

# Genome and language – two scripts of heredity

Edward N. Trifonov

University of Haifa, Israel

Beograd, 2013

Both biological sequences and human languages  
are represented by linear scripts on comparable size alphabets:

alphabet size

DNA and RNA	4
Proteins	20 (22)
Morse	4
Hawaiian	18
Hebrew	20
English	26
Polish	32
Thai	59

Apparently, the linear script on small to moderate alphabet size is most economical and least ambiguous way to communicate information.

This, perhaps, is the reason why **in biological evolution first the nucleotide sequences and amino-acid sequences appeared** (in an unknown yet process of initial competition and selection),

**and then the language scripts emerged** ,  
with development of speech apparatus and writing in *Homo sapiens*

**Both genetic sequences and language scripts are products of biological evolution,**

**both inherited in their own ways, and both are subjects of broadly understood natural selection**

According to developing theory of origin of genomes  
(Frenkel ZM, ENT, J Biomol Str & Dynamics 2012)

Genomes and genes emerged first as simple repeating sequences  
which gradually accumulated useful mutational changes,  
while new simple sequences continued to appear (and mutate)  
in the genomes

# “repetitive elements

(simple sequence repeats and transposable elements)

may comprize over two-thirds of the human genome”

(De Koning APJ *et al. PLoS Genet*, 2011)

# 15-mers of human genome (sorted)

1	1	198	780	TTTTTTTTTTTTTTTTT	$T_n$
2	1	190	667	AAAAAAAAAAAAAAAAA	$A_n$
3		366	285	TGTGTGTGTGTGTGT	$TG_n$
4		362	623	ACACACACACACACA	$AC_n$
5		348	215	GTGTGTGTGTGTGTG	$GT_n$
6		344	421	CACACACACACACAC	$CA_n$
7		223	424	GCTGGGATTACAGGC	Alu
8		223	011	GCCTGTAATCCCAGC	Alu
9		222	894	TATATATATATATAT	$TA_n$
10		222	730	ATATATATATATATA	$AT_n$
11-67					Alu
68		169	033	TTTTTTTTTTTTTTTG	$T_n$
69-72					Alu
73		167	889	CAAAAAAAAAAAAAAA	$A_n$
74		167	361	CTAAAATAACAAAAA	Alu
75		150	349	CTTTTTTTTTTTTTTT	$T_n$
76		149	748	AAAAAAAAAAAAAAG	$A_n$
77-82					Alu

-----

## Three known pathologically expanding (“aggressive”) classes of triplets

**GCU** (GCU, CUG, UGC, AGC, GCA, CAG) ,

**GCC** (GCC, CCG, CGC, GGC, GCG, CGG) and

**GAA** (AAG, AGA, GAA, CTT, TTC, TCT).

They cause neurodegenerative diseases and chromosome fragility

# According to the Theory of Early Molecular Evolution based on the Evolutionary Chart of Codons

(Trifonov, E. N., Consensus temporal order of amino acids and evolution of the triplet code. Gene 2000  
Trifonov, E. N. The triplet code from first principles. J Biomolec Str Dyn 2004)

the very first genes have been (aggressive) repeats

...GGC GGC GGC GGC GGC GGC...

and complementary

...GCC GCC GCC GCC GCC GCC...

encoding Gly<sub>n</sub> and Ala<sub>n</sub>, respectively



# Life is self-reproduction with variations

Trifonov, E. N., Origin of the genetic code and of the earliest oligopeptides, Res. Microbiol. 2009

Trifonov, E. N. Vocabulary of definitions of life suggests a definition, J Biomolec Str Dyn. 2011

Any system capable of  
replication and mutation  
is alive (Oparin 1961).

self-reproduction and variation

Could it be that protein-coding sequences,  
actually, are ALL originally made  
from the simple tandem repetitions?

We don't recognize all the original repeats  
just because they have  
extensively mutated.

If this view is correct, then we should see in mRNA sequences

1. Ideal repeats of some codons
2. Imperfect, mutated repeats. In particular, the codons “sandwiched” between two identical codons should be often their point mutation derivatives
3. Those codons which are more frequent in tandem repeats should be also of higher usage in non-repeats

We, thus, undertook analysis  
of the largest non-redundant database of mRNAs available,  
of total ~5 000 000 000 codons,  
from eukaryotes, prokaryotes, viruses, organelles together

Z. Frenkel, E. Trifonov, JBSD, 30, 201-210 (2012)

## Sorted occurrence of the triplet repeats for different groups ("aggressive" triplets)

	group of codons	Occurrence
1	<b>GCC, CCG, CGC, GGC, GCG, CGC</b>	<b>1 784302</b>
2	<b>GCA, CAG, AGC, UGC, GCU, CUG</b>	<b>1 436660</b>
3	<b>GAA, AAG, AGA, UUC, UCU, CUU</b>	<b>1 131214</b>
4	AAU, AUA, <b>uaa</b> , AUU, UUA, UAU	932105 (1 118526)
5	AUC, UCA, CAU, GAU, AUG, <b>uga</b>	735397 (882476)
6	ACC, CCA, CAC, GGU, GUG, UGG	726443
7	AGG, GGA, GAG, CCU, CUC, UCC	706484
8	AAC, ACA, CAA, GUU, UUG, UGU	694387
9	ACG, CGA, GAC, CGU, GUC, UCG	533888
10	ACU, CUA, UAC, AGU, GUA, <b>uag</b>	152747 (183296)

**1** . Tandem repeats of all 61 different codons are observed, strongest for aggressive groups, **as expected**

## 2. Middle codons abc

in “sandwiches” **GCU**abc**GCU**  
(total 3 168 933)

<b>GCU</b>	243706	
<b>GGU</b>	125946	
<b>GAU</b>	115500	
GAA	114278	the topmost in overall codon usage
<b>GUU</b>	102550	
<b>GCA</b>	95493	
<b>GCC</b>	92153	
AUU	89648	
UUU	87861	
AAA	84194	next topmost in codon usage
UUA	80660	
GGA	74934	
GGC	71770	
...		This also holds for most of other codons

## 2. The first derivatives between the identical codons in mRNA keep memory of initial tandem repetition of the codons

The sequences of the type

XYZ nnn nnn XYZ nnn nnn nnn nnn nnn XYZ

are likely descendants of

XYZ XYZ XYZ XYZ XYZ XYZ XYZ XYZ...



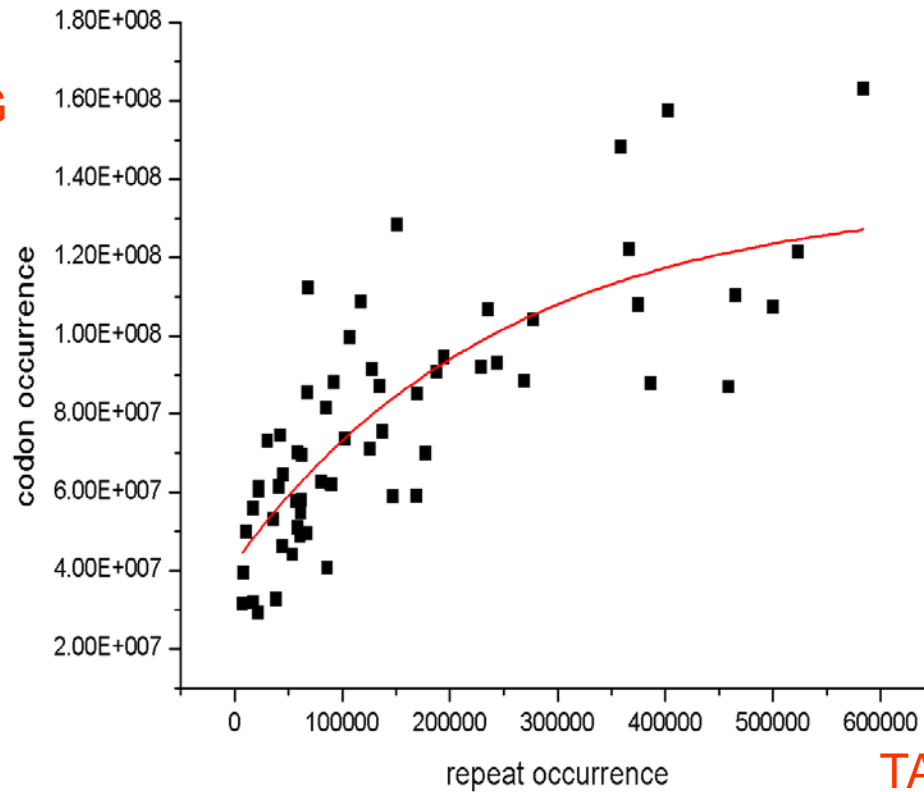
ATG GCT CTA ACC AAA GAA GAT ATT TTA AAC **GCA** ATT GCT **GAA** ATG CCA **GTA** ATG  
**GAC** CTT GTT **GAG** CTT ATC **GAA** GCT **GCA** **GAA** **GAA** AAA TTC GGT **GTA** ACA **GCT** ACT  
**GCT** **GCT** GTT **GCT** GCC **GCT** **GCT** CCT **GCT** **GCT** GGC GGT GAA **GCT** **GCT** GCA GAA CAA  
 ACT GAA TTT GAT GTT GTT TTG ACA TCT TTC GGT GGT AAC AAA GTT **GCT** **GTA** ATC  
**AAA** GCG **GTA** CGT GGC **GCA** ACT GGT CTT GGC TTG **AAA** **GAA** GCT AAA **GAA** **GTA** GTT  
**GAA** GCT **GCA** CCG AAA GCG ATT AAA **GAA** GGC GTT GCT **AAA** **GAA** **GAA** GCT **GAA** **GAA**  
 CTT AAG AAG ACG CTT GAA GAA GCT GGC GCT GAA GTT GAG CTT AAG

**GAA** and **GCT** “bricks” in mRNA of  
 ribosomal protein L12 of *Ps. Atlantica*

Frequent triplets make clusters,  
 remnants of original ideal repeats

3. The more frequently a given codon appears in tandem the more frequent it is also in non-repeating regions of mRNA

NON-REPEATING  
REGIONS



TANDEM

HALF

This result came as a surprize,  
considering **zillions of factors**  
**known to influence the codon usage**

The dominant codons are frequent because  
**they keep memory of**  
**their tandem repetition**  
in the past

**The triplet expansion of codons**  
**is the major single factor**  
**shaping the codon usage**

Thus, life started with the replication (and expansion) and subsequent mutations of tandemly repeating triplets GGC and GCC.

(self-reproduction with variation)

Life continued then to spontaneously emerge within the primitive early genomes and further on, in form of replication and expansion and subsequent mutations of other tandem repeats as well

(self-reproduction with variation)

**Life never stopped emerging**

Evolution of genetic sequences:

First simple repeats,

Then mutated forms of the repeats,

New repeats,

New mutations and insertions,

Self-reproduction all the way.

Self-reproduction with variations

Consonants, easily pronounceable by babies:

p, b, t, m, d, n, k, g, s, h, w, j,

Difficult to pronounce: f, v, th, sh, ch, l, r

This is why babies, 4-7 month old, babble

Pi-pi, Ba-ba, Ti-ti, Ma-ma etc.

even before they learn what these words may mean

These very sounds, most likely, have been babbled by babies of the earliest speaking hominids, because respective muscles of speech apparatus also appeared first in the developing human embryo, according to rule “Ontogeny recapitulates phylogeny”. (Ernst Haeckel).

Thus, babies of all hominids should have been babbling the same “words”, before the “words” acquired their meanings, obviously different at various times, places and evolutionary stages of the hominids.



Today European mother would enthusiastically respond to spontaneous “ma-ma-ma”, thus, establishing and further consolidating the first liaison of baby words with reality.

Georgian mother would react the same way to “da-da-da” (“dada” is mother in Georgian),

while Swahili speaking mother (“baba”) would respond to “ba-ba-ba”.

## The same baby words in different languages

Papa	mother (Jap)	father (Rus)	grandfather (Georg)
Baba	grandmother (Rus)	father (Bengali)	baby (Arabic)
Titi	breast (Rus)	father (Jap)	
Mama	food (Jap)	mother (Rus)	father (Georg)
Nana	mother (Fijian)	food (Arabic)	father (Telugu)
Deda (dada)	grandfather (Rus)	mother (Georg)	sleep (Arabic)
Kaka (caca)	feces (Intl)		
Gaga	geese (Rus)		
Sisi	breast (Rus)	bird (Arabic)	
Haha	mother (Jap)		
Weewee	pee, penis (Eng)		
Jojo	toy (Eng)		

## Sound imitations, mostly babies

Av-av (dog)  
Bi-bi (car)  
Cococo (chicken)  
Kva-kva (frog)  
Tik-tak (clock)  
Din'din' (ringbell)  
Ga-ga-ga (geese)  
Kria-kria (duck)  
Tuk-tuk-tuk (knocking)  
Kap-kap-kap (rain)  
Chmok-chmok (kisses)  
Top-top-top (walk)  
Skirly-skirly (wooden leg)

Rooster:

Ku ka re ku (Rus)  
Ku ke le ku (Dutch)  
Ki ke ri ki (German)  
Co co ri co (French)  
Cock-a-doodle-doo (English)

## Adult forms, perfect repeats:

O-o (warning)

Bebe

Da-da (come in)

Ja-ja (yes, German)

Ku-ku (crazy)

Ga-ga (crazy, English)

Hahaha

Nununu (warning to babies)

Tuktuk (Cambodia, Thailand, moto-rickshaw)

Tamtam (drum)

Tak-tak (all right)

Ks-ks-ks (calling cat)

Nuka-nuka (go ahead)

Chachacha

Leat-leat (slowly, Hebrew)

Tipa-tipa (little bit, Hebrew)

Tilki-tilki (barely fit, Ukrainian)

Trochi-trochi (little bit, Ukrainian)

Rock-rock-rock (Kenya, lullaby)

Langsam-langsam (slowly, Yiddish)

## **Adult forms, perfect repeats (mostly Russian):**

E-e (warning)

Ohoho (that much)

Mimimi (sweaty, cuty)

Bumbum (ignorant)

Lalala (empty talk)

Tsatsa (girl showing up)

Vot-vot (in a moment)

Idu-idu (coming)

Kto-kto? (who)

Gde-gde? (where)

Vas`-vas` (friends)

Tiny-tiny

Jele-jele (barely)

Kuda-kuda? (where)

Tolko-tolko (barely fit)

Chut`-chut` (little bit)

Hei-hei-hei (warning)

Chevo-chevo? (what)

Tsip-tsip-tsip (calling chicken)

Skolko-skolko? (how much)

Kak eto, kak eto? (why all of a sudden)

## **Mutated, imperfect repeats, babies and adults:**

Mamy (mother, English)

Baby

Bibika (car)

Mamaya (fruit, Brazil)

Papaya (similar fruit, Brazil)

O-la-la (surprize, French)

Cocook

To-to-je (Aliska, co to je, Czech)

Ta-ra-ram (mess)

Balalaika

Tarataika (type of a cart)

Yin`-yan` (Chinese)

Siusiukat` (imitate baby-talk)

Tsap-tsarap (catch, about cats)

Villi-nilli (against will, Latin)

Meli, Emelia (talking nonsense)

Olgoi-horhoi (Mongolian, ferrytale creature)

Volens-nolens (against will, Latin)

Naziuziukalsa (drunk)

Futy-nuty, lapti gnuty (mishap)

## **Mutated, imperfect repeats, babies and adults:**

Nu-i-nu (surprized)

Kukushka (coocook)

Coca-cola

Tra-ta-ta (thunder)

Futy-nuty (mishap)

Tiap-liap (lousy work)

Trali-vali (menstruation)

Dura duroi (stupid, her)

Figli-migli (flirt)

Shito-kryto (everything is fine)

Tram-tararam (mess)

Durak durakom (stupid, he)

Boogie-woogie

Trach-tararach (thunder)

Postolku-poskolku (as soon as)

Baiu-baiushki-baiu (lullaby)

Tiutelka v tiutelku (just exactly fit)

Martin Luther King, 1968:

"Yes, if you want to  
say that **I was a drum major**,  
say that **I was a drum major** for justice.  
Say that **I was a drum major** for peace.  
          **I was a drum major** for righteousness."

Criticized misquote:

"I was a drum major for justice,  
                                  for piece,  
                                  for righteousness."



...rhythm is an integral part of language.

(BBC Science, **TODAY**)

# Binary alphabet alternations

## In human languages

- alternation of consonants and vowels, like  
divinity (CVCVCVCV), wandering (CVCCVCCVCC), ammunition (VCCVCCVCCVCC)

for better “pronounceability”, and

## In protein sequences

– alternation of polar and non-polar aa residues,

like PPNNPNNNPNNPPNPPNPPNPPNPPN... with the period  $\sim 3.5$ ,  
in amphipathic alpha helices

Zemkova M., Trifonov E. N., Zahradnik D. One common structural feature of “words” in protein sequences and human texts. J. Biomol. Str. Dyn. 2013

Evolution of language scripts:

First simple repeats (baby-talk)

Then mutated forms of the repeats (advanced baby-talk)

New repeats (adult forms, consonant/vowel alternations),

New mutations and insertions (new texts),

(assisted) self-reproduction (rewriting, reprinting) all the way.

Self-reproduction with variations

## Biological sequences and languages evolve by the same scenario:

Evolution of **genetic sequences**:

First simple repeats,  
Then mutated forms of the repeats,  
New repeats,  
New mutations and insertions,  
Self-reproduction all the way.

Self-reproduction with variations

Evolution of **language scripts**:

First simple repeats,  
Then mutated forms of the repeats,  
New repeats,  
New mutations and insertions,  
Self-reproduction all the way.

Self-reproduction with variations

We are not the only possessors of the languages on the planet.

The biological succession from sequences to language scripts is not limited to *Homo sapiens* only

**Dolphins also talk to each other.**

For example, they have special phonogram words for personal names of the dolphins within the group.

Other words (and “letters”?) are still to be deciphered.

(SL King and VM Janik, PNAS July 22, 2013)

**The dolphin ultrasonic series are linear “script” as well**

Both genomes and languages are based on linear scripts,

Scripts of biological heredity,  
And scripts of cultural heritage.

Both are subjects of evolution and natural selection

THANKS TO

ZOHAR **KOREN,**

ZAKHARIA **FRENKEL,**

ALEXANDRA **RAPOPORT,**

THOMAS **BETTECKEN**

MISA **ZEMKOVA**



**Haifa**

**München**

**Prague**







# Baby talk words, perfect repeats

(Russian, if not specified)

Mama

Papa

Baba (grandma)

Pipi

Caca

Sisi (breast)

Bobo (pain)

Baibai (good night)

Tiatia (father)

Niania (nanny)

Ham-ham (eat, Vietnamese)

Ai-ai-ai (mishap)

Ne-ne-ne (no, Czech)

Wong-wong (drink, Vietnamese)

## **Baby talk words, perfect repeats**

Lala (doll, baby)

Kuku (from hiding)

Diadia (man)

Oi-oi-oi (mishap)

Ni-ni-ni (strictly no)

Niam-niam (eat)

Dai-dai-dai (give me)

# Mooring steamer to a pier

Sound imitations from “Adventures of Tom Sawyer” by Mark Twain:

He was boat and captain and engine-bells combined, so he had to imagine himself standing on his own hurricane-deck giving the orders and executing them:

"Stop her, sir! **Ting-a-ling-ling!**" The headway ran almost out, and he drew up slowly toward the sidewalk.

"Ship up to back! **Ting-a-ling-ling!**" His arms straightened and stiffened down his sides.

"Set her back on the stabboard! **Ting-a-ling-ling! Chow! ch-chow-wow! Chow!**"

His right hand, mean-time, describing stately circles—for it was representing a forty-foot wheel.

"Let her go back on the labboard! **Ting-a-ling-ling! Chow-ch-chow-chow!**"

The left hand began to describe circles.

"Stop the stabboard! **Ting-a-ling-ling!** Stop the labboard! Come ahead on the stabboard!

Stop her! Let your outside turn over slow! **Ting-a-ling-ling! Chow-ow-ow!**

Get out that head-line! *lively* now! Come—out with your spring-line—what're you about there! Take a turn round that stump with the bight of it! Stand by that stage,

now—let her go! Done with the engines, sir! **Ting-a-ling-ling! SH'T! S'H'T! SH'T!**"

(trying the gauge-cocks).

# Counting rhymes for various games

Ene bene rech  
Kenter menter zhech  
Ene bene raba  
Kenter menter zhaba

Eniki beniki  
Eli vareniki  
Eniki beniki klotz

Ine mine  
Minke tinke  
Fade rude  
Rolke tolke  
Wigel wagel weg (German)

# EVOLUTION OF THE TRIPLET CODE

E. N. Trifonov, December 2007, Chart 101

Consensus temporal order of amino acids:

	UCX	CUX	CGX	AGY	UGX	AGR	UYU	UAX														
<u>Gly Ala</u>	<u>Asp Val</u>	<u>Ser Pro</u>	<u>Glu Leu</u>	Thr	Arg	Ser	TRM	Arg	Ile	Gln	Leu	TRM	Asn	Lys	His	Phe	Cys	Met	Tyr	Trp	Sec	Pyl

1	GGC-GCC																					
2			GAC-GUC																			
3	GGA--	---	---	--UCC																		
4	GGG--	---	---	---	--CCC																	
5			(gag)-	---	---	--GAG-CUC																
6	GGU--	---	---	---	---	---	--ACC															
7		GCG--	---	---	---	---	---	--CGC														
8		GCU--	---	---	---	---	---	--AGC														
9		GCA--	---	---	---	---	---	--ugc									UGC					
10						CCG--	---	--CGG														
11						CCU--	---	---	--AGG													
12						CCA--	---	---	--ugg										UGG			
13				UCG--	---	---	---	--CGA														
14				UCU--	---	---	---	---	--AGA													
15				UCA--	---	---	---	---	--UGA													UGA
16								ACG-CGU														
17								ACU--	--AGU													
18								ACA--	---	--ugu								UGU				
19			GAU--	---	---	---	---	---	--AUC													
20				GUG--	---	---	---	---	---	--cac								CAC				
21								CUG--	---	---	--CAG											
22										aug-cau								CAU		AUG		
23						GAA--	---	---	---	--uuc								UUC				
24				GUA--	---	---	---	---	---	--uac										UAC		
25						CUA--	---	---	---	--UAG												UAG
26				GUU--	---	---	---	---	---	--AAC												
27						CUU--	---	---	---	--AAG												
28										CAA-UUG												
29										AUA--	---	--uau										UAU
30										AUU--	---	---	--AAU									
31												UUA-UAA										
32												uuu--	---	--AAA				UUU				

CONSECUTIVE ASSIGNMENT OF 64 TRIPLETS

CODON CAPTURE

aa "age":

17	17	16	16	15	14	13	13	12	11		10	9		8	7	6	5	4	3	2	1
----	----	----	----	----	----	----	----	----	----	--	----	---	--	---	---	---	---	---	---	---	---

"... if **variations** useful to any organic being ever do occur, assuredly individuals thus characterized will have the best chance of being preserved in the struggle for life; and from the strong principle of inheritance, these will tend to **produce offspring similarly characterized**"

*Charles Darwin, Origin of Species (1859)*

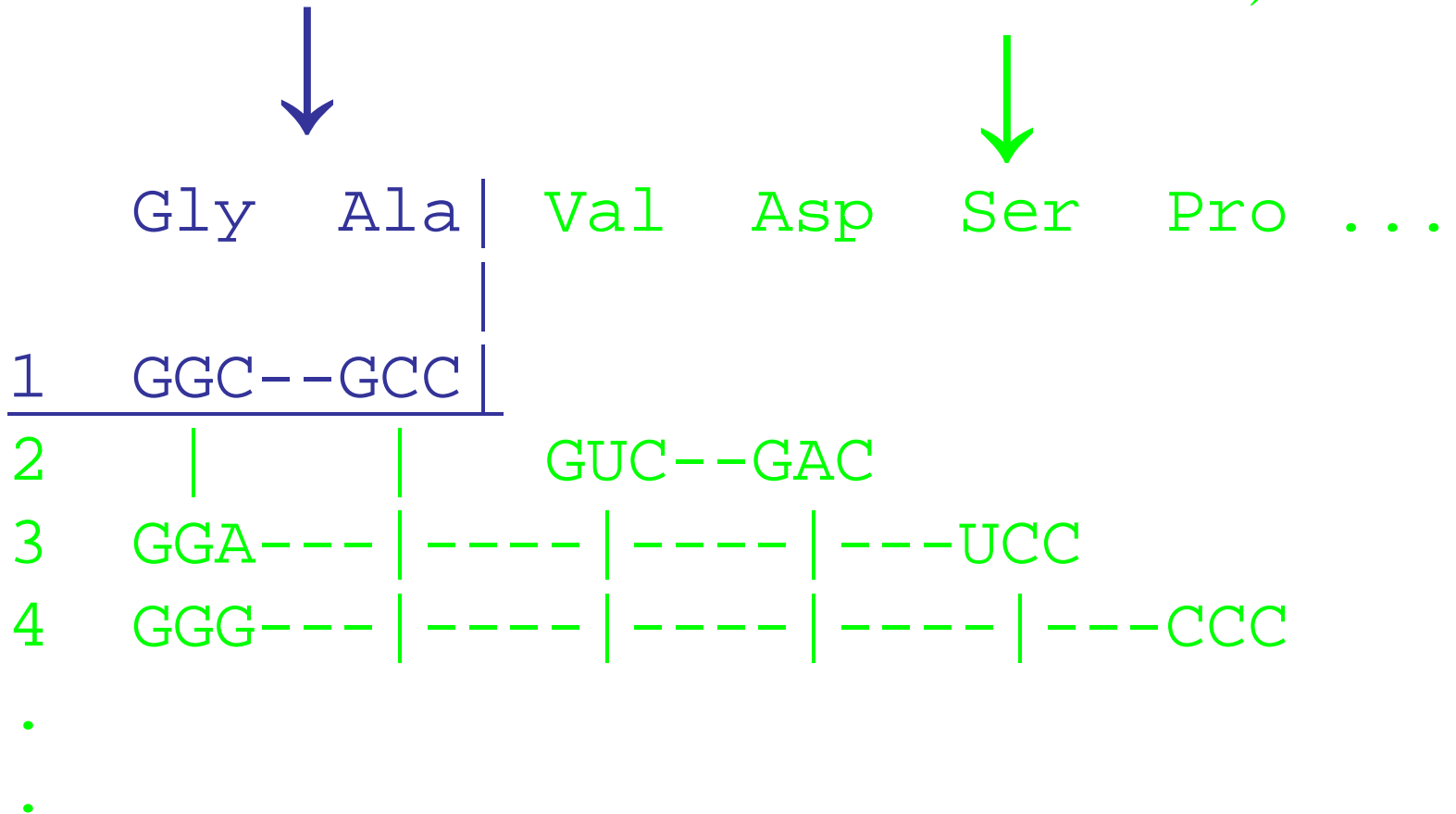
Rephrasing (ET):

Individuals with useful **variations** will **self-reproduce**

**self-reproduction and variations**

not Life yet  
(self-reproduction only)

Life  
(self-reproduction  
and variations)





# From vocabulary of 123 known definitions of life the following groups of meanings are revealed

<b>LIFE</b>	123
living	47
alive	10
being	6
biological	5
other related words	8
Sum	<b>199</b>

<b>SYSTEM</b>	43
systems	22
organization	14
organism	14
order	6
organisms	6
network	5
organized	5
other related words	40
Sum	<b>155</b>

<b>MATTER</b>	25
organic	11
materials	10
molecules	6
other related words	36
Sum	<b>88</b>

<b>CHEMICAL</b>	17
process	15
metabolism	14
processes	8
reactions	5
other related words	26
Sum	<b>85</b>

<b>COMPLEXITY</b>	13
information	8
complex	7
other related words	46
Sum	<b>74</b>

<b>REPRODUCTION</b>	10
reproduce	8
replication	7
self-reproduction	5
other related words	33
Sum	<b>63</b>

<b>EVOLUTION</b>	10
evolve	7
change	6
mutation	5
other related words	20
Sum	<b>48</b>

<b>ENVIRONMENT</b>	20
external	6
other related words	15
Sum	<b>41</b>

<b>ENERGY</b>	18
force	5
other related words	17
Sum	<b>40</b>

<b>ABILITY</b>	12
able	11
capable	11
capacity	5
other related words	1
Sum	<b>40</b>

# Life (*definiendum*)

## *Definientia:*

System

Matter

Chemical

Complexity

Reproduction

Evolution

Environment

Energy

Ability

These appear to be both  
necessary and sufficient  
for the definition of life

**We, thus, come again to the same definition:**

**Life is self-reproduction with variations**

## **Aggressive amino acids encoded by expanding triplets**

<b>Amino acid</b>	<b>Triplets</b>
<b>L</b> (leucine)	<b>CTG CTT</b>
<b>A</b> (alanine)	<b>GCT GCA GCC GCG</b>
<b>G</b> (glycine)	<b>GGC</b>
<b>P</b> (proline)	<b>CCG</b>
<b>S</b> (serine)	<b>AGC TCT</b>
<b>E</b> (glutamate)	<b>GAA</b>
<b>R</b> (arginine)	<b>CGG CGC AGA</b>
<b>Q</b> (glutamine)	<b>CAG</b>
<b>K</b> (lysine)	<b>AAG</b>
<b>F</b> (phenylalanine)	<b>UUC</b>
<b>C</b> (cysteine)	<b>UGC</b>

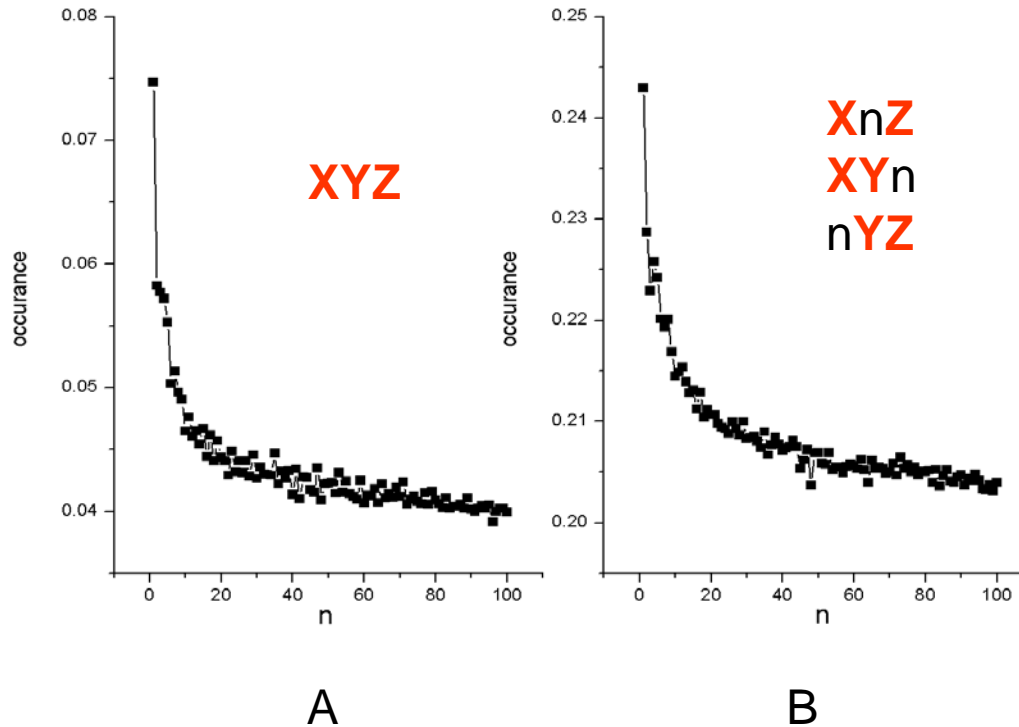
# Majority of homopeptides are built from aggressive amino acids

human tripeptides 1st exons	Score (tripept.)	eukar. (Faux et al.)	prokar. (Faux et al.)
<b>1. L3</b>	<b>4552</b>	<b>1446</b>	<b>70(5)</b>
<b>2. A3</b>	<b>4046</b>	<b>5465(3)</b>	<b>251(3)</b>
<b>3. G3</b>	<b>2972</b>	<b>5002(5)</b>	<b>310(2)</b>
<b>4. P3</b>	<b>2258</b>	<b>4157(7)</b>	<b>217(4)</b>
<b>5. S3</b>	<b>1981</b>	<b>5424(4)</b>	<b>378(1)</b>
<b>6. E3</b>	<b>1630</b>	<b>4334(6)</b>	<b>67(6)</b>
<b>7. R3</b>	<b>1145</b>	<b>462</b>	<b>60(8)</b>
<b>8. Q3</b>	<b>802</b>	<b>8022(1)</b>	<b>52(9)</b>
<b>9. K3</b>	<b>535</b>	<b>1920(9)</b>	<b>25</b>
-----			
10. V3	414	94	9
11. H3	273	1049	32
12. D3	269	1554	34
13. T3	267	2492(8)	63(7)
14. I3	109	34	3
<b>15. F3</b>	<b>103</b>	<b>175</b>	<b>1</b>
<b>16. C3</b>	<b>92</b>	<b>38</b>	<b>0</b>
17. N3	79	6962(2)	31
18. M3	34	19	0
19. Y3	32	39	4
20. W3	14	3	0
	<b>92%</b>	<b>75%</b>	<b>89% (Z. Koren, 2011)</b>

**22.5 min**

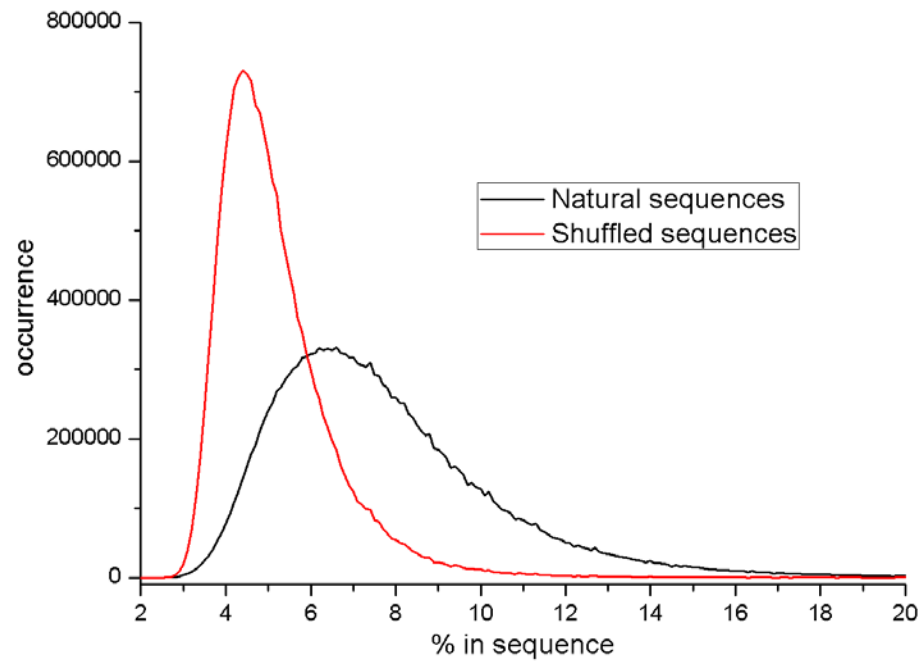
„Thick“ sandwiches

$XYZabc_1abc_2\dots abc_nXYZ$



Occurrence of the triplet  $XYZ$  (A) and its first derivatives (B)  
in the middle sequence  $abc_1abc_2\dots abc_n$

# Enrichment of mRNA sequences by one or another dominant codon





Ala	<b>GCC</b>	<b>110</b>	<b>465</b>	Arg	<b>CGC</b>	<b>70</b>	<b>177</b>	Arg	<b>AGA</b>	<b>55</b>	<b>62</b>
	GCA	94	195		CGU	46	45		AGG	29	22
	GCU	93	245		CGG	41	86				
	GCG	88	386		CGA	33	39				

1<sup>st</sup> columns - codons  
(millions)

Asn	<b>AAU</b>	<b>121</b>	<b>523</b>	Asp	<b>GAU</b>	<b>148</b>	<b>359</b>	Cys	<b>UGC</b>	<b>31.9</b>	<b>18</b>
	AAC	85	170		GAC	107	236		UGU	31.5	7

2<sup>nd</sup> columns - repeats  
(thousands)

Gln	<b>CAA</b>	<b>88</b>	<b>269</b>	Glu	<b>GAA</b>	<b>163</b>	<b>584</b>	<b>Gly</b>	<b>GGC</b>	<b>107</b>	<b>500</b>
	CAG	87	459		GAG	122	367		GGU	92	229
									GGA	87	135
									GGG	56	17

His	<b>CAU</b>	<b>58</b>	<b>62</b>	Ile	<b>AUU</b>	<b>128</b>	<b>151</b>	Leu	<b>UUA</b>	<b>91</b>	<b>127</b>
	CAC	49	61		AUC	100	107		UUG	73	30
					AUA	70	63				

Leu	<b>CUG</b>	<b>108</b>	<b>375</b>	Lys	<b>AAA</b>	<b>158</b>	<b>403</b>	Met	AUG	109	117
	CUU	75	43		AAG	104	277				
	CUC	70	59								
	CUA	40	8								

Phe	<b>UUU</b>	<b>112</b>	<b>68</b>	Pro	<b>CCA</b>	<b>62</b>	<b>89</b>	Ser	<b>UCU</b>	<b>63</b>	<b>81</b>
	UUC	82	85		CCG	59	169		UCA	62	90
					CCU	58	59		UCC	50	67
					CCC	50	11		UCG	44	54

Ser	<b>AGC</b>	<b>59</b>	<b>147</b>	Thr	<b>ACC</b>	<b>76</b>	<b>138</b>	Trp	UGG	60	22
	AGU	53	36		ACA	71	126				
					ACU	65	45				
					ACG	51	59				

Tyr	<b>UAU</b>	<b>86</b>	<b>68</b>	Val	<b>GUG</b>	<b>91</b>	<b>187</b>
	UAC	61	41		GUU	88	92
					GUC	74	103
					GUA	61	23

In 17 of 21 codon repertoires  
**the most frequent codon**  
**is also the most repetitive**

“... if (and oh what a big if) we could conceive in some warm little pond with all sort of ammonia and phosphoric salts, - light, heat, electricity etc., present, that a protein compound was chemically formed, ready to undergo still more complex changes, at the present day such matter would be **instantly devoured, or absorbed,** which would not have been the case before living creatures were formed.” (Darwin 1871)

With the new view on genome origin and evolution the emerging life **is not consumed** by the earlier life, **but rather protected** by the environment within the cell.

The tandem repeats have been considered as a class of “selfish DNA” (Orgel and Crick, 1980; Doolittle and Sapienza, 1980).

They are, actually, more than just parasites tolerated by genome.

They are even more than building material for the genome (Ohno, **Junk DNA**, 1972).

The tandem repeats represent constantly emerging life, and genomes are products of their everlasting domestication.

**Genomes are built by the expansion and mutational domestication of the tandem repeats**

**Genomes ARE the repeats  
(some already unrecognizable)**

# Painful symbiosis of repeats with genomes

## *For genomes*

accepted repeats are useful.

new repeats are dangerous.

## *For repeats*

genomes are natural habitats.

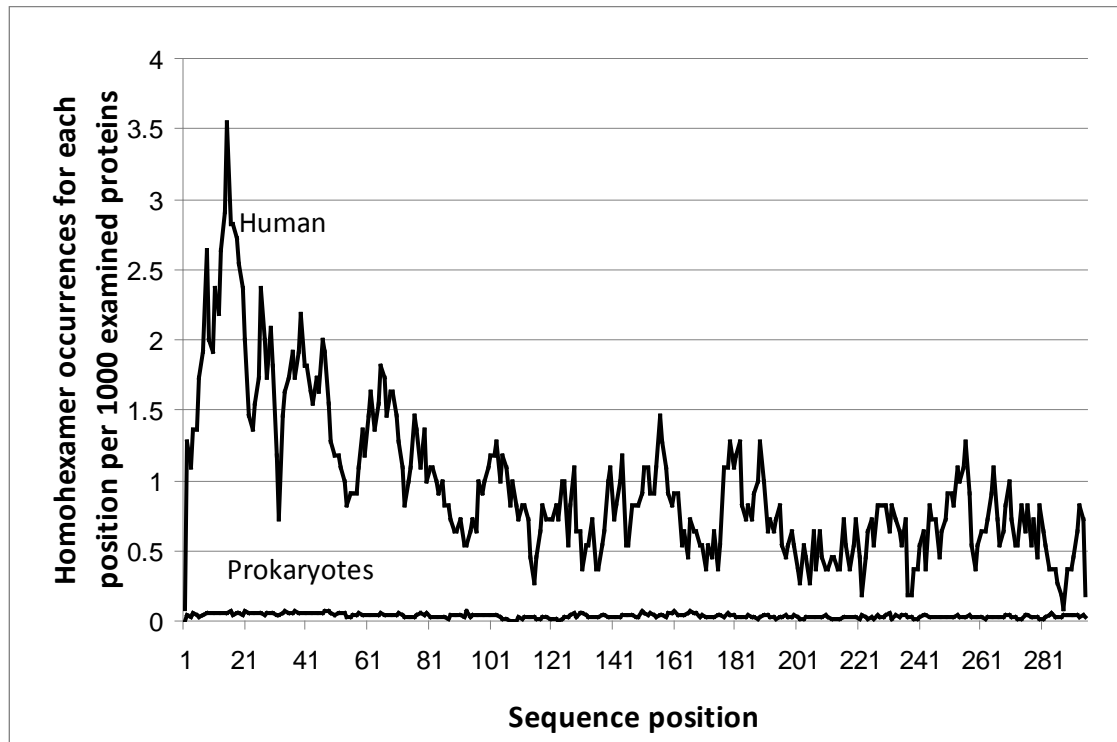
initiation is at high risk

PREDICTION:

**GENOMES SHOULD BE EQUIPPED BY**

**DEFENSE SYSTEMS**

**AGAINST CONSTANTLY EMERGING REPEATS**



The amino acid repeats in prokaryotes are far less frequent compared to eukaryotes.

Defense in prokaryotes:

Brutal negative selection,  
death of individuals contracting the repeats

Defense in eukaryotes:

Expulsion of the repeats into introns and intergenic sequences?  
(Alternative splicing as an intermediate stage)

Possible defense devices:

Prevention of slippage. Nucleosomes.

Excision of slippage loops.

Methylation of repeats.

Sequence-specific nucleases

.....



The simplest life forms – simple tandem repeats –  
represent a whole class of pathological agents,  
not considered as such up to now.

# Genomes evolve under constant attacks by various repeats.

Apparently, most of the attacks are normally stopped by the defense system.

Some of the new expansions or insertions are accommodated by the genomes.

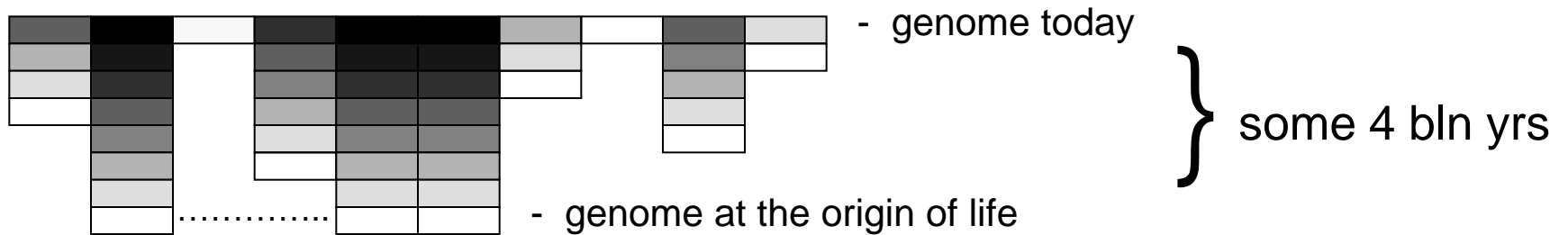
Some are neither stopped, nor accommodated, causing disaster.

A DIFFERENT VIEW ON CANCER, EXPANSION DISEASES  
AND DISEASES WITH UNKNOWN CAUSATIVE AGENT:

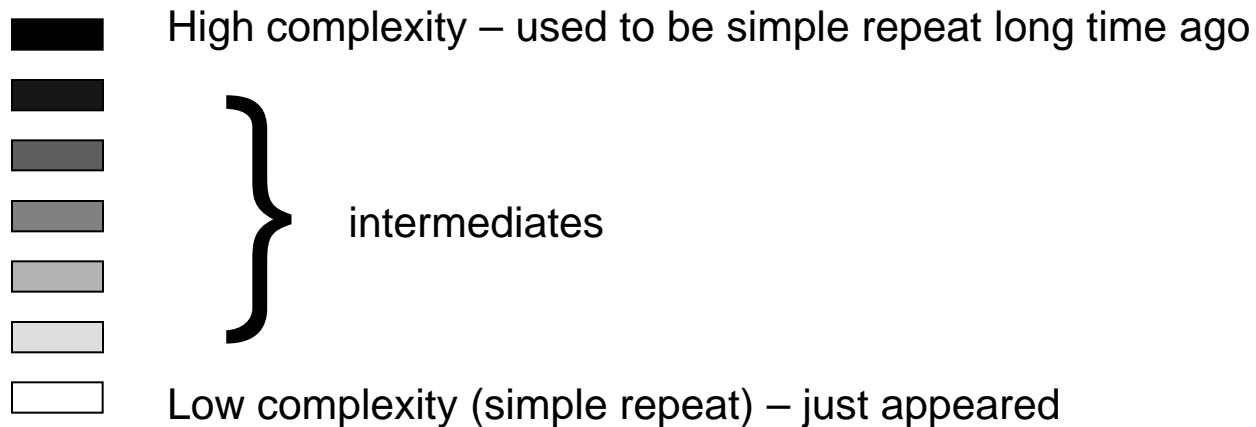
The repeats in the diseases are not **symptoms**.

They are **cause** of the diseases.





**Genomes are all built from simple repeats.  
Just many of them already unrecognizable**



**GAA GAA GAA GAA GAA GAA GAA GAA GAA GAA GAA GAA GAA**

GAA GAA CAA GAA GGA GAU GAA GAA UAC GAG GAA GAA AAA

CAA GAA CAA GGA GGA AAU GAA GCA UAC GAG GAA GGA AAU

CAG GUA CAG GGU GGA AAU GAA GCC UUC GGG GAA CGG ACU

CAG AUA CCG GGU GGG AAU UAC GCC UUC UGG AAA CGG ACU

CCG AUA CCG UGU GGG ACU UAC UCC UUC UGG AAC CGG ACU

**CCG AUC CCG UGU UGG ACU UCC UCC UUC UGG AGC CGG ACU**

83	138448	TTTTTTTTTTTTTTTGA	$T_n$
84	137643	TCAAAAAAAAAAAAAA	$A_n$
85	135070	TTTTTTTTTTTTTTGAG	$T_n$
86	134465	TTTTTTTTTTTTTTGAGA	$T_n$
87	134262	CTCAAAAAAAAAAAAAA	$A_n$
88	133917	TCTCAAAAAAAAAAAAAA	$A_n$
----- Alu and variants of the above			
185	85432	TTTATTTATTTATTT	$TTTA_n$
186	85142	AAATAAATAAATAAA	$AAAT_n$
-----			
293	70591	AGAGAGAGAGAGAGA	$AG_n$
-----			
298	70411	TCTCTCTCTCTCTCT	$TC_n$
-----			
945	33435	AATAATAATAATAAT	$AAT_n$
-----			
999	31742	CTTCCTTCCTTCCTT	$TTCC_n$
-----			

The list ends at line ~700 000 000

~300 000 000 15-mers do not appear at all  
(of total 1 073 741 824)

GCTGGGATTACAGGC

GCT RYY

GGG RRR

ATT RYY

ACA RYR

GGC RRY

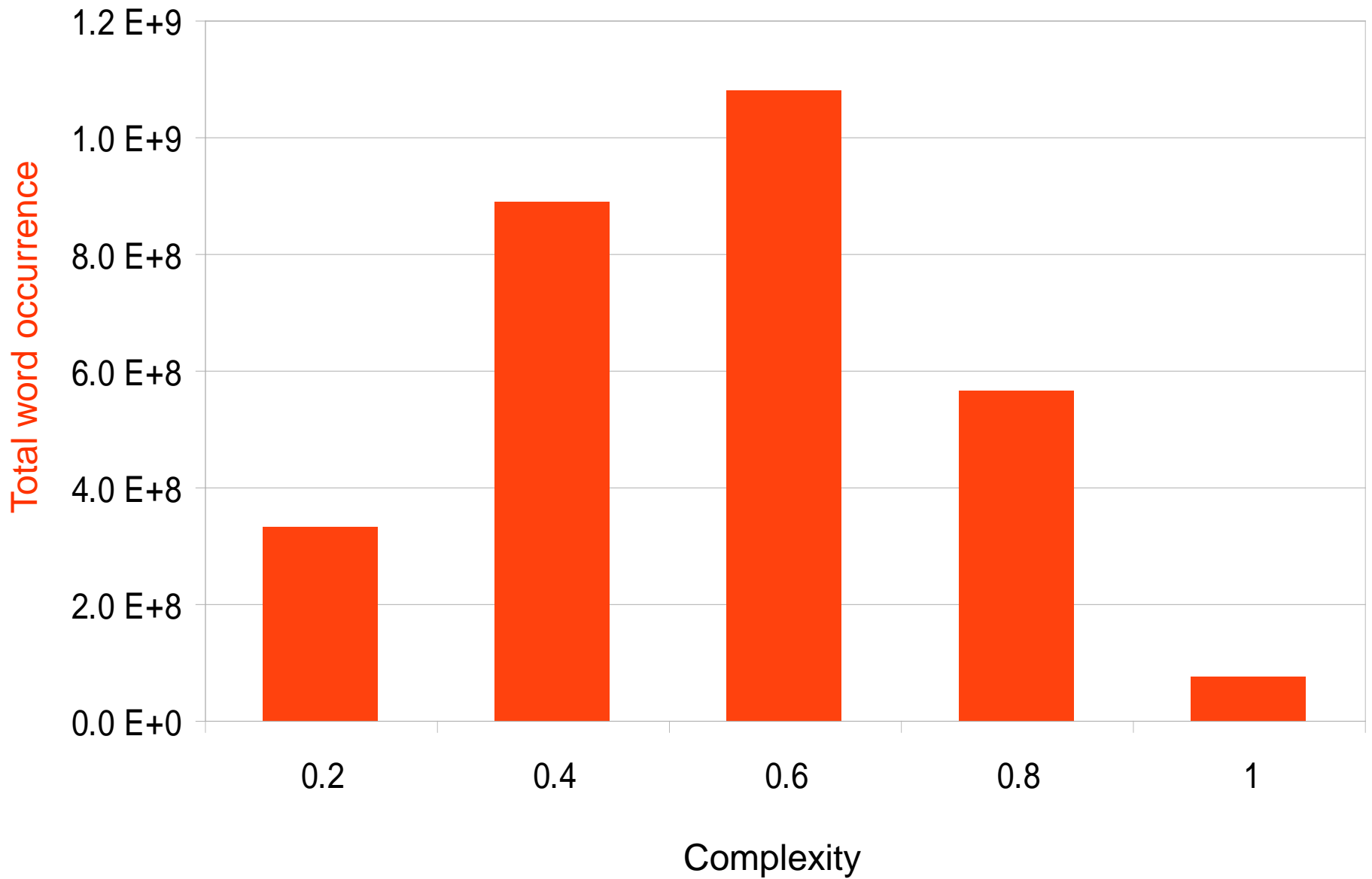
( Gct )<sub>n</sub> ( RYY )<sub>n</sub>

In the vocabulary of human genome 15-mers the simple repeats (low complexity words) dominate.

The high complexity words (of no repeat structure) are expected to be rather avoided.



# Occurrences of simple sequence 15-mers are anomalously high



**GCTGGGATTACAGGC** (Alu sequence)  
(complexity 0.68)

**GCT**

**GGG**

**ATT**

**ACA**

**GGC**

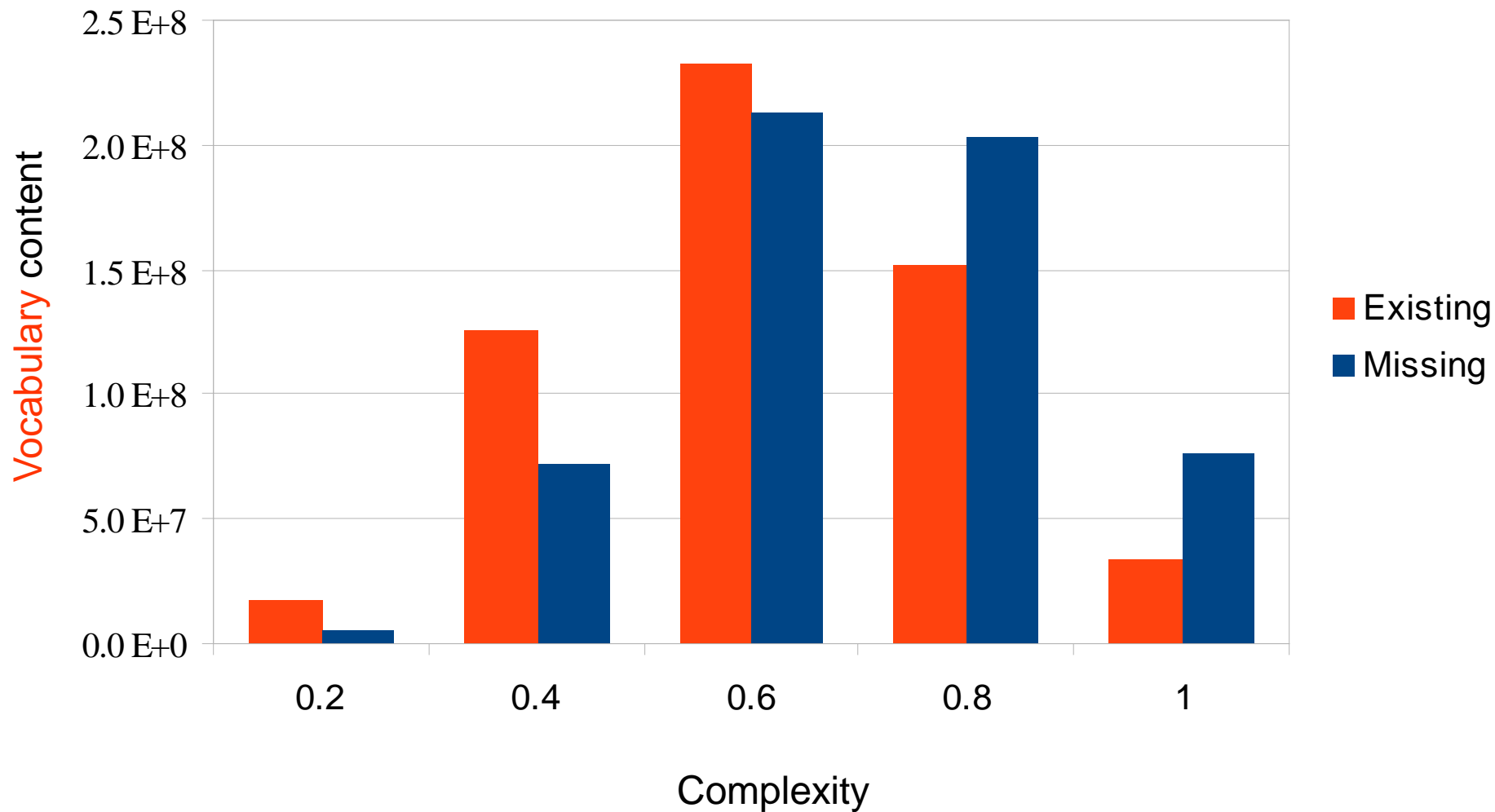
repeating

RYY<sub>5</sub>

**GCT**<sub>5</sub> aggressive triplet



15-mers of human genome are on low sequence complexity side.  
High complexity words are rather avoided



Genomes are simpler than we have thought  
They are dominated by simple sequences  
because they originate from simple sequences,  
as non-stop local births of new life