515 527–560. Available:
516 http://match.pmf.kg.ac.rs/electronic_versions/Match79/n3/match79n3_527-560.pdf
- 517 30. José M V, Morgado ER, Sánchez R, Govezensky T. The 24 Possible Algebraic
518 Representations of the Standard Genetic Code in Six or in Three Dimensions. *Adv Stud Biol.*
519 2012;4: 119–152. Available: <http://www.m-hikari.com/asb/asb2012/asb1-4-2012/joseASB1-4-2012-1.pdf>
- 520
- 521 31. Fuchs L. *Infinite Abelian Groups, Volume 1*. 1st Editio. Academic Press; 1970.
- 522 32. Stergachis AB, Haugen E, Shafer A, Fu W, Vernot B, Reynolds A, et al. Exonic transcription
523 factor binding directs codon choice and affects protein evolution. *Science* (80-). 2013;342:
524 1367–72. doi:10.1126/science.1243490
- 525 33. Reyna-Llorens I, Burgess SJ, Reeves G, Singh P, Stevenson SR, Williams BP, et al. Ancient
526 duons may underpin spatial patterning of gene expression in C 4 leaves. *Proc Natl Acad Sci*
527 *U S A*. 2018;115: 1931–1936. doi:10.1073/pnas.1720576115
- 528 34. Yadav VK, Smith KS, Flinders C, Mumenthaler SM, De S. Significance of duon mutations
529 in cancer genomes. *Sci Rep*. Nature Publishing Group; 2016;6: 1–9. doi:10.1038/srep27437
- 530 35. Pujato M, Kieken F, Skiles AA, Tapinos N, Fiser A. Prediction of DNA binding motifs from
531 3D models of transcription factors; identifying TLX3 regulated genes. *Nucleic Acids Res.*
532 2014;42: 13500–13512. doi:10.1093/nar/gku1228
- 533 36. Yang M, Derbyshire MK, Yamashita RA, Marchler-Bauer A. NCBI’s Conserved Domain
534 Database and Tools for Protein Domain Analysis. *Curr Protoc Bioinforma*. 2020;69: 1–25.
535 doi:10.1002/cpbi.90
- 536 37. Sanchez R, Mackenzie SA. Integrative Network Analysis of Differentially Methylated and
537 Expressed Genes for Biomarker Identification in Leukemia. *Sci Rep*. 2020;10: 2123.
538 doi:10.1038/s41598-020-58123-2
- 539 38. Schneider TD. Theory of molecular machines. II. Energy dissipation from molecular
540 machines. *J Theor Biol*. 1991;148: 125–137. Available:
541 <http://www.ncbi.nlm.nih.gov/pubmed/2016881>
- 542 39. Yu M, Ren B. The three-dimensional organization of mammalian genomes. *Annu Rev Cell*
543 *Dev Biol*. 2017;33: 265–289. doi:10.1146/annurev-cellbio-100616-060531
- 544 40. Salameh Y, Bejaoui Y, El Hajj N. DNA Methylation Biomarkers in Aging and Age-Related
545 Diseases. *Front Genet*. 2020;11: 1–11. doi:10.3389/fgene.2020.00171
- 546 41. Fuchs L. *Infinite Abelian Groups, Volume 2*. Academic Press; 1973.
- 547