# 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 DNA sequence domains suggesting a well-structured genomic architecture, where, for each species, genome populations are integrated by individual mutational variants. Herein, we show that, consistent with the fundamental theorem of Abelian finite groups, 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 of prime-power order defined on the genetic code and on the set of DNA bases, where populations can be stratified into subpopulations with the same canonical decomposition into $p$-groups. Through concrete examples we show that the architectures of current annotated genomic regions including (but not limited to) transcription factors binding-motif, promoter regulatory boxes, exon and intron arrangement associated to gene splicing are subjects for feasible modeling as decomposable Abelian $p$-groups. Moreover, we show that the epigenomic variations induced by diseases or environmental changes also can be represented as an Abelian group decomposable into homocyclic Abelian $p$-groups. The nexus between the direct sum of homocycle Abelian $p$-groups and the endomorphism ring paved the ways to unveil unsuspected stochastic-deterministic logical propositions ruling the ensemble of genomic regions. Our study aims to set the basis for concrete applications of the theory in computational biology and bioinformatics. Consistently with this goal, a computational tool designed for the analysis of fixed mutational events in gene/genome populations

33  represented as endomorphisms and automorphisms is provided. Results suggest that complex local

34  architectures and evolutionary features no evident through the direct experimentation can be unveiled

35  through the analysis of the endomorphism ring and the subsequent application of machine learning

36  approaches for the identification of stochastic-deterministic logical rules (reflecting the evolutionary

37  pressure on the region) constraining the set of possible mutational events (represented as

38  homomorphisms) and the evolutionary paths.

39

40

41  **Keywords**: Genomics, Genetic code, Abelian groups, genome algebra, automorphism, mutational

42  event

# 1  Introduction

43

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

45 genomics. Herein, with the term *genome architecture* we are adopting the definition given by Koonin

46 [1]: *Genome architecture can be defined as the totality of non-random arrangements of functional*

47 *elements (genes, regulatory regions, etc.) in the genome.*

48 Current bioinformatic tools make possible faster genome annotation process (identification of

49 locations for genes, regulatory regions, intron-exon boundaries, repeats, etc.) than some years ago

50 [2]. Current experimental genomic studies suggest that genome architectures must obey specific

51 mathematical biophysics rules [3–6]. Experimental results points to an injective relationship: *DNA*

52 *sequence → 3D chromatin architecture* [3,4,6], and failures of DNA repair mechanisms in preserving

53 the integrity of the DNA sequences lead to dysfunctional genomic rearrangements which frequently

54 are reported in several diseases [5]. Hence, **some hierarchical logic is inherent to the genetic**

55 **information system that makes it feasible for mathematical studies**. In particular, there exist

56 mathematical biology reasons to analyze the genetic information system as a communication system

57 [7–10].

58 We propose the study of genome architecture in the context of population genomics, where all

59 the variability constrained by the evolutionary pressure is expressed. Although the random nature of

60 the mutational process, only a small fraction of mutations is fixed in genomic populations. In

61 particular, fixation events, ultimately guided by random genetic drift and positive selection are

62 constrained by the genetic code, which permits a probabilistic estimation of the evolutionary

63 mutational cost by simulating the evolutionary process as an optimization process with genetic

64 algorithms [11].

## 1.1  The genetic code

65

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

67 genomic structure and robust coding apparatus, the genetic code is a fundamental link to the primeval

68    form of live, which played an essential role on the primordial architecture.  The genetic code is the

69    cornerstone of live on earth, the fundamental communication code from the genetic information

70    system [8,9]. The code-words from the genetic code are given in the alphabet of four DNA bases

71    $\mathfrak{B} = \{A, C, G, T\}$ and integrates a set of 64 DNA base-triplets $\{XYZ\}$ also named *codons*, where

72    $X, Y, Z \in \mathfrak{B}$. Each codon encodes the information for one aminoacids and each aminoacid is

73    encoded by one or more codons. Hence, at biomolecular level, the genetic code constitute a set of

74    biochemical rules (mathematically expressed as an injective mapping: *codon* → *aminoacid*) used by

75    living cells to translate information encoded within genetic material into proteins, which sets the basis

76    for our understanding of the mathematical logic inherent to the genetic information system [9,12].

77         The subjacent idea to impose a group structure on the set of codons resides on that the genetic

78    code is the code of a communication system, the genetic information system [8,13,14]. As suggested

79    by Andrews and Boss [15]: "*In codes used for electrical transmission of engineering signals,* **group**

80    **structure is imposed to increase efficiency and reduce error**. *Similarly, the group characteristics of*

81    *codon redundancy could serve to transmit additional information superimposed on the messages*

82    *directing amino acid order in protein synthesis*". As in the current human communication systems

83    [16], to impose a group structure (on biophysical basis) on the set of codons facilitate a better

84    understanding and evaluation of the error performance and efficiency of the genetic message carried

85    in the chromosomes across generations [15].

86    **1.2    The genetic code algebraic structures**

87    The basis of the current study are algebraic structures (specifically groups structures) defined on the

88    set of bases and on the codon sets.  We assume that readers are familiar with algebraic structures like

89    group, ring, and the classical mapping defined on them, homomorphisms, automorphism, and

90    translations. For readers not familiar with this subject, a brief basic introduction to these definitions

91    is given in the Appendix.

92    **The meaning of group operations.** Group operations are defined on the sets of DNA bases and

93    codons, are associated to physicochemical or/and biophysical relationships between DNA bases and
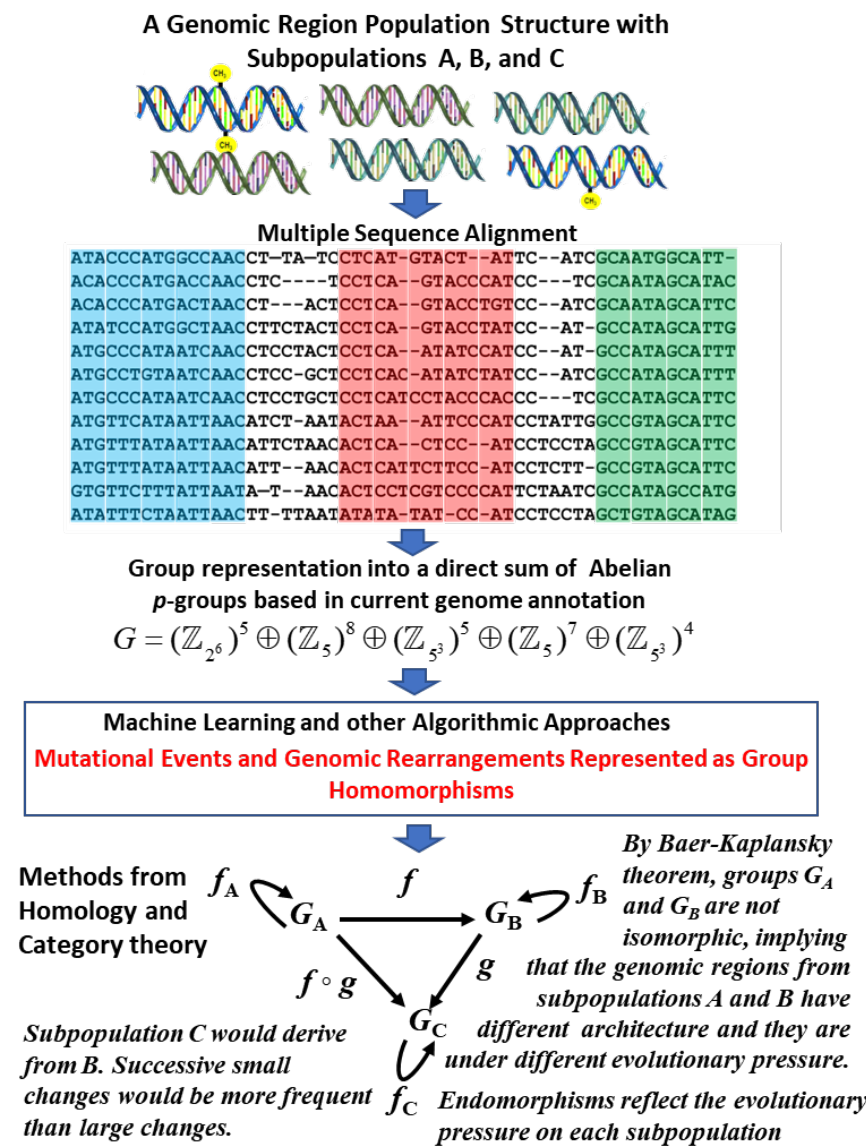
94    between codons and aminoacids. In other words, a proper definition of a group operation on the set

95    of bases or on the set of codons will encode the physicochemical or/and biophysical relationships

96    between the set's elements. Thus, by group operations defined on the set of bases or on the set of

97    codons, we understand an *encoding* applied to represent specified physicochemical or/and biophysical

98    relationships as group operations between the elements of the set. Then, we shall say that such an

99    encoding permits the *representation* of DNA bases, codons, genes, and genomic sequences as

100   elements from algebraic structures.

101        Obviously, depending on which physicochemical or biophysical relationship is under scrutiny,

102   different encodings of the group operations can be defined on the sets of bases and codons, as shown

103   in reference [17]. The meaning of the group operations has been subjects of the references where the

104   corresponding groups have been reported [11,17–20]. For example, in the DNA double helix,

105   nucleotide bases are paired following specific physicochemical relationships: 1) *the chemical type*

106   *sets the main rule for a paring*: *a purine base is paired with a pyrimidine*, 2) *paired bases must have*

107   *the same hydrogen-bonding capability*. These physicochemical relationships rule the DNA base

108   pairing: G:::C (*three hydrogen bonds*) and A::T (*two hydrogen bonds)*. In this scenario, the sum

109   operation is defined in [20], over the ordered set of bases $\mathfrak{B} = \left\{ D, \mathrm{A}, \mathrm{C}, \mathrm{G}, \mathrm{T} \right\}$, in such a way that the

110   DNA complementary bases are also complementary algebraic elements.

111   **Pioneering works on the genetic code algebraic structure**. Pioneering works were made in the 70s

112   [15,21–23], just few years after Nirenberg won the Nobel Prize in Physiology or Medicine (in 1968)

113   for his seminal work on the genetic code. Andrews and Boss proposed the cyclic groups of DNA

114   bases, which is isomorphic to the Abelian group defined on the set of integers modulo 4, $\mathbb{Z}_4 (\mathbb{Z}/4\mathbb{Z}$

115   ) [15]. Their approach also considered the base representation with cyclic group of complex numbers.

116   Further studies were focused on operational groups applied to transform bases and base-doublet into

117   each other. Dankworth and Neubert (1975) proposed the Klein-4-group structure (*K*) of doublet-

118   exchange operators and applied the direct product $K{\times}K$ to study the symmetries of genetic-code

119   doublets [22]. The four dimensional hypercube structure of the genetic-code doublets ($K{\times}K$ group)

120   was later studied by Bergman and Jungck (1979) [23].

121        Efforts with the application of *group representation theory* to study the origin and evolution of

122    the genetic code were made by Honors and Hornos [24,25], and extended to Lie superalgebras by

123    Forger and Sachse [26]. However, these efforts on the application of group representation theory are

124    heavily relying on physical interpretations disconnected from concrete molecular biology context,

125    which made hard a further application on concrete molecular biology or computational biology

126    studies, and on bioinformatic applications. Here, it is important to recall that the *representation* of

127    DNA bases, codons, genes, and genomic sequences as elements from algebraic structures must not

128    be confused with the term *group representation* typically used in algebra referring to the theory of

129    representations of algebraic structures or, particularly, the *group representation theory*. Nevertheless,

130    once a group structure has been defined, for example, in the set of codons, a further application of the

131    group representation theory can be developed.

132        In the current study, we aim to show that all possible genomic regions and, consequently, whole

133    chromosomes can be described by way of finite Abelian groups which can be split into the direct sum

134    of homocyclic 2-groups and 5-groups defined on the genetic code. Concepts and basic applications

135    are introduced step by step, sometimes with self-evident statements for a reader familiar with

136    molecular biology. However, it will be shown that the algebraic modeling is addressed to unveil more

137    complex relationships between molecular evolutionary process and the genomic architecture than

138    those eyes-visible relationships. This goal will be evidenced on section 3.2. Our algebraic model

139    approach is intended to set the theoretical basis for further studies addressed to unveil and to

140    understand the rules on how genomes are built. Concrete examples and an implementation in a R

141    package are provided to pave the way for future computational and bioinformatic applications. A

142    graphical summary of the modeling of DNA genomic regions proposed here is shown in Fig 1.

**Fig 1.** Graphical of the summary showing the bioinformatic and analytical steps followed in the algebraic modeling proposed in current work.

## 2    Materials and Methods

### 2.1    Preceding models applied in the current work

Of particular interest are the Abelian $p$-groups defined on the set of DNA bases $\mathfrak{B} = \{A, C, G, T\}$

and on the set of 64 codons $C_g = \{XYZ \,|\, X, Y, Z \in \mathfrak{B}\}$, which are applied to modeling the

physicochemical relationships between DNA bases in the codons [11,18]. Herein, for application

purposes in computational biology and bioinformatics addressed to the study of the genome

architecture, we focused our study on Abelian $p$-groups defined on $\mathfrak{B}$ and on $C_g$ isomorphic to the
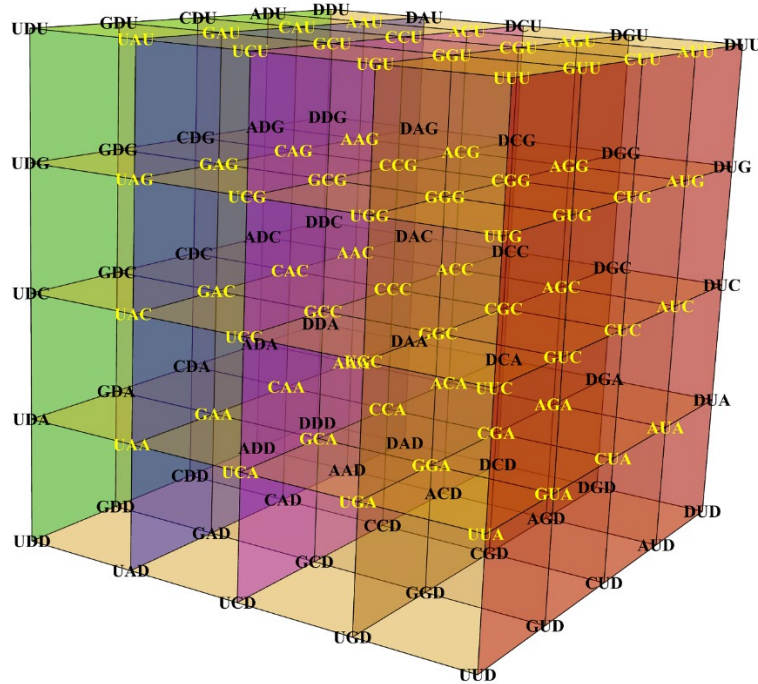
153    groups $\mathbb{Z}_{p_i^{\alpha_i}}$, $p_i^{\alpha_i} \in \{2^2, 2^6\}$, and on $\mathfrak{B}_+ = \{A, C, G, T, D\}$ and $C_{g+} = \{XYZ \mid X, Y, Z \in \mathfrak{B}_+\}$,

154    $p_i^{\alpha_i} \in \{5, 5^3\}$, as presented in references [11,17–20].

155    Setting different physicochemical restrictions on the definition of groups operations leads to

156    the 24 possible algebraic representations of the genetic code [17]. In particular, the Abelian $p$-group

157    representations on the set $C_G = \mathfrak{B} \times \mathfrak{B} \times \mathfrak{B}$ and $C_{G+} = \mathfrak{B}_+ \times \mathfrak{B}_+ \times \mathfrak{B}_+$ ($\mathfrak{B}_+ = \{A, C, G, T, D\}$,

158    where $D$ stands for an alternative base, see below) are isomorphic to Abelian groups defined on $\mathbb{Z}_{2^2}^3$

159    and $\mathbb{Z}_5^3$, respectively. These group structures lead to 24 (isomorphic) geometrical representations of

160    the genetic code as cubes inserted in three-dimensional space [11,17,19,20] (Fig 2 and SI Figs 1 and

161    3).

162    As shown in reference [11], a group structure isomorphic to the symmetric group of degree

163    four $S_4$ (preserving the group operations previously defined on the codon set) can be defined in set

164    the 24 genetic-code algebraic representations or in the set 24 cubes. Since the definition of a sum

165    operation over the base set is equivalent to define an order on it, *cubes are named according to the*

166    *base order on them.* For example, the cube shown in Fig 2 is denoted as ACGT, which correspond to

167    the group operation defined on the ordered set $\mathfrak{B} = \{A, C, G, T, D\}$ (the 'dual' cube TGCA is shown

168    in SI Fig 2 [11]). Simulation of the evolutionary mutational process with the application of genetic

169    algorithms indicates that fixed mutational events found in different protein populations are *very*

170    *restrictive* in the sense that the optimal evolutionary codon distances are reached for specific models

171    of genetic-code cube or for specific combination of genetic-code cube models [11]. In the present

172    work, it will be shown that codon mutational events represented in terms of automorphisms can be

173    also restrictive for specific genetic-code cube models (section 3.1).

174    All the Abelian $p$-group included in the current work are oriented to the study of the mutational

175    process [11,17–20]. That is, since we are interested in those structures that permit the analysis and

176    quantitative description of the mutational process in organismal populations, where mutational event

177    can be represented by means of endomorphisms, automorphisms, and translations on the defined

178    group, we do not include algebraic structures designed to study the origin and evolution of the genetic

179    code [11,18]. The genetic code is taken as currently is, without over-impose any evolutionary

180    hypothesis on it.



181

**Fig 2**. Geometrical representation of the genetic code as a cube inserted in three-dimensional space. The 2-group and 5-group representation defined on the sets $C_G = \mathfrak{B} \times \mathfrak{B} \times \mathfrak{B}$ and $C_{G+} = \mathfrak{B}_+ \times \mathfrak{B}_+ \times \mathfrak{B}_+$ isomorphic to the groups defined on $\mathbb{Z}_{2^2}^3$ and $\mathbb{Z}_5^3$, respectively, lead to the geometrical representations of the genetic code as a cube inserted in three-dimensional space. The cube corresponding to the base-triplets with coordinates on $\mathbb{Z}_{2^2}^3$ (yellow codons) is inserted in the cube with codon coordinates on $\mathbb{Z}_5^3$. The extended base-triplets including the alternative base D (in black) are located on the cartesian coordinate planes. Codons encoding for amino acids with similar physicochemical properties are located on the same vertical plane (for more details on the cube description see also SI Fig 1 and reference [11,17,19,20]).

192    A general model also consider Abelian 5-groups that includes a dummy variable (denoted by

193    letter D), which extends the DNA alphabet to five letters. The usefulness of including a fifth base in

194    the evolutionary analysis was shown in reference [20], where two evolutionary models, an algebraic

195    and a stationary Markov (process) models, were applied to phylogenetic analysis reaching (both

196    models) greater discriminatory power than the (now) classical Tamura-Neil evolutionary (Markov)

197    model based on four DNA alphabet [27]. Depending on the concrete application, letter "D" will take

198    a different value. The possible values in the context of the present modeling are: 1) the gap symbol

199  "-", which stands for insertion deletion/mutations in the multiple sequence alignment (MSA) of DNA

200  sequences, 2) alternative wobble base pairing (e.g., bases such as: inosine (in eukaryotes), agmatine

201  (in archaea), and lysidine (in bacteria) [17,21,22]), and 3) 5-methylcytosine ($C^m$) and N-6-

202  methyladenine ($A^m$) when intended for epigenetic studies.

203  A concrete application of the extended genetic-code cubes over the Galois field *GF*(5) to the

204  simulation of the mutational process proposed in reference [11] would be particularly relevant to

205  predict immunoescape epitope variants originated in populations of pathogenic microorganisms and

206  viruses. In addition, examples provided (here) on the application of the algebraic model to DNA

207  methylation (on 5-methylcytosine and on N-6-methyladenine) suggest its importance for epigenetic

208  studies. The analysis of the fixed mutational events on genes populations revealed that the mutational

209  process can be described by automorphisms on different cubes or sets of cubes [11]. The best genetic-

210  code cubes describing the mutational process on a given gene population are selected with the

211  application of an optimization algorithm (evolutionary (genetic) algorithms) using multiple sequence

212  alignment as raw data [11].

213  It is worthy to notice that, for all mentioned Abelian *p*-groups, the calculus can be

214  accomplished as symbolic computation on the set of DNA bases or on the set of codons (see e.g.,

215  [18]). However, for practical purposes, we take advantage of the group isomorphisms. That is, after

216  define group structures on the sets of bases and codons, for the sake of straightforward computation

217  it is convenient to take advantage of the group isomorphisms with the Abelian *p*-groups like: $\mathbb{Z}_{2^2}$,

218  $\mathbb{Z}_{2^2}^3$, $\mathbb{Z}_{2^6}$, $\mathbb{Z}_5$, $\mathbb{Z}_{5^3}$ and $\mathbb{Z}_5^3$, which will be used in our study instead of the original groups defined

219  on the sets of bases and codons (base-triplets). An introductory summary on the mentioned algebraic

220  structure defined on the set of codons is provided as supporting information in S1.

221  In the context of genetic-code algebraic structures, by the term "*representation*" of DNA bases,

222  codons, genes, and genomic sequences as elements from algebraic structures, we understand the

223  symbolic representation of the mentioned biomolecules and the physicochemical relationships

224  between them by means of group operations defined on the given set of biomolecules.

## 2.2 Aligned DNA sequences and data sets

All the DNA sequence alignments and data sets used in this work are available within the R package *GenomAutomorphism* (version 1.0.0) [28]. In addition, the pairwise sequence alignments of SARS coronaviruses used the analyses shown in Fig 8**a** and **b** are also available at GitHub in: https://github.com/genomaths/seqalignments/tree/master/COVID-19. The multiple sequence alignment (MSA) of primate somatic cytochrome c and data description are available on GitHub at: https://github.com/genomaths/seqalignments/tree/master/CYCS. This MSA includes DNA protein-coding sequences from: human, gorilla, silvery gibbon, white cheeked gibbon, Francois langur, olive baboon, golden monkey, rhesus monkeys, gelada baboon, and orangutan. The MSA of primate BRCA1 (transcript variant 4) DNA repair gene used to compute the automorphism shown Fig 8**d** is available on GitHub at https://github.com/genomaths/seqalignments/tree/master/BRCA1. The MSA, coordinates and R script to create the sequence-logo from Fig 4 are given in the Supporting Information.

## 2.3 Software applied for the mathematical and statistical analyses

Results shown in Fig 8 and Fig 9 were obtained applying the *GenomAutomorphism* R package [28] (version 1.0.0), which is available at Bioconductor (the open source software for Bioinformatics, version: 3.16) and, also, in GitHub at: https://github.com/genomaths/GenomAutomorphism. The whole R script pipeline applied in the estimation of automorphisms (Fig 8) and decision tree (Fig 9) are available as tutorials (vignettes) at the *Geno Automorphism* website: https://github.com/genomaths/GenomAutomorphismm.

The estimation of the best fitted probability distribution shown in Fig 8**f** was accomplished with R package *usefr* available at GitHub: https://github.com/genomaths/usefr, and the goodness-of-fit tests are reported in the mentioned tutorials.

The genetic-code cube shown in Fig 2 was obtained from the Wolfram Mathematica Notebook: *Introduction to $\mathbb{Z}_5$-Genetic-Code vector space*, free available at https://github.com/genomaths/GenomeAlgebra_SymmetricGroup.

251 ## 2.4 Theoretical Model

252 According to the fundamental theorem of Abelian finite groups (FTAG) [29,30], any finite Abelian

253 group can be decomposed into a direct sum of homocyclic $p$-groups [29], i.e., a group in which the

254 order of every element is a power of a primer number $p$. Herein, it will be showed that, in a general

255 scenario, genomic regions and, consequently, whole genome populations from any species or close

256 related species, can be algebraically represented as a direct sum of Abelian homocyclic groups or

257 more specifically Abelian $p$-groups of *prime-power order*. The multiple sequence alignments (MSA)

258 of a given genomic region of $N$ base-pair (bp) length can be represented as the direct sum:

$$G = \left( \mathbb{Z}_{p_1^{\alpha_1}} \right)^{n_1} \oplus \left( \mathbb{Z}_{p_2^{\alpha_2}} \right)^{n_2} \oplus \cdots \oplus \left( \mathbb{Z}_{p_k^{\alpha_k}} \right)^{n_k} \tag{1}$$

259

260 Where $p_i^{\alpha_i} \in \{2, 5, 2^6, 5^3\}$, $n_i$ stands for the number of cyclic groups $\mathbb{Z}_{p_i^{\alpha_i}}$ integrating the homocyclic

261 group $\left( \mathbb{Z}_{p_i^{\alpha_i}} \right)^{n_i} = \overset{n_i \ times}{\mathbb{Z}_{p_i^{\alpha_i}} \oplus \ldots \oplus \mathbb{Z}_{p_i^{\alpha_i}}}$. Here, we assume the usual definition of direct sum of groups

262 [30]. For $p_j^{\alpha_j} \in \{2^6, 5^3\}$ the cyclic group $\mathbb{Z}_{p_j^{\alpha_j}}$ will cover three bases, otherwise only one base (see

263 examples below). Considering such groups (not necessarily in the order given in Eq. 1) we have:

264 $N = n_1 + \ldots + n_j + n_{j+1} + \ldots + n_{j+m} + \ldots + n_k$. Throughout the exposition of the theory and

265 examples given in the next sections, it will be obvious that the group representations can be extended,

266 starting from small genomic regions till cover whole chromosomes and, consequently, the whole

267 genome, i.e., the set of all chromosomes.

268 Let $B_i$ $(i \in I = \{1, \ldots, n\})$ be a family of subgroups of $G$, subject to the following two
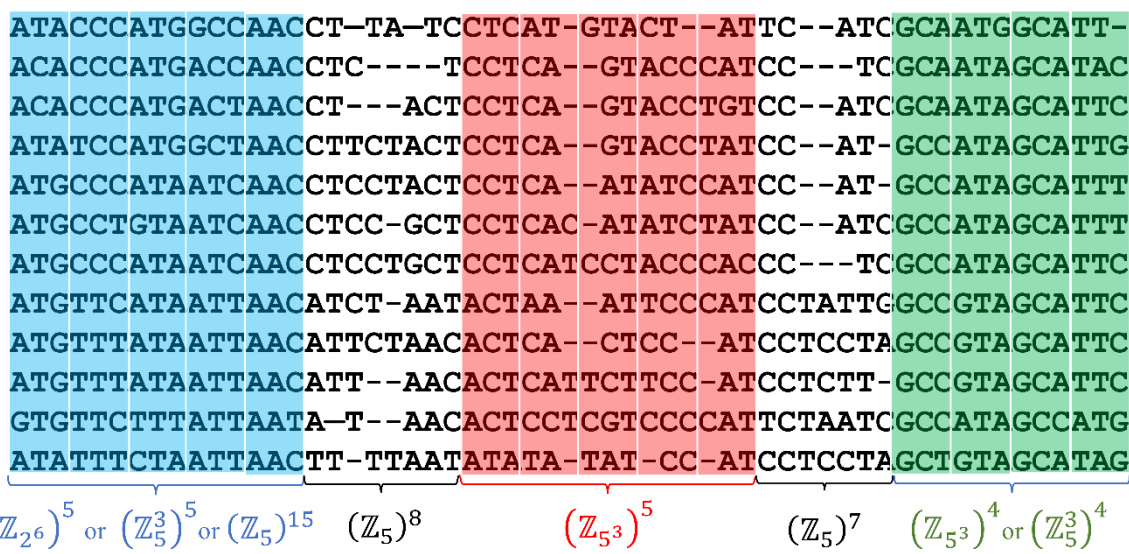
269 conditions:

270 1) $\sum B_i = G$. That is, $B_i$ together generate $G$.

271 2) For every $i \in I$ and $i \neq j$: $B_i \cap \sum B_j = 0$.

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

273     expression: $G = \bigoplus_i B_i$ or $G = B_1 \oplus \ldots \oplus B_n$.

274     Genomic DNA sequences from superior organisms are integrated by intergenic regions and

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

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

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

278     of substitutions, insertions, and deletion (*indel*) mutations. As a result, the alignment of homolog

279     genomic regions or whole chromosome DNA sequences from several individuals from the same or

280     close-related species can be split into well-defined subregions or domains, and each one of them can

281     be represented as homocyclic Abelian groups, i.e., as the direct sum of cyclic group of the same

282     *prime-power* order (Fig 3). As a result, each DNA sequence is represented as a *N*-dimensional vector

283     with numerical coordinates representing bases and codons.



284

**Fig 3**. An illustration of a typical DNA multiple sequence alignment (MSA) including segments of
protein-coding regions. A MSA would include the presence of substitution, insertion, and deletion
mutations (*indel* mutations). The aligned sequences can be grouped into blocks, which can be
algebraically represented by Abelian groups. A homocyclic group covering a MSA block corresponds
to a sub-classification of the protein-coding region into subregions and, consequently, leading to a
more accurate molecular taxonomy of species. In protein-coding regions cyclic groups $\mathbb{Z}_{2^6}$ and $\mathbb{Z}_5^3$
are appropriated to study exon regions, while $\mathbb{Z}_5$ for non-coding intron regions. As shown in section
1.4, the group representation leads us the analysis of the more frequent mutational events (represented
as endomorphisms and translations) observable in genes from organismal populations.

294

295   An intuitive mathematical representation of a MSA is implicit in Fig 3, with the following

296   observations:

a) Bases or codons can be represented as elements of an Abelian group defined on the set of bases or on the set of codons. In the second block (including gap symbol '–') each base from each sequence is represented as an element from the Abelian group defined on the set {A, C, G, T, $D$ } where $D = $ '–', which is isomorphic to the Abelian $p$-group defined on the set $\mathbb{Z}_5$. The extended base triplets (including gaps symbol '–') from each sequence in the third aligned block are represented as elements from the Abelian $p$-group defined on the set of extended base-triplets (125 element, see SI Table 1) which is isomorphic to the Abelian group defined on the set $\mathbb{Z}_{5^3}$, and so on.

b) Every DNA sequence from the MSA and every subsequence on it can be represented as a *numerical vector* with element coordinates defined in an Abelian group. For practical computational purposes we take advantage of the group isomorphism to work with numerical representations of DNA bases and codons. For example, codons from the first aligned block (in blue) can be represented as elements from an Abelian group defined on the set of codons, which can be isomorphic to $\mathbb{Z}_{2^6}$ or to $\mathbb{Z}_5^3$. That is, since $\left(C_g, +\right) \cong \left(\mathbb{Z}_{2^6}, +\right)$, the first five codons $\{ATA, CCC, ATG, GCC, AAC\} \in C_g$ from the first DNA sequence from Fig 3, can be represented by the vector of integers: $\{48, 21, 50, 25, 1\}$ where each coordinate is an element from group $\left(\mathbb{Z}_{2^6}, +\right)$ (see Table 1 from reference [18]).

c) Any MSA can be algebraically represented as a symbolic composition of Abelian groups each one of them is isomorphic to an Abelian group of integers module $n$. Such a composition can be algebraically represented as a direct sum of homocyclic Abelian $p$-groups. For example, the MSA from Fig 3 can be represented by the direct sum of five homocyclic Abelian $p$-groups:

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

320      Where the length of each region determines the number of cyclic $p$-groups in the

321      corresponding homocyclic Abelian $p$-group $\mathbb{Z}_{p_1^{\alpha_1}}$ representing each region. For example, in

322      Eq. 2 we have the homocyclic group: $\left(\mathbb{Z}_{5^3}\right)^4 = \oplus_{i=1}^4 \mathbb{Z}_{5^3}$ , which is a direct sum of 4 cyclic

323      5-groups $\left(\mathbb{Z}_{5^3}, +\right) \cong \left(C_{g+}, +\right)$. Since group $G$ is the direct sum of homocyclic Abelian $p$-

324      groups of different prime-order, we shall say that $G$ is a heterocyclic group.

325     In more specific scenario, the MSA from Fig 3 can be represented by only one homocyclic

326 Abelian 5-group:

$$G = \left(\mathbb{Z}_5\right)^{57} \tag{3}$$

327

328 But this representation ignores the local variability detected by the MSA algorithm. Hence, preserving

329 the highlighted features, the MSA can be represented as the direct sum of homocyclic Abelian 5-

330 groups:

$$G = \left(\mathbb{Z}_5^3\right)^5 \oplus \left(\mathbb{Z}_5\right)^8 \oplus \left(\mathbb{Z}_{5^3}\right)^5 \oplus \left(\mathbb{Z}_5\right)^7 \oplus \left(\mathbb{Z}_5^3\right)^4 \tag{4}$$

331

332      Although the above *direct sums* of Abelian $p$-groups provides a useful compact representation

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

334 product (*cartesian sum or complete direct sums*) [30]. Next, let $S$ be a set of Abelian cyclic groups

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

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

337 sequence on the $M$ can be represented by a cartesian product $\left(b_1, \ldots, b_n\right)$ where $b_i \in S_i$ $\left(i = 1, \ldots, n\right)$

338 and $n = |S|$. Let $G_i$ be a group defined on the set of all elements $\left(0, \ldots, 0, b_i, 0, \ldots 0\right)$ where $b_i \in S_i$

339 stands on the $i^{th}$ place and 0 everywhere else. It is clear that $S_i \cong G_i$. In this context, the set of all

340 vectors $\left(b_1, \ldots, b_n\right)$ with equality and addition of vectors defined coordinate-wise becomes a group (

341 $\mathcal{G}$ ) named direct product (cartesian sum) of groups $S_i$ ( $G_i$ ), i.e.:

342
$$G = \otimes_i S_i = \oplus_i G_i \qquad (5)$$

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

## 3   Results

345 Results essentially comprise an application of the fundamental theorem of Abelian finite groups

346 [29,30]. By this theorem, every finite Abelian group $G$ is isomorphic to a direct sum of cyclic groups

347 of prime-power order of the form:

348
$$G = \mathbb{Z}_{p_1^{\alpha_1}} \oplus \mathbb{Z}_{p_2^{\alpha_2}} \oplus \cdots \oplus \mathbb{Z}_{p_n^{\alpha_n}} \qquad (6)$$
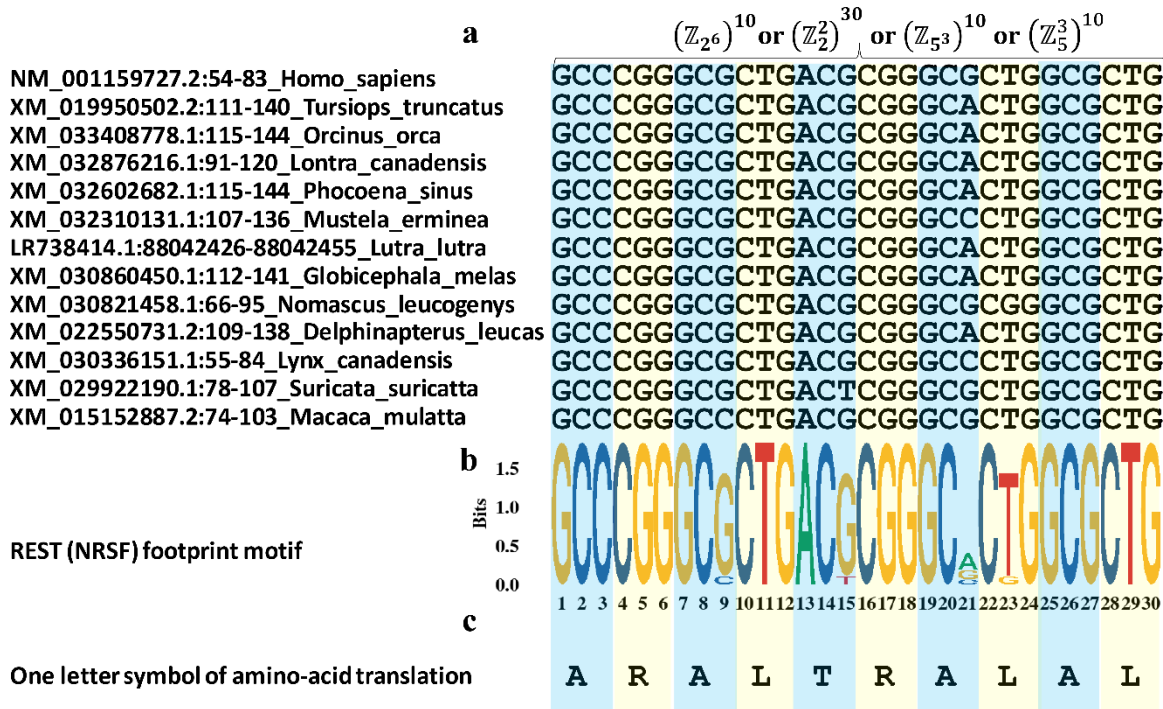
349 Or (in short) $G = \oplus_{i=1}^n \mathbb{Z}_{p_i^{\alpha_i}}$, where the $p_i$'s are primes (not necessarily distinct), $\alpha_i \in \mathbb{N}$ and $\mathbb{Z}_{p_i^{\alpha_i}}$

350 is the group of integer module $p_i^{\alpha_i}$. The Abelian group representation of the MSA from Fig 3 given

351 by Eq. 2 correspond to a heterocyclic group that split into a direct sum of homocyclic Abelian 2-

352 groups and 5-groups, each one of them split into the direct sum of cyclic $p$-groups with same order;

353 while in Eqs. 3 and 4, the Abelian group $G$ is decomposed into a direct sum of homocyclic Abelian

354 5-groups [29,30].

355        Notice that for a large enough genomic region of fixed length $N$ we can build a *manifold of (a*

356 *set of various) heterocyclic groups $S_i$*, where each one of them can have different decomposition into

357 $p$-groups. The set $S$ of all possible Abelian $p$-group representations $S_i$ of a large genomic region of

358 fixed length (having numerous different parts, elements, features, forms, etc.) that split into the direct

359 sum of several heterocyclic groups $G_k$ ( $S_i = \oplus_{k=1}^n G_k$ ) shall be called a *heterocyclic-group manifold*.

360 So, each genomic region can be characterized by means of their corresponding *heterocyclic-group*

361 *manifold*.

### 3.1   Examples of genomic regions group representations

363 A group representation is particularly interesting for the analysis of DNA sequence motifs, which

364 typically are highly conserved across the species. As suggested in Fig 3 and 4, there are subregions

365 of DNA or protein sequences where there are few or not gaps introduced and mostly substitution

366 mutations are found. Such subregions conform blocks that can cover complete DNA sequence motifs

367 targeted by DNA biding proteins like transcription factors (TFs, Fig 4), which are identifiable

368 applying bioinformatic algorithms like BLAST [31].

369
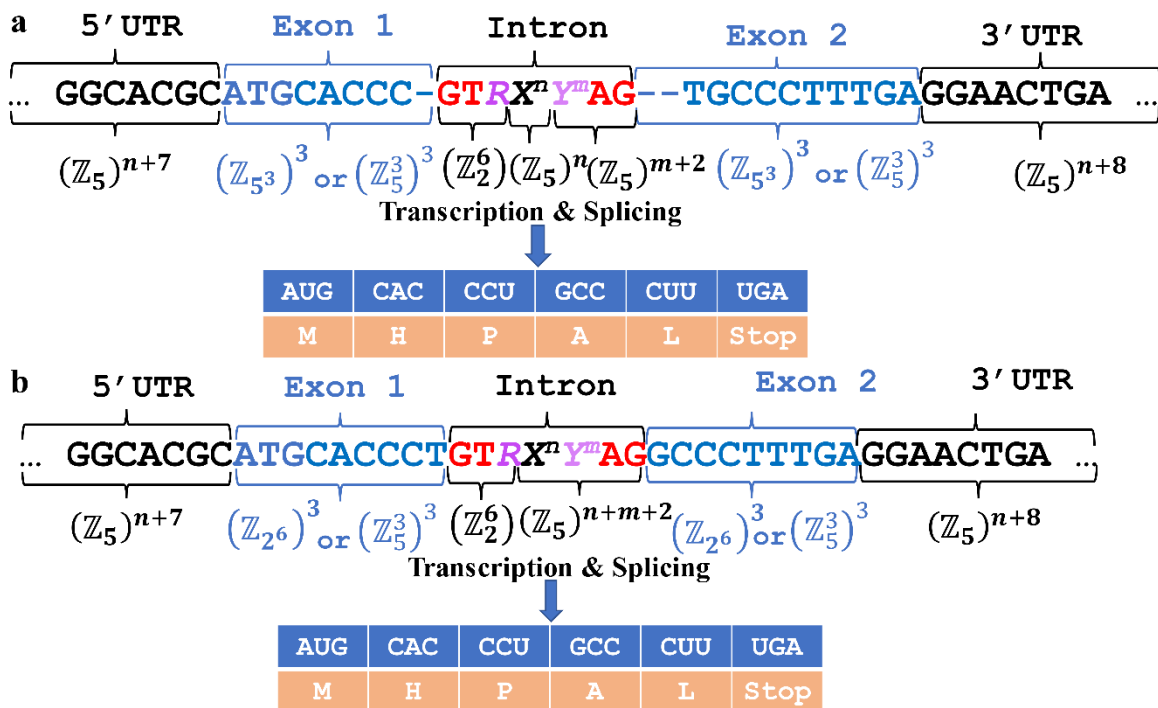


**Fig 4**. The DNA sequence motifs targeted by transcription factors usually integrate genomic building block across several mammal species. **a**, DNA sequence alignment of the protein-coding sequences from phospholipase B domain containing-2 (PLBD2) carrying the footprint sequence motif recognized (targeted) by the Silencing Transcription factor (REST), also known as Neuron-Restrictive Silencer Factor (NRSF) REST (NRSF). **b**, Sequence logo of the footprint motif recognized REST (NRSF) on the exons. **c**, Translation of the codon sequences using the one-letter symbol of the aminoacids.

379 The case of group representation on a TF binding motif is exemplified in Fig 4, where an exon

380 region from the enzyme *phospholipase B domain containing-2* (PLBD2) simultaneously encodes

381 information for several aminoacids and carries the footprint to be targeted by the transcription factor

382 REST. Herein, the case of double encoding called our attention, where the DNA sequence

383 simultaneously encodes the information for transcription enhancer target motif and for a codon

384 sequence (base-triplets) encoding for aminoacids. These types of double-coding regions are also

385 called *duons* [32–34].

386   Four group representations for this exon subregion are suggested in the top of the Fig 4 (panel

387   **a**). However, the MSA's sequence logo (panel **b**) suggests that this transcription factor binding-motif

388   is a highly conserved codon sequence in mammals (with no indel mutations on it) and, in this case,

389   the Abelian group $\left(C_g, +\right) \cong \left(\mathbb{Z}_{2^6}, +\right)$ defined on the standard genetic code is the appropriated model

390   to represent these motifs (Fig 4). The homocyclic group representation of conserved and biological

391   relevant DNA sequence motifs, illustrated in Figs. 3 and 4, stablish the basis for the study of the

392   molecular evolutionary process in the framework of group endomorphisms and automorphisms as

393   suggested in [18,20] (section 1.4).

394   In Fig 5, two different protein-coding (gene) models from two different genome populations

395   can lead to the same direct sum of Abelian $p$-groups and to the same final aminoacids sequence

396   (protein).



397

**Fig 5**. 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 $(Y^m AG)$ 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

405    exons, 1 and 2, carries a complete set of three codons (base-triplets). Both gene models, from panels
406    **a** and **b**, share a common group representation as direct sum of Abelian 5-groups.
407

408      The respective exon regions have different lengths and gaps ("-", representing base D in the

409    extended genetic code) were added to exons 1 and 2 (from panel **a**) to preserve the reading frame in

410    the group representation (after transcription and splicing gaps are removed). Both gene models, from

411    panel **a** and **b**, share a common direct sum of Abelian 2-groups and 5-groups:

412    $\left(\mathbb{Z}_5\right)^{n+7} \oplus \left(\mathbb{Z}_5^3\right)^3 \oplus \left(\mathbb{Z}_2^6\right) \oplus \left(\mathbb{Z}_5\right)^{n+m+2} \oplus \left(\mathbb{Z}_5^3\right)^3 \oplus \left(\mathbb{Z}_5\right)^{n+8}$ . The analysis of theses gene

413    models suggests that *DNA sequences sharing a common group representation as direct sum of*

414    *Abelian p-groups would carry the same or similar, or close related biological information*. However,

415    it does not imply that the architecture of these protein-coding regions is the same. The gene model in

416    panel      **b**      permits      the      direct      sum      representation:
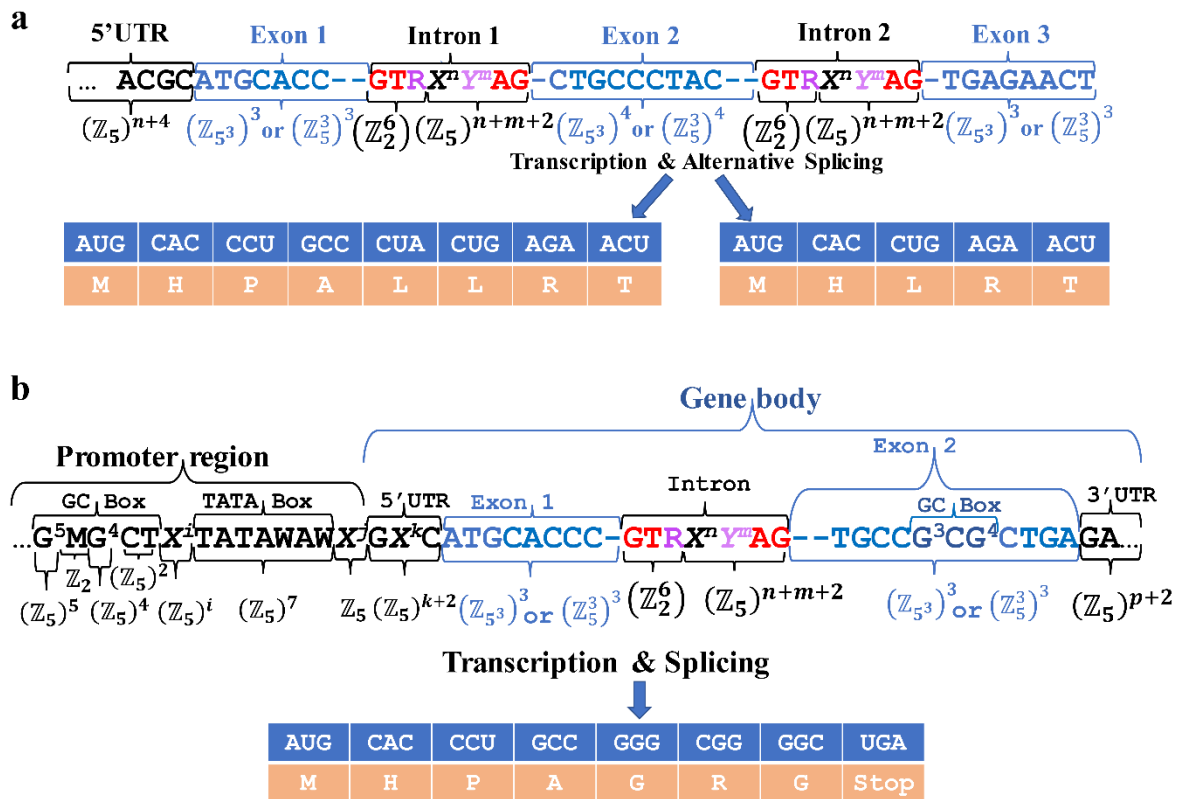
417    $\left(\mathbb{Z}_5\right)^{n+7} \oplus \left(\mathbb{Z}_{2^6}\right)^3 \oplus \left(\mathbb{Z}_2^6\right) \oplus \left(\mathbb{Z}_5\right)^{n+m+2} \oplus \left(\mathbb{Z}_{2^6}\right)^3 \oplus \left(\mathbb{Z}_5\right)^{n+8}$ , which is no possible for the

418    gene model from panel **a**. That is, the *heterocyclic-group manifold* from the gene model in panel **a** is

419    different from the one in panel **b**. The difference of group representation just captures the obvious

420    fact that these gene models are different and, consequently, their gene architectures are different.

421      At this point we shall introduce the concept of *equivalent class of genomic region*. We shall

422    say that two genomic regions belong to same *equivalent class of genomic region* if they hold the same

423    heterocyclic-group manifold (and, consequently, they hold same architecture). Under this definition,

424    the region architecture of the protein-coding regions from Fig 5**a** and **b** are not equivalent. The

425    concept of *equivalent class of genomic region* is relevant for further applications of the group

426    representation on the taxonomy study of organismal populations.

427      Taxonomy is the study of the scientific classification of biological organisms into groups based

428    on shared characteristics. Mathematically, this is a way to split biological organisms into classes of

429    equivalences. Numerical taxonomy is a well-established application of multivariate statistics on the

430    analysis of plant germplasm banks. The group representations of genomic regions will lead to a higher

431    accuracy in the taxonomy study of organismal populations.

432     No matter how complex a genomic region might be, it has an Abelian group representation. A

433     further application of group theory would unveil more specific decomposition of small genomic

434     regions into Abelian groups. For example, the set of base-triplets found in a typical sequence motif

435     targeted by the spliceosome donor, GT$R$ (Figs. 5 and 6), is in the vertical line GTZ (GUZ) of the

436     vertical plane $XTZ$ ($XUZ$) from the cube ACGU shown in Fig 2 (see also SI Fig 3).



437

**Fig 6**. The Abelian group representation of a given genome only depend on our current knowledge on its annotation. **a,** the alternative splicing specified for an annotated gene model does not alter the Abelian group representation and only would add information for the decomposition of the existing cyclic groups into subgroups. **b,** a more complex gene model including detailed information on the promoter regions. A GC box (G5MG4CU) motif is located upstream of a TATA box (TATAWAW) motif in the promoter region. The GC box is commonly the binding site for Zinc finger proteins, particularly, Sp1 transcription factors. A putative GC box was included in exon 2, which is an atypical scenario, but it can be found, e.g., in the second exon from the gene encoding for sphingosine kinase 1 (SPHK1), transcript variant 2 (NM_182965, CCDS11744.1). In this group representation, the spliceosome donor GT$R$ can be represented by the elements from a quotient group (see main text).

448

449     Since purine bases (R: A and G) are the only accepted variants at the third codon position, it is

450     convenient to model these base-triples with the group defined on the cube AGCU [11] (SI Fig 3).

451     Next, following analogous reasoning as in [19], it turns out that the set of base-triplets GT$R$ is a coset

452    from the quotient group $\left(C_G,+\right)/G_{\mathrm{AAG}}$, where here $\left(C_G,+\right)\cong\left(\mathbb{Z}_2^6,+\right)$ is the additive group from

453    the genetic code Galois field $GF(64)$ reported in reference [35] and $G_{\mathrm{AAG}}=\left(\{\mathrm{AAA},\mathrm{AAG}\},+\right)$ is a

454    subgroup from the Klein four group defined on the set $\{\mathrm{AAA},\mathrm{AAG},\mathrm{AAC},\mathrm{AAU}\}$ (see operation

455    table in the SI Table 2), i.e., $\mathrm{GT}R=\mathrm{GTA}+G_{\mathrm{AAG}}$ (SI Fig 3).
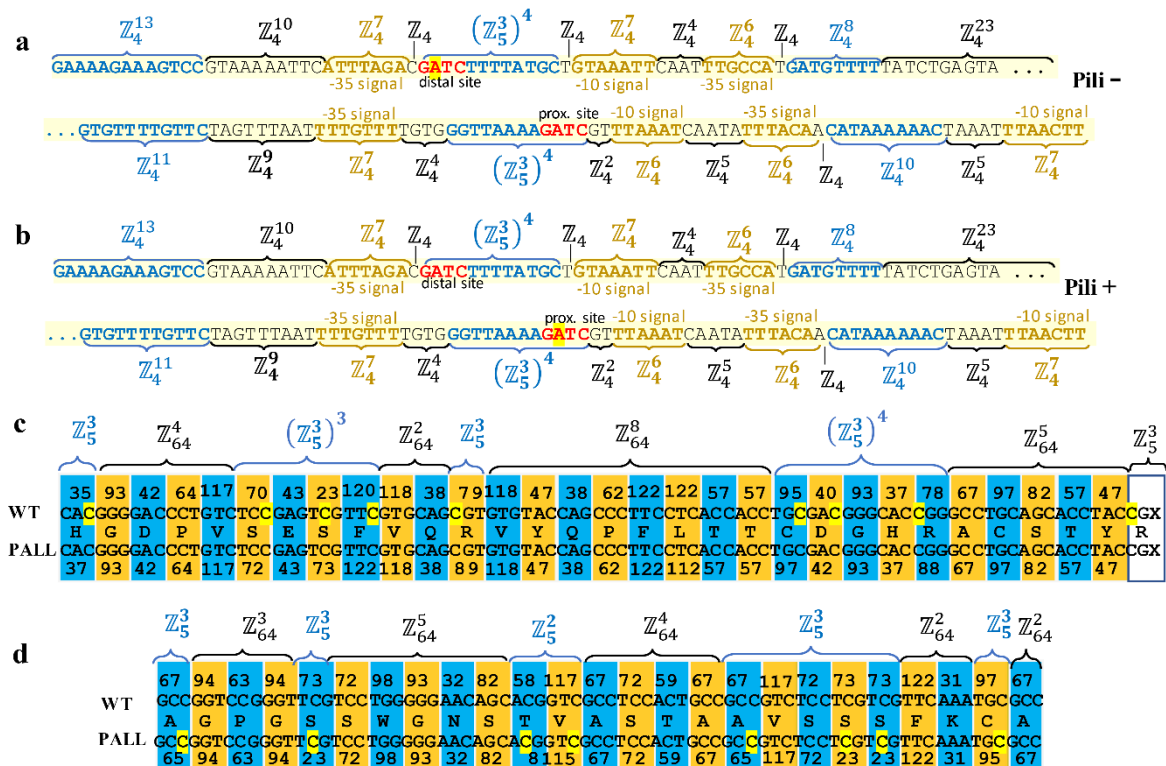
456    There exists strong evolutionary pressure on splicing donor site to keep the base-triplet GTR

457    in the vertical line GTZ (GUZ) vertical line (coset). As shown in the clinical report [36] mutational

458    variants, located in different cube's vertical lines (different cosets, SI Fig 3) GCZ and CTZ (CUZ),

459    within intron 3 have led to four aberrant RNAs transcripts that causes rare X-chromosome-linked

460    congenital deafness. As will be shown below (in section 3.1) the strong connection between DNA

461    sequences and non-disrupting mutational events is mathematically (and accurately) modeled by the

462    strong relationship between a group representation and the endomorphism ring on it.

463    An example considering changes on the gene-body reading frames as those observed in

464    alternative splicing is shown in Fig 6. Gene-bodies with annotated alternative splicing can easily be

465    represented by any of the groups $\left(\mathbb{Z}_5^3\right)^n$ or $\left(\mathbb{Z}_{5^3}\right)^n$ (Fig 6**a**). The splicing can include enhancer

466    regions as well (Fig 6**b**) [37]. Enhancers are key regulator of differential gene expression programs.

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

468    extended base-triple group representation. Typically, the analysis of methylome data is addressed to

469    identify methylation changes induced by, for example, environmental changes, lifestyles, age, or

470    diseases. So, in this case the letter $D$ stands for methylated adenine and cytosine ($D=C^m$), since

471    only epigenetic changes are evaluated.

472    Concrete examples of adenine in bacteria linked to the regulation of pyelonephritis-associated

473    pilus (pap) expression by DNA methylation on the *Escherichia coli* operon (locus X14471) and

474    cytosine methylation in two (humans) genes from patients with pediatric acute lymphoblastic

475    leukemia (PALL) are presented in Fig 7. On protein-coding regions methylation change can be

476    analyzed on the homocyclic groups composed by the cyclic group $\mathbb{Z}_5^3$ or $\mathbb{Z}_{5^3}$ (Fig 7c and **d**). Notice

477    that adenine methylation is found in humans as well and, usually, it plays a very specific regulatory

478    role [38,39].



**Fig 7**. Vector representation of differentially methylated gene regions. **a** and **b**, regulation of pyelonephritis-associated pilus (pap) expression by DNA methylation on the Escherichia coli operon (locus X14471). **c** and **d**, exons regions from genes EGEL7 and P2RY1 from patients with pediatric acute lymphoblastic leukemia (PALL). In panel **a**, two 5'-GATC-3' DNA adenine methyltransferase (Dam) methylation sites in the middle of each set of the leucine-responsive regulatory protein (Lrp) binding sites (in blue). In the inactive state, panel **b**, a Lrp octamer is bound to the three proximal Lrp 3' sites, while the GATC$^{dist}$ site in Lrp site 5 is fully methylated, and the system remains in phase OFF (Pili -) with regard to pilus expression. In the active state, the adenine from the GATC$^{prox}$ is methylated permitting to bend the DNA to recruit CRP to activate transcription of papBA genes (Pili +). Pap pili are multisubunit fibers essential for the attachment of uropathogenic Escherichia coli to the kidney (see [40]). In panel **c**, a segment of exon-6 from gene EGFL7 located at chromosome 9: 139,563,008-139,563,124 is shown. On average, this gene is hypo-methylated in the control group with respect to PALL group. **d**. Segment of exon-1 from gene P2RY1. Methylated cytosines are highlighted in yellow background. In PALL patients, gene EGEL7 mostly hypomethylated and gene P2RY1 mostly hypermethylated in respect to healthy individuals (WT). The encoded aminoacid sequence is given using the one letter symbols. Both genes, EGEL7 and P2RY1, were identified in the top ranked list of differentially methylated genes integrating clusters of hubs in the protein-protein interaction networks from PALL reported in reference [41]. The integer number at the top and bottom of panel **c** and **d** stand for the codon coordinates in $\mathbb{Z}_{5^3}$ (see SI Table 1).

500        It is obvious that the MSA from a whole genome derives from the MSA of every genomic

501    region, from the same or closed related species. At this point, it is worthy to recall that there is not,

502    for example, just one human genome or just one from any other species, but populations of human

503    genomes and genomes populations from other species. Since every genomic region can be represented

504    by the direct sum of Abelian homocyclic groups of prime-power order, then the whole genome

505    population from individuals from the same or closed related species can be represented as an Abelian

506    group, which will be, in turns, the direct sum of Abelian homocyclic groups of prime-power order.

507    Hence, results lead us to the representation of genomic regions from organismal populations from the

508    same species or close related species (as suggested in Fig 3 to 7) by means of direct sum of their

509    group representation into Abelian cyclic groups. A general illustration of this modelling would be,

510    for example:

511
$$G = (\mathbb{Z}_{5^3})^{n_1} \oplus \overbrace{(\mathbb{Z}_{2^6})^{m_1}}^{motif} \oplus (\mathbb{Z}_{5^3})^{n_2} \oplus ... \oplus \overbrace{(\mathbb{Z}_2^2)^{m_2}}^{domain} \oplus ... \oplus \overbrace{(\mathbb{Z}_{5^3})^{n_p}}^{domain} \oplus \overbrace{(\mathbb{Z}_{2^6})^{m_p}}^{motif} \quad (7)$$

512    That is, Eq. 7 expresses that any large enough genomic region can be represented as direct sum of

513    homocyclic Abelian groups of prime-power order. In other words, the fundamental theorem of

514    Abelian finite groups (FTAG) has an equivalent in genomics.

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

516    as an Abelian group isomorphic to a direct sum of homocyclic Abelian groups of prime-power order.

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

518    Basically, group representations of the genetic code lead to group representations of local genomic

519    domains in terms of cyclic groups of prime-power order, for example, $\left(C_g, +\right) \cong \left(\mathbb{Z}_{2^6}, +\right)$,

520    $\left(C_{G+}, +\right) \cong \left(\mathbb{Z}_5^3, +\right)$ or $\left(C_{g+}, +\right) \cong \left(\mathbb{Z}_{5^3}, +\right)$, till covering the whole genome. As for any finite Abelian

521    group, the Abelian group representation of genome populations can be expressed in terms of a direct

522    sum of Abelian homocyclic groups of prime-power order. Any new discovering on the annotation of

523    a given genome population will only split an Abelian group, already defined on some genomic

524    domain/region, into the direct sum of Abelian subgroups ■.

525         The application of the FTAG in terms of the group representation of genomic regions $G$, as

526    given in Eq. 7, establishes the basis to the study the molecular evolutionary process in terms of

527    endomorphisms. That is, fixed mutational events in the organismal population can be modeled as

528    homomorphism: endomorphisms and automorphisms, all elements of the endomorphism ring $\mathfrak{R}(G)$

529    on $G$ (see next section). In the context of comparative evolutionary genomics, the analysis of the

530    endomorphism ring $\mathfrak{R}(G)$ is an intermediate step for the further application of methods from

531    Category theory, which has the potential to unveil unsuspected features of the genome architecture,

532    hard to be inferred from the direct experimentation.

533    **3.2   The endomorphism ring**

534    A biologically relevant application of the theory presented here relies on the fact that if a finite group

535    $G$ is written as a direct sum of subgroups $G_i$, as given in Eq. 7, then endomorphism ring $End(G)$ is

536    isomorphic to the ring matrices $(A_{ij})$, where $A_{ij} \in Homo(G_i, G_j)$ (homomorphism between $G_i$ and

537    $G_j$ ), with the usual matrix operations [30]. In the case of genomic regions from the species or closed

538    related genomic regions from distinct species, the endomorphism that transform the DNA aligned

539    sequence $\alpha$ into $\beta$ ($\alpha, \beta \in G$) is represented by a matrix with only non-zero elements in the principal

540    diagonal. These diagonal elements are sub-matrices $A_{ii} \in End(G_i)$ or $A_{ii} \in Aut(G_i)$. In other

541    words, mutational events fixed in gene/genome populations can be quantitatively described as

542    endomorphisms and automorphisms.

543    In the Abelian $p$-group defined on $\mathbb{Z}_{p_i^{\alpha_i}}$, the endomorphisms $\eta_i \in End\left(\mathbb{Z}_{p_i^{\alpha_i}}\right)$ are described

544    as functions $f(x) = k\,x \bmod p_i^{\alpha_i}$, where $k$ and $x$ are elements from the set of integers modulo $p_i^{\alpha_i}$.

545    For example, in the cube ACGT the sequence ATACCCATGGCCAAC (blue block in Fig. 3)

546    represented by the vector $(48, 21, 50, 25, 1) \in \left(\mathbb{Z}_{2^6}\right)^5$ is transformed into the sequence

547 ACACCCATGACCAAC, represented by the vector $(16, 21, 50, 17, 1) \in \mathbb{Z}_{2^6}$, by the automorphism:

548
$$\begin{bmatrix} 3 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 57 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}, \text{ i.e.: } (48, 21, 50, 25, 1) \begin{bmatrix} 3 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 57 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} \mod 64 = \begin{pmatrix} 16 \\ 21 \\ 50 \\ 17 \\ 1 \end{pmatrix}.$$

549 Now, it is not difficult to realize that the set of all endomorphisms $\eta_i \in End\left(\mathbb{Z}_{p_i^{\alpha_i}}\right)$ hold the ring

550 axioms mentioned in the Introduction. That is, the set of all endomorphisms $\eta_i \in End\left(\mathbb{Z}_{p_i^{\alpha_i}}\right)$ forms

551 a ring on $\mathbb{Z}_{p_i^{\alpha_i}}$ that we shall denote as $\mathfrak{R}\left(\mathbb{Z}_{p_i^{\alpha_i}}\right)$.

552 As shown in reference [30], if $G = G_1 \oplus G_2 \ldots \oplus G_n$ *is a direct decomposition with fully invariant*

553 *summands*, then :

554
$$End(G) = End(G_1) \oplus End(G_2) \ldots \oplus End(G_n) \tag{8}$$

555 In this modeling, mutational events are represented as endomorphisms $\eta_i \in End\left(\mathbb{Z}_{p_i^{\alpha_i}}\right)$ on $\mathbb{Z}_{p_i^{\alpha_i}}$

556 . This fact permits the study of the genome architecture through the study of the evolutionary

557 (mutational) process in a genome population. Moreover, the decomposition of the endomorphism ring

558 into subgroups, quotient groups, and cosets can lead to a deterministic algebraic taxonomy of the

559 species based on their genome architecture, which is not limited by our current biological knowledge.

560 Particularly relevant for the evolutionary comparative genomics is Baer-Kaplansky theorem: *If G*

561 *and H are p-groups such that* $\mathfrak{R}(G) \cong \mathfrak{R}(H)$, then $G \cong H$ ([29,42]). That is, two Abelian finite

562 groups are isomorphic if, and only if, their endomorphism rings are isomorphic [42]. In other words,

563 genomic regions experiencing mutational events representable by isomorphic rings are algebraically

564 represented by isomorphic Abelian groups and, consequently, have similar genome architecture.

565 Application of Baer-Kaplansky theorem implies that two gene-body regions encoding exactly for

566 the same polypeptide but with different region architecture (Fig 5) are under different evolutionary

567 pressure. That is, if the group representations of two gene-body regions are not isomorphic, then their
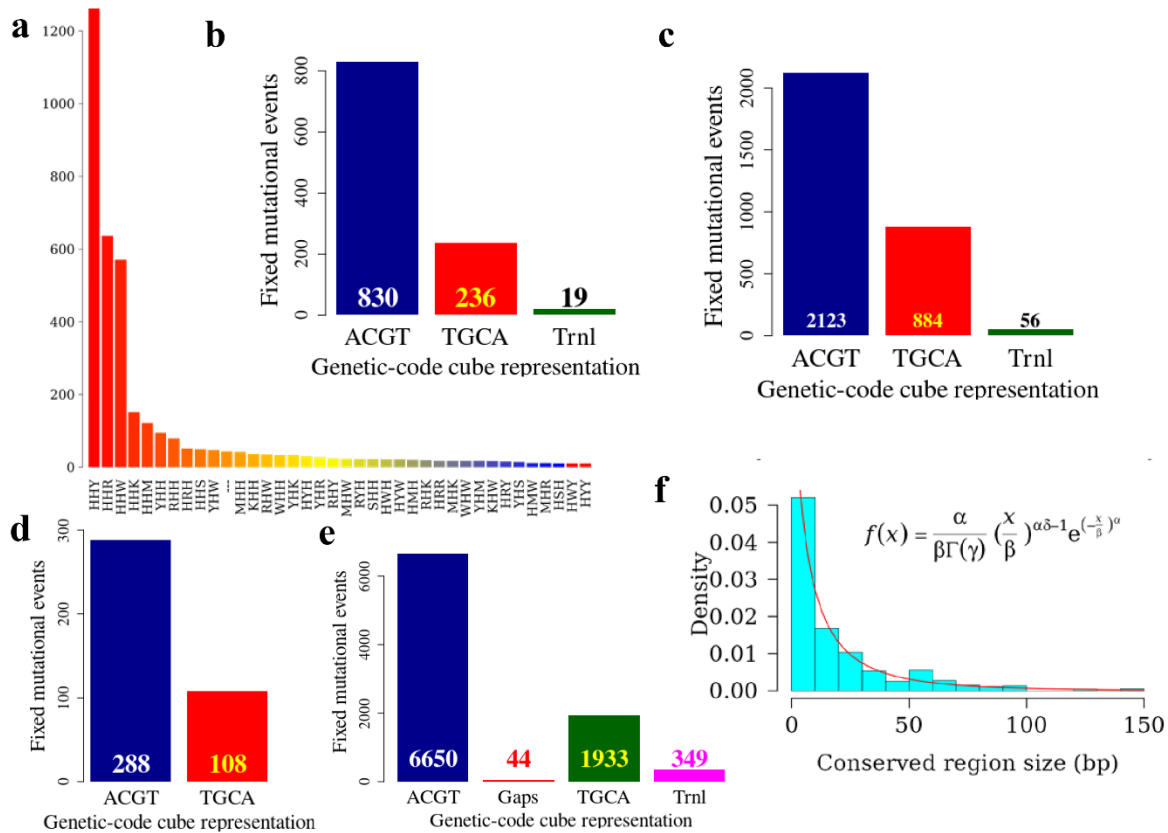
568     endomorphism rings are not isomorphic either and, consequently, they will be under different

569     evolutionary pressure, experiencing different subsets of mutational events, which are represented as

570     endomorphisms from their corresponding endomorphism ring. This scenario is typically found in

571     some isoforms, which are proteins that are similar to each other and perform similar roles within cells

572     [43]. This is the case where two or more closely related genes are responsible for the same translated

573     protein, illustrated in Fig 5. They can be simply duplicated, or paralogous genes, where both paralogs

574     can remain similar (paralog isoforms) if an increased production of the protein is advantageous or if

575     a dosage balance occurs in conjunction with other gene products or where different transcripts can

576     lead to different subcellular localization [44].

577     A screening of mutational events on subsets of aligned genes suggests that the decomposition

578     of protein-coding regions is tractable, conforming Eq. 8. Results with the alignments of several

579     protein-coding regions are shown in Fig 8. In this example, we searched for automorphisms on the

580     *groups of dual cubes* [11]: ACGT – TGCA and CATG – GTAC on $\mathbb{Z}_{2^6}$, which comprise four of the

581     24 possible algebraic representations of the standard genetic code [17] isomorphic to $\mathbb{Z}_{2^6}$.

582     The analysis of the frequency of mutational events (automorphisms, COVID: human *vs* bat

583     strains) by mutation types is shown in Fig 8**a**. Results are consistent with the well-known observation

584     highlighted by Crick: *the highest mutational rate is found in the third base of the codon, followed by*

585     *the first base, and the lowest rate is found in the second one* [45]. However, estimations on different

586     gene sets suggest that the evolutionary pressure on each codon position depends on the

587     physicochemical properties (annotated according to IUPAC nomenclature [36]) of DNA bases. For

588     example, in Fig 8**a** pyrimidine (Y) transitions on the third codon position (HHY) are, by far, the most

589     frequent observed mutational events. While, in BRCA1 gene (SI Fig 2), the frequency of purine

590     (HHR) transitions is followed by pyrimidine (HHY) transitions.

591     The analysis on the pairwise alignment of protein-coding regions of SARS and Bat SARS-like

592     coronaviruses is presented in Fig 8**b** an **c**. Most of the mutational events distinguishing human SARS

593     from Bat SARS-like coronaviruses can be described by automorphism on cube ACGT. This

594     observation was confirmed in primate somatic cytochrome c (Fig 8**c**) and BRCA1 DNA repair gene

595  (Fig 8**d**). Since automorphisms transform the null element (gap-triplet DDD/---) into itself, insertion-

596  deletion mutational events cannot be described by automorphisms but as translations on the groups

597  (denoted as *Trnl* in Fig 8). The representation of conserved genomic regions with homocyclic *p*-group

598  is straightforward. However, their frequency in the genome architecture exponentially decreases with

599  the size of the region (Fig 8**f** and SI Fig 4).



600

601  **Fig 8**. Analysis of mutational events in terms of automorphisms on DNA protein-coding regions
602  represented as homocyclic groups on $\mathbb{Z}_{64}$. In the Abelian group defined on $\mathbb{Z}_{64}$, automorphisms are
603  described as functions $f(x) = k\,x \bmod 64$, where $k$ and $x$ are elements from the set of integers modulo
604  64.  **a**, Frequency of mutational events (automorphisms) according to their mutation type. That is,
605  every single base mutational event across the MSA was classified according IUPAC nomenclature
606  [46]: 1) According to the number of hydrogen bonds (on DNA/RNA double helix): strong S={C, G}
607  (three hydrogen bonds) and weak W={A, U} (two hydrogen bonds). According to the chemical type:
608  purines R= {A, G} and pyrimidines Y= {C, U}. 3). According to the presence of amino or keto groups
609  on the base rings: amino M= {C, A} and keto K= {G, T}. Constant (hold) base positions were labeled
610  with letter H. So, codon positions labeled as HKH means that the first and third bases remains constant
611  and mutational events between bases G and T were found in the MSA. **b** and **c**, Bar plots showing the
612  frequency of automorphisms found on the *group of dual cubes* (see [11]): ACGT – TGCA  and CATG
613  – GTAC  on $\mathbb{Z}_{64}$ between SARS coronavirus GZ02 and bat SARS-like coronaviruses: **a**, isolate
614  Rs7327 (GenBank: KY417151.1, protein-coding regions) and **c**, isolate bat-SL-CoVZC45 (GenBank:
615  MG772933.1:265-1345513455-21542, nonstructural polyprotein). **d**, frequency of automorphisms
616  between human somatic cytochrome c and other nine primates (monkeys). **e**, frequency of
617  automorphisms between human BRCA1 DNA repair gene and other seven primates (see Material and
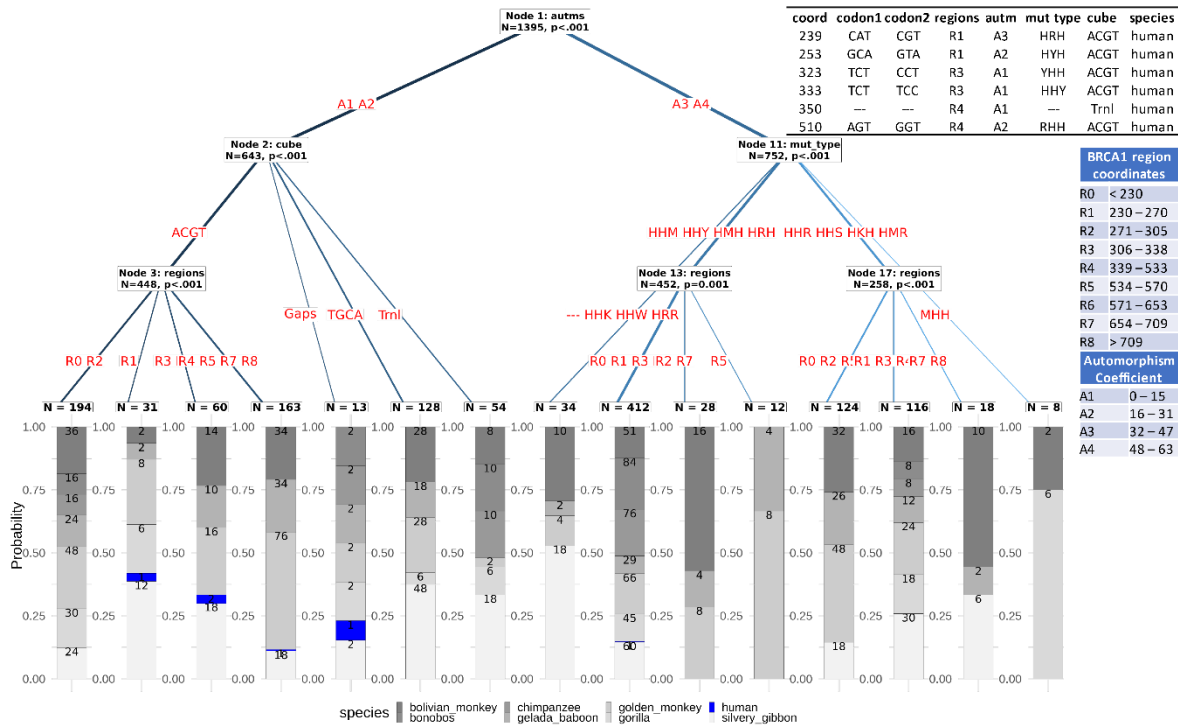
618    Method section). **f**, Distribution of the conserved COVID-19 genomic regions according to their size.
619    The graphics result from the analysis SARS coronavirus GZ02 versus the two mentioned bat strains.
620    The best fitted probability distribution turned out to be the generalized gamma distribution.
621

622        Next, under the assumption that Eq. 8 holds, different protein-coding regions must experience

623    "*preference*" for specific type of automorphisms. To illustrating the concept, an analysis based on the

624    application of Theorem 1 and Eq. 8 on gene/genome population studies, an application of decision

625    tree algorithms was conducted on primate BRCA1 genes. Results for the analysis with Chi-squared

626    Automated Interaction Detection (CHAID) is presented in Fig 9. It is important to keep in mind that

627    this is only an illustrative example with small sample size, and that definite conclusions related to

628    BRCA1 genes can only be derived with larger sample size from humans and non-human primate

629    sequences. In this algorithmic approach, for each compound category consisting of three or more of

630    the original categories, the algorithm finds the most significant binary split for a node (split-variable)

631    based on a chi-squared test [47].

632        For a given MSA of protein-coding regions, the resulting decision tree leads to stochastic-

633    deterministic logical rules (propositions) permitting a probabilistic estimation of the best model

634    approach holding Eq. 8. For example, since only one mutational event human-to-human from class

635    A3 is reported in the right side of the tree (Fig 9), with high probability the proposition: "(A4 $\vee$ (A3

636    $\wedge \neg$ HRH) $\rightarrow \neg$ human" is true. That is, with high probability only non-humans hold the last rule.

637    Due to graphic printing limitations not all tree details are shown in Fig 9 (calculations details are

638    given in the tutorials links provided at SI).

639        Results shown in Fig 9 are only for the purpose to illustrate the application of the theory, since

640    for the sake of visualization and simplicity, were limited to small sample data set and to the

641    application of a relatively "modest" (unsupervised) machine-learning approach which, however, is

642    sufficient to illustrate the concepts. Next, let us suppose that the decision tree from Fig 9 holds on a

643    large enough sample-size (to minimize the classification error) of primate BRCA1-gene populations.

644    Then, with high probability the logical rule: "A1 $\wedge$ R3 $\wedge$ (YHH $\vee$ HHY) $\rightarrow$ human" is true. That is,

645    with high probability transitions mutations (T $\leftrightarrow$ C) on region R3 from BRCA1 gene (specifically at

646    positions 323 and 333, Fig 9) in the first and third codon positions, represented by automorphisms

647    with coefficient between 0 and 15, are not observed in primates other than humans.



648

**Fig 9**. Decision tree based on automorphisms estimated on primate BRCA1 genes. Symbols R0 to
R8 denote the protein regions as given in UniProt plus inter regions segments (see
https://www.uniprot.org/uniprot/P38398#family_and_domains). Only regions experiencing fixed
mutational events are included in the analysis. The range of automorphism coefficients $k$ ($f(x) = kx$
mod 64) are denoted after the isomorphism between the genetic-code cyclic group defined in the set
of codons and the Abelian group defined on $\mathbb{Z}_{64}$. For the sake of graphic comprehension, the
coordinates of human-to-human mutations were added. Every branch (path) from the top to the leaf
node is equivalent to a stochastic-determinist logical rule defining the automorphism preference for
each protein region in the subset of analyzed primate BRCA1 genes. For example, with high
probability the rule: "(A4 ∨ (A3 ∧ ¬ R1)) → ¬ human" is true (see Supporting Information).

660        Obviously, the predictive power of the stochastic rules depends on the size of the samples from

661    the populations under scrutiny. A larger data set including 41 variants of the BRCA1 gene and a rough

662    estimation of the (encoded) *mutational cost* given in the term of a quasichemical energy of aminoacid

663    interactions in an average buried environment [11,48] (data included in the GenomAutomorphism R

664    package [28]) allow reach more robust rules after the application of decision tree algorithms.

665    Likewise, an estimation of *mutational cost* can be given in terms of distances between aminoacids

666    based on codon distances defined on a specific genetic-code cube model or on a combination of two

667    models [11,49]. Examples of stochastic some mutational rules are given in Table 1.

668

669 **Table 1**. Examples of stochastic mutational rules found in aligned DNA sequences from primate
670 BRCA1 genes.

| Mutational cost (MC) | Stochastic Rule[3] |
|---|---|
| Aminoacid contact potential[1] | MC(0.03) → ¬human |
| | MC(-0.47) ∧ R4 ∧ A4 → human |
| | MC(-0.47) ∧ (R0 ∨ R0. ∨ R3 ∨ R5 ) → bonobos |
| | MC(0.08) → bolivian_monkey |
| Aminoacid distance based on genetic-code codon distances[2] | MC(1.34) ∧ R0 → ¬human |
| | MC(1.36) → gorilla |
| | MC(0.28) ∧ R4 ∧ A4 → human |
| | (MC(0.12) ∨ MC(0.12)) ∧ (R1 ∨ R5) → silvery gibbon |
| | MC(0.26) ∧ (A1 ∨ A2) ∧ ¬ R4 ∧ ¬ HHW → human |
| | MC(0.99) ∧ HHS ∧ (R7 ∨ R4) → golden monkey |

671 [1]Aminoacid contact potentials are given in reference [48]. [2]Aminoacid distance based on the codon distance are given in
672 reference [49] and applied (together with the concept of encoded mutational cost) in reference [11]. 3 The decision trees
673 using CHAID algorithm are given in the Supporting Information (also available in the tutorials at
674 https://genomaths.github.io/genomautomorphism).
675

676       Our results provides supporting evidence to the previous finding reported in [11] about that the

677 selection of the genetic-code cube model cannot be arbitrary, since the automorphisms and the

678 estimation of mutational costs (as defined in [11]) on different local DNA protein-coding regions

679 shows clear "preference" for specific models. Obviously, the mathematical model is only a tool (a

680 representation of the physicochemical relationships given between molecules) applied to uncovering

681 the existence of specific evolutionary constraints.

682 **3.3   Future theoretical developments**

683 In this section we want to highlight a direction of future theoretical development. A full coverage of

684 this topic is out of the limits of the current work. Nevertheless, a sketch on a future direction is

685 presented here. Our goal will be the description of mutational process on protein-coding regions in

686 terms of homomorphisms of different algebraic structures.

687       Genomic regions represented as an Abelian group decomposable into homocyclic Abelian $p$-

688 groups, e.g. $\mathbb{Z}_{2^6} \oplus \ldots \oplus \mathbb{Z}_{2^6}$ (n times), can be studied  as $R$-algebras [18], which in particular is a $R$-module

689 and after considering only the sum operation of the ring $\mathbb{Z}_{2^6}$, it is also a $G$-module. Recall that our

690 modeling just takes advantage of the group isomorphism: $(\mathbb{Z}_{64},+) \cong (C_g,+)$ (for the sake of

691     simplicity we are using the same sum operation symbol in both groups, $\mathbb{Z}_{64}$ and $C_g$). Thus, the $\mathbb{Z}_{64}$-

692     algebra of the group $S = \left( \left( C_g \right)^n, + \right) = \left( C_g, + \right) \oplus \overset{n \; times}{\ldots} \oplus \left( C_g, + \right)$ over the ring $\mathbb{Z}_{64}$ can be defined

693     [18].

694     In our current case (considering the codon coordinate level), we are interested on heterocycle

695     groups $S = \oplus G_i$ of $C_g$ and $C_{g+}$ $(G_i \in \{ C_g, C_{g+} \})$, as suggested in Fig. 1, which permits the analysis

696     of multiple sequence alignments including insertion-deletion (indel) mutations. It is not hard to notice

697     that the collection of all the $R$-Module of groups $S$ over the ring $R = \otimes R_i$, $(R_i \in \{ \mathbb{Z}_{64}, \mathbb{Z}_{125} \})$ together

698     with $R$-Module homomorphisms conform to a category of **$R$-Modules**, also denoted as **$R$-Mod**. Let

699     $\mathcal{C}_N$ be the category **Ab** with the Abelian groups of the DNA sequences of length $\leq N$ as objects and

700     group homomorphisms as morphisms (see Appendix A). Fredy's theorem states that every Abelian

701     category is a subcategory of some category of modules over a ring [50]. Mitchell has reinforced

702     Fredy's result, proving that every Abelian category is a full subcategory of a category of modules

703     over a ring [51].

704     At codon coordinate level, the group defined on the set of codon is a subgroup of the group

705     defined on the set of extended base-triplets ( $C_g \subset C_{g+}$ ) and the $\mathbb{Z}_{125}$-Module of group $C_g$ is a

706     submodule of the $\mathbb{Z}_{125}$-Module of group $C_{g+}$ over the ring $\mathbb{Z}_{125}$. The triplet of gaps '---' corresponds

707     to the identity element of group $C_{g+}$, which is mapped into $0 \in \mathbb{Z}_{5^3}$ by $Hom\left( C_{g+}, \mathbb{Z}_{5^3} \right)$. A

708     homomorphism always maps the identity element from the domain of group, say $\mathbf{0}_{C_g}$, into the identity

709     element from the codomain $\mathbf{0}_{C_{g+}}$, which in $C_{g+}$ is $0_{C_{g+}} = $ '---'.

710     The following example illustrates a possible sequence of attainable analytical steps with

711     concrete computational biology application. Let $A = $ GACAGAGCAGTATTAGCTTCACAC and $B$

712     $= $ GAAAACGTATTATCAAAG DNA sequence segments represented as elements from the groups:

713     $G_A = C_g^{\text{ACGT}} \oplus C_g^{\text{TGCA}} \oplus (C_g^{\text{ACGT}})^6$ and $G_B = C_g^{\text{ACGT}} \oplus C_g^{\text{TGCA}} \oplus (C_g^{\text{ACGT}})^4$, respectively, where $C_g^X$

714     is the Abelian $p$-group defined on the set of 64 codons and base orders (cubes): $X = \{ \text{ACGT}, \text{TGCA} \}$.

715    Groups $G_A$ and $G_B$ are elements of the **Ab** category $\mathcal{C}_N$ defined on the collection of heterocyclic

716    group $\left(C_g^X\right)^N$ defined on the set of DNA sequences (of codons) with length $N \le 8$.

717         Since the triplet of gaps cannot be arbitrary allocated in the sequence, the alignment of DNA

718    sequence is an essential step required for the application of this modeling preserving the biological

719    meaning. The pairwise alignment of the corresponding aminoacid sequences from $A$ and $B$ yields:

720    DRAVLASQ ,        which    corresponds    to    the    DNA    sequence    alignment:    *aln*
       EN-VL-SN

721    $= \begin{pmatrix} \text{GACAGAGCAGTATTAGCTTCACAC} \\ \text{GAAAAC---GTATTA---TCAAAG} \end{pmatrix}$. That is, to preserve the reading frame, a robust alignment is

722    accomplished translating the codon sequence into aminoacid sequence alignment.

723         Sequences $A$ and $B' = $ GAAAAC---GTATTA---TCAAAG can also be represented as elements

724    from group:

725    $$G_{A'} = C_g^{\text{ACGT}} \oplus C_g^{\text{TGCA}} \oplus C_{g+}^{\text{ACGT}} \oplus (C_g^{\text{ACGT}})^2 \oplus C_{g+}^{\text{ACGT}} \oplus (C_g^{\text{ACGT}})^2$$

726    This group is an element of the **Ab** category $\mathcal{C}_{A'}$, which is a subcategory of the $\boldsymbol{R_{A'}}$-**Mod** category

727    over the ring $R_{A'} = \left(\mathbb{Z}_{2^6}\right)^2 \otimes \mathbb{Z}_{5^3} \otimes \left(\mathbb{Z}_{2^6}\right)^2 \otimes \mathbb{Z}_{5^3} \otimes \left(\mathbb{Z}_{2^6}\right)^2$. The group isomorphism $F_B : G_B \to G_{B'}$

728    is the functor that maps DNA sequences from group $G_B \in \mathcal{C}_N$ into an element from group $G_{B'} \in \mathcal{C}_R$

729    (see Appendix B). That is, for all element $b = (X_1, X_2, X_3, X_4, X_5, X_6)$ ($b \in G_B$) there is a unique

730    element $b' = (X_1', X_2', 0, X_3', X_4', 0, X_5', X_6')$ ($X_i' = X_i$ and $b' \in G_{B'}$).

731         Also, there is an injective morphism $F_A : G_A \to G_{A'}$ that transforms each element

732    $a = (X_1, X_2, X_3, X_4, X_5, X_6, X_7, X_8)$        ($a \in G_A$)        into        a        unique        element

733    $a' = (X_1', X_2', X_3', X_4', X_5', X_6', X_7', X_8')$ ($a' \in G_{A'}$), which is evident since $C_g \subset C_{g+}$ and,

734    consequently, codons are preserved, i.e., $X_i' = X_i$ and $G_A$ is isomorphic to the image $F_A(G_A)$. The

735    homomorphism $F_{A'} : G_A \to G_{A'}$ is also a functor which maps elements from the $\boldsymbol{R_A}$-**Mod** category

736    over the ring $R_A = \otimes_8 \mathbb{Z}_{2^6}$ into the $\boldsymbol{R_{A'}}$-**Mod** category. Notice that $F_B(B)$ is a subgroup of $F_A(A)$.

737    In practice, for the sake of computational genomics implementations, the aligned DNA

738    sequences $A$ and $B$ can be represented by the numerical vectors a $a = (9,32,24,56,60,27,28,5)$ and

739    $b = (8,1,56,60,28,1)$, respectively, with coordinates on $\mathbb{Z}_{2^6}$. The application of the morphisms $F_A$

740    and $F_B$ permits the new representations: $a' = \left((9,32) \in \mathbb{Z}_{2^6}, 66 \in \mathbb{Z}_{5^3}, (56,60) \in \mathbb{Z}_{2^6}, 69 \in\right.$

741    $\left.\mathbb{Z}_{5^3}, (28,5) \in \mathbb{Z}_{2^6}\right)$    and    $b' = \left((8,1) \in \mathbb{Z}_{2^6}, 0 \in \mathbb{Z}_{5^3}, (56,60) \in \mathbb{Z}_{2^6}, 0 \in \mathbb{Z}_{5^3}, (28,1) \in \mathbb{Z}_{2^6}\right)$,

742    respectively. The group homomorphism $\varphi$ with matrix representation with diagonal elements

743    $\left((8,2) \in \mathbb{Z}_{2^6}, 0 \in \mathbb{Z}_{5^3}, (1,1) \in \mathbb{Z}_{2^6}, 0 \in \mathbb{Z}_{5^3}, (1,1) \in \mathbb{Z}_{2^6}\right)$ maps sequence $a'$ into $b'$, *i.e.*, $\varphi(a') = b'$

744    :

$$
(9,32,66,56,60,69,28,5)
\begin{bmatrix}
8 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 2 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 13
\end{bmatrix}
=
\begin{pmatrix}
8 \\ 1 \\ 0 \\ 56 \\ 60 \\ 0 \\ 28 \\ 1
\end{pmatrix}
$$

746    Where the third and sixth rows are computed modulo 125 and the rest modulo 64. The group

747    homomorphism $h : G_{B'} \to G_A$ that accomplish the mapping $h(b') = a'$ is computed as:

$$
(8,1,0,56,60,0,28,1)
\begin{bmatrix}
58 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 32 & 66 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\
57 & 0 & 0 & 0 & 0 & 69 & 0 & 5
\end{bmatrix}
=
\begin{pmatrix}
9 \\ 32 \\ 66 \\ 56 \\ 60 \\ 69 \\ 28 \\ 5
\end{pmatrix}
$$

749    Or by means of the affine transformation:

$$
\left(
(8,1,0,56,60,0,28,1)
\begin{bmatrix}
58 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 32 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 5
\end{bmatrix}
\bmod 64 +
\begin{pmatrix}
118 \\ 0 \\ 66 \\ 0 \\ 0 \\ 69 \\ 0 \\ 0
\end{pmatrix}
\right)
\bmod 125 =
\begin{pmatrix}
9 \\ 32 \\ 66 \\ 56 \\ 60 \\ 69 \\ 28 \\ 5
\end{pmatrix}
$$

751    In summary, a future theoretical development in the framework of category theory opens new

752    horizons for the analysis of the mutational process in a wider computational genomic scenario not

753    previously studies in molecular evolutionary biology.

## 4    Discussions

755    The encoding of the physicochemical relationships between nucleotides (nitrogenous bases) in the

756    DNA double helix in terms of group operations permits a mathematical representation of genome

757    architecture interpretable in a molecular evolutionary context. The group representations of the

758    genetic code are logically extended from protein-coding DNA regions to the entire genome. As shown

759    in Fig 1, the Abelian group representation of genomic regions into the direct sum of Abelian $p$-group

760    is only one of several steps addressed to get better understanding on how genomes are built.

761    The advantage on using group representations is that there exists a well-established

762    mathematical development that leads to an objective study of the genome architecture in a molecular

763    evolutionary context, through the analysis of mutational events in terms of group homomorphisms:

764    endomorphisms, automorphisms, and translations. On this scenario the analysis of group

765    homomorphisms permits us the uncovering of stochastic rules constraining the local architecture on

766    genes and genomic regions. The goal is unveiling hidden genomic architecture and rules hard to be

767    detected by current experimental approaches. All the information required can be retrieved from the

768    MSA of DNA sequences, which is particularly relevant for poorly annotated genomes.

769    Examples shown in Figs 3 to 4 indicates that whatever would be the genomic architecture for

770    given species, the observed variations in the individual populations and in populations from closed

771    related species, it can be quantitatively described as the direct sum of Abelian cyclic groups. The

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

773    known Abelian homocyclic or cyclic groups representing a genomic subregion into direct sums of

774    subgroups. In such algebraic representation DNA sequence motifs for which only substitution

775    mutations happened are specifically represented by the Abelian group $\left( C_g, + \right) \cong \left( \mathbb{Z}_{64}, + \right)$, in protein

776    coding regions, and by any combination of groups $\left(\mathfrak{B},+\right)\cong\left(\mathbb{Z}_2^2,+\right)$, $\left(\mathfrak{B}^2,+\right)\cong\left(\mathbb{Z}_2^4,+\right)$ or some

777    quotient group like $C_G/G_{\mathrm{GGA}}\cong\left(\mathbb{Z}_2,+\right)$ in non-protein coding regions.

778        Notably, the genetic code Abelian group $\left(C_{G+},+\right)\cong\left(\mathbb{Z}_5^3,+\right)$ is enough for an algebraic

779    representation of the genome population from the same species or close related species. However,

780    such a decomposition leads to a poor description of local architecture that, as suggested in Figs. 3 to

781    6, can mask relevant biological features. Figure 3 to 6 illustrate the basic Abelian group

782    representations for further analysis of genome architecture through the study of the mutational events,

783    as essential transformations inherent to the molecular evolutionary process, in terms of

784    endomorphisms and automorphisms, elements of the endomorphism ring.

785        For the sake of reader's comprehension, the examples on the group representation of genomic

786    regions presented here are simple. However, the analysis demands for the development of novel

787    computational algebraic approaches to study the genomic architecture. Unlike to traditional

788    computational algebra, we can take advantage of the group isomorphisms, which permits decreasing

789    the computational complexity by avoiding symbolic computation. Nevertheless, results presented

790    here show that the architecture of genome region in an entire population can be quantitively studied

791    in the framework of Abelian group theory.

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

793    architecture unveiled when represented it in terms of sums of finite Abelian groups, which can be

794    further studied with the application of methods from category theory, the potential success of which

795    has been proven in programming languages and in linguistic [52]. The future developments of

796    genome annotation from several species can certainly lead to the discovery of logical rules from such

797    a language, finding the viable variations in the populations. The identification of quotient groups (at

798    larger scale) can permit the stratification of large genome population into equivalence classes

799    (quotient subgroups) corresponding to individual subpopulations, each one of them carrying

800    particular viable variations of species genome architecture.

801      As indicated in reference [18], natural genomic rearrangement like DNA recombination and

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

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

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

805      of living organisms [5,53]. The analysis and discussion on the application of the endomorphism ring

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

807      Particularly promising is the application of the genomic Abelian groups on epigenomic studies,

808      which results from the model where base $D$ stands for the methylated cytosine and adenine. As

809      suggested in Fig 7**a** and **b**, a precise decomposition of methylation motif into the direct sum of Abelian

810      finite groups can lead to their classification into unambiguous equivalence classes. The group

811      structure of the methylation regulatory regions: **GATC**TTTTATGC and GGTTAAAA**GATC,** both

812      represented by the homocyclic group on $\left(\mathbb{Z}_5^3\right)^4$, breaks from the monotone homocyclic group

813      representation of the region in terms of cyclic groups on $\mathbb{Z}_4$ (Fig 7**a** and **b**). The group representation

814      of protein-coding regions (or base-triplet sequences) as numerical vectors with coordinates on $\mathbb{Z}_{5^3}$

815      (Fig 7**c** and **d**) facilitates the analysis of methylation changes represented as group

816      endomorphism/automorphisms of the cyclic group on $\mathbb{Z}_{5^3}$.

817      Results indicate that, as a consequence of the genetic code constraints and the evolutionary

818      pressure on protein-coding regions, stochastic-deterministic logical rules can be inferred on a large

819      enough sample-size from a gene/genomic-region population. Such a stochastic-deterministic rules

820      lead to specific applications of Theorem 1 and Eq.8, consequently, the analysis of mutational process

821      on each group, subgroup, and coset. For example, mutational events on a MSA column (identified)

822      from class YHH (with discriminatory classification power as shown in Fig 8) where the second and

823      third DNA bases remain invariant (H) and the first base are pyrimidines (Y) experiencing transition

824      mutations (across individuals sequences) are represented by automorphisms on a subgroup (from the

825      genetic code Abelian subgroup $\left(C_G,+\right)$) defined on the set {THH, CHH} isomorphic to $\left(\mathbb{Z}_2,+\right)$ [20].

826      Figure 9 provides illustrative example that motivates further applications of based machine-

827      learning bioinformatic approaches to unveil the subjacent logic to the genome architecture and its

828      association with the DNA cytosine/adenine methylation patterning found on individual populations

829      and the changes (repatterning) induced by, e.g., environmental changes, aging process and diseases,

830      which is of particular interest in genomic medicine [54]. Machine learning applications on MSA

831      involving large sample size of genomic regions from populations of different species can unveil

832      further decompositions into the direct sum of Abelian groups, which do not depend on our current

833      knowledge of the annotated genomes. As suggested in Fig 9, we can expect that most of the hiding

834      genomic DNA sequence motif can be unveiled by studying the molecular evolutionary (mutational)

835      process in a genome population through the lens of the endomorphism ring. In other words, as a

836      consequence of the injective relationship: *DNA sequence → 3D chromatin architecture* [3,4,6], fixed

837      mutational events (in organismal populations) on DNA sequence motifs involved in the 3D chromatin

838      architecture are under evolutionary pressure, biophysically and biochemically constrained to preserve

839      the chromatin biological functions.

840      Results shown in Figs. 8 and 9 also suggest deep implications of Baer-Kaplansky theorem on

841      the genome architecture unknown by the current knowledge and understanding of genome annotation,

842      which currently relies on the DNA sequence itself. Concretely, on an evolutionary context, the fact

843      that two genomic regions from two different species are almost identical and, event would encode for

844      the same functional protein, does not necessarily imply that they hold to the same genome

845      architecture. The evolutionary pressure in both such hypothetical regions must be same, which

846      implies that the regions experience the same type of mutational events in terms of

847      automorphism/endomorphism representations.

848      For example, let's suppose that the results shown in Fig 9 were derived from a large sample

849      size (large enough to derive statistically significant rules), then the rule "A1 ∧ R3 ∧ (YHH ∨ HHY)

850      → human" (Fig 9) implies that the gene regions of BRCA1 from human and non-human primates do

851      not belong to the same equivalent class of genomic region. In particular, since the endomorphism

852      rings $\mathfrak{R}\left(G_{human}^{\text{BRCA1}}\right)$ and $\mathfrak{R}\left(G_{non-human}^{\text{BRCA1}}\right)$ on the Abelian groups $G_{human}^{\text{BRCA1}}$ and $G_{non-human}^{\text{BRCA1}}$ defined on the

853     human and non-human primates BRCA1 genes, respectively, are not isomorphic, then according to

854     the Baer-Kaplansky theorem groups $G_{human}^{\text{BRCA1}}$ and $G_{non-human}^{\text{BRCA1}}$ are also not isomorphic. Hence, region

855     architectures of BRCA1 gene in human and non-human primates are (in this hypothetical scenario)

856     implicitly different, which is not obvious to human eyes from their MSA (see supporting information).

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

858     evolutionary biology, which heavily relies on subjective evolutionary null hypotheses about the past.

859     As suggested in reference [11], the genomic Abelian groups open new horizons for the study of the

860     molecular evolutionary stochastic processes (at genomic scale) with relevant biomedical applications,

861     founded on a deterministic ground, which only depends on the physicochemical properties of DNA

862     bases and aminoacids. In this scenario, the only molecular evolutionary hypothesis needed about the

863     past is a fact, the existence of the genetic code.

864     Remarkably, further studies applying the theory presented here do not require for special

865     experimental datasets but for the DNA sequences of the genomic regions under scrutiny. Although

866     the accuracy of the predictions depends on the sample size, the number of sequenced genomes stored

867     in the databases grows year-after-year. Large samples of DNA sequences (from homolog genomic

868     regions) from at least two or more species facilitate application of Baer-Kaplansky theorem and

869     further studies applying methods of Categorical theory to unveil the grammar embedded in the DNA

870     sequences.

871     The theory and concretes examples provided here make explicit the basic foundation for a

872     further unprecedented application of the last advances in Abelian group theory incorporating methods

873     from Category theory, where groups and group homomorphisms (in our context: mutational events)

874     are the main players, which have the potential to discover unsuspected features of the genome

875     architecture, opening new horizon to the genomic taxonomy of species in accordance with the state-

876     of-the-art in mathematics, logic, and computational sciences. In other words, these applications have

877     the potential to elevate the genomic studies from the current descriptive level to the vanguard level

878     marked in the frontier of science by mathematics, physics, and computational sciences.

## 5   Conclusions

Results to date indicate that the genetic code and the physicochemical properties of DNA bases on which the genetic code algebraic structure are defined, has a deterministic effect, or at least partially rules on the current genome architectures, in such a way that the Abelian group representations of the genetic code are logically extended to the whole genome. In consequence, the fundamental theorem of Abelian finite groups can be applied starting from genomic regions till cover whole chromosomes. This result opens new horizons for further genomics studies with the application of the Abelian group theory, which currently is well developed and well documented [30,55].

Results suggest that the architecture of current population genomes is quite far from randomness and obeys stochastic-deterministic rules. The nexus between the Abelian finite group decomposition into homocycle Abelian $p$-groups and the endomorphism ring paved the ways to unveil unsuspected stochastic-deterministic logical propositions ruling the ensemble of genomic regions and sets the basis for a novel algebraic taxonomy of the species, which is not limited by our current biological knowledge.

In the context of evolutionary comparative genomics, the theory presented here open new horizons for the application of Group theory including methods of Category theory, which have the potential to unveil hidden features and rules inherent to the genome architecture, leading to an unprecedented understanding on how genomes are built.

We believe that the mathematical formalism proposed here sets the theoretical ground for a further development in genomics, transitioning the field from a fully empirical science to a predictive science, where the theoretical and empirical research coexist in a tight positive feedback loop, a development stage only reached so far in the field of physics.

All the above claims are feasible, only limited by our computational power and the availability of samples of sequenced genomes from the same species and from multiple species.

At this point we emphasize that *an accurate understanding of the genome architecture and population's structure, on a formal mathematical framework, is as essential for the future of genetic*

905     *engineering and genome editing as the physics of architecture is to the design of sturdy and stable*

906     *energy-efficient building.*

## 907   6   Appendix A. Genetic code algebraic structures defined on the base and
## 908       codon sets

909     An Abelian group structure $(B,+)$ is a set **B** together with a binary operation '+' that combines any

910     two elements $a \in B$ and $b \in B$ to form another element of $c \in B$, denoted $a + b = c$, which satisfy

911     the following axioms:

912         1)   Associativity. For all $a, b, c \in B$, the equality $(a + b) + c = a + (b + c)$ holds.

913         2)   Identity. There exists an element $e \in B$ named identity element of $B$, such that for any

914             $a \in B$, the equality $a + e = a$ holds.

915         3)   Commutativity. For all $a, b \in B$, the equality $a + b = b + a$

916         The Abelian groups considered here are finite cyclic groups $(G, +)$ isomorphic to the Abelian

917     group defined on the set of integers modulo $n$, denoted as $\mathbb{Z}_n (\mathbb{Z}/n\mathbb{Z})$. That is, the integers

918     $1, 2, 3, \ldots, n-1$ form a cyclic group of order $n$ under addition (modulo $n$) and 0 as the identity element.

919     This group will be denoted as $(\mathbb{Z}_n, +)$. However, for the sake of simplicity in the figures it will be

920     denoted simply as $\mathbb{Z}_n$, i.e., without making distinction between the set $\mathbb{Z}_n$ and group structure

921     defined on it. The particular interest for the current work is the Abelian $p$-group derived when $n = p^\alpha$

922     where $p$ is a prime number and $\alpha$ an integer. The group operations defined on the set of bases or on

923     the codon set are associated to physicochemical properties of DNA bases (see the next sections).

924

925     ***Homomorphisms and isomorphisms***

926     In modern algebra, a group homomorphism is a map $f : A \rightarrow B$ between two group structures $(A, \bullet)$

927     and $(B, \circ)$ such that for all $a, b \in A$ holds: $f(\alpha_1 \bullet \alpha_2) = f(\alpha_1) \circ f(\alpha_2) = \beta_1 \circ \beta_2$, where

928 $\beta_1, \beta_2 \in B$. A group isomorphism is a one-to-one correspondence (mapping) between two sets that

929 preserves binary relationships between elements of the sets. That is, an isomorphism is a

930 homomorphism holding the inverse mapping: $f^{-1}(\beta_1 \circ \beta_2) = f^{-1}(\beta_1) \bullet f^{-1}(\beta_2) = \alpha_1 \bullet \alpha_2$. For

931 example, since there exists only one cyclic group with four elements up to isomorphism, for each one

932 of the 24 cyclic group $(\mathfrak{B}, +_b)$ defined on the set of bases $\mathfrak{B} = \{A, C, G, T\}$ ([17,18]) there exists a

933 one-to-one mapping $f$ such that for each base $\beta \in \mathfrak{B}$ there is an integer $\iota \in \mathbb{Z}_4$ such that $f(\beta) = \iota$

934 and:

935     1. $f(\beta_1 +_b \beta_2) = f(\beta_1) + f(\beta_2) = \iota_1 + \iota_2$, $\beta_1, \beta_2 \in \mathfrak{B}$ and $\iota_1, \iota_2 \in \mathbb{Z}_4$.

936     2. The inverse mapping $f^{-1}(\iota_1 + \iota_2) = f^{-1}(\iota_1) +_b f^{-1}(\iota_2) = \beta_1 +_b \beta_2$

937 To highlight the fact that the sum operations are defined on different ways on the sets $\mathfrak{B}$ and $\mathbb{Z}_4$, we

938 have used the symbols ' $+_b$ ' and '+', respectively. However, for sake of brevity of the symbolic

939 notation, such knowledge will be considered implicit, writing simply '+'. Then, we said that groups

940 $(\mathfrak{B}, +_b)$ and $(\mathbb{Z}_4, +)$ are isomorphic; in symbols $(\mathfrak{B}, +_b) \cong (\mathbb{Z}_4, +)$. $f$ and its inverse $f^{-1}$ are

941 called isomorphisms. If $f$ (and its inverse $f^{-1}$) is a mapping from a group into itself, then $f$ is called

942 an automorphism. A mapping $g$, not necessarily one-to-one, of the elements from a group into itself

943 is called a group endomorphism. An endomorphism that is also an isomorphism is an automorphism.

944       A ring algebraic structure is obtained when together with the sum operation "+" (as defined

945 above) a new operation "·" is defined on the set $B$ holding the properties:

946     1. Associativity: $(a \cdot b) \cdot c = a \cdot (b \cdot c)$ for all $a, b, c \in B$

947     2. Multiplication is distributive with respect to addition:

948       a. $(a + b) \cdot c = (a \cdot c) + (b \cdot c)$ for all $a, b, c \in B$ (right distributivity).

949       b. $c \cdot (a + b) = c \cdot a + c \cdot b$ for all $a, b, c \in B$ (left distributivity).

950　　As it is shown in the next section, these algebraic structures have been defined on the genetic code.

951　　In particular, the ring $\left(\mathbb{Z}_{2^6}, +. \cdot\right)$ and endomorphism ring (section 3.1) has been defined and studied

952　　on the genetic code [18].

### *Appendix B. Category*

954　　Category theory is a general mathematical theory of structures and of systems of structures that

955　　occupy a central position in contemporary mathematics, theoretical computer science, and linguistics

956　　[56].

957　　**Definition**: A category $\mathcal{C}$ can be described as a collection of objects $\mathcal{O}$ satisfying the

958　　following three conditions:

959　　　1) *Morphism*: For every pair $X$, $Y$ of objects, there is a set $Hom(X,Y)$ called the

960　　　　 *morphisms* from $X$ to $Y$ in $\mathcal{C}$. If $f$ is a morphism from, we write $f: X \rightarrow Y$.

961　　　2) *Identity*: For every object $X$, there exists a morphism $id_X$ in $Hom(X,Y)$, called the

962　　　　 *identity* on $X$ (also denoted as $1_X$).

963　　　3) *Composition*: For every triple $X$, $Y$, and $Z$ of objects, there exists a partial binary

964　　　　 operation from $Hom(X,Y) \times Hom(Y,Z)$ to $Hom(X,Z)$, called the composition of

965　　　　 morphisms in $\mathcal{C}$. If $f: X \rightarrow Y$ and $g: Y \rightarrow Z$, this composition is written as the

966　　　　 mapping $(g \circ f): X \rightarrow Z$.

967　　Identity, morphisms, and composition satisfy two axioms:

968　　*Associativity*: If $f: X \rightarrow Y$, $g: Y \rightarrow Z$, and $h: Z \rightarrow W$, then $h \circ (g \circ f) = (h \circ g) \circ f$.

969　　*Identity*: If $f: X \rightarrow Y$, then $f_X \circ f = f$ and $f \circ f_X = f$.

970

971　　**Definition**: A functor $F$ is a function between two categories $\mathcal{C}$ and $D$ which maps objects to

972　　objects and morphisms to morphisms. That is:

973　　　• For each $X \in \mathcal{C}$ there is an object $F(Y) \in D$

974     •    For each morphism $f: X \rightarrow Y$ in $\mathcal{C}$ there is morphism $F(f): F(X) \rightarrow F(Y)$ in $D$ such

975        that the following conditions hold:

976        i.      $F(g \circ f) = F(g) \circ F(f)$ for all morphisms $f: X \rightarrow Y$ and g: $X \rightarrow Y$ in $\mathcal{C}$

977        ii.      $F(id_X) = id_{F(X)}$ for all $X \in \mathcal{C}$.

## 7   Supporting Information

979 A summary with the reported genetic code Abelian groups relevant for the current study is provided

980 as supporting information in a file named: *Supporting_Information.docx*.

981        All the data, computational and statistical analyses can be reproduced following the R scripts

982 provided in tutorials available at the GenomAutomorphism R package website

983 https://genomaths.github.io/genomautomorphism/. In particular, data and R scripts used in the

984 computation of automorphisms and the decision tree from Fig 9 are available within

985 GenomAutomorphism R package and in a tutorial at:

986 https://genomaths.github.io/genomautomorphism/articles/automorphism_and_decision_tree.html.

## 8   References

988 [1]    E. V. Koonin, Evolution of genome architecture, Int. J. Biochem. Cell Biol. 41 (2009) 298–
989       306. https://doi.org/10.1016/j.biocel.2008.09.015.

990 [2]    S. Whalen, R.M. Truty, K.S. Pollard, Enhancer-promoter interactions are encoded by complex
991       genomic signatures on looping chromatin, Nat. Genet. 48 (2016) 488–496.
992       https://doi.org/10.1038/ng.3539.

993 [3]    J. Nuebler, G. Fudenberg, M. Imakaev, N. Abdennur, L.A. Mirny, Chromatin organization by
994       an interplay of loop extrusion and compartmental segregation, Proc. Natl. Acad. Sci. U. S. A.
995       115 (2018) E6697–E6706. https://doi.org/10.1073/pnas.1717730115.

996 [4]    M.J. Rowley, V.G. Corces, Organizational principles of 3D genome architecture, Nat. Rev.
997       Genet. 19 (2018) 789–800. https://doi.org/10.1038/s41576-018-0060-8.

998 [5]    A. Piazza, W.D. Heyer, Homologous Recombination and the Formation of Complex Genomic
999       Rearrangements, Trends Cell Biol. 29 (2019) 135–149.
1000     https://doi.org/10.1016/j.tcb.2018.10.006.

1001 [6]    H. Zheng, W. Xie, The role of 3D genome organization in development and cell
1002     differentiation, Nat. Rev. Mol. Cell Biol. 20 (2019) 535–550. https://doi.org/10.1038/s41580-

1003    019-0132-4.

1004    [7]    T.D. Schneider, Evolution of biological information, Nucleic Acids Res. 28 (2000) 2794–9.
1005           http://www.pubmedcentral.nih.gov/articlerender.fcgi?artid=102656&tool=pmcentrez&rende
1006           rtype=abstract.

1007    [8]    H.P. Yockey, Origin of life on earth and Shannon's theory of communication, Comput. Chem.
1008           24 (2000) 105–123. https://doi.org/10.1016/S0097-8485(99)00050-9.

1009    [9]    R. Sanchez, R. Grau, A genetic code Boolean structure. II. The genetic information system as
1010           a    Boolean    information    system,    Bull.    Math.    Biol.    67    (2005)    1017–1029.
1011           https://doi.org/10.1016/j.bulm.2004.12.004.

1012    [10]   R. Sanchez, S.A. Mackenzie, Information thermodynamics of cytosine DNA methylation,
1013           PLoS One. 11 (2016) e0150427. https://doi.org/10.1371/journal.pone.0150427.

1014    [11]   R. Sanchez, Symmetric Group of the Genetic-Code Cubes. Effect of the Genetic-Code
1015           Architecture on the Evolutionary Process, MATCH Commun. Math. Comput. Chem. 79
1016           (2018) 527–560. http://match.pmf.kg.ac.rs/electronic_versions/Match79/n3/match79n3_527-
1017           560.pdf.

1018    [12]   R. Sánchez, E. Morgado, R. Grau, A genetic code Boolean structure. I. The meaning of
1019           Boolean    deductions,    Bull.    Math.    Biol.    67    (2005)    1–14.
1020           https://doi.org/10.1016/j.bulm.2004.05.005.

1021    [13]   H.P. Yockey, Information theory, evolution and the origin of life, Inf. Sci. (Ny). 141 (2002)
1022           219–225. https://doi.org/10.1016/S0020-0255(02)00173-1.

1023    [14]   R. Sanchez, R. Grau, A genetic code Boolean structure. II. The genetic information system as
1024           a    Boolean    information    system,    Bull.    Math.    Biol.    67    (2005)    1017–1029.
1025           https://doi.org/10.1016/j.bulm.2004.12.004.

1026    [15]   D.. Andrews, M.L. Boss, DNA code may transmit "Superinformation," Yale Sci. Mag. 45
1027           (1971) 50–51.

1028    [16]   G.S. Chirikjian, Algebraic and Geometric Coding Theory, in: Stoch. Model. Inf. Theory, Lie
1029           Groups, Vol. 2 Anal. Methods Mod. Appl., Birkhäuser Boston, Boston, 2012: pp. 313–336.
1030           https://doi.org/10.1007/978-0-8176-4944-9_9.

1031    [17]   M. V José, E.R. Morgado, R. Sánchez, T. Govezensky, The 24 Possible Algebraic
1032           Representations of the Standard Genetic Code in Six or in Three Dimensions, Adv. Stud. Biol.
1033           4 (2012) 119–152. http://www.m-hikari.com/asb/asb2012/asb1-4-2012/joseASB1-4-2012-
1034           1.pdf.

1035    [18]   R. Sanchez, E. Morgado, R. Grau, Gene algebra from a genetic code algebraic structure, J.
1036           Math. Biol. 51 (2005) 431–457. https://doi.org/10.1007/s00285-005-0332-8.

1037    [19]   R. Sanchez, R. Grau, E. Morgado, A novel Lie algebra of the genetic code over the Galois
1038           field    of    four    DNA    bases,    Math.    Biosci.    202    (2006)    156–174.

1039          https://doi.org/10.1016/j.mbs.2006.03.017.

1040   [20]   R. Sánchez, R. Grau, An algebraic hypothesis about the primeval genetic code architecture,
1041          Math. Biosci. 221 (2009) 60–76. https://doi.org/S0025-5564(09)00114-X [pii]
1042          10.1016/j.mbs.2009.07.001.

1043   [21]   D.. Andrew, M.L. Boss, DNA code may transmit "Superinformation," Chem. Eng. News. 48
1044          (1970) 50–51.

1045   [22]   H.J. Danckwerts, D. Neubert, Symmetries of genetic code-doublets, J Mol Evol . 5 (1975)
1046          327–332. https://doi.org/10.1007/BF01732219.

1047   [23]   M.O. Bertman, J.R. Jungck, Group graph of the genetic code, J. Hered. 70 (1979) 379–384.
1048          https://doi.org/10.1093/oxfordjournals.jhered.a109281.

1049   [24]   J.E.M. Hornos, Y.M.M. Hornos, Algebraic model for the evolution of the genetic code, Phys.
1050          Rev. Lett. 71 (1993) 4401–4404. https://doi.org/10.1103/PhysRevLett.71.4401.

1051   [25]   J.E.M. Hornos, Y.M.M. Hornos, M. Forger, Symmetry and Symmetry Breaking: an Algebraic
1052          Approach To the Genetic Code, Int. J. Mod. Phys. B. 13 (1999) 2795–2885.
1053          https://doi.org/10.1142/S021797929900268X.

1054   [26]   M. Forger, S. Sachse, Lie superalgebras and the multiplet structure of the genetic code. I.
1055          Codon representations, J. Math. Phys. 41 (2000) 5407–5422.
1056          https://doi.org/10.1063/1.533417.

1057   [27]   K. Tamura, M. Nei, Estimation of the number of nucleotide substitutions in the control region
1058          of mitochondrial DNA in humans and chimpanzees, Mol. Biol. Evol. 10 (1993) 512–26.
1059          http://www.ncbi.nlm.nih.gov/pubmed/8336541.

1060   [28]   R. Sanchez, GenomAutomorphism: Compute the automorphisms between DNA's Abelian
1061          group representations. R package version 1.0.0, (2020).
1062          https://doi.org/10.18129/B9.bioc.GenomAutomorphism.

1063   [29]   L. Fuchs, Abelian groups, Publishing House of the Hungarian Academy of Sciences, 1958.

1064   [30]   L. Fuchs, Infinite Abelian Groups - Volume I, 1st Editio, Academic Press, 1970.
1065          https://doi.org/10.1007/978-3-319-19422-6.

1066   [31]   M. Yang, M.K. Derbyshire, R.A. Yamashita, A. Marchler-Bauer, NCBI's Conserved Domain
1067          Database and Tools for Protein Domain Analysis, Curr. Protoc. Bioinforma. 69 (2020) 1–25.
1068          https://doi.org/10.1002/cpbi.90.

1069   [32]   A.B. Stergachis, E. Haugen, A. Shafer, W. Fu, B. Vernot, A. Reynolds, A. Raubitschek, S.
1070          Ziegler, E.M. LeProust, J.M. Akey, J.A. Stamatoyannopoulos, Exonic transcription factor
1071          binding directs codon choice and affects protein evolution., Science (80-. ). 342 (2013) 1367–
1072          72. https://doi.org/10.1126/science.1243490.

1073   [33]   I. Reyna-Llorens, S.J. Burgess, G. Reeves, P. Singh, S.R. Stevenson, B.P. Williams, S.
1074          Stanley, J.M. Hibberd, Ancient duons may underpin spatial patterning of gene expression in

1075    C 4 leaves, Proc. Natl. Acad. Sci. U. S. A. 115 (2018) 1931–1936.
1076    https://doi.org/10.1073/pnas.1720576115.

1077 [34]  V.K. Yadav, K.S. Smith, C. Flinders, S.M. Mumenthaler, S. De, Significance of duon
1078    mutations in cancer genomes, Sci. Rep. 6 (2016) 1–9. https://doi.org/10.1038/srep27437.

1079 [35]  R. Sanchez, R. Grau, E.R. Morgado, R. Sánchez, L.A. Perfetti, R. Grau, E.R.M. Morales, A
1080    New DNA Sequence Vector Space on a Genetic Code Galois Field, MATCH Commun. Math.
1081    Comput. Chem. 54 (2005) 3–28.
1082    http://match.pmf.kg.ac.rs/electronic_versions/Match54/n1/match54n1_3-28.pdf.

1083 [36]  Y. Lv, J. Gu, H. Qiu, H. Li, Z. Zhang, S. Yin, Y. Mao, L. Kong, B. Liang, H. Jiang, C. Liu,
1084    Whole-exome sequencing identifies a donor splice-site variant in SMPX that causes rare X-
1085    linked congenital deafness, Mol. Genet. Genomic Med. 7 (2019) e967.
1086    https://doi.org/10.1002/mgg3.967.

1087 [37]  A. Panigrahi, B.W. O'Malley, Mechanisms of enhancer action: the known and the unknown,
1088    Genome Biol. 22 (2021) 1–30. https://doi.org/10.1186/s13059-021-02322-1.

1089 [38]  X. Sheng, J. Wang, Y. Guo, J. Zhang, J. Luo, DNA N6-Methyladenine (6mA) Modification
1090    Regulates Drug Resistance in Triple Negative Breast Cancer, Front. Oncol. 10 (2021) 1–6.
1091    https://doi.org/10.3389/fonc.2020.616098.

1092 [39]  Q. Lin, J. wei Chen, H. Yin, M. an Li, C. ren Zhou, T. fang Hao, T. Pan, C. Wu, Z. ran Li, D.
1093    Zhu, H. fan Wang, M. sheng Huang, DNA N6-methyladenine involvement and regulation of
1094    hepatocellular carcinoma development, Genomics. 114 (2022) 110265.
1095    https://doi.org/10.1016/j.ygeno.2022.01.002.

1096 [40]  M. Zamora, C.A. Ziegler, P.L. Freddolino, A.J. Wolfe, A Thermosensitive, Phase-Variable
1097    Epigenetic Switch: pap Revisited, Microbiol. Mol. Biol. Rev. 84 (2020).
1098    https://doi.org/10.1128/mmbr.00030-17.

1099 [41]  R. Sanchez, S.A. Mackenzie, Integrative Network Analysis of Differentially Methylated and
1100    Expressed Genes for Biomarker Identification in Leukemia, Sci. Rep. 10 (2020) 2123.
1101    https://doi.org/10.1038/s41598-020-58123-2.

1102 [42]  G. Ivanov, Generalizing the Baer-Kaplansky Theorem, J. Pure Appl. Algebr. 133 (1998) 107–
1103    115. https://doi.org/10.1016/S0022-4049(97)00187-4.

1104 [43]  P.W. Gunning, E.C. Hardeman, Isoforms: Fundamental differences, Elife. 7 (2018) 1–3.
1105    https://doi.org/10.7554/eLife.34477.

1106 [44]  L.P. Iñiguez, G. Hernández, The evolutionary relationship between alternative splicing and
1107    gene duplication, Front. Genet. 8 (2017) 1–7. https://doi.org/10.3389/fgene.2017.00014.

1108 [45]  F.H.C. Crick, The Origin of the Genetic Code, J. Mol. Biol. 38 (1968) 367–379.

1109 [46]  A. Comnish-Bowden, Nomenclature for incompletely specified bases in nucleic acid
1110    sequences: recommendations 1984, Nucleic Acids Res. 13 (1985) 3021–3030.

1111 [47] G.V.Kass, An Exploratory Technique for Investigating Large Quantities of Categorical Data,
1112 J. Roral Stat. Soc. 29 (1980) 119–127.

1113 [48] S. Miyazawa, R.L. Jernigan, Residue-Residue Potentials with a Favorable Contact Pair Term
1114 and an Unfavorable High Packing Density Term, for Simulation and Threading, J Mol Biol.
1115 256 (1996) 623–644. https://doi.org/10.1006/jmbi.1996.0114.

1116 [49] R. Sanchez, Evolutionary Analysis of DNA-Protein-Coding Regions Based on a Genetic Code
1117 Cube Metric, Curr. Top. Med. Chem. 14 (2014) 407–417.
1118 https://doi.org/10.2174/1568026613666131204110022.

1119 [50] P. Freyd, Abelian Categories: An Introduction to the Theory of Functors, Harper & Row, New
1120 York, 1964.

1121 [51] B. Mitchell, The Full Imbedding Theorem, Amer. J. Math. 86 (1964) 619.
1122 https://doi.org/10.2307/2373027.

1123 [52] Y. Maruyama, Category theory and foundations of life science: A structuralist perspective on
1124 cognition, BioSystems. 203 (2021) 1–20. https://doi.org/10.1016/j.biosystems.2021.104376.

1125 [53] M. Yu, B. Ren, The three-dimensional organization of mammalian genomes, Annu. Rev. Cell
1126 Dev. Biol. 33 (2017) 265–289. https://doi.org/10.1146/annurev-cellbio-100616-060531.

1127 [54] Y. Salameh, Y. Bejaoui, N. El Hajj, DNA Methylation Biomarkers in Aging and Age-Related
1128 Diseases, Front. Genet. 11 (2020) 1–11. https://doi.org/10.3389/fgene.2020.00171.

1129 [55] L. Fuchs, Infinite Abelian Groups, Volume 2, Academic Press, 1973.

1130 [56] J.-P. Marquis, Category Theory, in: E.N. Zalta (Ed.), {Stanford} Encycl. Philos., {F}all 202,
1131 Metaphysics Research Lab, Stanford University, 2021.

1132