

Accepted Manuscript

The “Wow! signal” of the terrestrial genetic code

Vladimir I. shCherbak, Maxim A. Makukov

PII: S0019-1035(13)00079-1

DOI: <http://dx.doi.org/10.1016/j.icarus.2013.02.017>

Reference: YICAR 10549

To appear in: *Icarus*

Received Date: 26 June 2012

Revised Date: 31 January 2013

Accepted Date: 12 February 2013



Please cite this article as: V.I. shCherbak, M.A. Makukov, The “Wow! signal” of the terrestrial genetic code, *Icarus* (2013), doi: <http://dx.doi.org/10.1016/j.icarus.2013.02.017>

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

Title:

The “Wow! signal” of the terrestrial genetic code

Authors:

Vladimir I. *sh*Cherbak ^a, Maxim A. Makukov ^b

Author affiliations:

^a Department of Mathematics, al-Farabi Kazakh National University, al-Farabi Avenue 71, 050038 Almaty,

Republic of Kazakhstan

e-mail: Vladimir.shCherbak@kaznu.kz

^b Fesenkov Astrophysical Institute, Observatory 23, 050020 Almaty, Republic of Kazakhstan

e-mail: makukov@aphi.kz

Corresponding author:

Maxim A. Makukov

e-mail: makukov@aphi.kz

Post address:

Fesenkov Astrophysical Institute, Observatory 23, 050020 Almaty, Republic of Kazakhstan

Phone:

Office: +7 727 260 7590

Mobile: +7 777 130 2255

Abstract

It has been repeatedly proposed to expand the scope for SETI, and one of the suggested alternatives to radio is the biological media. Genomic DNA is already used on Earth to store non-biological information. Though smaller in capacity, but stronger in noise immunity is the genetic code. The code is a flexible mapping between codons and amino acids, and this flexibility allows modifying the code artificially. But once fixed, the code might stay unchanged over cosmological timescales; in fact, it is the most durable construct known. Therefore it represents an exceptionally reliable storage for an intelligent signature, if that conforms to biological and thermodynamic requirements. As the actual scenario for the origin of terrestrial life is far from being settled, the proposal that it might have been seeded intentionally cannot be ruled out. A statistically strong intelligent-like “signal” in the genetic code is then a testable consequence of such scenario. Here we show that the terrestrial code displays a thorough precision-type orderliness matching the criteria to be considered an informational signal. Simple arrangements of the code reveal an ensemble of arithmetical and ideographical patterns of the same symbolic language. Accurate and systematic, these underlying patterns appear as a product of precision logic and nontrivial computing rather than of stochastic processes (the null hypothesis that they are due to chance coupled with presumable evolutionary pathways is rejected with P -value $< 10^{-13}$). The patterns display readily recognizable hallmarks of artificiality, among which are the symbol of zero, the privileged decimal syntax and semantical symmetries. Besides, extraction of the signal involves logically straightforward but abstract operations, making the patterns essentially irreducible to natural origin. Plausible ways of embedding the signal into the code and possible interpretation of its content are discussed. Overall, while the code is nearly optimized biologically, its limited capacity is used extremely efficiently to pass non-biological information.

1. Introduction

1 Recent biotech achievements make it possible to employ genomic DNA as data storage more durable
2 than any media currently used (Bancroft et al., 2001; Yachie et al., 2008; Ailenberg and Rotstein, 2009).
3 Perhaps the most direct application for that was proposed even before the advent of synthetic biology.
4 Considering alternative informational channels for SETI, Marx (1979) noted that genomes of living cells
5 may provide a good instance for that. He also noted that even more durable is the genetic code. Exposed to
6 strong negative selection, the code stays unchanged for billions of years, except for minor variations
7 (Knight et al., 2001) and context-dependent expansions (Yuan et al., 2010). And yet, the mapping between
8 codons and amino acids is malleable, as they interact via modifiable molecules of tRNAs and aminoacyl-
9 tRNA synthetases (Giegé et al., 1998; Ibba and Söll, 2000; see also *Appendix A*). This ability to reassign
10 codons, thought to underlie the evolution of the code to multilevel optimization (Bollenbach et al., 2007),
11 also allows to modify the code artificially (McClain and Foss, 1988; Budisa, 2006; Chin, 2012). It is
12 possible, at least in principle, to arrange a mapping that both conforms to functional requirements and
13 harbors a small message or a signature, allowed by 384 bits of informational capacity of the code. Once
14 genome is appropriately rewritten (Gibson et al., 2010), the new code with a signature will stay frozen in
15 the cell and its progeny, which might then be delivered through space and time to putative recipients. Being
16 energy-efficient (Rose and Wright, 2004) and self-replicating, the biological channel is also free from
17 problems peculiar to radio signals: there is no need to rely on time of arrival, frequency and direction.
18 Thus, due to these restrictions the origin of the famous “Wow!” signal received in 1977 remains uncertain
19 (Ehman, 2011). The biological channel has been given serious considerations for its merits in SETI, though
20 with the focus on genomes (Yokoo and Oshima, 1979; Freitas, 1983; Nakamura, 1986;
21 Davies, 2010, 2012).

22 Meanwhile, it has been proposed to secure terrestrial life by seeding exoplanets with living cells
23 (Mautner, 2000; Tepfer, 2008), and that seems to be a matter of time. The biological channel suggests itself
24 in this enterprise. To avoid anthropocentric bias, it might be admitted that terrestrial life is not the starting
25 point in the series of cosmic colonization (Crick and Orgel, 1973; Crick, 1981). If so, it is natural to expect

26 a statistically strong intelligent-like “signal” in the terrestrial genetic code (Marx, 1979). Such possibility is
27 incited further by the fact that how the code came to be apparently non-random and nearly optimized still
28 remains disputable and highly speculative (for reviews of traditional models on evolution of the code see
29 Knight et al., 1999; Gusev and Schulze-Makuch, 2004; Di Giulio, 2005; Koonin and Novozhilov, 2009).

30 The only way to extract a signal, if any, from the code is to arrange its elements – codons, amino acids
31 and syntactic signs – by their parameters using some straightforward logic. These arrangements are then
32 analyzed for patterns or grammar-like structures of some sort. The choice of arrangements and parameters
33 should exclude arbitrariness. For example, only those parameters should be considered which do not
34 depend on systems of physical units. However, even in this case *a priori* it is unknown exactly what kind
35 of patterns one might expect. So there is a risk of false positives, as with a data set like the genetic code it
36 is easy to find various patterns of one kind or another.

37 Nonetheless, the task might be somewhat alleviated. First, it is possible to predict some general aspects
38 of a putative signal and its “language”, especially if one takes advantage of active SETI experience. For
39 example, it is generally accepted that numerical language of arithmetic is the same for the entire universe
40 (Freudenthal, 1960; Minsky, 1985). Besides, symbols and grammar of this language, such as positional
41 numeral systems with zero conception, are hallmarks of intelligence. Thus, interstellar messages sent from
42 the Earth usually began with natural sequence of numbers in binary or decimal notation. To reinforce the
43 artificiality, a symbol of zero was placed in the abstract position preceding the sequence. Those messages
44 also included symbols of arithmetical operations, Egyptian triangle, DNA and other notions of human
45 consciousness (The Staff at the NAIC, 1975; Sagan et al., 1978; Dumas and Dutil, 2004).

46 Second, to minimize the risk of false positives one can impose requirements as restrictive as possible on
47 a putative signal. For example, it is reasonable to expect that a genuinely intelligent message would
48 represent not just a collection of patterns of various sorts, but patterns of the same “linguistic style”. In this
49 case, if a potential pattern is noticed, further search might be narrowed down to the same sort of patterns.
50 Another stringent requirement might be that patterns should involve each element of the code in each

51 arrangement, whereas the entire signal should occupy most, if not all, of the code's informational capacity.
52 By and large, given the nature of the task, specifics of the strategy are defined en route.

53 Following these lines, we show that the terrestrial code harbors an ensemble of precision-type patterns
54 matching the requirements mentioned above. Simple systematization of the code reveals a strong
55 informational signal comprising arithmetical and ideographical components. Remarkably, independent
56 patterns of the signal are all expressed in a common symbolic language. We show that the signal is
57 statistically significant, employs informational capacity of the code entirely, and is untraceable to natural
58 origin. The models of emergence of primordial life with original signal-free genetic code are beyond the
59 scope of this paper; whatever it was, the earlier state of the code is erased by palimpsest of the signal.

60 2. Background

61 Should there be a signal in the code, it would likely have manifested itself somehow during the half-
62 century history of traditional analysis of the code organization. So it is of use to summarize briefly what
63 has been learned about that up to date. Also, for the sake of simplicity in data presentation, we will mention
64 in advance some *a posteriori* information concerning the signal to be described, with fuller discussion in
65 due course. We suggest to a reader unfamiliar with molecular mechanisms behind the genetic code first to
66 refer to *Appendix A*, where it is also explained why the code is amenable to intentional "modulation" (to
67 use the language of radio-oriented SETI) and, at the same time, is highly protected from casual
68 "modulation" (i.e., has strong noise immunity).

69 2.1. The code at a glance

70 As soon as the genetic code was biochemically cracked (Nirenberg et al., 1965), its non-random
71 structure became evident (Woese, 1965; Crick, 1968). The most obvious pattern that emerged in the code
72 was its regular redundancy. The code comprises 16 codon families beginning with the same pair of bases,
73 and these families generally consist of either one or two equal series of codons mapped to one amino acid
74 or to *Stop* (Fig. 1a). In effect, the standard code is nearly symmetric in redundancy. There are only two
75 families split unequally: those beginning with TG and AT. The minimum action to restore the symmetry is

76 to match TG-family against AT-family by reassigning TGA from *Stop* to cysteine. Incidentally, this
77 symmetric version is not just a theoretical guess but is also found in nature as the nuclear code of euplotid
78 ciliates (Meyer et al., 1991). While the standard code stores the arithmetical component of the signal, the
79 symmetrical euplotid version keeps the ideographical one (the interrelation between these two code
80 versions is discussed later, see Section 4.2). Regular redundancy leads also to the block structure of the
81 genetic code. This makes it possible to depict the code in a contracted form, where each amino acid
82 corresponds to a single block, or a contracted series (Fig. 1b). The three exceptions are Arg, Leu and Ser,
83 which have one IV-series and one II-series each.

84 Apart from regular redundancy, a wealth of other features were reported afterwards, among which are
85 robustness to errors (Alff-Steinberger, 1969), correlation between thermostability and redundancy of codon
86 families (Lagerkvist, 1978), non-random distribution of amino acids among codons if judged by their
87 polarity and bulkiness (Jungck, 1978), biosynthetic pathways (Taylor and Coates, 1989), reactivity
88 (Siemion and Stefanowicz, 1992), and even taste (Zhuravlev, 2002). The code was also shown to be
89 effective at handling additional information in DNA (Baisnée et al., 2001; Itzkovitz and Alon, 2007).
90 Apparently, these features are related, if anything, to the direct biological function of the code. There are
91 also a number of abstract approaches to the code, such as those based on topology (Karasev and Stefanov,
92 2001), information science (Alvager et al., 1989), and number theory (Gonzalez, 2004). However, the main
93 focus of these approaches is in constructing theoretical model descriptions of known features in the code,
94 rather than dealing with new ones.

95 All in all, only two intrinsic regularities, observed early on in the study of the code, might suggest
96 possible relation to a putative signal due to their conspicuous and unambiguous character. They also
97 suggest two dimensionless integer parameters for signal extraction. These are quantity of codons in a series
98 mapped to one amino acid (redundancy) and quantity of nucleons in amino acid molecules. These
99 parameters might be called “ostensive numerals” by analogy with the quantity of radio beeps in *Lingua*
100 *Cosmica* (Freudenthal, 1960).

101 2.2. Rumer's bisection

102 Rumer (1966) bisected the code by redundancy – the first “ostensive numeral”. There are 8 whole
103 families and 8 split families in the code (Fig. 2a). Rumer found that codons in these families are mapped to
104 each other in a one-to-one fashion with a simple relation $T \leftrightarrow G$, $C \leftrightarrow A$, now known as Rumer’s
105 transformation. There are two more transformations of such type: $T \leftrightarrow C$, $A \leftrightarrow G$ and $T \leftrightarrow A$, $C \leftrightarrow G$. They
106 also appear in Rumer’s bisection and each makes half of what Rumer’s transformation makes alone.

107 Arbitrary bisection of the code has small chances to produce a transformation, and still less – their
108 ordered set (see *Appendix B*). Rumer’s finding was rediscovered by Danckwerts and Neubert (1975), who
109 also noted that this set might be described with a structure known in mathematics as the Klein-4 group.
110 That triggered a series of yet other models involving group theory to describe the code (Bertman and
111 Jungck, 1979; Hornos and Hornos, 1993; Bashford et al., 1998), which, admittedly, did not gain decisive
112 insights. Meanwhile, in traditional theories of the code evolution this accurate feature was ignored
113 altogether, though it was repeatedly rediscovered again (e.g., see Wilhelm and Nikolajewa, 2004).
114 Noteworthy, this regularity – which turns out to be a small portion of the signal – was first noticed
115 immediately after codon assignments were elucidated. Together with the fact of rediscoveries, this speaks
116 for the anticryptographic nature of the signal inside the code.

117 2.3. Amino acid nucleons

118 Hasegawa and Miyata (1980) arranged amino acids in order of increasing nucleon number – the second
119 “ostensive numeral” which, unlike other amino acid properties, does not rely on arbitrarily chosen system
120 of units. Such arrangement reveals a rough anticorrelation: the greater the redundancy the smaller the
121 nucleon number (Fig. 2b). This promoted speculations that prevailing small amino acids occupied the
122 series of higher redundancy during the code evolution. As shown below, this anticorrelation is a derivative
123 of the signal. Moreover, exactly this observation suggests simple systematization for both “ostensive
124 numerals”: monotonous arraying of nucleon and redundancy numbers in opposite directions.

125 On the whole, Hasegawa and Miyata dealt with amino acids whereas Rumer dealt with codons.
126 Combined, these approaches yield assignments between codons and amino acid nucleon numbers

127 convenient for systematization. *Stop*-codons code for no amino acid; therefore, to include them into the
128 systematization, they are assigned a zero nucleon number.

129 2.4. *The activation key*

130 All arithmetical patterns considered further appear with the differentiation between blocks and chains in
131 all 20 amino acids and with the subsequent transfer of one nucleon from side chain to block in proline
132 (Fig. 2b). Proline is the only exception from the general structure of amino acids: it holds its side chain
133 with two bonds and has one hydrogen less in its block. The mentioned transfer in proline “standardizes” its
134 block nucleon number to $73 + 1$ and reduces its chain nucleons to $42 - 1$. In itself, the distinction between
135 blocks and chains is purely formal: there is no stage in protein synthesis where amino acid side chains are
136 detached from standard blocks. Therefore, there is no any natural reason for nucleon transfer in proline; it
137 can be simulated only in the mind of a recipient to achieve the array of amino acids with uniform structure.
138 Such nucleon transfer thus appears artificial. However, exactly this seems to be its destination: it protects
139 the patterns from any natural explanation. Minimizing the chances for appealing to natural origin is a
140 distinct concern in messaging of such kind, and this problem seems to be solved perfectly for the signal in
141 the genetic code. Applied systematically without exceptions, the artificial transfer in proline enables
142 holistic and arithmetically precise order in the code. Thus, it acts as an “activation key”. While nature deals
143 with the actual proline which does not produce the signal in the code, an intelligent recipient easily finds
144 the key and reads messages in arithmetical language (see also Section 4.1).

145 2.5. *Decimalism*

146 The arithmetical patterns to be described hold true in any numeral system. However, as it turned out,
147 expressed in positional decimal system, they all acquire conspicuously distinctive notation. Therefore, here
148 we briefly provide some relevant information.

149 Nature is indifferent to numerical languages contrived by intelligence to represent quantities, including
150 zero. A privileged numeral system is therefore a reliable sign of artificiality. Intentionally embedded in an
151 object, a privileged system might then demonstrate itself through distinctive notation to any recipient

152 dealing with enumerable elements of that object. For example, digital symmetries of numbers divisible by
153 prime 037 exist only in the positional decimal system with zero conception (Fig. 3). Thus, distinctive
154 decimals 111, 222 and 333 look ordinarily 157, 336 and 515 in the octal system. This notational feature
155 was marked by Pacioli (1508) soon after the decimal system came to Europe. Analogous three-digit feature
156 exists in some other systems, including the quaternary one (see *Appendix C*).

157 3. Results

158 3.1. General structure of the signal

159 The overall structure of the signal is shown in Fig. 4, which might be used as guidance in further
160 description. The signal is composed of arithmetical and ideographical patterns, where arithmetical units are
161 represented by amino acid nucleons, whereas codon bases serve as ideographical entities. The patterns of
162 the signal are displayed in distinct logical arrangements of the code, thereby increasing both the
163 informational content of the signal and its statistical significance. Remarkably, all of the patterns bare the
164 same general style reflected in Fig. 4 with identical symbols in each signal component (represented by
165 boxes). Namely, distinct logical arrangements of the code and activation key produce exact equalities of
166 nucleon sums, which furthermore display decimalism and are accompanied by Rumer's and/or half-
167 transformations. One of these arrangements furthermore leads to ideography and semantical symmetries.
168 All elements of the code – 64 codons, 20 amino acids, *Start* and *Stop* syntactic signs – are involved in each
169 arrangement.

170 Unlike radio signals which unfold in time and thus have sequential structure, the signal in the genetic
171 code has no beginning and ending, similar to the pictorial message of Pioneer plaques (Sagan et al., 1972).
172 However, instead of providing pictograms the signal in the genetic code provides patterns that do not
173 depend on visual symbols chosen to represent them (be it symbols for nucleotide bases or for the notation
174 of “ostensive numerals”). These patterns make up the organic whole, so there is no unique order in
175 presenting them. We will begin with arithmetical component and then move on to ideography.

176 3.2. The arithmetical component

177 3.2.1. Full-size standard code

178 One logically plain arrangement of the code was proposed by George Gamow in his attempt to guess the
179 coding assignments theoretically before the code was cracked (see Hayes, 1998). He could not have known
180 in the fifties about the signal inside the code but one of his models, though it did not predict the actual
181 mapping correctly, coincided remarkably with one of the signal component. Gamow arranged codons
182 according to their composition, since 20 combinations of four bases taken three at a time could account for
183 20 amino acids (Gamow and Yčas, 1955). Application of the activation key and few “freezing” conditions
184 to this arrangement reveals total nucleon balancing ornate with decimal syntax.

185 Codons with identical and unique bases comprise two smaller sets (Fig. 5a). Halved, both sets show the
186 balance of side chains with $703 = 037 \times 19$ nucleons in each half as well as the balance of whole molecules
187 with $1665 = 666 + 999 \times 1$ nucleons. Importantly, the halving is not arbitrary. Codons are opposed by
188 Rumer’s transformation along with the half-transformation $T \leftrightarrow C$, $A \leftrightarrow G$ in the first set and $T \leftrightarrow A$, $C \leftrightarrow G$
189 in the second set. The *Spin* \rightarrow *Antispin* transformation does not affect the first set but finally freezes
190 elements of the second one. There is only one degree of freedom left since there are no reversible
191 transformations that might connect both sets, so one of them is free to swap around the axis. The balance
192 appears in one of the two alternative states.

193 The third set includes codons with two identical bases. When halved according to whether they are
194 purines or pyrimidines, regardless of the unique base type, this set shows the balance $999 = 999$ of side
195 chains (Fig. 5b). Besides, such halving keeps Rumer’s and one of the half-transformations again in place.
196 In its turn, the right half of the set is threefold balanced. Codons with adenine side by side, guanine side by
197 side and palindromic codons make up three equal parts with 333 nucleons each.

198 In Fig. 5c the same set is halved according to whether unique bases are purines or pyrimidines, this time
199 regardless of the identical bases type. Though not balanced, these halves again show distinctive decimal
200 syntax with 888 and $1110 = 111 + 999 \times 1$ nucleons. Decimalism of one of these sums is algebraically
201 dependent, as from the previous case (Fig. 5b) the sum of the whole set is known to be divisible by 037; if
202 a part of this set is decimally distinctive, the other one will be such automatically. Notably, an independent

203 pattern nonetheless stands out here. Namely, a part of the previous threefold balance has an equivalent in
 204 one half here, where the same amino acids are represented by synonymous codons (Fig. 5b and c). Whole
 205 molecules of this equivalent – 333 side chain and 444 standard block nucleons – are balanced with 777
 206 chain nucleons in the rest of the subset.

207 Note that all those distinctive notations of nucleon sums appear only in positional decimal system. The
 208 decimal notation is so customary in our culture that most of its users hardly remember a fairly complex rule
 209 behind it that encodes numbers as $a_{n-1} \times q^{n-1} + \dots + a_1 \times q^1 + a_0 \times q^0$, where $q = 10$, n is the quantity of digits
 210 in the notation, and a_i – digits 0-9 that are left in the final notation.

211 3.2.2. Decomposed standard code

212 Another arrangement of the code is brought about by decomposition of its 64 full-size codons. This
 213 yields 192 separate bases and reveals a pattern of the same type as in full-size format. Identical bases make
 214 up four sets of 48 bases in each. Each base retains the amino acid or *Stop* of its original codon (Fig. 6a).
 215 Thus, the four sets get their individual chain and block nucleon sums.

216 In total, there are $222 + 999 \times 10$ side chain nucleons in the decomposed code – obviously, thrice as
 217 much as the total sum in the previous full-size case (with the activation key still applied). Only one
 218 combination of the four sets displays distinctive decimalism of side chain nucleon sums. These are $666 +$
 219 999×2 nucleons in the T-set and $555 + 999 \times 7$ nucleons in the joint CGA-set (Fig. 6b). Meanwhile, there
 220 are exactly $222 + 999 \times 10$ block nucleons in the CGA-set (note that the sets have unequal block sums due
 221 to different accumulation of *Stops*). Thus, while chain nucleons are outnumbered by block nucleons overall
 222 the code, they are neatly balanced with their CGA-part.

223 3.2.3. Contracted code and the systematization rule

224 In a sense, contraction of codon series (see Fig. 1b) is an operation logically opposite to decomposition.
 225 Besides displaying new arithmetical patterns, contracted code also reveals ideographical component of the
 226 signal. The systematization rule leading to the ideography combines findings of Rumer (1966) and of
 227 Hasegawa and Miyata (1980) and is ostensive by itself (*shCherbak*, 1993). Contracted series are sorted into

228 four sets according to their redundancy; within those sets they are aligned side-by-side in order of
229 monotonously changing (e.g., increasing) nucleon number. The sets themselves are then arranged in
230 antisymmetrical fashion (e.g., in order of decreasing redundancy number). *Stop*-series is placed at the
231 beginning of its set representing zero in its special position. Finally, Rumer's bisection opposes the IV-set
232 to III, II, I sets. The resulting arrangement is shown in Fig. 7 for the euplotid code, with ideography of
233 codon bases (see next section) in Fig. 7a and arithmetical patterns of amino acids (shared by both code
234 versions) in Fig. 7b.

235 A new balance is found in the joint III, II, I set. Side chain nucleons of all its amino acids are equalized
236 with their standard blocks: $111 + 999 \times 1 = 111 + 999 \times 1$ (Fig. 7b). This pattern manifests as the
237 anticorrelation mentioned by Hasegawa and Miyata (1980). Chain nucleon sum of all series in the code is
238 less than the sum of all blocks. Only a subset of series coding mainly bigger amino acids may equalize its
239 own blocks. Exactly this happens in the joint III, II, I set. As a consequence, smaller amino acids are left in
240 the set of redundancy IV.

241 Meanwhile, there are 333 chain and 592 block nucleons and $333 + 592 = 925$ nucleons of whole
242 molecules in the IV-set. With 037 cancelled out, this leads to $3^2 + 4^2 = 5^2$ – numerical representation of the
243 Egyptian triangle, possibly as a symbol of two-dimensional space. Incidentally, codon series in the
244 ideogram (Fig. 7a) are arranged in the plane rather than linearly in a genomic fashion.

245 Rumer's bisection is based on redundancy and thus makes use of third positions in codon series.
246 Divisions of the contracted code based on first and center positions also reveal similar patterns (Fig. 8).
247 Another arithmetical phenomenon presumably related to the signal – the cytoplasmic balance – is described
248 in *Appendix D*.

249 Thus, the standard code reveals same-style and yet algebraically independent patterns simultaneously in
250 decomposed, full-size, and contracted representations (see Fig. 4). It is a highly nontrivial algebraic task to
251 find the solution that maps amino acids and syntactic signs to codons in a similar fashion. Normally this
252 would require considerable computational power.

253 3.3. The ideographical component

254 3.3.1. Upper strings

255 We refer to the product of systematization in Fig. 7a as the ideogram. The ideogram of the genetic code
256 is based on symmetries of its strings (shCherbak, 1988). The strings are read across contracted series.

257 The upper short string demonstrates *mirror*, *translation* and *inversion* symmetries (Fig. 9a). Its bases are
258 invariant under combined operation of the *mirror* symmetry and *inversion* of the type
259 base→complementary base. A minimum pattern of the *translation* symmetry is represented by *RRYY*
260 quadruplet.

261 The same three symmetries arrange the long upper string (Fig. 9b). The pair of flanking TATAT
262 sequences is *mirror* symmetrical. The pair of central AGC codons forms a minimum pattern of the
263 *translation* symmetry. First and third bases in the set of redundancy II are interconnected in an
264 axisymmetric manner with purine↔pyrimidine *inversion* and its opposite operation – the unit
265 transformation producing no exchange.

266 3.3.2. Center strings

267 Placed coaxially, the short and the long center strings appear interconnected with purine↔pyrimidine
268 *inversion* (Fig. 10a). Both strings exhibit purine-pyrimidine *mirror* symmetry. The long string keeps the
269 mirror symmetry even for ordinary bases.

270 Codons of the short string CCC and TCT break the mirror symmetry of ordinary bases, but they share a
271 palindromic feature, i.e. direction of reading invariance. This feature restores the mirror symmetry, this
272 time of the *semantical* type (Fig. 10b). As in the previous case, two center strings are expected to share the
273 same set of symmetries. Therefore, the semantical symmetry of palindromic codons flanked by G-bases
274 may indicate a similar feature in the long string. Indeed, semantical symmetry is found there in the triplet
275 reading frame starting after flanking G-base (Fig. 10c). This reading frame is remarkable with the regular
276 arrangement of all syntactic signs of the euplotid code – both *Stop*-codons and the *Start*-codon repeated
277 twice. The reading frame displays the *semantical mirror* symmetry of antonyms with homogeneous AAA-
278 codon in the center.

279 The codons of this reading frame are purely abstract symbols, given that they are read across contracted
280 series. However, they are regularly crossed with the same codons in the ideogram, thereby reinforcing the
281 semantical symmetry and making the current frame unique (Fig 10c). Besides, direction of reading now
282 becomes distinguished since such “crossword” disappears if read in opposite way, though the palindrome
283 itself remains the same.

284 Remarkably, the triplet string in Fig. 10c is written with the code symbols within the code itself. This
285 implies that the signal-harboring mapping had to be projected preliminarily (see Section 4.3 in *Discussion*).
286 Besides, translation of this string with the code itself reveals the balance $222 = 222$ of chains and blocks
287 (Fig. 10d). Additional palindrome in the frame shifted by one position (Fig. 10e) reproduces the chain sum
288 of 222, confirming that the ideogram is properly “tuned in” to the euploid version: TGA stands for Cys
289 here, not for *Stop* of the standard code.

290 4. Discussion

291 4.1. Artificiality

292 To be considered unambiguously as an intelligent signal, any patterns in the code must satisfy the
293 following two criteria: (1) they must be highly significant statistically and (2) not only must they possess
294 intelligent-like features (Elliott, 2010), but they should be inconsistent in principle with any natural
295 process, be it Darwinian (Freeland, 2002) or Lamarckian (Vetsigian et al., 2006) evolution, driven by
296 amino acid biosynthesis (Wong, 2005), genomic changes (Sella and Ardell, 2006), affinities between
297 (anti)codons and amino acids (Yarus et al., 2009), selection for the increased diversity of proteins (Higgs,
298 2009), energetics of codon-anticodon interactions (Klump, 2006; Travers, 2006), or various pre-
299 translational mechanisms (Wolf and Koonin, 2007; Rodin et al., 2011).

300 The statistical test for the first criterion is outlined in *Appendix B*, showing that the described patterns
301 are highly significant. The second criterion might seem unverifiable, as the patterns may result from a
302 natural process currently unknown. But this criterion is equivalent to asking if it is possible at all to embed
303 informational patterns into the code so that they could be unequivocally interpreted as an intelligent
304 signature. The answer seems to be yes, and one way to do so is to make patterns virtual, not actual. Exactly

305 that is observed in the genetic code. Strict balances and their decimal syntax appear only with the
306 application of the “activation key”. Physically, there are no strict balances in the code (e.g., in Fig. 5b one
307 would have $1002 \neq 999$ instead of $999 = 999$). Artificial transfer of a nucleon in proline turns the balances
308 on and thereby makes them virtual. This is also the reason why we interpret distinctive notation as an
309 indication of decimalism, rather than as a physical requirement (yet unknown) for nucleon sums to be
310 multiples of 037: in general, physically there is no such multiplicity in the code. In its turn, notationally
311 preferred numeral system is by itself a strong sign of artificiality. It is also worth noting that all three-digit
312 decimals – 111, 222, 333, 444, 555, 666, 777, 888, 999 (as well as zero, see below) – are represented at
313 least once in the signal, which also looks like an intentional feature.

314 However, it might be hypothesized that amino acid mass is driven by selection (or any other natural
315 process) to be distributed in the code in a particular way leading to approximate mass equalities and thus
316 making strict nucleon balances just a likely epiphenomenon. But it is hardly imaginable how a natural
317 process can drive mass distribution in abstract representations of the code where codons are decomposed
318 into bases or contracted by redundancy. Besides, nucleon equalities hold true for free amino acids, and yet
319 in these free molecules side chains and standard blocks had to be treated by that process separately.
320 Furthermore, no natural process can drive mass distribution to produce the balance in Fig. 10d: amino acids
321 and syntactic signs that make up this balance are entirely abstract since they are produced by translation of
322 a string read across codons.

323 Another way to make patterns irreducible to natural events is to involve semantics, since no natural
324 process is capable of interpreting abstract symbols. It should be noted that notions of symbols and
325 meanings are used sometimes in a natural sense (Eigen and Winkler, 1983), especially in the context of
326 biosemiotics (Barbieri, 2008) and molecular codes (Tlustý, 2010). The genetic code itself is regarded there
327 as a “natural convention” that relates symbols (codons) to their meanings (amino acids). However, these
328 approaches make distinction between organic semantics of molecular codes and interpretive or linguistic
329 semantics peculiar to intelligence (Barbieri, 2008). Exactly the latter type of semantics is revealed in the
330 signal of the genetic code. It is displayed there not only in the symmetry of antonymous syntactic signs

(Fig. 10c), but also in the symbol of zero. For genetic molecular machinery there is no zero, there are nucleotide triplets recognized sterically by release factors at the ribosome. Zero – the supreme abstraction of arithmetic – is the interpretive meaning assigned to *Stop*-codons, and its correctness is confirmed by the fact that, being placed in its proper front position, zero maintains all ideogram symmetries. Thus, a trivial summand in balances, zero, however, appears as an *ordinal* number in the ideogram. In other words, besides being an integral part of the decimal system, zero acts also as an individual symbol in the code.

In total, not only the signal itself reveals intelligent-like features – strict nucleon equalities, their distinctive decimal notation, logical transformations accompanying the equalities, the symbol of zero and semantical symmetries, but the very method of its extraction involves abstract operations – consideration of idealized (free and unmodified) molecules, distinction between their blocks and chains, the activation key, contraction and decomposition of codons. We find that taken together all these aspects point at artificial nature of the patterns.

Though the decimal system in the signal might seem a serendipitous coincidence, there are few possible explanations, from ten-digit anatomy as an evolutionary near-optimum for bilateral beings (Dennett, 1996) to the fact that there are conveniently $74 = 2 \times 37$ nucleons in the standard blocks of α -amino acids. Besides, the decimal system shares the triplet digital symmetry with the quaternary one (see *Appendix C*), establishing a link to the “native” language of DNA. After all, some of the messages sent from the Earth included the decimal system as well (Sagan et al., 1978; Dumas and Dutil, 2004), though they were not supposed to be received necessarily by ten-digit extraterrestrials. Whatever the actual reason behind the decimal system in the code, it appears that it was invented outside the solar system already several billions years ago.

4.2. Two versions of the code

The nearly symmetric code version with arithmetical patterns acts as the universal standard code. With this code at hand it is intuitively easy to infer the symmetric version with its ideography. Vice versa, if the symmetric version were the universal one, it would be hardly possible to infer the nearly symmetric code with all its arithmetical patterns. Therefore, with the standard version alone it is possible to “receive” both

357 arithmetical and ideographical components of the signal, even if the symmetric version was not found in
358 nature. There are two possible reasons why it is actually found in euplotid ciliates: either originally when
359 Earth was seeded there were both versions of the code with one of them remaining currently in euplotid
360 ciliates, or originally there was only the standard version, and later casual modification in euplotid lineage
361 coincided with the symmetric version.

362 What concerns other known versions of the code, they seem neither to have profound pattern ensembles,
363 nor to be easily inferable from the standard code. Most probably they represent casual deviations caused by
364 ambiguous intermediates or codon captures (Moura et al., 2010).

365 4.3. *Embedding the signal*

366 To obtain a code with a signature one might search through all variant mappings and select the “most
367 interesting” one. However, this method is unpractical (at least with the present-day terrestrial computing
368 facilities), given the astronomically huge number of variant codes. In a more realistic alternative, the
369 pattern ensemble of the signal is projected preliminarily as a system of algebraic expressions which is then
370 solved relatively easily to deduce the mapping of the code. Thus, all described patterns might be
371 represented *post factum* as a system of Diophantine equations and inequalities, and numerical analysis of
372 this system shows that it uniquely determines the mapping between codon series and nucleon numbers,
373 including zeros for *Stop*-codons (see *Appendix E*). Though some amino acids have equal nucleon numbers,
374 as the case for Leu and Ile, or Lys and Gln, even they are not interchangeable, as suggested by distinctive
375 notation of nucleon sums in β , γ and other positional levels of side chains in the contracted code (Figs. 7b
376 and 8a). The activation key applies here as well (note that β - and δ -carbons in proline are positionally
377 equivalent). The standard chemical nomenclature of carbon atoms is extended here to denote positions of
378 other nodal atoms. Decimalism in different combinations of levels circumvents algebraic dependence and
379 defines chemical structure of amino acids more rigidly.

380 These patterns within side chains go even deeper into chemical structure. Some of the canonical amino
381 acids – His, Arg and Trp – might exist in alternative neutral tautomeric forms differing in the position of
382 one hydrogen atom in their side chains (Taniguchi and Hino, 1981; Rak et al., 2001; Li and Hong, 2011).

383 Though some of these tautomers occur very rarely at cytoplasmic pH (as the case for indolenine tautomer
384 of Trp shown in Fig. 7b), all neutral tautomers are legitimate if idealized free molecules are considered, and
385 taking only one of them would introduce arbitrariness. Notably, however, that while one Trp tautomer
386 maintains the patterns in Fig. 7b, another one does the job in Fig. 8a, whereas any neutral tautomer of His
387 and Arg might be taken in both cases without affecting the patterns at all (which is easily checked; to this
388 end, both Arg tautomers are shown in Fig. 8a and both His tautomers are shown in Figs. 7b and 8a).

389 Importantly, preliminary projecting of a signal admits imposition of functional requirements as extra
390 formal conditions. The terrestrial code is known to be conservative with respect to polar requirement
391 (Freeland and Hurst, 1998), but not to molecular size (Haig and Hurst, 1991). The signal in the code does
392 not involve polar requirement as such, so it might be used in a parallel formal condition to reduce effect of
393 misreadings. However, the signal does involve nucleon numbers which correlate with molecular volume.
394 That interferes with an attempt to make the code conservative with respect to size of amino acids as well.

395 4.4. Possible interpretation

396 Besides having the function of an intelligent signature as such, the signal in the genetic code might also
397 admit sensible interpretations of its content. Without claim to be correct, here we propose our own version.
398 It is now tempting to think that the main body of the message might reside in genomes (Marx, 1979; see
399 also Hoch and Losick, 1997). Though the idea of genomic SETI (Davies, 2010) might seem naïve in view
400 of random mutations, things are not so obvious. For example, a locus with a message might be exposed to
401 purifying selection through coupling to essential genes, and there is even possible evidence for that (*ibid.*).
402 Whatever the case, the ideogram does seem to provide a reference to genomes. Thus, complementary
403 mirror-symmetrical bases of the short upper string (Fig. 9a) resemble Watson-Crick pairs; the four central
404 bases TC|GA and the central axis therefore possibly represent the symbol of the genomic DNA itself.
405 Flanking TATAT bases (Fig. 9b) might symbolize consensus sequence found in promoters of most genes.
406 Coding sequences of genes are located between *Start*- and *Stop*-codons. Vice versa, nontranslated regions
407 are found between *Stop*- and *Start*-codons of neighbor genes. Therefore the triplet string in Fig. 10c might
408 symbolize intergenic regions, and may be interpreted as the address of the genomic message.

409 The privileged numeral system in the code might also be interpreted as an indication of a similar feature
410 in genomes. It is often said that genomes store hereditary information in quaternary digital format. There
411 are 24 possible numberings of DNA nucleotides with digits 0, 1, 2, 3. The ideogram seems to suggest the
412 proper one: T \equiv 0, C \equiv 1, G \equiv 2, A \equiv 3. In this case the TCGA quadruplet (Fig. 9a), read in the
413 distinguished direction, represents the natural sequence preceded by zero. Palindromic codons CCC and
414 TCT (Fig. 10b) become a symbol of the quaternary digital symmetry 111_4 and the radix of the
415 corresponding system $010_4 = 4$, respectively. Translationally related AGC, or 321_4 , codons (Fig. 9b)
416 possibly indicate positions in quaternary place-value notation, with higher orders coming first. The sum of
417 digital triplets in the string TAG + TAA + AAA + ATG + ATG (Fig. 10c) equals to the number of
418 nucleotides in the code $3000_4 = 192$. Besides, T as zero is opposed to the other three “digits” in the
419 decomposed code (Fig. 6). Finally, each complementary base pair in DNA sums to 3, so the double helix
420 looks numerically as $333\dots_4$, and the central AAA codon in Fig. 10c becomes the symbol of duplex DNA
421 located between genes. Should this particular numbering have relation to the genomic message, if any, is a
422 matter of further research.

423 It is worth mentioning that all genomes, despite their huge size and diversity, do possess a feature as
424 universal as the genetic code itself. It is known as the second Chargaff’s rule. In almost all genomes – from
425 viral to human – the quantities of complementary nucleotides, dinucleotides and higher oligonucleotides up
426 to the length of ~ 9 are balanced to a good precision within a *single* DNA strand (Okamura et al., 2007).
427 Unlike the first Chargaff’s rule which quickly found its physicochemical basis, the second rule with its
428 total orderliness still has no obvious explanation.

429 **Appendix A. Molecular implementation of the genetic code**

430 Here we outline molecular workings behind the genetic code which explain why it stays unchanged for
431 billions of years and, at the same time, might be readily modified artificially, e.g., for embedding a signal.
432 For simplicity, we skip the details such as U instead of T in RNA, ATP energetics, wobble pairing, etc.,
433 that do not affect understanding of the main point (for details see, e.g., Alberts et al., 2008).

434 The first type of molecule behind the genetic code is transfer RNA (tRNA). They deliver amino acids
435 into ribosomes, where protein synthesis takes place. tRNAs are transcribed as a final product from tRNA
436 genes in genomes by RNA polymerase (Fig. A.1a; for definiteness, the mechanism is shown for amino acid
437 Ser and its TCC codon). With the length varying around 80 nucleotides, tRNA transcripts fold in a specific
438 spatial configuration due to base-pairing between different sections of the same RNA strand, similar to as it
439 occurs between two strands of DNA helix (Fig. A.1b). At its opposite sides the folded tRNA molecule has
440 an unpaired anticodon and the acceptor end to which amino acid is to be bound. tRNAs with differing
441 anticodons specifying the same amino acid (remember the code is redundant) are identical in their overall
442 configuration. tRNAs specifying distinct amino acids differ from each other in anticodons as well as other
443 spots, so they have slightly different overall configurations. However, acceptor ends are identical in all
444 tRNAs, so for tRNA itself it makes no difference which amino acid is bound to it, no matter which
445 anticodon it has at the opposite side. The process of binding amino acids to tRNAs is performed by protein
446 enzymes called aminoacyl-tRNA synthetases (aaRSs, Fig A.1b, bottom). Normally, there are 20 types of
447 aaRSs, one for each amino acid, and they themselves are translated from appropriate genes in genome.
448 Each of these enzymes recognizes with great specificity both its cognate amino acid and all tRNAs with
449 anticodons specifying that amino acid; however, tRNAs are recognized primarily by their overall
450 configuration (Fig. A.1c). After binding and additional checking, aaRS releases tRNA charged with amino
451 acid to be delivered to ribosome (Fig. A.1d). In its turn, the ribosome does not care if tRNA carries an
452 amino acid specified by its anticodon; it only checks if the anticodon of tRNA matches complementarily
453 the current codon in messenger RNA (mRNA; Fig. A.1e). If so, the amino acid is transferred from tRNA to
454 the growing peptide chain and tRNA is released to be recycled. If codon and anticodon do not match,
455 tRNA with its amino acid is dislodged from the ribosome to be used later until it matches codon on mRNA
456 (even with this overshoot the bacterial ribosome manages to add ~20 amino acids per second to a peptide
457 chain). The described mechanism results in relationships between mRNA codons and amino acids (Fig.
458 A.1f) which, collected together in any convenient form (one possibility is shown in Fig. 1a), constitute the
459 genetic code.

460 The key point in terms of changeability of the genetic code is that there is no direct chemical interaction
461 between mRNA codons and amino acids at any stage. They interact via molecules of tRNA and aaRS both
462 of which might be modified so that a codon is reassigned to another amino acid. As an example,
463 Figures A.1g-k show a simple way of changing the code where two amino acids – Ser and Ala –
464 interchange two of their codons. It is known that in most organisms tRNA anticodons are not involved in
465 recognition by aaRSs cognate for these amino acids (Giegé et al., 1998; the fact reflected in Fig. A.1c with
466 SARS not touching the anticodon). Therefore, the three nucleotides in tRNA^{Ser} gene corresponding to
467 anticodon might be replaced (Fig. A.1g), in particular, to get GGC anticodon corresponding to GCC codon
468 in mRNA, which normally codes Ala. (To get anticodon for a codon, or vice versa, one has to apply
469 complementarity rule and reverse the resulting triplet, since complementary DNA/RNA strands have
470 opposite directionalities). After that, SARS will still bind Ser to tRNA^{Ser}, even though it now has new GGC
471 anticodon (Fig. A.1h). If analogous procedure is performed with tRNA^{Ala} genes to produce tRNA^{Ala} with
472 GGA anticodon, the genetic code would be modified: Ser and Ala would have interchanged some of their
473 codons (actually, two codons, due to wobble pairing). However, the cell will not survive such surgery,
474 since all coding genes in genome remain “written” with the previous code and after translation with the
475 new code they all produce non- or at best semi-functional proteins, with Ala occasionally replaced by Ser
476 and vice versa. To fix the new code in a cell lineage, one also has to change coding mRNAs appropriately
477 to leave amino acid sequences of coded proteins unaltered (Fig. A.1i). That would be automatically
478 fulfilled if all coding genes are rewritten all over the genome so that TCC codons are replaced with GCC
479 and vice versa (Fig. A.1j); such operation is possible when genomes are even rewritten from scratch
480 (Gibson et al., 2010). Now, amino acid sequences of proteins stay unaltered and a cell proliferates with the
481 new genetic code (Fig. A.1k).

482 It must be clear now why the genetic code is highly protected from casual modifications. If a mutation
483 occurs in tRNA or aaRS leading to codon reassignment, all genes in genome remain written with the
484 previous code, and a cell quickly goes off the scene without progeny. The chances that such mutation in
485 tRNA/aaRS is accompanied by corresponding mutations in coding genes all over the genome resulting in

unaltered proteins are vanishingly small, given that there are dozens of such codons in thousands of genes in a genome. Thus, the machinery of the genetic code experiences exceptionally strong purifying selection that keeps it unchanged over billions of years.

It should be reminded that in reality the process of intentional modification of the code is more complicated. For example, details of tRNA recognition by aaRSs vary depending on tRNA species and organism, and in some cases anticodon is involved, partially or entirely, in that process. However, this is avoidable, in principle, with appropriate methods of molecular engineering. Another issue is that modifications in the code that leave proteins unaltered still might affect the level of gene expression (Kudla et al., 2009). Therefore, additional measures might have to be taken to restore the expression pattern with the new genetic code. These are surmountable technical issues; the point is that there are no principal restrictions for changing the code artificially in any desired way. In effect, elaborate methods of modifying the overall tRNA configuration and/or aaRS recognition sites might allow not only interchanging two amino acids, but introducing new ones.

Appendix B. Statistical test

It is appropriate to ask if the presented patterns are merely an artifact of data fishing. To assess that, one might compare information volumes of the data set itself (V_0) and of the pattern ensemble within that set (V_p). The artifact of data fishing might then be defined as the case when $V_p \ll V_0$. As shown in *Appendix E*, the presented ensemble of patterns might be described with a system of Diophantine equations, where nucleon numbers of amino acids serve as unknowns. Given the set of canonical amino acids (the range of possible values for the unknowns), this system is completely defined: it has a single solution and that turns out to be the actual mapping of the code (this also implies that there are no more algebraically independent patterns of the same sort in the code). Hence, $V_p = V_0$, so the pattern ensemble employs informational capacity of the code entirely, making the assumption of data fishing artifact irrelevant.

One might ask then how likely such pattern ensemble is to appear in the genetic code by chance. Since this question implies that the current mapping of the code has been shaped by natural processes, it is more appropriate to ask how likely such pattern ensemble is to appear by chance under certain conditions

512 reflecting presumable evolutionary pathways. We tested both versions of the null hypothesis (“the patterns
513 are due to chance alone” and “the patterns are due to chance coupled with presumable evolutionary
514 pathways”). The results are of the same order of magnitude; we describe only the version with presumable
515 natural conditions. Three such conditions reflecting predominant speculations on the code evolution were
516 imposed on computer-generated codes in this test:

517 (1) Redundancy must be on average similar to that of the real code. This is thought to be due to the
518 specifics of interaction between the ribosome, mRNA and tRNA (Novozhilov et al., 2007). Besides, we
519 took into account possible dependence of the probability for a codon family to stay whole or to be split on
520 the type of its first two bases. This follows from the difference in thermostability between codon-anticodon
521 pairs enriched with strong (G and C) bases and those enriched with weak (A and T) bases (Lagerkvist,
522 1978). For that, the probability for a family of four codons with leading strong doublets to specify a single
523 amino acid is adopted to be 0.9, for those with weak doublets – 0.1, and for mixed doublets it is 0.5. Each
524 of the 20 amino acids and *Stop* is recruited at least once; therefore codes with less than 21 generated blocks
525 are discarded. After that blocks are populated randomly with amino acids and *Stop*.

526 (2) Reduced effect of mutations/mistranslations due to natural selection. The cost function for polar
527 requirement was adopted from Freeland and Hurst (1998), taking into account transversion-transition and
528 mistranslation biases (see also Novozhilov et al., 2007). Only those codes are passed further which have
529 cost function value smaller than $\varphi_0 + \sigma$, where φ_0 is the value for the universal code, and σ is the standard
530 deviation for all random codes filtered through the previous condition.

531 (3) Small departure from the cytoplasmic balance (see *Appendix D*). As argued by Downes and
532 Richardson (2002), this balance might reflect evolutionary pathways optimizing the distribution of mass in
533 proteins. With C standing for all side chain nucleons in the code and B for all nucleons in block residues,
534 the value $\delta = (C - B)/(C + B)$ is distributed approximately normally with $\mu = 0.043$ and $\sigma = 0.024$ (under
535 the first condition described above). Only those codes were considered which had δ in the range $0 \pm \sigma$,
536 centered on the value of the standard code. As that range corresponds to codes with smaller (“early”) amino

537 acids predominating, this condition also reflects presumable history of the code expansion (Trifonov, 2000;
538 Wong, 2005).

539 The random variable in question is the number of independent patterns of the same sort in a code.
540 Obviously, the more such patterns are observed in a code, the less likely such observation is. Probably, a
541 good approximation here would be a binomial distribution since, for example, a nucleon balance might be
542 regarded as a Bernoulli trial: in a given arrangement the balance is either “on” or “off”, where probability
543 for “on” is much smaller than for “off”. However, probabilities for balances in distinct arrangements might
544 differ, especially under conditions imposed. Situation is even more complex with ideogram symmetries:
545 symmetry is not just “on” or “off”, it is also characterized by the length of the string and the number of
546 nucleotide types involved. Therefore, we do not apply any approximations but use brute-force approach to
547 find distributions for appropriately defined scores for the patterns. Proline was considered with one nucleon
548 transferred from its side chain to its block (note that since the activation key is applied universally, the
549 actual code and the code with the key applied are equivalent statistically).

550 ***Nucleon balances.*** Arithmetical patterns of the standard code are all of the same style: equality of nucleon
551 sums + their distinctive decimal notation + at least one of the three transformations (except the decomposed
552 case). The search for a random code with a few patterns of this sort (not to mention that they should form
553 an algebraically defined system) turned out to be time-consuming, so the requirements were simplified.
554 Only nucleon equalities were considered, without requirement of distinctive notation in any numeral
555 system. Presence of transformations was required only in Gamow’s arrangement for codons with identical
556 and unique bases, since transformations act there in the first place, not as companions of another sorting
557 logic. Also for simplicity, only global patterns were considered; “local” features like the threefold balance
558 in Fig. 5b were not checked.

559 Alternative codes might have balances in arrangements and combinations different from those in the real
560 code. Contrary to as it might seem, there are not so many ways of arranging the code that fix its elements
561 unambiguously. For example, along with Gamow’s sorting, several other arrangements were proposed
562 during early attempts to deduce the code theoretically (see Hayes, 1998). One of them is known as the

563 “code without commas” (Crick et al., 1957). However, unlike Gamow’s sorting, this and other proposed
564 arrangements do not allow “freezing” the code elements completely, leaving a large degree of arbitrariness.

565 Ultimately, the following arrangements were considered in the test:

- 566 - divisions based on redundancy;
- 567 - divisions based on positions in codons (alternating all combinations such as S or W in the first position,
568 R or Y in the second position, etc.);
- 569 - sortings based on nucleotide composition of codons (alternating all combinations of “freezing”
570 conditions and division logic);
- 571 - arrangements based on decomposition of codons into bases (alternating all combinations of the four
572 nucleotide sets).

573 Besides, the first two types might be arranged with full-size or contracted codons. The only possible
574 balance of the peptide representation (*Appendix D*) was also checked. In total, 160 potential balances (of
575 both chain-chain and block-chain types) were checked in all these arrangements. Precautions were made to
576 ignore arithmetical dependencies, as for certain code versions some balances are trivially fulfilled if few
577 others occur. A simple scoring scheme was adopted: the score of a code is the number of algebraically
578 independent nucleon equalities it happens to possess in all arrangements. In this scheme the simplified
579 version of arithmetical patterns in the standard code has the score 7. Computer estimation shows that
580 probability for a code to have the score not less than 7 by chance under imposed conditions is $p_1 = 1.5 \times 10^{-8}$
581 (Fig. B.1a).

582 **Ideogram symmetries.** An ideogram might be built for each variant code in the same way as shown in Fig.
583 7 (however, no requirement is made for whole and split families to be linked with any transformation).
584 There are a few more conceivable ways to build an ideogram using contracted codon series (ideograms
585 based on full-size codons suffer with ambiguities). For example, nucleon and redundancy numbers might
586 be arranged in the same direction, rather than antisymmetrically. Another way is to divide the code by
587 positions in codons (e.g., R or Y in the first position; though these ideograms are simpler as two of their
588 four upper strings are always binary, whereas in ideograms based on redundancy all strings are, in general,

589 quaternary). In total, 9 ideogram versions were built for each code and checked for symmetries. Namely,
 590 each of the four strings was checked for \mathcal{M} , $\mathcal{M} + I$, \mathcal{T} , $\mathcal{T} + I$, where \mathcal{M} and \mathcal{T} stand for mirror and translation
 591 symmetries and I denotes pair inversions of all three types. For each symmetry a string of length L gets the
 592 score $L/2$, if it contains only two types of bases (or if the symmetry holds only in binary representation RY,
 593 SW or KM), and L , if it contains three or all four types of bases. Only whole-string symmetries were
 594 considered (in this case multiple symmetries organizing different parts of a string such as in Fig. 9b are not
 595 detected; the whole string in Fig. 9b, however, is mirror symmetrical in KM representation). For each
 596 ambiguous position (two neighboring series with equal nucleon numbers) the penalty $L/3$ was introduced.
 597 Semantical symmetries and balances of translated amino acids were not checked. Finally, if at least one of
 598 the four strings has none of the symmetries, the score is divided by 2. The euplotid code has the score 35 in
 599 this scheme: 8 for $\mathcal{M} + I_{(\text{T} \leftrightarrow \text{A}, \text{C} \leftrightarrow \text{G})}$ and 4 for \mathcal{T}_{RY} in the upper short string, 4 for \mathcal{M}_{RY} in the center short
 600 string, 8 for \mathcal{M}_{KM} in the upper long string, 16 for \mathcal{M} in the center long string, penalty $-16/3 \approx -5$ for Lys and
 601 Gln (though in this case their interchange affects neither \mathcal{M}_{KM} in the upper string, nor \mathcal{M} in the center one).
 602 Computer estimation shows that probability for a code to have the score not less than 35 by chance under
 603 imposed conditions is $p_2 = 9.4 \times 10^{-5}$ (Fig. B.1b).

604 We also checked transformations in Rumer's bisections of generated codes, since these transformations
 605 served as the guiding principle for signal extraction in the real code. Under the conditions imposed,
 606 probability for a random code to have equal numbers of whole and split families which are furthermore
 607 linked with any of the three possible transformations was found to be 4.6×10^{-2} . Given that one
 608 transformation takes places, the other two might be distributed among codons in the ratios 8:0 ($p = 0.125$),
 609 4:4 ($p = 0.375$), or 2:6 ($p = 0.5$). For the real code this ratio is 4:4 (see Fig. 2a), so finally $p_3 = 1.7 \times 10^{-2}$.

610 As suggested by a separate computational study, mutual influence of the three types of patterns is
 611 negligible, so the total probability for a (very simplified) signal to occur by chance in a single code under
 612 imposed conditions is $p_1 p_2 p_3 = 2.4 \times 10^{-14}$. Since the redundancy-symmetric code is not even needed to be
 613 found in nature to reveal the ideogram, the final P -value will not differ much from that value.

614 This result gives probabilities for the specific type of patterns – nucleon equalities and ideogram
 615 symmetries. However, testing the hypothesis of an intelligent signal should take into account patterns of
 616 other sorts as well, as long as they meet the requirements outlined in *Introduction*. After analysis of the
 617 literature on the genetic code our opinion is still that nucleon and redundancy numbers are the best
 618 candidates for “ostensive numerals”. We admit though that there could be other possibilities and that the
 619 obtained *P*-value should be regarded as a very rough approximation (keep in mind simplifications in the
 620 test as well). But admittedly, there just cannot be enough candidates for “ostensive numerals” and
 621 corresponding logical arrangements to compensate for the small *P*-value obtained and to raise it close to
 622 the significance level.

623 **Appendix C. Digital symmetries of positional numeral systems**

624 The digital symmetry described in the main text for the decimal system is related to a divisibility
 625 criterion that might be used to effectively perform checksums. Consider the number 27014319417 as an
 626 example. Triplet reading frame splits this number into digital triplets 270, 143, 194, 170 (any of the three
 627 reading frames might be chosen; zeros are added at flanks to form complete triplets). The sum of these
 628 triplets equals to 777. Its distinctive notation indicates that the original number is divisible by 037. In four-
 629 digit numbers that appear during summations thousand’s digits are transferred to unit’s digits. If notation of
 630 the resulting sum is not distinctive, add or subtract 037 once. Subsequent distinctive notation will confirm
 631 the divisibility of the original number by 037 while its absence will disprove it. Thus, the other two frames
 632 for the exemplary number yield:

$$633 \quad 002 + 701 + 431 + 941 + 700 = 2775 \rightarrow 002 + 775 = 777;$$

$$634 \quad 027 + 014 + 319 + 417 = 777.$$

635 This criterion applies to numbers of any length and requires a register with only three positions. Moving
 636 along a linear notation, such register adds digital triplets together and transfers thousand’s digits to unit’s
 637 digits.

638 The same triplet digital symmetry and related divisibility criterion exist in all numeral systems with
639 radix q that meets the requirement $(q - 1)/3 = Integer$. The symmetry-related prime number in those
640 systems is found as $111_{q/3}$. Thus, the feature exists in the quaternary system ($q = 4$) with prime number 7
641 (013_4), septenary system ($q = 7$) with prime number 19 (025_7), decimal system ($q = 10$) with prime number
642 037 , the system with $q = 13$ and prime number 61 (049_{13}), and so on. The digital symmetry of the
643 quaternary system is shown in Fig. C.1.

644 **Appendix D. The cytoplasmic balance**

645 Fig. D.1 represents the entire genetic code as a peptide. Each amino acid is inserted into this peptide as
646 many times as it appears in the standard code. Amino acid block residues make up the peptide backbone.
647 The resulting polymer is 61 amino acids long. If its N- and C-termini are eliminated by closing the peptide
648 into a ring, its backbone and side chains appear precisely balanced. Notably, this feature is common to
649 natural proteins: their mass is distributed approximately equally between peptide backbone and side chains
650 (Downes and Richardson, 2002). This also automatically implies that frequencies of amino acids in natural
651 proteins correlate with their abundance in the genetic code (see data in Gilis et al., 2001).

652 Not only the activation key is discarded in this balance, but amino acid molecules are considered as they
653 appear in cytoplasmic environment (where side chains of some of them are ionized). For these reasons the
654 balance shown in Fig. D.1 is referred to as natural or cytoplasmic. Nevertheless, unusual peptide form
655 (though circular peptides do occur rarely in nature, see Conlan et al., 2010) and distinction between amino
656 acid blocks and chains suggest that the cytoplasmic balance and the “virtual” balances shown in the main
657 text are likely to be related phenomena. Possibly, this balance is intended to validate the artificial nature of
658 the activation key, showing that only actual proline could maintain patterns in natural environment. This
659 balance was found by Downes and Richardson (2002) from biological aspect. Simultaneously,
660 Kashkarov et al. (2002) found it with a formal arithmetical approach.

661 **Appendix E. Algebraic representation of the signal**

662 Here we describe a possible way the signal-harboring mapping might have been obtained. As initial
663 data, one has a set of 64 codons and another set of 20 canonical amino acids plus *Stop*. Suppose, the

664 mapping between those two sets is unknown and it has to be deduced from the given pattern ensemble of
 665 the signal. There are $\sim 10^{83}$ possible mappings between the two sets, provided that each element from both
 666 sets is represented at least once. Knowing the ideogram (without knowing nucleon numbers mapped to
 667 individual codons) is equivalent to knowing the block structure of the code. From this follows the first
 668 portion of equations $ggt = ggc = gga = ggg = ggn$, $ttt = ttc = tty$, etc., where codons are used to denote
 669 variables – unknown nucleon numbers of amino acid side chains. Thus, the number of elements in the first
 670 set is essentially reduced from 64 to 24. But there are still $\sim 10^{30}$ possible mappings left. Now one might
 671 write down the nucleon sums from Figs. 5-8 and 10 (leaving out algebraically dependent parts and standard
 672 block sums, as we are provided with the set of canonical amino acids):

$$673 \quad ggn + gcn + tcn + ccn + gtn + acn + ctg + cgn = 333 \text{ (Fig. 7b);}$$

$$674 \quad tgy + tga + ath + tar + agy + ttr + aay + gay + car + aar + gar + cay + tty + agr + tay + atg + tgg =$$

$$675 \quad 111 + 999 \text{ (Fig. 7b);}$$

$$676 \quad tty + ttr + tcn + tay + tar + tgy + tga + tgg + ctg + ccn + cay + car + cgn = 814 \text{ (Fig. 8a);}$$

$$677 \quad tty + ttr + tcn + tay + tar + tgy + tga + tgg + gtn + gcn + gay + gar + ggn = 654 \text{ (Fig. 8b);}$$

$$678 \quad tty + ttr + ctg + ath + atg + gtn + tgy + tga + tgg + cgn + agy + agr + ggn = 789 \text{ (Fig. 8b);}$$

$$679 \quad tty + aar + ath + tcn + cay + 2gcn + ctg + tgy + tga + gay + atg + car + agy = 703 \text{ (Fig. 5a);}$$

$$680 \quad ggn + ccn + ctg + 2acn + tay + tcn + 2gtn + 2cgn + agy + tar + gay = 703 \text{ (Fig. 5a);}$$

$$681 \quad tty + 2ttr + 3ccn + 2ctg + ath + gtn + 2tcn + acn + gcn + tay + tgy + cay + cgn = 999 \text{ (Fig. 5b);}$$

$$682 \quad 2aay + aar + tar + car + gar = 333 \text{ (Fig. 5b);}$$

$$683 \quad 3ggn + tgg + cgn + agr = 333 \text{ (Fig. 5b);}$$

$$684 \quad ath + acn + agr + gtn + gcn + gar = 333 \text{ (Fig. 5b);}$$

$$685 \quad tty + 2ctg + 2tcn + ccn + 2aay + tar + ath + car + acn + 2ggn + tgg + gtn + cgn + gcn = 888 \text{ (Fig. 5c);}$$

$$686 \quad 5tty + 4ttr + 5ctg + 4ath + atg + 5gtn + 5tcn + ccn + acn + gcn + 3tay + 2tar + cay + aay + gay + 3tgy$$

$$687 \quad + tga + tgg + cgn + agy + ggn = 666 + 999 \times 2 \text{ (Fig. 6b);}$$

$$688 \quad 2tar + aar + 2atg = 222 \text{ (Fig. 10d);}$$

$$689 \quad agy + 2aar + tgh = 222 \text{ (Fig. 10e).}$$

690 There are also additional inequalities provided by the ideogram (Fig. 7a):

$$691 \quad ggn \leq gcn \leq tcn \leq ccn \leq gtn \leq can \leq ctn \leq cgn;$$

$$692 \quad tgh \leq ath;$$

$$693 \quad tar \leq agy \leq ttr \leq aay \leq gay \leq car \leq aar \leq gar \leq cay \leq tty \leq agr \leq tay;$$

$$694 \quad atg \leq tgg.$$

695 Finally, $tgh = tgy$ to account for two code versions. In total, there are 26 unknowns, 16 equations and 4
 696 inequalities (the cytoplasmic balance is not accounted here as it has no algebraic connection to this system
 697 due to the activation key). Generally, such systems of Diophantine equations have multiple solutions. Since
 698 we are interested here in deducing the mapping of the code from available patterns and the fixed set of
 699 canonical amino acids, the solution is to be searched within the fragmentary domain $\{0, 1, 15, 31, 41, 43,$
 700 $45, 47, 57, 58, 59, 72, 73, 75, 81, 91, 100, 107, 130\}$. In this case, analysis of the system with any software
 701 algebraic solver shows that this system has a single solution coinciding with the actual mapping of nucleon
 702 numbers onto codons: $tty = 91$, $ggn = 1$, $tga = 0$, $ath = 57$, etc. That still leaves us with several mappings
 703 for amino acids though, since two of the roots – 57 and 72 – represent two amino acids each. This
 704 ambiguity is eliminated when side chain patterns (Figs. 7b and 8a) are also taken into account. After that
 705 the actual mapping of the code is deduced completely and unambiguously from the algebraic system of
 706 patterns, given the set of canonical amino acids and *Stop*. In fact, analysis shows that unambiguous solution
 707 is achieved even if the restriction of the fragmentary domain is applied only to some of the unknowns. In
 708 another approach (*shCherbak*, 2003) unambiguous solution is achieved even without fixed set of amino
 709 acids and with only small assumptions concerning the amino acid set itself.

710 Acknowledgements

711 Partially the study was financed by the Ministry of Education and Science of the Republic of
 712 Kazakhstan. The research was promoted by Professor Bakytzhan Zhumagulov from the National
 713 Engineering Academy of the Republic of Kazakhstan. Part of the research was made during V.I.S.' stay at
 714 Max-Planck-Institut für biophysikalische Chemie (Göttingen, Germany) on kind invitation of Professor
 715 Manfred Eigen. V.I.S. expresses special thanks to Ruthild Winkler-Oswatitsch for her valuable help.

M.A.M. acknowledges the support by the administration of Fesenkov Astrophysical Institute. The authors are grateful to Professor Paul Davies, Vladimir Kashkarov, Artem Novozhilov, Denis Tulinov, Artem Yermilov and Denis Yurin for objective criticism and fruitful discussions of the manuscript.

Authors' contributions

V.I.S. conceived of and performed the research, developed graphic arts. V.I.S. and M.A.M. analyzed data, introduced interpretation of the activation key, outlined structure of the paper. M.A.M. performed statistical test and algebraic analysis, wrote the manuscript.

Authors' information

V.I.S. received his Master degree in Physics from Al-Farabi Kazakh National University in 1970. Since 1980 he is a senior researcher in the Laboratory of computer science in biology at the Department of Mathematics, where he received his PhD and became the chief of the Laboratory in 1995. Since 1989 he is a full member of the International Society for the Study of the Origin of Life.

M.A.M. received his Master degree in Physics from Lomonosov Moscow State University in 2004. Since 2006 he is in Fesenkov Astrophysical Institute, where he has been involved in astrophysical and cosmological computer simulations. His another research activity is concerned with astrobiology and cell biology.

References

- Ailenberg, M., Rotstein, O.D., 2009. An improved Huffman coding method for archiving text, images, and music characters in DNA. *BioTechniques* 47, 747-754.
- Alberts, B., Johnson, A., Lewis, J., Raff, M., Roberts, K., Walter, P., 2008. *Molecular biology of the cell*, 5th edition. Garland Science, New York.
- Alff-Steinberger, C., 1969. The genetic code and error transmission. *Proc. Natl. Acad. Sci. USA* 64, 584-591.

- 739 Alvager, T., Graham, G., Hilleke, R., Hutchison, D., Westgard, J., 1989. On the information content of the
740 genetic code. *BioSystems* 22, 189-196.
- 741 Baisnée, P.-F., Baldi, P., Brunak, S., Pedersen, A.G., 2001. Flexibility of the genetic code with respect to
742 DNA structure. *Bioinformatics* 17, 237-248.
- 743 Bancroft, C., Bowler, T., Bloom, B., Clelland, C.T., 2001. Long-term storage of information in DNA.
744 *Science* 293, 1763-1765.
- 745 Barbieri, M., 2008. Biosemiotics: a new understanding of life. *Naturwissenschaften* 95, 577-599.
- 746 Bashford, J.D., Tsohantjis, I., Jarvis, P.D., 1998. A supersymmetric model for the evolution of the genetic
747 code. *Proc. Natl. Acad. Sci. USA* 95, 987-992.
- 748 Bertman, M.O., Jungck, J.R., 1979. Group graph of the genetic code. *J. Hered.* 70, 379-384.
- 749 Bollenbach, T., Vetsigian, K., Kishony, R., 2007. Evolution and multilevel optimization of the genetic
750 code. *Genome Res.* 17, 401-404.
- 751 Budisa, N., 2006. *Engineering the Genetic Code: Expanding the Amino Acid Repertoire for the Design of*
752 *Novel Proteins.* Wiley-VCH, Weinheim.
- 753 Chin, J.W., 2012. Reprogramming the genetic code. *Science* 336, 428-429.
- 754 Conlan, B.F., Gillon, A.D., Craik, D.J., Anderson, M.A., 2010. Circular proteins and mechanisms of
755 cyclization. *Biopolymers* 94, 573-583.
- 756 Crick, F.H.C., 1968. The origin of the genetic code. *J. Mol. Biol.* 38, 367-379.
- 757 Crick, F.H.C., 1981. *Life Itself: Its Origin and Nature.* Simon and Schuster, New York.
- 758 Crick, F.H.C., Griffith, J.S., Orgel, L.E., 1957. Codes without commas. *Proc. Natl. Acad. Sci. USA* 43,
759 416-421.
- 760 Crick, F.H.C., Orgel, L.E., 1973. Directed panspermia. *Icarus* 19, 341-346.
- 761 Danckwerts, H.-J., Neubert, D., 1975. Symmetries of genetic code-doublets. *J. Mol. Evol.* 5, 327-332.
- 762 Davies, P.C.W., 2010. *The Eerie Silence: Are We Alone in the Universe?* Penguin, London.
- 763 Davies, P.C.W., 2012. Footprints of alien technology. *Acta Astronaut.* 73, 250-257.

- 764 Dennett, D.C., 1996. *Darwin's Dangerous Idea: Evolution and the Meanings of Life*. Penguin, London, p.
765 131.
- 766 Di Giulio, M., 2005. The origin of the genetic code: theories and their relationships, a review. *BioSystems*
767 80:175-184.
- 768 Downes, A.M., Richardson, B.J., 2002. Relationships between genomic base content and distribution of
769 mass in coded proteins. *J. Mol. Evol.* 55, 476–490.
- 770 Dumas, S., Dutil, Y., 2004. The Eypatoria messages. <http://www.activeseti.org>, 'papers' section.
- 771 Ehman, J.R., 2011. "Wow!" – a tantalizing candidate. In: Shuch, H.P. (Ed.), *Searching for Extraterrestrial*
772 *Intelligence: SETI Past, Present, and Future*. Springer, Berlin, Heidelberg, pp. 47-63.
- 773 Eigen, M., Winkler, R., 1983. *Laws of the Game: How the Principles of Nature Govern Chance*. Princeton
774 Univ. Press, Princeton.
- 775 Elliott, J.R., 2010. Detecting the signature of intelligent life. *Acta Astronaut.* 67, 1419-1426.
- 776 Freeland, S.J., 2002. The Darwinian genetic code: an adaptation for adapting? *Genet. Programm. Evolvable*
777 *Mach.* 3, 113-127.
- 778 Freeland, S.J., Hurst, L.D., 1998. The genetic code is one in a million. *J. Mol. Evol.* 47, 238-248.
- 779 Freitas, R.A., 1983. The search for extraterrestrial artifacts (SETA). *J. Brit. Interplanet. Soc.* 36, 501-506.
- 780 Freudenthal, H., 1960. *LINCOS: Design of a Language for Cosmic Intercourse*. North-Holland Publishing
781 Company, Amsterdam.
- 782 Gamow, G., Yčas, M., 1955. Statistical correlation of protein and ribonucleic acid composition. *Proc. Natl.*
783 *Acad. Sci. USA* 41, 1011-1019.
- 784 Gibson, D.G., Glass, J.L., Lartigue, C., Noskov, V.N., Chuang, R.Y., Algire, M.A., Benders, G.A.,
785 Montague, M.G., Ma, L., Moodie, M.M., Merryman, C., Vashee, S., Krishnakumar, R., Assad-Garcia,
786 N., Andrews-Pfannkoch, C., Denisova, E.A., Young, L., Qi, Z.Q., Segall-Shapiro, T.H., Calvey, C.H.,
787 Parmar, P.P., Hutchison, C.A. III, Smith, H.O., Venter, J.C., 2010. Creation of a bacterial cell
788 controlled by a chemically synthesized genome. *Science* 329, 52-56.

- 789 Giegé, R., Sissler, M., Florentz, C., 1998. Universal rules and idiosyncratic features in tRNA identity.
790 Nucleic Acids Res. 26, 5017-5035.
- 791 Gilis, D., Massar, S., Cerf, N.J., Rooman, M., 2001. Optimality of the genetic code with respect to protein
792 stability and amino-acid frequencies. Genome Biol. 2, 49.1–49.12.
- 793 Gonzalez, D.L., 2004. Can the genetic code be mathematically described? Med. Sci. Monit. 10, HY11-17.
- 794 Gusev, V.A., Schulze-Makuch, D., 2004. Genetic code: Lucky chance or fundamental law of nature? Phys.
795 Life Rev. 1, 202-229.
- 796 Haig, D., Hurst, L.D., 1991. A quantitative measure of error minimization in the genetic code. J. Mol. Evol.
797 33, 412-417.
- 798 Hasegawa, M., Miyata, T., 1980. On the antisymmetry of the amino acid code table. Orig. Life 10, 265-
799 270.
- 800 Hayes, B., 1998. The invention of the genetic code. Am. Sci. 86, 8-14.
- 801 Higgs, P.G., 2009. A four-column theory for the origin of the genetic code: tracing the evolutionary
802 pathways that gave rise to an optimized code. Biol. Dir. 4, 16.
- 803 Hoch, A.J., Losick, R., 1997. Panspermia, spores and the *Bacillus subtilis* genome. Nature 390, 237-238.
- 804 Hornos, J.E.M., Hornos, Y.M.M., 1993. Algebraic model for the evolution of the genetic code. Phys. Rev.
805 Lett. 71, 4401-4404.
- 806 Ibba, M., Söll, D., 2000. Aminoacyl-tRNA synthesis. Annu. Rev. Biochem. 69, 617-650.
- 807 Itzkovitz, S., Alon, U., 2007. The genetic code is nearly optimal for allowing additional information within
808 protein-coding sequences. Genome Res. 17, 405-412.
- 809 Jungck, J.R., 1978. The genetic code as a periodic table. J. Mol. Evol. 11, 211-224.
- 810 Karasev, V.A., Stefanov, V.E., 2001. Topological nature of the genetic code. J. Theor. Biol. 209, 303-317.
- 811 Kashkarov, V.V., Krassovitskiy, A.M., Mamleev, V.S., *sh*Cherbak, V.I., 2002. Random sequences of
812 proteins are exactly balanced like the canonical base pairs of DNA. 10th ISSOL Meeting and 13th
813 International Conference on the Origin of Life, 121-122 (abstract).

- 814 Klump, H.H., 2006. Exploring the energy landscape of the genetic code. *Arch. Biochem. Biophys.* 453, 87-
815 92.
- 816 Knight, R.D., Freeland, S.J., Landweber, L.F., 1999. Selection, history and chemistry: the three faces of the
817 genetic code. *Trends Biochem. Sci.* 24, 241-247.
- 818 Knight, R.D., Freeland, S.J., Landweber, L.F., 2001. Rewiring the keyboard: evolvability of the genetic
819 code. *Nat. Rev. Genet.* 2, 49–58.
- 820 Koonin, E.V., Novozhilov, A.S., 2009. Origin and evolution of the genetic code: the universal enigma.
821 *IUBMB Life* 61, 99–111.
- 822 Kudla, G., Murray, A.W., Tollervey, D., Plotkin, J.B., 2009. Coding-sequence determinants of gene
823 expression in *Escherichia coli*. *Science* 324, 255-258.
- 824 Lagerkvist, U., 1978. “Two out of three”: an alternative method for codon reading. *Proc. Natl. Acad. Sci.*
825 *USA* 75, 1759-1762.
- 826 Li, S., Hong, M., 2011. Protonation, tautomerization, and rotameric structure of histidine: a comprehensive
827 study by magic-angle-spinning solid-state NMR. *J. Am. Chem. Soc.* 133, 1534-1544.
- 828 Marx, G., 1979. Message through time. *Acta astronaut.* 6, 221-225.
- 829 Mautner, M.N., 2000. *Seeding the Universe with Life: Securing Our Cosmological Future*. Legacy Books,
830 Christchurch.
- 831 McClain, W.H., Foss, K., 1988. Changing the acceptor identity of a transfer RNA by altering nucleotides in
832 a “variable pocket”. *Science* 241, 1804-1807.
- 833 Meyer, F., Schmidt, H.I., Plümper, E., Hasilik, A., Mersmann, G., Meyer, H.E., Engstörn, A., Heckmann,
834 K., 1991. UGA is translated as cysteine in pheromone 3 of *Euplotes octocarinatus*. *Proc. Natl. Acad.*
835 *Sci. USA* 88, 3758-3761.
- 836 Minsky, M., 1985. Why intelligent aliens will be intelligible. In: Regis, E. (Ed.), *Extraterrestrials: Science*
837 *and Alien Intelligence*. Cambridge Univ. Press, Cambridge, pp. 117-128.
- 838 Moura, G.R., Paredes, J.A., Santos, M.A.S., 2010. Development of the genetic code: Insights from a fungal
839 codon reassignment. *FEBS Lett.* 584, 334–341.

- 840 Nakamura, H., 1986. SV40 DNA – A message from ϵ Eri? *Acta Astronaut.* 13, 573-578.
- 841 Nirenberg, M., Leder, P., Bernfield, M., Brimacombe, R., Trupin, J., Rottman, F., O’Neal, C., 1965. RNA
842 codewords and protein synthesis, VII. On the general nature of the RNA code. *Proc. Natl. Acad. Sci.*
843 USA 53, 1161-1168.
- 844 Novozhilov, A.S., Wolf, Y.I., Koonin, E.V., 2007. Evolution of the genetic code: partial optimization of a
845 random code for robustness to translation error in a rugged fitness landscape. *Biol. Dir.* 2, 24.
- 846 Okamura, K., Wei, J., Scherer, S.W., 2007. Evolutionary implications of inversions that have caused intra-
847 strand parity in DNA. *BMC Genomics* 8, 160.
- 848 Pacioli, L., 1508. *De Viribus Quantitatis*, manuscript, Library of the University of Bologna, code number
849 250.
- 850 Rak, J., Skurski, P., Simons, J., Gutowski, M., 2001. Low-energy tautomers and conformers of neutral and
851 protonated arginine. *J. Am. Chem. Soc.* 123, 11695-11707.
- 852 Rodin, A.S., Szathmáry, E., Rodin, S.N., 2011. On origin of genetic code and tRNA before translation.
853 *Biol. Dir.* 6, 14.
- 854 Rose, C., Wright, G., 2004. Inscribed matter as an energy-efficient means of communication with an
855 extraterrestrial civilization. *Nature* 431, 47-49.
- 856 Rumer, Yu.B., 1966. Codon systematization in the genetic code. *Dokl. Acad. Nauk SSSR* 167, 1393-1394
857 (in Russian).
- 858 Sagan, C., Drake, F.D., Druyan, A., Ferris, T., Lomberg, J., Sagan, L.S., 1978. *Murmurs of Earth: The*
859 *Voyager Interstellar Record*. Random House, New York.
- 860 Sagan, C., Sagan, L.S., Drake, F., 1972. A Message from Earth. *Science* 175, 881-884.
- 861 Sella, G., Ardell, D.H., 2006. The coevolution of genes and genetic codes: Crick’s frozen accident
862 revisited. *J. Mol. Evol.* 63, 297-313.
- 863 *shCherbak*, V.I., 1988. The co-operative symmetry of the genetic code. *J. Theor. Biol.* 132, 121-124.
- 864 *shCherbak*, V.I., 1993. The symmetrical architecture of the genetic code systematization principle. *J.*
865 *Theor. Biol.* 162, 395-398.

- 866 *sh*Cherbak, V.I., 2003. Arithmetic inside the universal genetic code. *BioSystems* 70, 187-209.
- 867 Siemion, I.Z., Stefanowicz, P., 1992. Periodical change of amino acid reactivity within the genetic code.
868 *BioSystems* 27, 77-84.
- 869 Taniguchi, M., Hino, T., 1981. Cyclic tautomers of tryptophans and tryptamines – 4. *Tetrahedron* 37, 1487-
870 1494.
- 871 Taylor, F.J.R., Coates, D., 1989. The code within the codons. *BioSystems* 22, 177-187.
- 872 Tepfer, D., 2008. The origin of life, panspermia and a proposal to seed the Universe. *Plant Science* 175,
873 756-760.
- 874 The Staff at the National Astronomy and Ionosphere Center, 1975. The Arecibo message of November,
875 1974. *Icarus* 26, 462-466.
- 876 Tlusty, T., 2010. A colorful origin for the genetic code: Information theory, statistical mechanics and the
877 emergence of molecular codes. *Phys. Life Rev.* 7, 362-376.
- 878 Travers, A., 2006. The evolution of the genetic code revisited. *Orig. Life Evol. Biosph.* 36, 549-555.
- 879 Trifonov, E.N., 2000. Consensus temporal order of amino acids and evolution of the triplet code. *Gene*
880 261, 139-151.
- 881 Vetsigian, K., Woese, C., Goldenfeld, N., 2006. Collective evolution and the genetic code. *Proc. Natl.*
882 *Acad. Sci. USA* 103, 10696-10701.
- 883 Wilhelm, T., Nikolajewa, S., 2004. A new classification scheme of the genetic code. *J. Mol. Evol.* 59, 598-
884 605.
- 885 Woese, C.R., 1965. Order in the genetic code. *Proc. Natl. Acad. Sci. USA* 54, 71-75.
- 886 Wolf, Y.I., Koonin, E.V., 2007. On the origin of the translation system and the genetic code in the RNA
887 world by means of natural selection, exaptation, and subfunctionalization. *Biol. Dir.* 2, 14.
- 888 Wong, J.T.-F., 2005. Coevolution theory of the genetic code at age thirty. *BioEssays* 27, 416-425.
- 889 Yachie, N., Ohashi, Y., Tomita, M., 2008. Stabilizing synthetic data in the DNA of living organisms. *Syst.*
890 *Synth. Biol.* 2, 19-25.

- 891 Yarus, M., Widmann, J.J., Knight, R., 2009. RNA-amino acid binding: a stereochemical era for the genetic
892 code. *J. Mol. Evol.* 69, 406-429.
- 893 Yokoo, H., Oshima, T., 1979. Is bacteriophage ϕ X174 DNA a message from an extraterrestrial
894 intelligence? *Icarus* 38, 148-153.
- 895 Yuan, J., O'Donoghue, P., Ambrogelly, A., Gundllapalli, S., Sherrer, R.L., Palioura, S., Simonović, M.,
896 Söll, D., 2010. Distinct genetic code expansion strategies for selenocysteine and pyrrolysine are
897 reflected in different aminoacyl-tRNA formation systems. *FEBS Lett.* 584, 342–349.
- 898 Zhuravlev, Yu.N., 2002. Two rules of distribution of amino acids in the code table indicate chimeric nature
899 of the genetic code. *Dokl. Biochem. Biophys.* 383, 85-87.

900

901 **Fig. 1.** The genetic code. (a) Traditional representation of the standard, or universal, code. Codons coding the same amino acid
 902 form synonymic series denoted with opening braces. Number of codons in a series defines its redundancy (degeneracy). Whole
 903 codon families consist of one series of redundancy IV. Other families are split. Most split families are halved into two series of
 904 redundancy II each, one ending with pyrimidines {T, C} and another with purines {A, G}. Three codons in the standard code are
 905 not mapped to any amino acid and are used as *Stop* in translation. The *Start* is usually signified by ATG which codes Met.
 906 Closing brace shows the only difference between the euplotid and the standard code. (b) Contracted representation of the
 907 euplotid version. Synonymous full-size codons are replaced by a single contracted series with combined third base. FASTA
 908 designations are used: R and Y stand for purines and pyrimidines, respectively, N stands for all four bases and H stands for {T,
 909 C, A}. Series are placed vertically for further convenience. The pictogram on the left helps in figures below. Filled elements
 910 denote whole families here.

911 **Fig. 2.** Preceding observations. (a) Rumer's bisection. Whole families are opposed to split ones, thereby bisecting the code.
 912 Codons in opposed families are mapped to each other with the ordered set of Rumer's transformation and two half-
 913 transformations. Transformation of third bases is trivial as they are the same in any family; therefore contracted representation is
 914 adequate to show this regularity. The regularity is valid both for the standard and the euplotid (shown here) version.
 915 (b) Categorization of amino acids by nucleon numbers. Free molecules unmodified by cytoplasmic environment are shown.
 916 Each of them is formed of the standard block and a side chain. Blocks are identical in all amino acids except proline. Chains are
 917 unique for each amino acid. Numbers of nucleons, i.e. protons and neutrons, are shown for both blocks and chains. To avoid
 918 ambiguity, it is judicious to consider only most common and stable isotopes: ^1H , ^{12}C , ^{14}N , ^{16}O , ^{32}S . The bar at the bottom shows
 919 the redundancy of amino acids in the code. Cross-cut bonds symbolize the distinction between standard blocks and unique side
 920 chains of amino acids. The arrow in proline denotes hereafter the "activation key" (see Section 2.4).

921 **Fig. 3.** Digital symmetry of decimals divisible by 037. Leading zero emphasizes its equal participation in the symmetry. All
 922 three-digit decimals with identical digits 111, ..., 999 are divisible by 037. The sum of three identical digits gives the quotient of
 923 the number divided by 037. Analogous sum for numbers with unique digits gives the central quotient in the column. Digits in
 924 these numbers are interconnected with cyclic permutations that are mirror symmetrical in neighbor columns. Addition instead of
 925 division provides an efficient way to perform checksums (see *Appendix C*). The scheme extends to decimals with more than
 926 three digits, if they are represented as $a + 999 \times n$, where n is the quotient of the number divided by 999 and a is the remainder, to
 927 which the same symmetry then applies (for three-digit decimals $n = 0$). Numbers divisible by 037 and larger than 999 will be
 928 shown in this way.

929 **Fig. 4.** The structure of the signal. All details are discussed sequentially in the text. The image of scales represents precise
 930 nucleon equalities. DEC stands for distinctive decimal notation of nucleon sums. The dotted box denotes the cytoplasmic

balance (see *Appendix D*), the only pattern maintained by actual proline and cellular milieu. All other patterns are enabled by the “activation key” and are valid for free amino acids. K stands for {T, G}, M stands for {A, C}. Though all three types of transformations act in the patterns, only Rumer’s transformation is indicated for simplicity.

Fig. 5. Gamow’s sorting of codons according to their nucleotide base composition. Base combinations (shown on triangular frames) produce three sets: 4 codons with three identical bases, 24 codons with unique bases and 36 codons with two identical bases. (a) The first and the second sets halved by vertical axis with Rumer’s and half-transformations along with *Spin*→*Antispin* transformation denoted with circular arrows. Applied to triangular frames, these arrows define the sequence of bases in codons. Note that while any block sum (with the activation key applied) is divisible by 037 as each block has $74 = 2 \times 037$ nucleons, chain sums are not restricted in this way. (b) The third set halved according to whether identical bases are purines or pyrimidines. (c) The third set halved with horizontal axis according to whether unique bases are purines or pyrimidines.

Fig. 6. The decomposed standard code. (a) Decomposition shown for one family of codons. Three T-bases contribute three Cys molecules into T-set; one A-base contributes one *Stop* to A-set and so on for the entire code. (b) Identical bases are sorted into four sets regardless of their position in codons. The sets are shown twice for convenience.

Fig. 7. The contracted euplotid code with the systematization rule applied (compare with Fig. 2). (a) The resulting arrangement of contracted codon series forming the ideogram. Side-by-side alignment of vertical series produces three horizontal strings of peer-positioned bases. Gln and Lys have the same nucleon number; ambiguity in their positioning is eliminated by the symmetries considered further. (b) The arithmetical background of the ideogram (valid for the standard version as well, as it contributes another zero to the III, II, I set). For β and γ side chain levels see Section 4.3 in *Discussion*.

Fig. 8. Additional arithmetical patterns of the contracted code (shared by both code versions). (a) The code is divided according to whether first bases are purines or pyrimidines. This gives two sets with equal numbers of series. The halve with pyrimidines in first positions reveals a new balance of chains and blocks analogous to that in Fig. 7b. Another halve is algebraically dependent except the decimal sum of its β , δ , ζ levels, see *Discussion*, Section 4.3. (b) The code is divided according to whether first bases are K or M (left) or whether central bases are K or M (center). Both divisions produce halves with identical chain nucleon sums. As algebraic consequence of these divisions, series with K in first and central positions and series with M in first and central positions are chain-balanced (right). Each of the three divisions is accompanied by half-transformations and, remarkably, also produces equal numbers of series in each half. This pattern is the only one that shows no divisibility by 037. However, all three numbers – 654, 789 and 369 – are again specific in decimal notation where digits in each of them appear as arithmetic progressions.

959 **Fig. 9.** Patterns of the short (a) and the long (b) upper strings. The strings are arranged with the same set of symmetries: *mirror*
 960 symmetry (denoted with the central vertical axis), *translation* symmetry (denoted with italicized letters and skewed frames) and
 961 purine \leftrightarrow pyrimidine *inversion* (denoted with color gradient, where black and white stand for pyrimidines and purines,
 962 respectively). The image of DNA at the top illustrates possible interpretation of the short string (see Section 4.4 in *Discussion*).

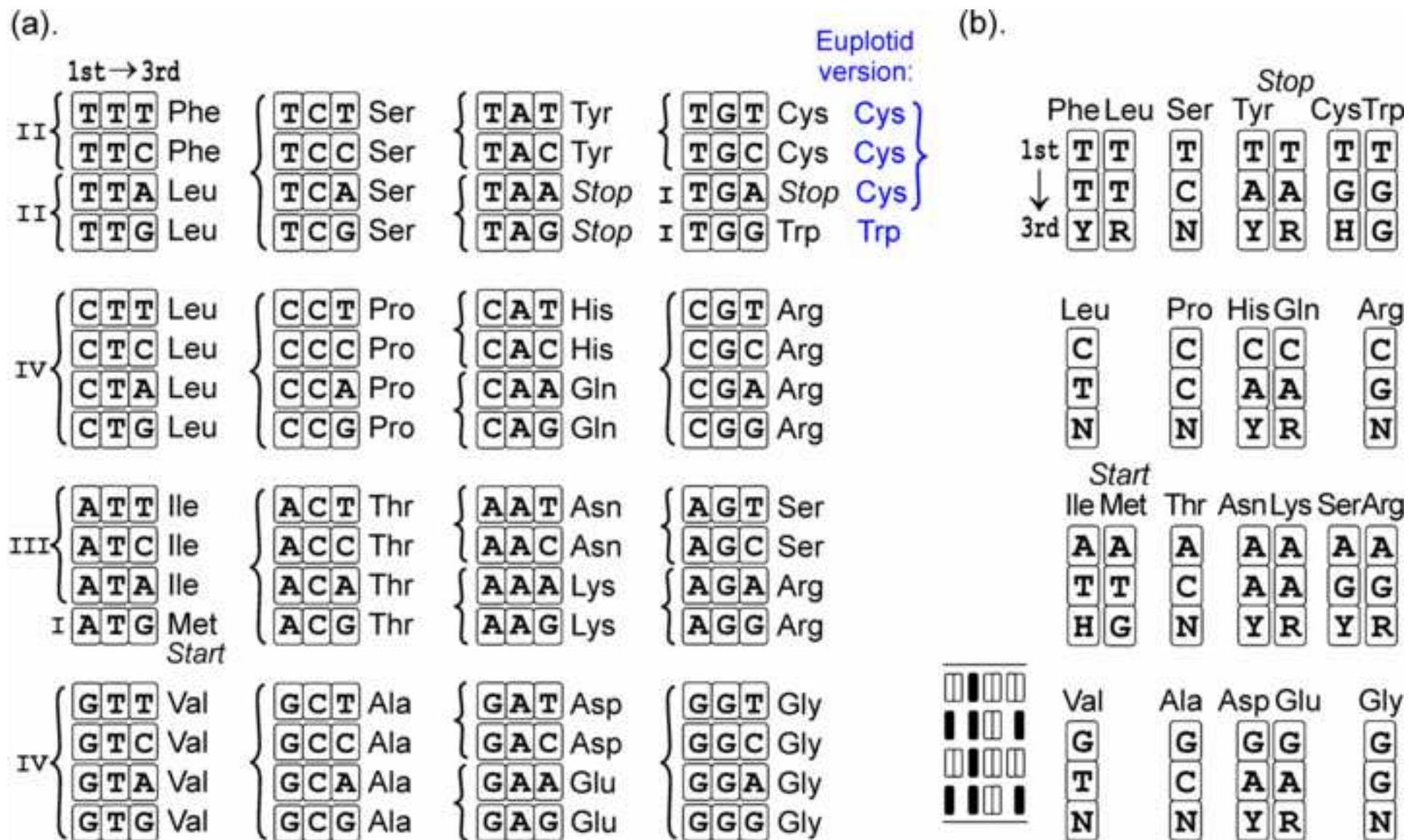
963 **Fig. 10.** Patterns of the short (a, b) and the long (a, c, d, e) center strings. Both strings are arranged with purine-pyrimidine
 964 *mirror* symmetry, purine \leftrightarrow pyrimidine *inversion* and *semantical* symmetry. The first two are denoted in the same way as in
 965 Fig. 9, π denotes palindrome.

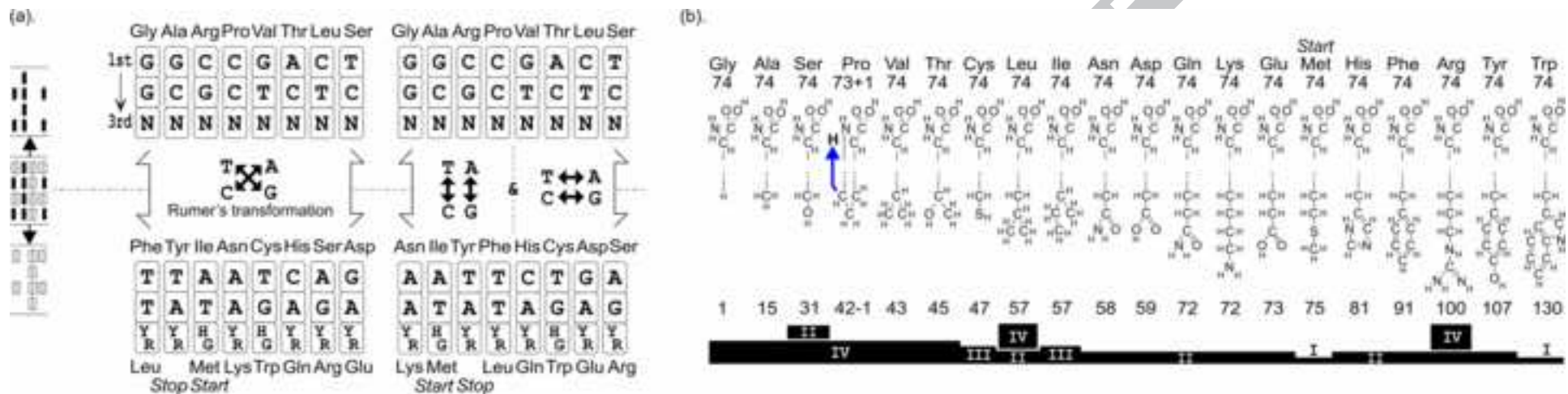
966 **Fig. A.1.** Molecular mechanisms of the genetic code (shown for the case of serine amino acid) and a simple example of its
 967 artificial modification. The contour arrows indicate directionality of DNA and RNA strands as defined by orientation of their
 968 subunits (designated in biochemistry as 5'→3' orientation; replication, transcription and translation occur only in that direction).
 969 (a) tRNA^{Ser} gene (the gene of tRNA that specifies Ser in the standard code) is transcribed by RNA polymerase from genomic
 970 DNA. (b) The folded tRNA^{Ser} molecule (top), serine molecule (middle) and seryl-tRNA synthetase (SARS, an aaRS cognate for
 971 amino acid serine; bottom). (c) SARS recognizes both serine and tRNA^{Ser} and binds them together. (d) Ser-tRNA^{Ser} released
 972 from SARS and ready to be delivered to ribosome. (e) The process of peptide synthesis at the ribosome (as an example, the
 973 mRNA with the gene fragment of the SARS itself is shown). (f) The resulting fragment of the genetic code (also shown is Ala
 974 group, which will be used in an example below). (g)-(k). A simple way of genetic code modification. The shaded sequence in (j)
 975 corresponds to the region shown in (e).

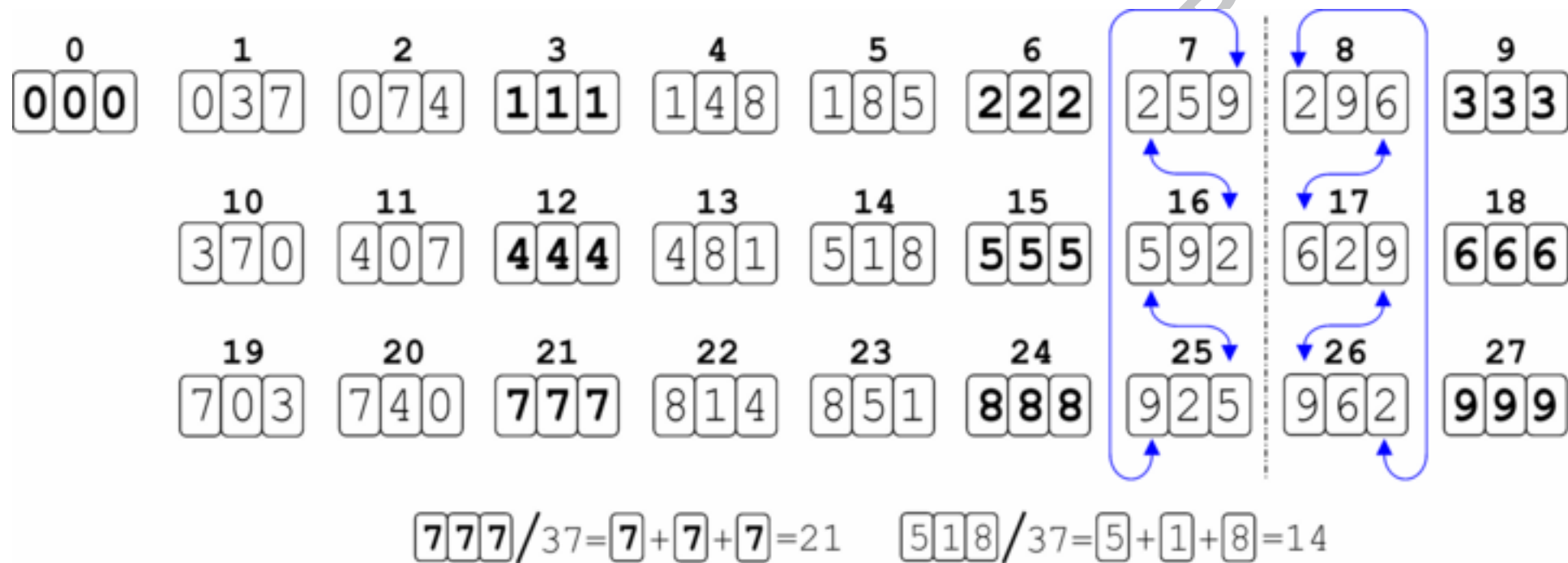
976 **Fig. B.1.** Distribution of variant codes by their scores for (a) nucleon equalities and (b) ideogram symmetries. The size of the
 977 sample in both cases is one billion codes.

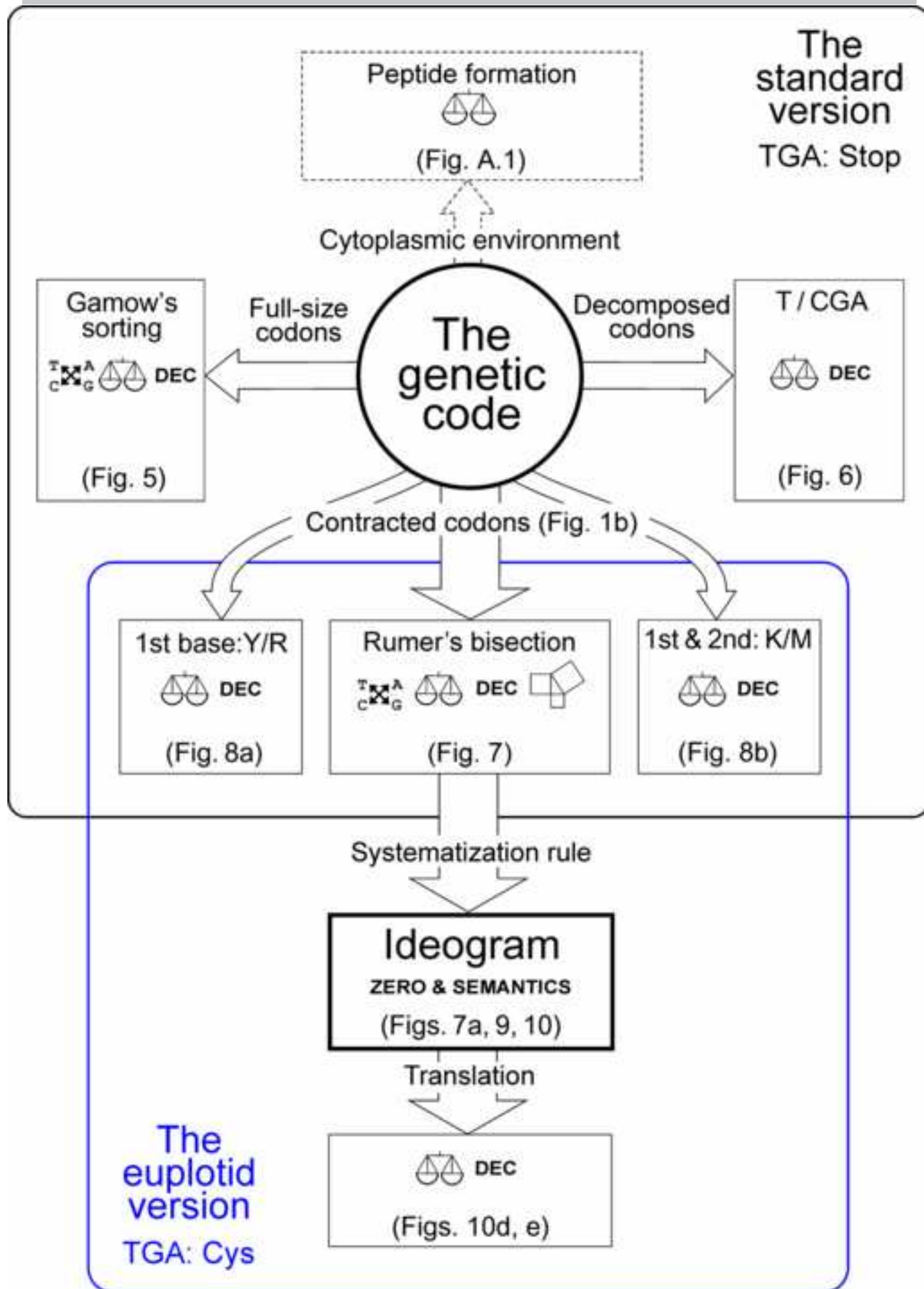
978 **Fig. C.1.** Similar to the decimal system, the quaternary system also displays symmetry of digital triplets, where 7 (013_4) acts
 979 instead of 037 .

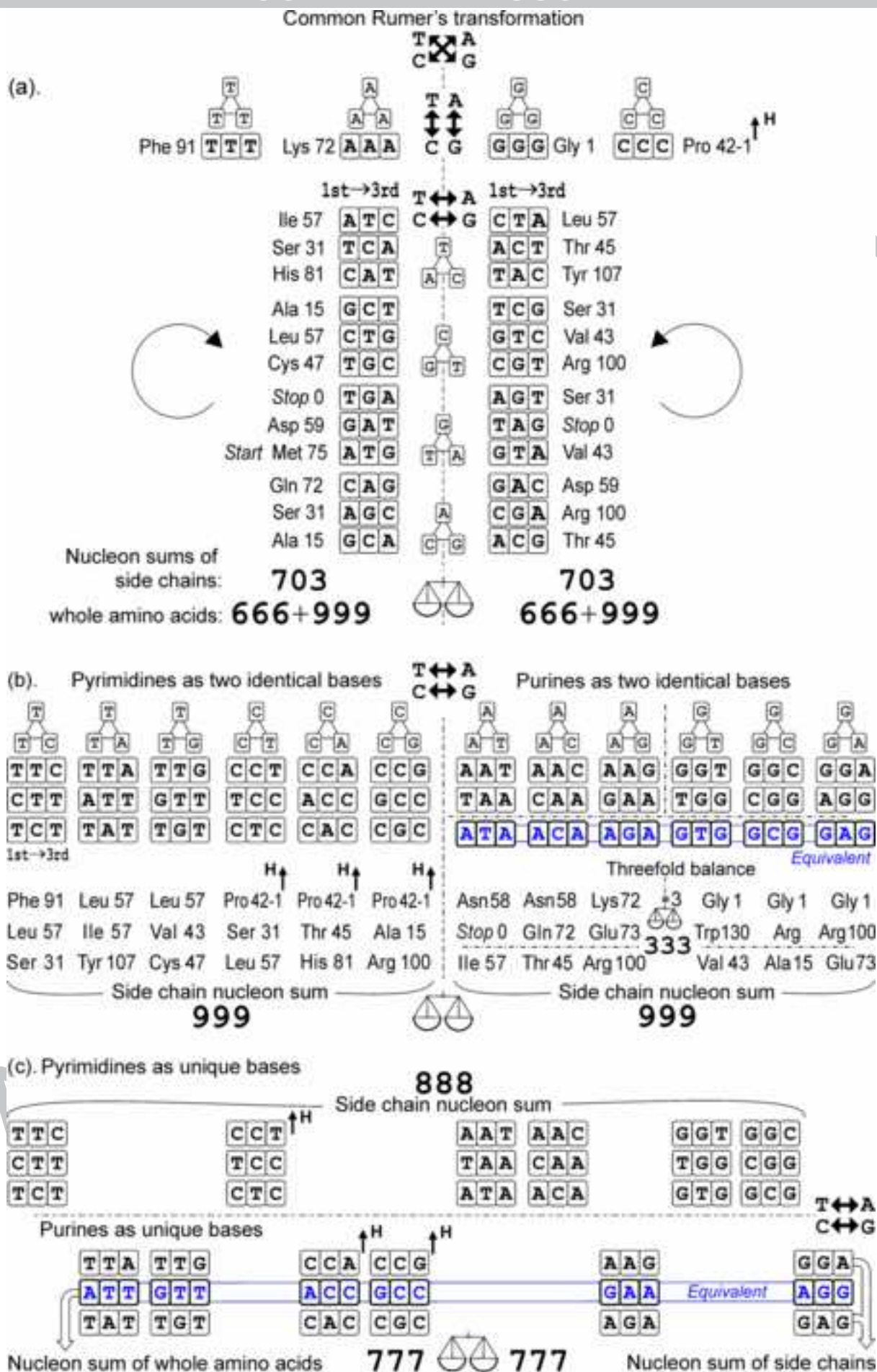
980 **Fig. D.1.** Amino acids of the standard genetic code in the form of a circular peptide (sequence order does not matter). The
 981 peptide is formed by aggregating standard blocks of amino acids into polymer backbone. Formation of each peptide bond
 982 releases a water molecule reducing each amino acid block to 56 nucleons (55 in proline). Asp and Glu lose one proton each from
 983 their side chains at cytoplasmic pH, while Arg and Lys gain one proton each (denoted with -1 and +1, respectively). Other
 984 amino acids are predominantly neutral in cytoplasmic environment (Alberts et al., 2008). As a result, nucleon sum of the peptide
 985 backbone is exactly equal to that of all its side chains.

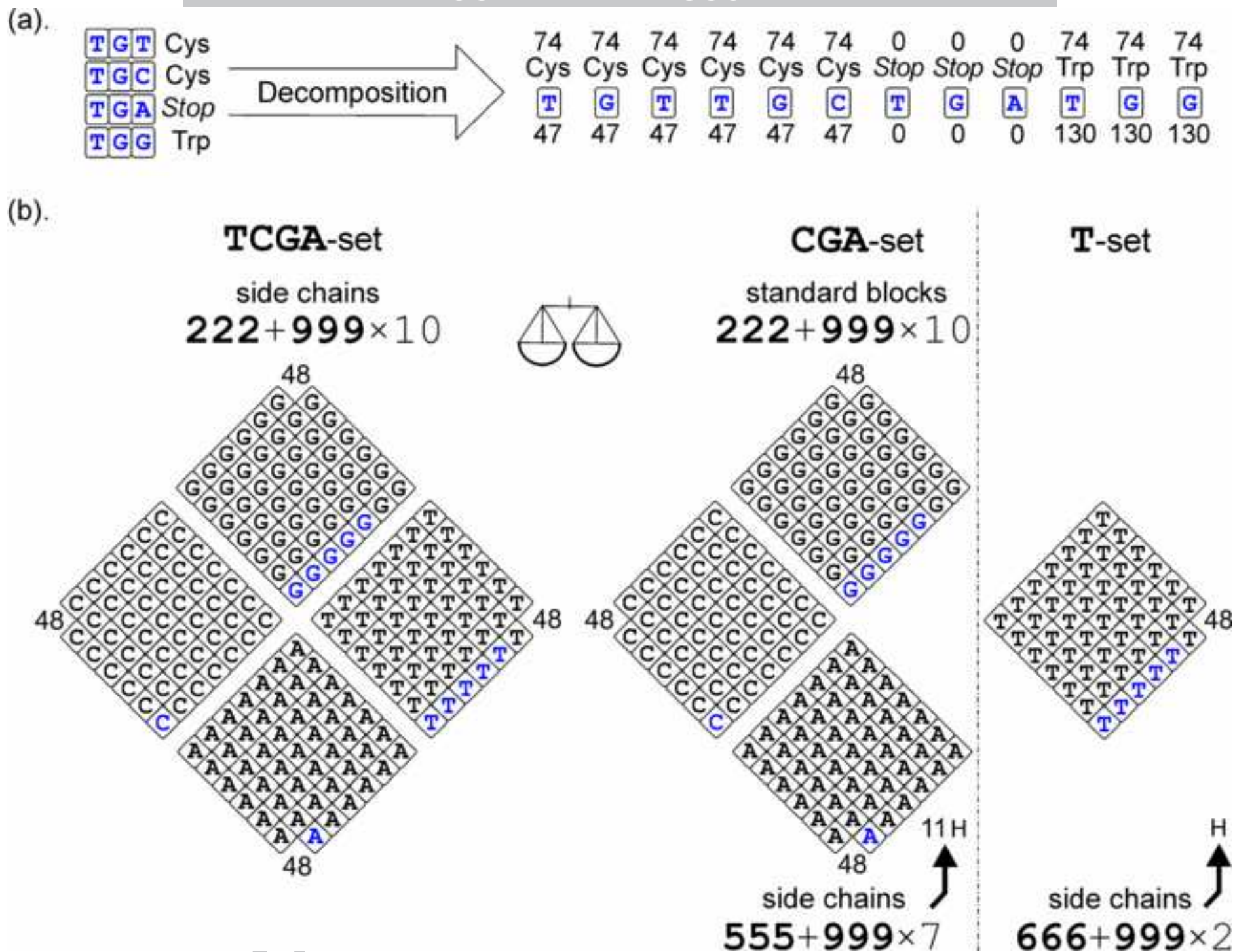


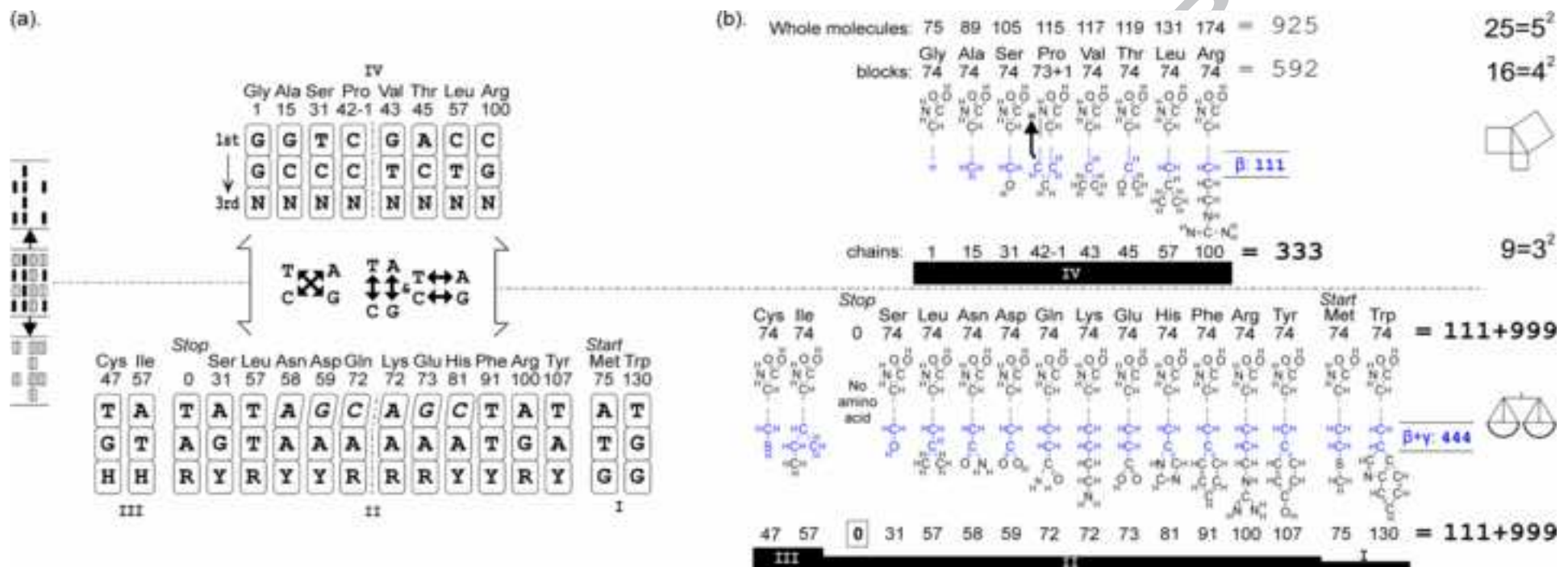






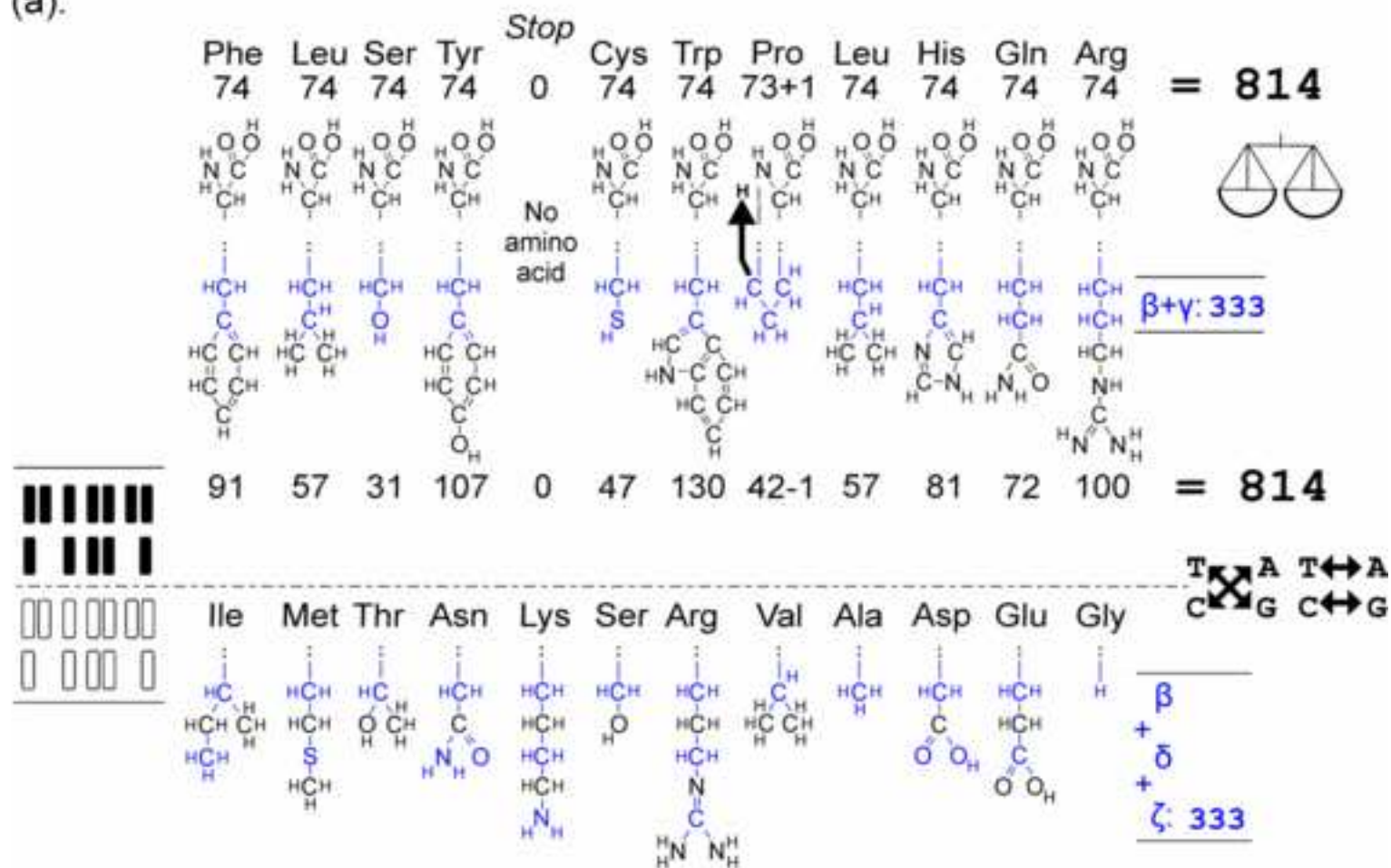




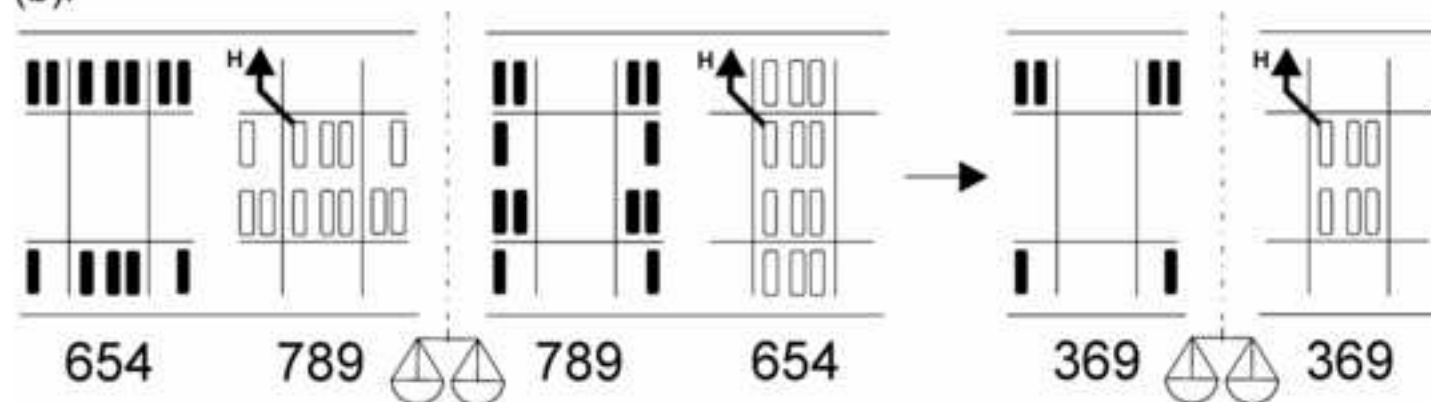


ACCEPTED

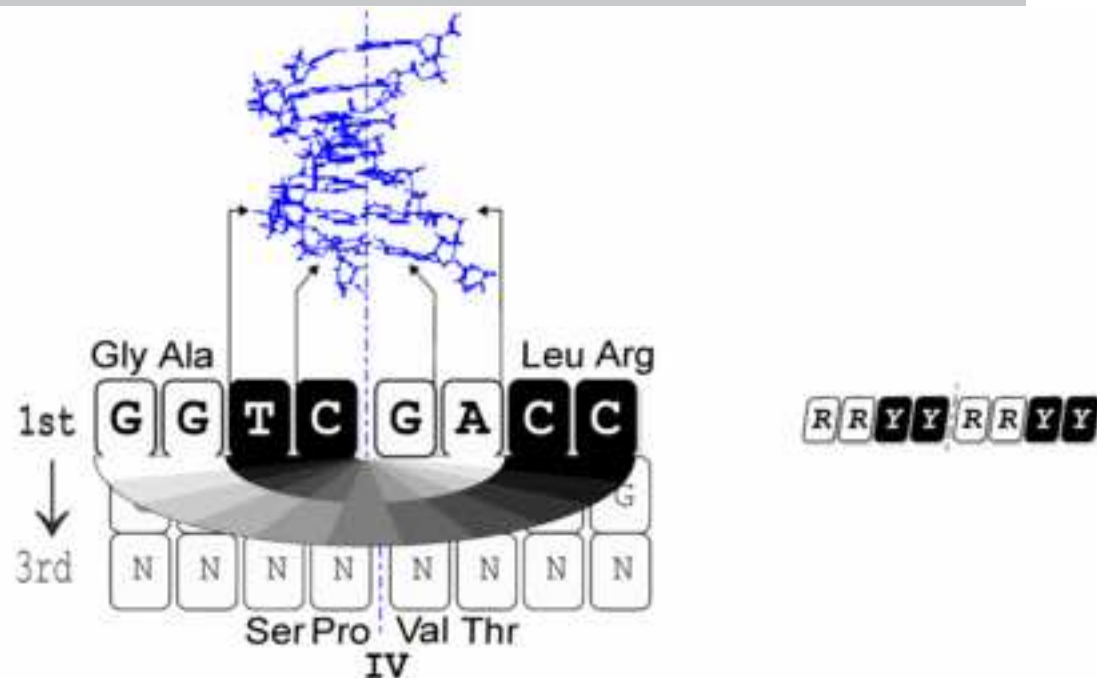
(a).



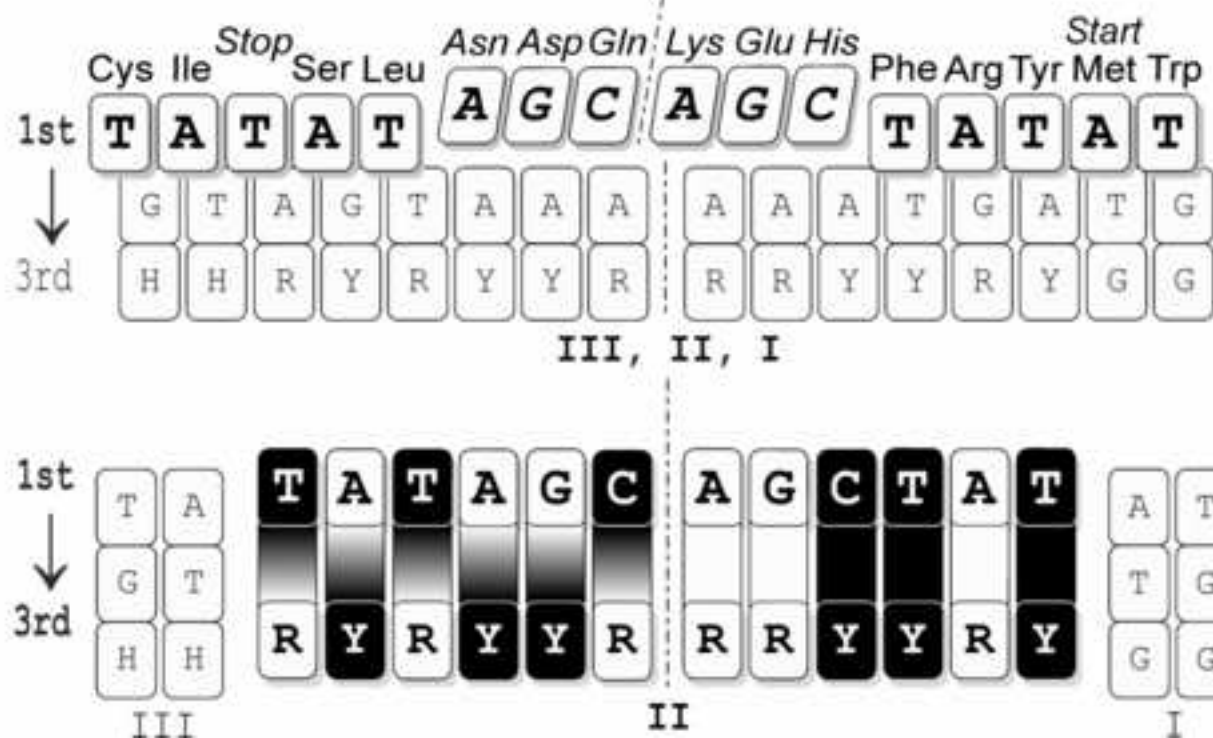
(b).

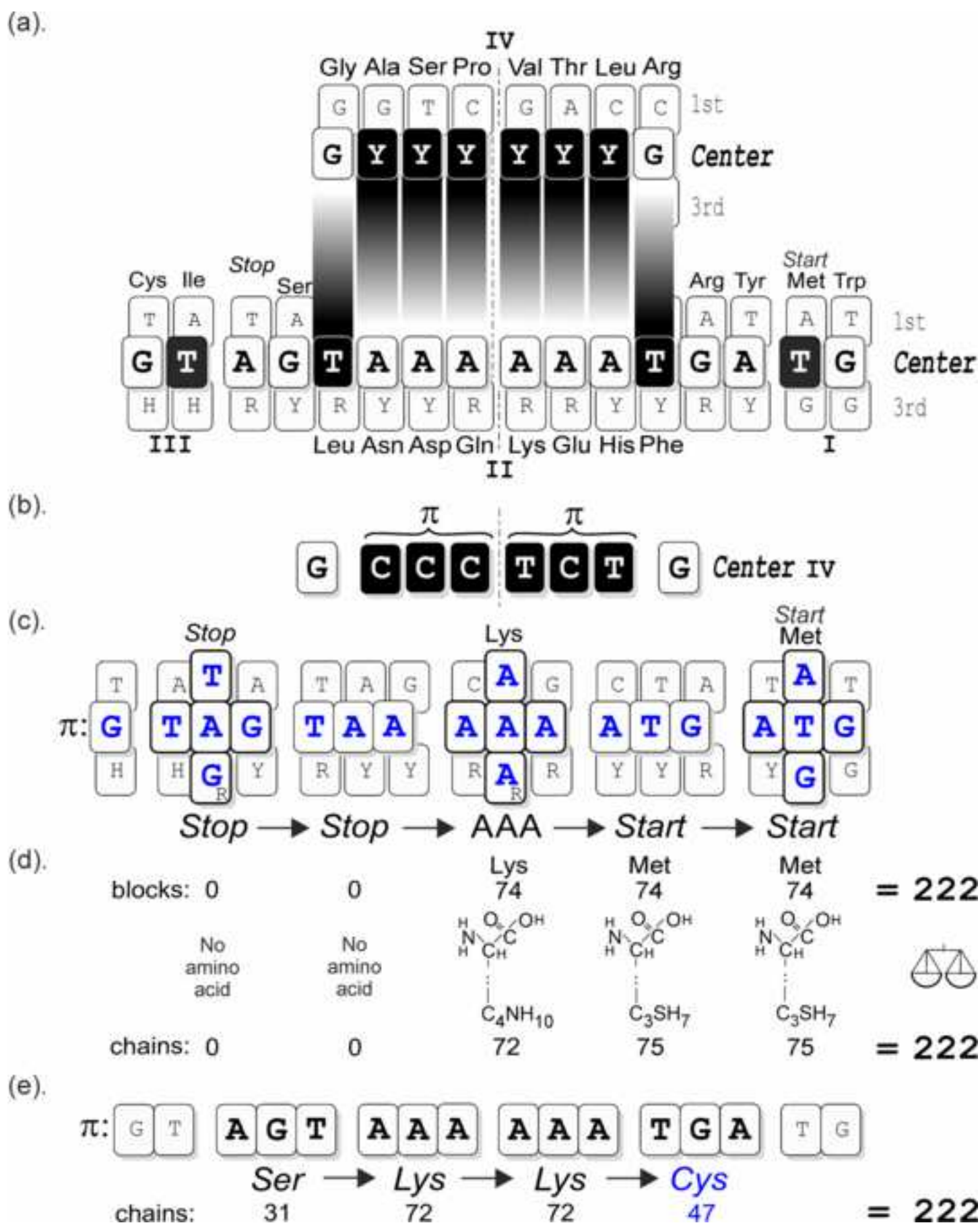


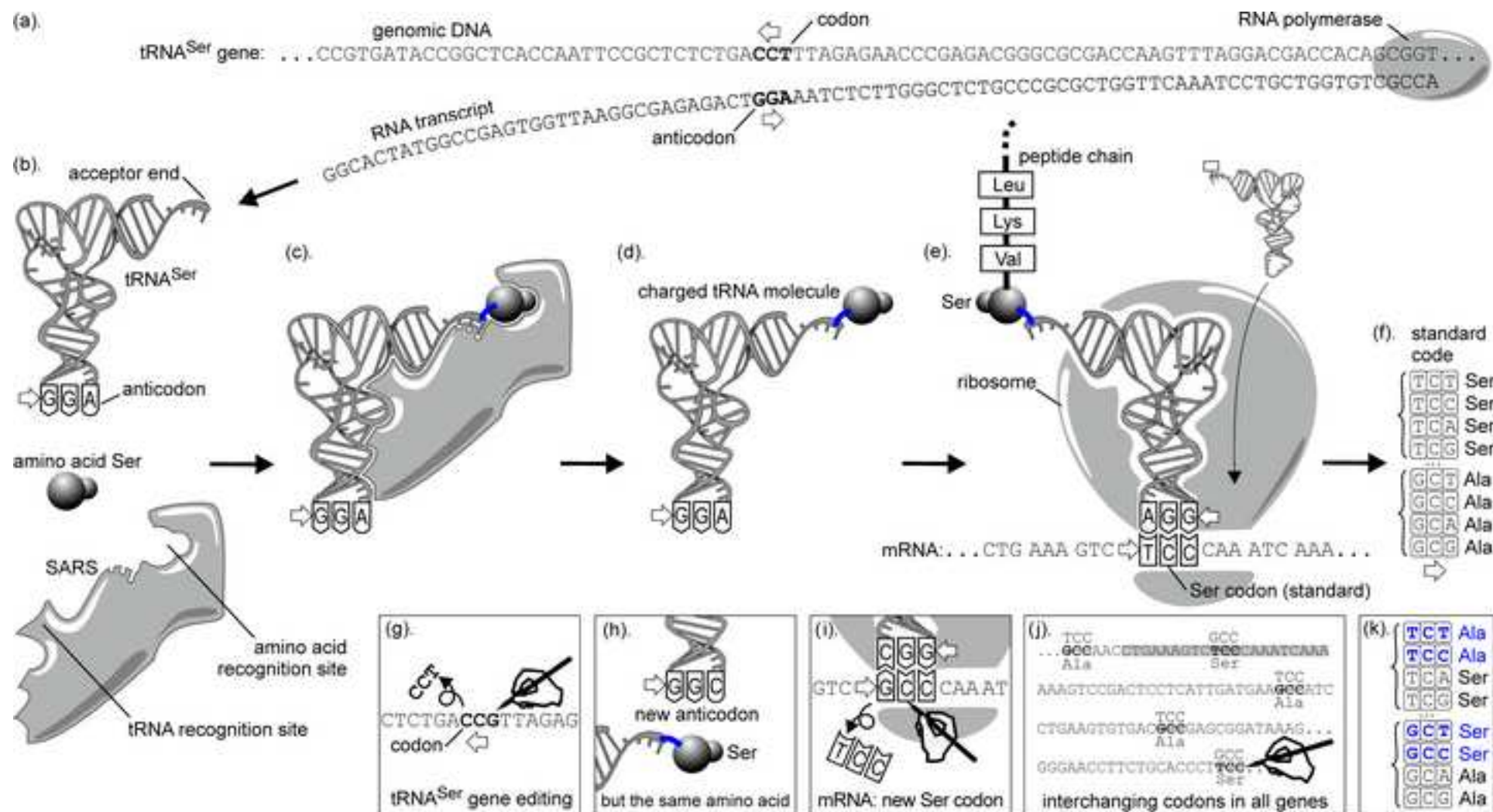
(a).

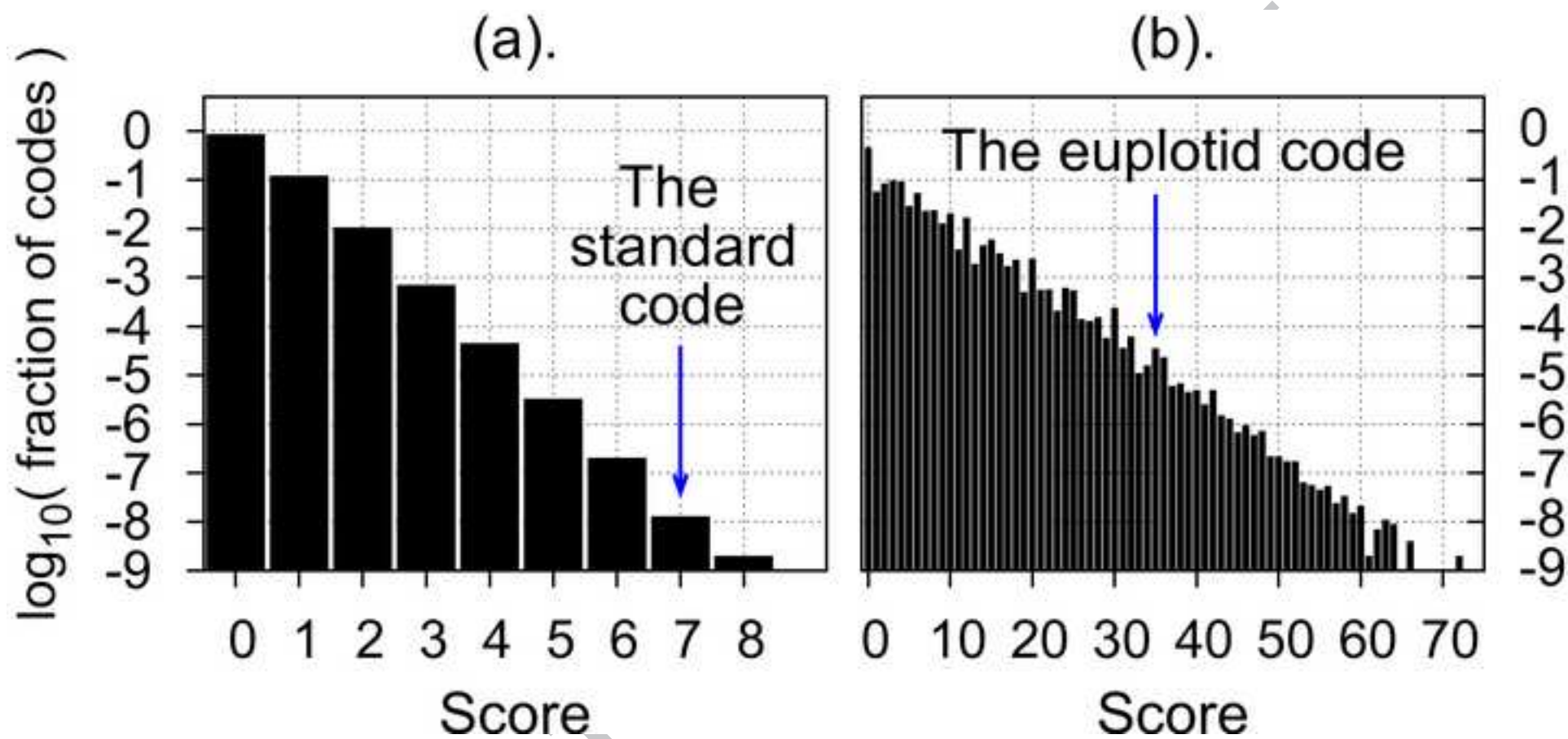


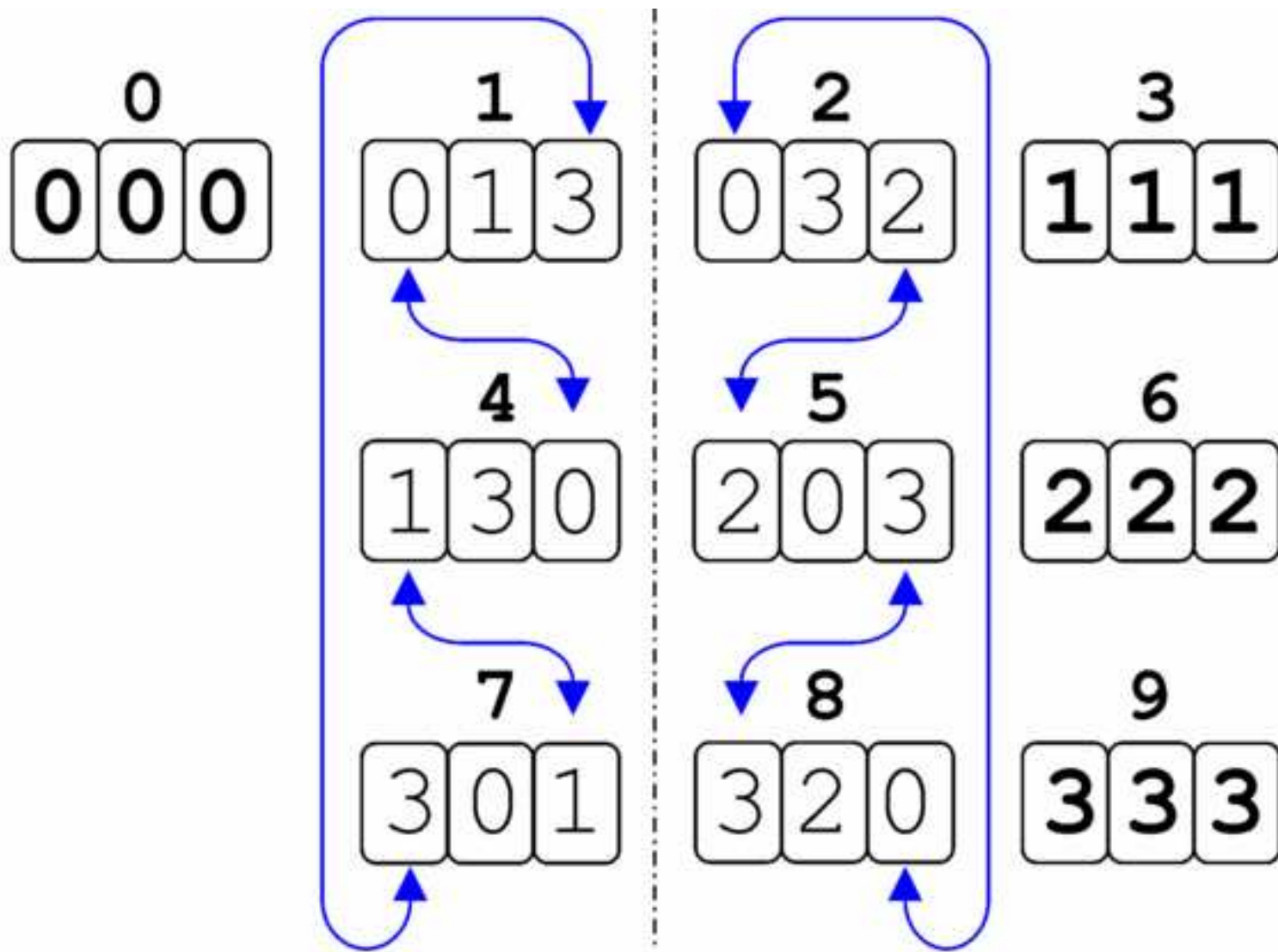
(b).



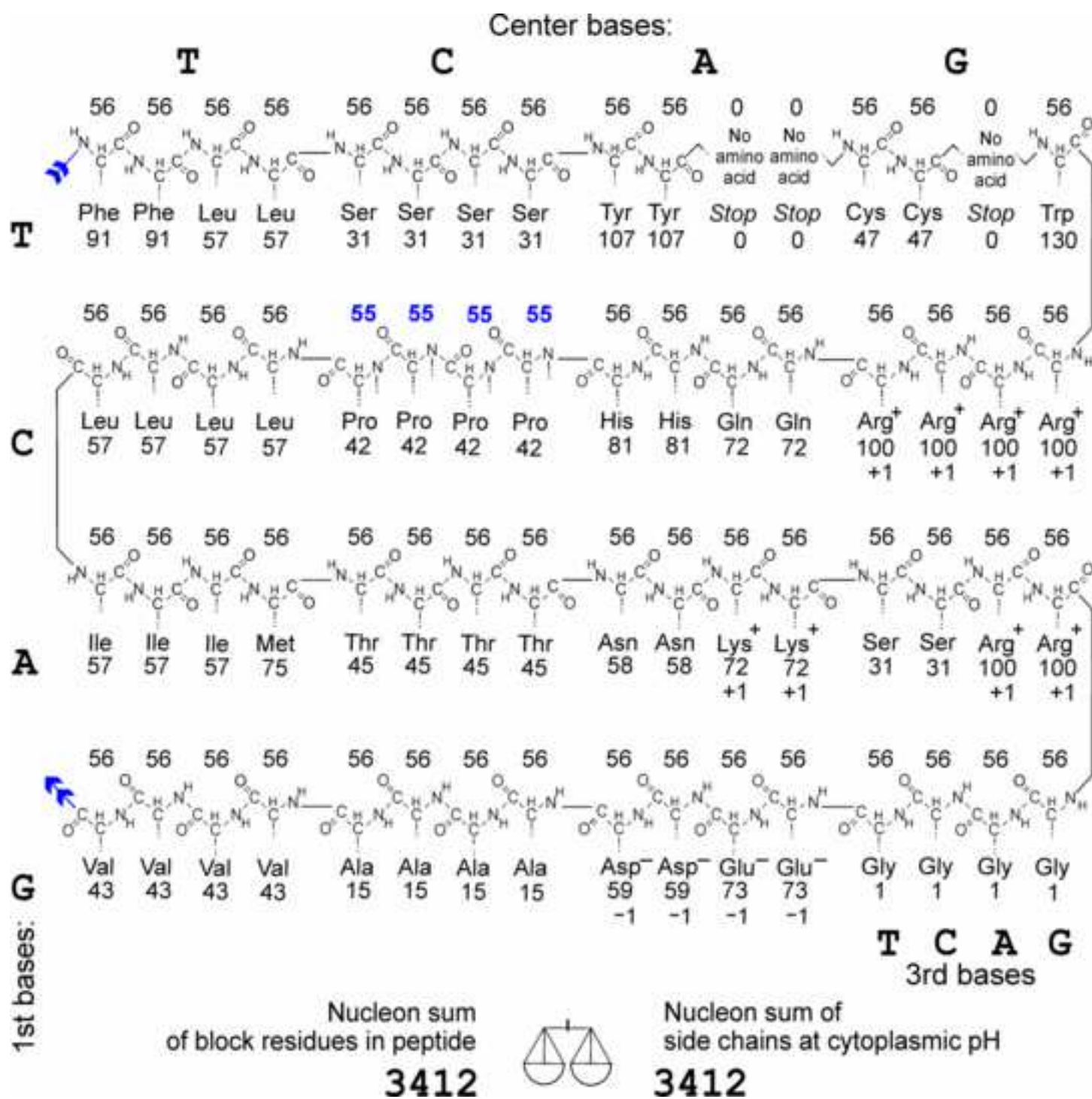








Quaternary $\begin{bmatrix} 0 & 1 & 3 \end{bmatrix}$ and $\begin{bmatrix} 3 & 3 & 3 \end{bmatrix}$ = Decimal 7 and 63



- The SETI hypothesis of an intelligent signal in the genetic code is tested
- The code is shown to possess an ensemble of same-style precision-type patterns
- The patterns are shown to match the criteria of an intelligent signal

ACCEPTED MANUSCRIPT