# Why are frameshift homologs widespread within and across species?

*Xiaolong Wang[*1], Quanjiang Dong[2], Gang Chen[1], Jianye Zhang[1], Yongqiang Liu[1], Jinqiao Zhao[1], Haibo Peng1, Yalei Wang1, Yujia Cai1, Xuxiang Wang1, Chao Yang1*

1. *College of Life Sciences, Ocean University of China, Qingdao, 266003, P. R. China*

    2. *Qingdao Municipal Hospital, Qingdao, Shandong, 266003, P. R. China*

# Abstract

Frameshifted coding genes yield truncated and dysfunctional proteins, frameshift mutations have been therefore considered as utterly harmful and of little importance for the evolution of novel proteins. However, frameshifted yet functional proteins and coding genes have been frequently observed. Here we report that frameshift homologs are widespread within a genome and across species. We showed that protein coding genes have a *ca*-0.5 quasi-constant shiftability: given any protein coding sequence, at least 50% of the amino acids remain conserved in a frameshifted protein sequence. In the natural genetic code, amino acid pairs assigned to frameshift codon substitutions are more conservative than those to random codon substitutions, and the frameshift tolerability of the natural genetic code ranks among the best 6.3% of all compatible genetic codes. Hence, the shiftability of coding genes was predefined by the genetic code, while additional sequence-level shiftability was achieved through biased usages of codons and codon pairs. We concluded that during early evolution the genetic code was optimized to tolerate frameshifting.

[1] To whom correspondence should be addressed: *Xiaolong Wang, Ph.D., Department of Biotechnology, Ocean University of China, No. 5 Yushan Road, Qingdao, 266003, Shandong, P. R. China*, Tel: *0086-139-6969-3150*, E-mail: *Xiaolong@ouc.edu.cn*.

1

## 1. Introduction

The genetic code was discovered in the early 1960s [1]. It consists of 64 triplet codons: 61 sense codons for the twenty amino acids and the remaining three nonsense codons for stop signals. The natural genetic code has several important properties: (1) The genetic code is universal for all organisms, with only a few variations found in some organelles or organisms, such as mitochondrion, archaea and yeast; For details, see the webpage: https://www.ncbi.nlm.nih.gov/Taxonomy/Utils/wprintgc.cgi. (2) The triplet codons are redundant, degenerative and wobble (the third base tends to be interchangeable); (3) In an open reading frame, an insertion/deletion (InDel) causes a frameshift unless the size of the InDel is a multiple of three.

The natural genetic code was optimized for translational error minimization [2], which is extremely efficient at minimizing the effects of mutation or mistranslation errors [3], and optimization for kinetic energy conservation in polypeptide chains [4]. Moreover, it was presumed that the natural genetic code resists frameshift errors by increasing the probability that a stop signal is encountered upon frameshifts, because frameshifted codons for abundant amino acids overlap with stop codons [5].

Presumably, most frameshifted coding DNA sequences (CDSs) yield truncated, non-functional, potentially cytotoxic products, lead to waste of cell energy, resources and the activity of the biosynthetic machinery [6, 7]. Therefore, frameshift mutations were considered as utterly harmful and of little importance for the evolution of novel proteins [8, 9]. However, frameshifted yet functional proteins and coding genes have been frequently observed [10-13]. For example, in yeast, a frameshifted coding gene for mitochondrial cytochrome c oxidase subunit II (COXII), the sequence is translated in an alternative frame [13]. Moreover, it was reported that frameshift mutations can be retained for millions of years and enable the acquisition of new gene functions [14], shed light into the role of frameshift mutation in molecular evolution.

A protein can be dysfunctioned even by changing a few residues, it is therefore a puzzle how the frameshift proteins kept their structures and functionalities while their

amino acid sequences has been changed substantially. Here we report that frameshift homologs are widespread within a genome and across species, and this is because the natural genetic code was optimized to tolerate frameshifting in early evolution.

## 2. Materials and Methods

### 2.1 Protein and coding DNA sequences

All available protein sequences in all species (Release 2016_04 of 13-Apr-2016 of UniProtKB/TrEMBL, contains 63686057 sequence entries) were downloaded from the UniprotKB protein database. All available reference protein sequences and their coding DNA sequences (CDSs) in nine model organisms, including *Escherichia coli*, *Saccharomyces cerevisiae*, *Arabidopsis thaliana*, *Caenorhabditis elegans*, *Drosophila melanogaster*, *Danio rerio*, *Xenopus tropicalis*, *Mus musculus* and *Homo sapiens*, were retrieved from *UCSC, Ensembl* and/or *NCBI* Genome Databases. Ten thousand CDSs each containing 500 random sense codons were simulated by *Recodon* 1.6.0 using default settings [15]. The human/simian immunodeficiency virus (HIV/SIV) strains were derived from the seed alignment in Pfam (pf00516). The CDSs of their envelop glycoprotein (GP120) were retrieved from the HIV sequence database [16].

### 2.2 Aligning and computing the similarity of the frameshifted protein sequences

A java program, *Frameshift-Align*, was written to translate CDSs in three reading frames, align the three translations and compute their similarities. Every CDS was translated into three protein sequences in its three reading frames in the same strand using the standard genetic code, while all internal nonsense codons were *readthrough* according to the above *readthrough rules* (Table 1). Each protein sequence and the two frameshifted amino acid sequences were aligned by ClustalW2 using default parameters. The pairwise similarity between a protein sequence and its frameshifted protein sequence is given by the percent of sites in which the matched amino acids are conserved (having a positive or zero amino acid substitution score in a scoring matrix, BLOSSUM62, PAM250 or GON250).

### 2.3 Blastp searching for frameshift homologs

1   A java program, *Frameshift-Translate*, was written and used to translate CDSs in

2   the alternative reading frames, and the frameshift translations were used as queries to

3   search against the UniprotKB protein database by local blastp, and the Blast hits were

4   filtered with a stringent cutoff criterion (*E-value*≤1e-5, *identity*≥30%, and *alignment*

5   *length*≥20 AAs).

6   Given a coding gene, its alternative reading frames often contain a certain number

7   of off-frame stop codons. Therefore, frameshifted coding sequences are commonly

8   translated into inconsecutive protein sequences interrupted by some stop signals (*).

9   To find frameshift homologs by blastp, the query sequences is better to be consecutive

10   sequences devoid of stop signals. Therefore, in *Frameshift-Translate*, when the CDSs

11   were translated into protein sequences in the alternative reading frames, every internal

12   nonsense codon was translated into an amino acid according to a set of *readthrough*

13   *rules* (Table 1).

14   **2.4 Computational analysis of frameshift codon substitutions**

15   A protein sequence consisting of $n$ amino acids is written as, $A_1 A_2 \ldots A_i A_{i+1} \ldots$

16   $A_n$, where $A_i = \{ A, C, D, E, F, G, H, I, K, L, M, N, P, Q, R, S, T, V, W, Y \}$, $i = 1 \ldots n$; its

17   coding DNA sequence consists of $n$ triplet codons, which is written as,

18   $\boldsymbol{B_1 B_2 B_3} \mid B_4 B_5 B_6 \mid \boldsymbol{B_7 B_8 B_9} \mid \ldots \mid \boldsymbol{B_{3i+1} B_{3i+2} B_{3i+3}} \mid B_{3i+4} B_{3i+5} B_{3i+6} \mid \ldots \mid B_{3n-2} B_{3n-1} B_{3n}$

19   Where $B_k = \{ A, G, U, C \}$, $k = 1 \ldots 3n$. Without loss of generality, let a frameshift

20   be caused by deleting or inserting one or two bases in the start codon:

21   (1) *Delete one:*   $\boldsymbol{B_2 B_3} B_4 \mid B_5 B_6 \boldsymbol{B_7} \mid \ldots \mid \boldsymbol{B_{3i+2} B_{3i+3}} B_{3i+4} \mid B_{3i+5} B_{3i+6} \boldsymbol{B_{3i+7}} \mid \ldots$

22   (2) *Delete two:* $\boldsymbol{B_3} B_4 B_5 \mid B_6 \boldsymbol{B_7 B_8} \mid \ldots \mid \boldsymbol{B_{3i+3}} B_{3i+4} B_{3i+5} \mid B_{3i+6} \boldsymbol{B_{3i+7} B_{3i+8}} \mid \ldots$

23   (3) *Insert one:*   $B_0 B_1 B_2 \mid \boldsymbol{B_3} B_4 B_5 \mid B_6 \boldsymbol{B_7 B_8} \mid \ldots \mid \boldsymbol{B_{3i+3}} B_{3i+4} B_{3i+5} \mid B_{3i+6} \boldsymbol{B_{3i+7} B_{3i+8}} \mid \ldots$

24   (4) *Insert two:* $B_{-1} B_0 \boldsymbol{B_1} \mid \boldsymbol{B_2 B_3} B_4 \mid B_5 B_6 \boldsymbol{B_7} \mid \ldots \mid \boldsymbol{B_{3i+2} B_{3i+3}} B_{3i+4} \mid B_{3i+5} B_{3i+6} \boldsymbol{B_{3i+7}} \mid \ldots$

25   So, if a frameshift mutation occurred in the first codon, the second codon $B_4 B_5 B_6$

26   and its encoded amino acid $A_2$ has two and only two possible changes:

27   (1) *Forward frameshifting (FF)*: $\boldsymbol{B_3} B_4 B_5$ ($\rightarrow A_{21}$)

28   (2) *Backward frameshifting (BF)*: $B_5 B_6 \boldsymbol{B_7}$ ($\rightarrow A_{22}$)

1    And so forth for each of the downstream codons. The results are two frameshifted

2    protein sequences, which were denoted as *FF* and *BF*. In either case, in every codon

3    all three bases are changed when compared base by base with the original codon.

4    According to whether the encoded amino acid is changed or not, codon substitutions

5    have been classified into two main types: (1) *Synonymous substitution* (SS); (2)

6    *Nonsynonymous substitution* (NSS). Based on the above analysis, we classified codon

7    substitutions into three subtypes: (1) *Random substitution*; (2) *Wobble substitution*; (3)

8    *Frameshift substitution*.

9    The amino acid substitution score of a frameshift codon substitution is defined as

10   frameshift substitution score (FSS). A java program, *Frameshift-CODON,* was written

11   to compute the average substitution scores for distinct kinds of codon substitutions by

12   using a scoring matrix (BLOSSUM62, PAM250 or GON250).

13   **2.5  Computational analysis of alternative codon tables**

14   A java program, *Frameshift-GC*, was written to produce "compatible" alternative

15   codon tables according to the method used in reference [3], by changing amino acids

16   assigned to sense codons and keeping degenerative codons synonymous. One million

17   alternative genetic codes were randomly selected from all $(20! = 2.43290201 \times 10^{18})$

18   "compatible" genetic codes. The sum and average FSSs for each genetic code were

19   computed and sorted, and compared with that of the natural genetic code.

20   **2.6  Analysis of codon pairs and their frameshift substitution scores**

21   For a given pair of amino acids, written as, $A_1 A_2$, where $A_i = \{$ *A, C, D, E, F, G, H,*

22   *I, K, L, M, N, P, Q, R, S, T, V, W, Y* $\}$ , *i = 1, 2*; its encoding codon pair is written as, $\boldsymbol{B_1}$

23   $\boldsymbol{B_2 B_3} / B_4 B_5 B_6$ , where $B_k = \{$ *A, G, U, C* $\}$ , *k = 1...6*. There are 400 different amino

24   acid pairs and 4096 different codon pairs.

25   Without loss of generality, let a frameshift be caused by inserting or deleting one

26   base in the first codon, the codon pair and its encoded amino acids has two and only

27   two types of changes:

28   (1) *Forward frameshifting:*    $B_0 \boldsymbol{B_1 B_2} / \boldsymbol{B_3} B_4 B_5 (\to A_{11} A_{21})$

29   (2) *Backward frameshifting:*   $\boldsymbol{B_2 B_3} B_4 / B_5 B_6 \boldsymbol{B_7} (\to A_{12} A_{22})$

1    A java program, *Frameshift-CODONPAIR*, was written to compute the average

2    amino acid substitution scores for each codon pair. The result of these calculations is a

3    list of 4096 codon pairs with their corresponding FSSs.

4    *2.7 Computational analysis of the usage of codon and codon pairs*

5    The usage of codons and codon pairs was analyzed on the above dataset using the

6    same method used in reference [17]. The program *CODPAIR* was rewritten in java as

7    the original program is not available. For each sequence, it enumerates the total

8    number of codons, and the number of occurrences for each codon and codon pair. The

9    observed and expected frequencies were then calculated for each codon and codon

10   pair. The result of these calculations is a list of 64 codons and 4096 codon pairs, each

11   with an expected (*E*) and observed (*O*) number of occurrences, usage frequency,

12   together with a value for $\chi_I^2 = (O \ \square \ E)^2/E$. The codons and dicodons whose *O-value*

13   is greater/smaller than their *E-value* were identified as *over-/under-represented*, their

14   average FSSs and the total weighted average FSSs were computed and compared.

# 3.  Results and Analysis

## 3.1 *Frameshift homologs are widespread within and across species*

17   Frameshift mutations disrupt the function of proteins, as every codon is changed,

18   and often many nonsense codons emerge in a frameshifted CDS. However, in the

19   development of codon and amino acid unified sequence alignment (CAUSA) [18, 19],

20   we noticed that protein sequences encoded by frameshifted CDSs are highly similar to

21   the wild-type protein sequences when they were aligned with each other. For example,

22   in different HIV/SIV strains, including HIV, SIVCZ and SIVGB, HIV was originated

23   from SIVCZ, and SIVCZ was from SIVGB [20-22]. As shown in Fig 1A, the envelop

24   glycoprotein gene (*gp120*) underwent a series of evolutionary events, including base

25   substitution, insertion, deletion, frameshifting and recombination. Especially, several

26   whole or partial, forward or backward frameshifting events occurred in *gp120,* but

27   their encoded protein sequences remain highly similar to each other (Fig 1B), and

28   these frameshifted proteins (GP120) are surely all functional, as the infection of these

29   virus into their host cells relies on these proteins.

1    As we know, a frameshift mutation is caused by one or more InDels in a protein

2    coding gene whose length is not a multiple of three. Consequently, the reading frame

3    is altered, either fully or partially. In this study, a *frameshift homolog* is defined as a

4    blastp hit using an artificially frameshifted protein sequence as a query. A frameshift

5    homolog is not a frameshift pseudogene, which often contains a certain number of

6    internal nonsense codons and is usually considered as dysfunctional. A frameshift

7    homolog, however, does not necessarily contain internal stop codons, and is usually a

8    frameshifted coding gene that encodes a functional protein.

9    By searching Uniprot database using blastp with artificially frameshifted protein

10   sequences as queries, we found that frameshift homologs are widespread within a

11   genome and across species. These frameshift homologs were classified into two types:

12   (1) **Frameshift ortholog***: given a coding gene *A* in a species (*sp*. 1), a frameshift

13       homolog (gene *a*) exists in another species (*sp*. 2)*, which was evolved from a

14       common ancestral gene via speciation and frameshifting (Fig 2A).

15   (2) **Frameshift paralog**: given a coding gene *A* in a species, a frameshift homolog

16       (gene *B*) exists in the same species*, which was evolved from a common

17       ancestral gene via duplication and frameshifting (Fig 2B).

18   As shown in Supplementary Dataset 1(Frameshift homologs.xlsx), large numbers

19   of frameshift paralogs and orthologs were found in the genomes of all species tested.

20   For example, in *Homo sapiens*, using frameshifted protein sequences translated from

21   the alternative reading frames of human reference CDSs (hg38, GRCh38) as queries,

22   blastp detected 3974 frameshift paralogs in the human genome and 23224 frameshift

23   homologs (including frameshift orthologs and frameshift paralogs) in all species.

24   These frameshift homologs were mapped onto the human genome and displayed in

25   the UCSC genome browser in two custom tracks, *frameshift homologs* and *frameshift*

26   *paralogs* (Fig 1C), respectively. The supplementary dataset, the source code of the

27   programs, and the custom track files for the UCSC genome browser are available in

28   the supporting information listed in the end of this article.

29   A modified blastp method for searching frameshift homologs was first established

30   by Claverie in 1993 [8] and then by Pellegrini and Yeates in 1999 [9]. Both studies

1  relied on more sophisticated use of amino acid scoring tables as a way to account for

2  models of protein sequence divergence, and both provided more robust statistical

3  treatments. Claverie suggested that, setting aside cases of accepted overlapping genes

4  in certain microbes and viruses, and cases of likely sequencing errors, there were only

5  a very small number of detectable cases of frameshift homologs. Pellegrini and Yeates

6  performed more careful sequence shuffling experiments to establish a baseline for

7  random expectations, and concluded that some weak signal existed in the databases to

8  suggest frameshifting as an evolutionary mechanism, concluded that strong inferences

9  about frameshift relationships between specific modern sequences was not possible.

10      Their method is more sophisticated and could be better than ours in the matter of

11  specificity and accuracy. However, their results published in the 1990s were based on

12  small datasets: Claverie used only 28,154 protein sequences from UniProt in 1993.

13  The method of Pellegrini and Yeates requires consensus sequences and there were

14  only 8,823 entries of consensus sequences available in Prodom in 1999.

15      The size of the UniProtKB database has been growing exponentially. The UniProt

16  data we used (Release 2016_04) contains 63,686,057 entries, which is 2262-fold

17  greater than the size of database used by Claverie. The blastp hits were filtered with a

18  very rigorous cutoff criteria (E-value≤1e-5, identity≥30% and alignment length≥20

19  AAs), but it might be not sufficient to filter out all false positives. Although the

20  method we used is rudimentary, it is based on the well-established blastp program and

21  we can adjust the cutoff criterion to raise the specificity and ensure that most of the

22  hits are true frameshift homologs. In human, when the cutoff criteria was raise to

23  E-value≤1e-5, identity≥50% and alignment length≥20 AAs, there were still 1120

24  frameshift paralogs in human genome and 6371 frameshift homologs in all species.

25      Moreover, hundreds of frameshifted queries are ≥95% identical to other known

26  protein sequences, and hundreds of them have a match length ≥100 AAs. Clearly,

27  they are recently derived from frameshifting of other coding genes. Despite there may

28  still be some false positives (e.g., frameshifts caused by sequencing errors), most of

1  them were considered as true frameshift homologs evolved from a common ancestral

2  gene via frameshifting rather than random similarities or artifacts. Finally, in the last

3  decade, a few studies have already been reported that frameshift homologs are widely

4  exist in many species [14, 23]. Therefore, we concluded that: *frameshift homologs are*

5  *widespread within and across species*.

6  ### 3.2 Frameshift proteins are always highly similar to their wild-types

7      As mentioned above, we noticed that frameshifted protein sequences are always

8  highly similar to their wild-types. To further validate this, the coding sequences were

9  translated each into three protein sequences in their three reading frames, the three

10  translations were aligned by ClustalW, and their pairwise similarities were computed.

11  For a given CDS, let $\delta_{ij} = \delta_{ji}$ be the similarity between a pair of protein sequences

12  encoded in reading frame *i* and *j* (*i, j=1,2,3, i ≠ j*), the average pairwise similarity

13  among the three protein sequences translated from the three different reading frames

14  on the same strand is defined as *the shiftability of the protein coding gene* (*δ*),

$$\delta = \frac{1}{3}(\delta_{12} + \delta_{13} + \delta_{23})$$

15      By analyzing all available reference CDSs in nine major model organisms, we

16  show that *δ* was centered approximately at 0.5 in all CDSs, in all species, as well as in

17  the simulated CDSs (Table 2 and Supplementary Dataset 2). In other words, *in most*

18  *coding genes*, the three protein sequences encoded in their three reading frames are

19  always highly similar to each other, with an average similarity of ~50%. Therefore,

20  we proposed that *protein coding genes have ca-0.5 quasi-constant shiftability, i.e., in*

21  *a protein coding gene, approximately 50% of its amino acids remain conserved in a*

22  *completely frameshifted protein sequence*.

23      For partially frameshifted coding genes, obviously, site conservation is inversely

24  proportional to the numbers of frameshifted sites, therefore, partial frameshifts are all

25  highly similar to the wild-type. Hence, it is guaranteed that in a frameshifted protein

26  at least half of its aa sites are conserved when compared to the wide-type. This does

27  not mean that frameshifted variants are all functional, however, quite many of them

28  could maintain their structure and function, forming the basis of frameshift tolerating.

1  In addition, the wild type is not necessarily the "best" form. In a frameshifted protein,

2  the other half of sites change into dissimilar amino acids, provides a fast and effective

3  means of molecular evolution for improving or altering the structure and function of

4  proteins, or developing the overlapping genes.

5  ### *2.8 Explanation of the readthrough rules and their impact on computation*

6  The *readthrough rules* were summarized from nonsense suppression tRNAs

7  reported in *E. coli*. The suppressor tRNAs are expressed *in vivo* to correct nonsense

8  mutations, including *amber suppressors* (*supD* [24], *supE* [25], *supF* [26]), *ochre*

9  *suppressor*s (*supG* [27]) and *opal suppressors* (*supU* [26], *su9* [28]). These suppressor

10  tRNAs are taken as *readthrough rules*, because *translational readthrough* occurs upon

11  activity of a suppressor tRNA with an anticodon matching a stop codon. The

12  suppressor tRNAs frequently occur in the negative strand of a regular tRNA [29-31].

13  It was found that translational readthrough occurred by using these suppressor tRNAs

14  allows the translation of off-frame peptides [32-35]. There have been a lot of reports

15  that translational readthrough functions not only in *E. coli*, but also in yeast and many

16  eukaryotes species (including human), while the readthrough rules may vary [36, 37].

17  In addition, there have been increasing evidences show that translational readthrough

18  is related to frameshift tolerating, ribosomal frameshifting or frameshift repair. For

19  example, interaction of eRF3 with RNase L leads to increased readthrough efficiency

20  at premature termination codons and +1 frameshift efficiency [38].

21  However, in this study, the readthrough-rules are taken simply as 'computational

22  rules' borrowed from biology to obtain consecutive frameshifted protein sequences,

23  without the interruption of stop signals. Therefore, the artificial frameshifting and *in*

24  *silicon* readthrough operations performed on the coding sequences are distinct from *in*

25  *vivo* translational readthrough, since the frameshifted amino acid sequences translated

26  from the artificially frameshifted CDSs were used as inputs to ClustalW for multiple

27  sequence alignment (MSA). The purpose of MSA is only to compute the similarities

28  of the protein sequences encoded in the three reading frames.

29  The artificially frameshifted protein sequences were also used as query for blastp

30  to search for frameshift homologs in the Uniprot database. Although the frameshifts

1    themselves are not really exist in biology, a blastp hit found in the Uniprot database is

2    a true biological protein sequences in most cases (unless the hit itself contains an

3    artificial sequencing error), and the hits are the homologs (ancestors or descendants)

4    of the corresponding frameshifted query, called *frameshift homologs*.

5         We performed ClustalW aligning and blastp searching by using both readthrough

6    and non-readthrough frameshifted protein sequences. For example, as shown in Fig

7    2D, in the MSA of wild-type zebrafish VEGFAA with their frameshifted translations,

8    the alignment for readthrough and non-readthrough frameshifted protein sequences

9    are same to each other, except for the stop signals presented in the alignments. As

10   shown in Fig 2D, 62.2% (117/188) of their sites are kept conserved in physiochemical

11   properties. The shiftability of vegfaa computed by readthrough and non-readthrough

12   is 0.5354 and 0.5573, respectively. So, *in silicon* readthrough has a negligible impact

13   on the computation of the shiftability.

14        The blastp results, however, were slightly better (higher score, more positives and

15   better E-value) in readthrough than in non-readthrough queries. As shown in Fig 2E,

16   in the blastp result of a frameshifted query, the stop signals of the query match each

17   with an amino acid in the subject, suggesting that the corresponding stop codons were

18   substituted each by a sense codon in evolution. So, we translated the frameshifted

19   coding sequences by using the readthrough rules, but, it does not require or imply that

20   these *in-silicon* readthrough rules must function in *E. coli* or any other species, but

21   simply a computational method to obtain consecutive frameshifted protein sequences.

22   **3.3  The genetic code was optimized for frameshift tolerating**

23        In Table 2, the shiftability of the protein coding genes is similar in all species, and

24   all genes, and the standard deviation is very small, suggesting that the shiftability is

25   largely species- and sequence-independent. This implies that the shiftability is defined

26   mainly by the genetic code rather than by the coding sequences themselves. This is

27   also suggested by the simulated coding sequences, whose shiftability is comparable

28   with those of the real coding genes.

29        As described above in the method section, we computed the average amino acid

30   substitution scores respectively for random, wobble and forward/backward frameshift

1  codon substitutions. As shown in Table 3 and Supplementary Dataset 3, in all 4096

2  possible codon substitutions, most (192/230=83%) of the synonymous substitutions

3  are wobble, and most (192/256=75%) wobble substitutions are synonymous, thus the

4  average substitution score of the wobble substitutions is the highest. For frameshift

5  codon substitutions, including the four triplet codons (TTT, AAA, GGG, CCC) which

6  are kept unchanged in frameshifting, only a small proportion (28/512=5.5%) of the

7  frameshift codon substitutions are synonymous (Table 4), and the others (95.9%) are

8  all nonsynonymous. However, a substantial proportion (35.9%) of them are positive

9  (including SSs and positive NSSs), which is significantly higher than the proportion

10  of positive substitutions in random codon substitutions (25.7%). In summary, in the

11  natural genetic code, SSs are assigned mainly to wobble codon substitutions, while

12  positive NSSs are assigned mainly to frameshift substitutions.

13      In addition, no matter which substitution scoring matrix (BLOSSUM62, PAM250

14  or GON250) was used for computation, the average FSSs are significantly higher than

15  the substitution scores of the random codon substitutions (t-test P << 0.01), suggesting

16  that the amino acid substitutions assigned to the frameshift substitutions are more

17  conservative than those to the random substitutions.

18      The amino acid substitution scoring matrix is widely used to determine similarity

19  and conservation in sequence alignment and blast searching, which forms the basis of

20  most bioinformatics analysis. In commonly used scoring matrix, either BLOSSUM62,

21  PAM250 or GON250, most of the amino acid substitution scores are negative and the

22  percent of positive scores is less than 30%. So, percent of positive scores for random

23  codon substitutions is ~30%. However, as shown in Table 2, the frameshifted protein

24  sequences are always ~50% similar to the wild types: the ~35% similarity derived

25  from the frameshift codon substitutions, combined with the ~25% similarity from the

26  random codon substitutions, deduct their ~10% intersection, well explains the ~50%

27  similarities observed among the frameshifted protein sequences and the wild types.

28  Therefore, it is suggested that the shiftability of coding genes was predefined mainly

29  by the genetic code, and is largely independent on the coding sequences themselves.

1    To further investigate optimization for frameshift tolerance of the natural genetic

2    code, one million alternative genetic codes were randomly selected from all (20! =

3    $2.43290201 \times 10^{18}$) "compatible" genetic codes by changing the amino acids assigned

4    to the sense codons randomly, while keeping all degenerative codons synonymous. By

5    computing and sorting the average FSSs for these alternative genetic codes (Table 5),

6    the FSSs of the natural genetic code ranks in the best 6.3% of all compatible genetic

7    codes. Hence the genetic code was indeed optimized for tolerating frameshifts, clearly

8    demonstrating that the shiftability of coding genes is defined by the genetic code.

9    *3.4  The genetic code is symmetric in frameshift tolerating*

10   The genetic code shows the characteristics of symmetry in many aspects [39-41],

11   and it evolved probably through progressive symmetry breaking [42-44]. Here in all

12   CDSs both forward and backward frameshift proteins have comparable similarities

13   with the wild-type (Table 2). In addition, in the natural genetic code both forward and

14   backward frameshift substitutions have the same number of SSs/NSSs and roughly

15   equal FSSs (Table 3). These data suggested that the genetic code is also symmetric in

16   terms of shiftability and frameshift tolerating, so that a coding gene has an ability to

17   tolerate frameshifting in both forward and backward directions at the same time (Fig

18   2). This could also explain why the codons in the natural genetic code are not tetrad

19   but triplet: triplet codon could be kept symmetric for both forward and backward

20   frameshifting, while for tetrad codons the situation will be much more complicated in

21   frameshifting.

22   *3.5  The shiftability at sequence level*

23   Although the shiftability of a coding sequence is predefined mainly by the genetic

24   code, shiftability may also exist at the sequence level. Functionally important coding

25   genes, such as housekeeping genes, which are more conserved, may also have greater

26   shiftability when compared with other genes. At first, we thought that a biased usage

27   of codons may contribute to the sequence-level shiftability. However, as shown in

28   Table 6 and Supplementary Dataset 4, it is somewhat surprising that in *E. coli* and *C.*

29   *elegans* the average FSSs weighted by their codon usages are even lower than for

30   unweighted calculations (equal usage of codons). In the other species tested, although

1 the weighted average FSSs are higher than for unweighted analyses, the difference is

2 not statistically significant in all species tested (P>0.05), suggesting that the usage of

3 codons has little or no direct impact on the shiftability. However, the usage of codons

4 may influence the shiftability indirectly, *e.g.*, by shaping the pattern of codon pairs.

5 Given a pair of amino acids, $A_1 A_2$, if $A_1$ and $A_2$ have $m_1$ and $m_2$ degenerative

6 codons, respectively, their encoding dicodons, $B_1 B_2 B_3 | B_4 B_5 B_6$, has $m_1 \times m_2$ possible

7 combinations, called *degenerative codon pairs* (DCPs). It has been reported that

8 codon pair usages are highly biased in various species, such as bacteria, human and

9 animals [17, 45-50]. As shown in Table 7, and Supplementary Dataset 5, in all species

10 tested, the average FSSs of the over-represented codon pairs are all positive, while

11 those of the under-represented codon pairs are all negative; in addition, the weighted

12 average FSSs of codon pairs are all positive, while that of the equal usage of codon

13 pairs is negative, suggesting that a selective pressure was working on the codon pairs,

14 so that frameshift-tolerable DCPs are present more frequently in these genomes than

15 non-frameshift-tolerable DCPs. Therefore, sequence-level shiftability does exist, and

16 was achieved through a biased usage of codons and codon pairs. There have been

17 many studies on the causes and consequences of the usage of codons, such as gene

18 expression level [51-56], mRNA structure [57-64], protein abundance [61, 65-67], and

19 stability [68-70]. The above analysis suggested that the usages of codon pairs is either

20 a cause or a consequence of the shiftability of the protein-coding genes.

## 4. Discussion

### 4.1 The genetic code was optimized for frameshift tolerating

23 The natural genetic code results from selection during early evolution, and it was

24 optimized along several properties when compared with other possible genetic codes

25 [71-82]. It was reported that the natural genetic code was optimized for translational

26 error minimization, because the amino acids whose codons differed by a single base

27 in the first and third positions were similar with respect to polarity and hydropathy,

28 and the differences between amino acids were specified by the second position is

29 explained by selection to minimize the deleterious effects of translation errors during

the early evolution of the genetic code [2]. In addition, it was reported that only one in every million alternative genetic codes is more efficient than the natural genetic code, which is extremely efficient at minimizing the effects of point mutation or translation errors [3]. It was demonstrated that the natural genetic code is nearly optimal for allowing additional information within coding sequences, such as out-of-frame hidden stop codons (HSCs) and secondary structure formation (self-hybridization) [5].

In the above, we showed that the code- and sequence-level shiftability of coding genes guaranteed at least half of the sites are kept conserved in a frameshifted protein when compared with the wild-type protein. This is the basis for frameshift tolerating, and explains why the usage of codons and codon pairs are biased and why frameshift homologs are widespread within and across species.

The sequence-level shiftability caused by the biased usages of codon pairs are probably relevant to the circular code. The circular code is a set of 20 codons that are overrepresented in the regular coding frame of genes as compared to frameshifted frames [83-85]. The mechanism by which the circular code maintains the translation frame is unknown [85-89], but *in silico* frame detection was made possible by using the empirical circular code [88-93]. However, the relationship among the shiftability, the biased usages of codons and codon pairs, and the circular code is unknown.

## 4.2 The universality of the shiftability

Here we analyzed the shiftability of protein-coding genes only in some model organisms, thus it is interesting to further validate this mechanism in other species. It has been reported that in some animal species frameshift mutations are tolerated by the translation systems in mitochondrial genes [94-96]. For example, a (+1) frameshift insertion is tolerated in the *nad3* in birds and reptiles [94]. Moreover, frameshifted overlapping coding genes have been found in mitochondria genes in fruit fly and turtles [97, 98]. It was reported that the levels of translational readthrough and frameshifting in *E. coli* are both high and growth phase dependent [99]. Meanwhile, translational readthrough has been widely observed in various species [100-107]. Frameshift tolerating has also been explained by *ribosomal frameshifting* [108-111].

1  However, the shiftability of protein coding genes may also contribute, at least partially,

2  to the functioning, repairing and evolution of the frameshifted protein coding genes.

## 5.  Conclusion

4      The above analysis conclude that frameshift homologs are widespread within

5  and across species, and this is because the genetic code was optimized for frameshift

6  tolerating. The shiftability of coding genes guarantees a near-half conservation after a

7  frameshifting event, endows coding genes an inherent ability to tolerate frameshifting.

8  The natural genetic code, which exists since the origin of life, was optimized by

9  competition with other alternative genetic codes during early evolution [112-115]. As

10  the *bottom design* for all genes and genomes for all species, the natural genetic code

11  allows coding genes to tolerate both forward and backward frameshifting, could have

12  a better fitness in the early evolution. Thanks to this ingenious property of the genetic

13  code, the shiftability serves an innate mechanism for protein-coding genes to deal

14  with frameshift mutations, by which the disastrous frameshifting events were utilized,

15  becoming a driving force for molecular evolution.

### Author Contributions

17      Xiaolong Wang conceived the study, coded the programs, analyzed the data, prepared the

18  figures, tables and wrote the paper; Quanjiang Dong proofread the paper and gave conceptual

19  advices. Gang Chen and Jianye Zhang provided materials and supports. Yujia Cai, Yongqiang

20  Liu and Jinqiao Zhao analyzed the data for alternative genetic codes.

### Acknowledgements

### Figure Legends

28      **Fig 1. The alignment of the coding and the protein sequences of HIV/SIV GP120.** (A)

29  The alignment of coding sequences, with highlights showing that the coding genes contain several

30  frameshifting events. In other words, the coding gene is expressed in different reading frames in

1    different virus strains. (B) The alignment of protein sequences, showing that the GP120 sequences

2    for different virus, which are encoded in different reading frames of *gp120*, are highly similar.

3    **Fig 2. Diagram of different frameshift homologs.** (A) Frameshift orthologs; (B) Frameshift

4    paralog; (C) Custom tracks for the frameshift homologs displayed in the UCSC genome browser;

5    (D) the ClustalW alignment of the wild-type VEGFAA and its readthrough or non-readthrough

6    frameshifts; (E) The outputs of blastp searching for frameshift homologs, the query is an artificial

7    protein sequence translated from a frameshifted CDS by readthrough or non-readthrough.

8    **Additional Information**

9    We declare that the authors have no competing interests.

1    Table 1. The natural suppressor tRNAs (*readthrough rules*) for nonsense mutations.

| Site | tRNA (AA) | Wild type | | Correction | |
|---|---|---|---|---|---|
| | | **Code** | **Anti-code** | **Code** | **Anti-code** |
| *supD* | Ser (S) | → UCG | CGA← | → UAG | CUA← |
| *supE* | Gln (Q) | → CAG | CUG← | → UAG | CUA← |
| *supF* | Tyr (Y) | → UAC | GUA← | → UAG | CUA← |
| *supG* | Lys (K) | → AAA | UUU← | → UAA | UUA← |
| *supU* | Trp (W) | → UGG | CCA← | → UGA | UCA← |

2

3

1

2        Table 2. The similarities of natural and simulated proteins and their frameshift forms.

| No. | Species | Number of CDSs | Average Similarity | | | | | |
|---|---|---|---|---|---|---|---|---|
| | | | $\delta_{12}$ | $\delta_{13}$ | $\delta_{23}$ | $\delta$ | MAX | MIN |
| 1 | H. sapiens | 71853 | 0.5217±0.0114 | 0.5044±0.0122 | 0.4825±0.0147 | 0.5028±0.0128 | 0.5948 | 0.4357 |
| 2 | M. musculus | 27208 | 0.5292±0.042 | 0.5058±0.0437 | 0.4869±0.0418 | 0.5073±0.0425 | 0.8523 | 0.1000[*] |
| 3 | X. tropicalis | 7706 | 0.5190±0.0013 | 0.4987±0.0013 | 0.4855±0.0008 | 0.5010±0.0008 | 0.5962 | 0.4790 |
| 4 | D. rerio | 14151 | 0.5234±0.0007 | 0.5022±0.0008 | 0.4921±0.0005 | 0.5059±0.0004 | 0.5240 | 0.4784 |
| 5 | D. melanogaster | 23936 | 0.5162±0.0015 | 0.4921±0.001 | 0.4901±0.0013 | 0.4995±0.0008 | 0.6444 | 0.4667 |
| 6 | C. elegans | 29227 | 0.5306±0.0007 | 0.5035±0.0008 | 0.5002±0.001 | 0.5115±0.0006 | 0.6044 | 0.4864 |
| 7 | A. thaliana | 35378 | 0.5389±0.0508 | 0.5078±0.0481 | 0.5062±0.048 | 0.5176±0.0388 | 0.9540 | 0.2162[*] |
| 8 | S. cerevisiae | 5889 | 0.5174±0.0011 | 0.4811±0.001 | 0.5072±0.0006 | 0.502±0.0007 | 0.5246 | 0.4577 |
| 9 | E. coli | 4140 | 0.5138±0.0019 | 0.4871±0.0046 | 0.481±0.0015 | 0.494±0.0012 | 0.7778 | 0.4074 |
| 10 | Simulated | 10000 | 0.5165±0.0282 | 0.4745±0.0272 | 0.4773±0.0263 | 0.4894±0.0013 | 0.6489 | 0.3539 |

3        * Very large and small similarity values were observed in a few very short or repetitive

4    peptides.

5

1          Table 3. The amino acid substitution scores for different kinds of codon substitutions.

| Codon Substitution | | ALL (Random) | Frameshift | | Wobble |
|---|---|---|---|---|---|
| | | | FF | BF | |
| | All | 4096 | 256 | 256 | 256 |
| Type of Codon Substitution | Unchanged (%) | 64 (1.6%) | 4 (1.6%) | 4 (1.6%) | 64 (25%) |
| | Changed (%) | 4032 (98.4%) | 252 (98.4%) | 252 (98.4%) | 192 (75%) |
| | SS (%) | 230 (5.6%) | 14 (5.5%) | 14 (5.5%) | 192 (75%) |
| | NSS-Positive (%) | 859 (20.1%) | 76 (29.7%) | 76 (29.7%) | 40 (15.6%) |
| | NSS-Negative (%) | 3007 (73.4%) | 166 (64.8%) | 166 (64.8%) | 24 (9.4%) |
| Average Substitution Score | BLOSSUM62 | -1.29 | -0.61 | -0.65 | 3.77 |
| | PAM250 | -4.26 | -0.84 | -0.84 | 3.68 |
| | GON250 | -10.81 | -1.78 | -1.78 | 35.60 |

2   SS/NSS: synonymous/nonsynonymous substitution; FF/BF: forward/backward frameshift codon
3   substitution.

4

1

2

Table 4. The synonymous frameshift substitutions

| Forward Frameshifting | | | | Backward Frameshifting | | | |
|---|---|---|---|---|---|---|---|
| From | | To | | From | | To | |
| 1 | AAA K | AAA K | | 1 | AAA K | AAA K | |
| 2 | AAA K | AAG K | | 2 | AAG K | AAA K | |
| 3 | GGG G | GGA G | | 3 | GGA G | GGG G | |
| 4 | GGG G | GGG G | | 4 | GGG G | GGG G | |
| 5 | GGG G | GGC G | | 5 | GGC G | GGG G | |
| 6 | GGG G | GGT G | | 6 | GGT G | GGG G | |
| 7 | CCC P | CCA P | | 7 | CCA P | CCC P | |
| 8 | CCC P | CCG P | | 8 | CCG P | CCC P | |
| 9 | CCC P | CCC P | | 9 | CCC P | CCC P | |
| 10 | CCC P | CCT P | | 10 | CCT P | CCC P | |
| 11 | CTT L | TTA L | | 11 | TTA L | CTT L | |
| 12 | CTT L | TTG L | | 12 | TTG L | CTT L | |
| 13 | TTT F | TTC F | | 13 | TTC F | TTT F | |
| 14 | TTT F | TTT F | | 14 | TTT F | TTT F | |

3

1    Table 5. The frameshift substitution score of the natural and alternative genetic codes (computed
2                              by using the amino acid scoring matrix BLOSSUM62).

| Number of alternative genetic codes Sampled | The natural genetic code | | FSS of the alternative genetic codes | | | | |
|---|---|---|---|---|---|---|---|
| | FSS Score | Rank | MAX | MIN | Average A* | Average B** | Average |
| 1,000,000 | -294 | 62007 | -43 | -814 | -256.842 | -438.930 | -427.375 |

3    * Average A: the average FSS of the genetic codes ranks above (better than) the natural genetic
4    code;
5    ** Average B: the average FSS of the genetic codes ranks below (worse than) the natural genetic
6    code;

7

8

1

2             Table 6. The usage of codons and their weighed average FSSs (Gon250)

| NO | Species (Codon Usage) | Weighted Average FSS |
|----|-----------------------|----------------------|
| 1 | H. sapiens | -9.82 |
| 2 | M. musculus | -13.47 |
| 3 | X. tropicalis | -12.75 |
| 4 | D. rerio | -20.58 |
| 5 | D. melanogaster | -19.43 |
| 6 | C. elegans | -23.38 |
| 7 | A. thaliana | -22.52 |
| 8 | S. cerevisiae | -14.08 |
| 9 | E. coli | -28.59 |
| 10 | Equal usage | -22.27 |

3

4

1

Table 7. The usage of codon pairs and their weighed average FSSs (Gon250)

| NO | Species (Codon Usage) | Average FSS of over-represented Codon pairs | Average FSS of under-represented Codon pairs | Weighted Average FSS of All Codon pairs |
|---|---|---|---|---|
| 1 | H. sapiens | 41.30 | -25.94 | 102.41 |
| 2 | M. musculus | 41.09 | -26.09 | 98.55 |
| 3 | X. tropicalis | 42.20 | -25.81 | 98.24 |
| 4 | D. rerio | 40.91 | -26.17 | 87.38 |
| 5 | D. melanogaster | 39.77 | -25.95 | 79.51 |
| 6 | C. elegans | 40.85 | -26.18 | 81.48 |
| 7 | A. thaliana | 40.54 | -26.09 | 90.64 |
| 8 | S. cerevisiae | 40.85 | -26.18 | 99.21 |
| 9 | E. coli | 39.27 | -30.75 | 77.03 |
| 10 | Equal Usage | N/A | N/A | -28.50 |

2

3

# References

1. Nirenberg, M.W. and J.H. Matthaei, *The dependence of cell-free protein synthesis in E. coli upon naturally occurring or synthetic polyribonucleotides.* Proc Natl Acad Sci U S A, 1961. **47**: p. 1588-602.

2. Haig, D. and L.D. Hurst, *A quantitative measure of error minimization in the genetic code.* J Mol Evol, 1991. **33**(5): p. 412-7.

3. Freeland, S.J. and L.D. Hurst, *The genetic code is one in a million.* Journal of Molecular Evolution, 1998. **47**(3): p. 238-248.

4. Guilloux, A. and J.L. Jestin, *The genetic code and its optimization for kinetic energy conservation in polypeptide chains.* Biosystems, 2012. **109**(2): p. 141-4.

5. Itzkovitz, S. and U. Alon, *The genetic code is nearly optimal for allowing additional information within protein-coding sequences.* Genome Research, 2007. **17**(4): p. 405-412.

6. Seligmann, H. and D.D. Pollock, *The ambush hypothesis: hidden stop codons prevent off-frame gene reading.* DNA Cell Biol, 2004. **23**(10): p. 701-5.

7. Tse, H., et al., *Natural selection retains overrepresented out-of-frame stop codons against frameshift peptides in prokaryotes.* BMC Genomics, 2010. **11**: p. 491.

8. Claverie, J.M., *Detecting frame shifts by amino acid sequence comparison.* J Mol Biol, 1993. **234**(4): p. 1140-57.

9. Pellegrini, M. and T.O. Yeates, *Searching for frameshift evolutionary relationships between protein sequence families.* Proteins, 1999. **37**(2): p. 278-83.

10. Xu, J., R.W. Hendrix, and R.L. Duda, *Conserved translational frameshift in dsDNA bacteriophage tail assembly genes.* Molecular Cell, 2004. **16**(1): p. 11-21.

11. Pai, H.V., et al., *A frameshift mutation and alternate splicing in human brain generate a functional form of the pseudogene cytochrome P4502D7 that demethylates codeine to morphine.* Journal of Biological Chemistry, 2004. **279**(26): p. 27383-27389.

12. Baykal, U., A.L. Moyne, and S. Tuzun, *A frameshift in the coding region of a novel tomato class I basic chitinase gene makes it a pseudogene with a functional wound-responsive promoter.* Gene, 2006. **376**(1): p. 37-46.

13. Fox, T.D., *Five TGA "stop" codons occur within the translated sequence of the yeast mitochondrial gene for cytochrome c oxidase subunit II.* Proc Natl Acad Sci U S A, 1979. **76**(12): p. 6534-8.

14. Raes, J. and Y. Van de Peer, *Functional divergence of proteins through frameshift mutations.* Trends Genet, 2005. **21**(8): p. 428-31.

15. Arenas, M. and D. Posada, *Recodon: coalescent simulation of coding DNA sequences with recombination, migration and demography.* BMC Bioinformatics, 2007. **8**: p. 458.

16. Abecasis, A.B., A.M. Vandamme, and P. Lemey, *Sequence Alignment in HIV Computational Analysis*, in *HIV Sequence Compendium*, T. Thomas, et al., Editors. 2007, Theoretical Biology and Biophysics Group, Los Alamos National Laboratory,: Los Alamos, NM. LA-UR 07-4826. p. 2-16.

17. Gutman, G.A. and G.W. Hatfield, *Nonrandom utilization of codon pairs in Escherichia coli.* Proceedings of the National Academy of Sciences of the United States of America, 1989. **86**(10): p. 3699-3703.

1 18. C, W.X.Y., *CAUSA 2.0: accurate and consistent evolutionary analysis of proteins*
2 *using codon and amino acid unified sequence alignments.* PeerJ PrePrints **3**.
3 19. Wang, X., Fu, Yu , Zhao, Yue , Wang, Qi , Pedamallu, Chandra Sekhar , Xu,
4 Shuang-yong , Niu, Yingbo , and Hu, Jingjie . , *Accurate Reconstruction of Molecular*
5 *Phylogenies for Proteins Using Codon and Amino Acid Unified Sequence Alignments (CAUSA). .*
6 Nature Precedings
7 20. Holmes, E.C., *On the origin and evolution of the human immunodeficiency virus*
8 *(HIV).* Biol Rev Camb Philos Soc, 2001. **76**(2): p. 239-54.
9 21. Rambaut, A., et al., *Human immunodeficiency virus. Phylogeny and the origin of*
10 *HIV-1.* Nature, 2001. **410**(6832): p. 1047-8.
11 22. Paraskevis, D., et al., *Analysis of the evolutionary relationships of HIV-1 and SIVcpz*
12 *sequences using bayesian inference: implications for the origin of HIV-1.* Mol Biol Evol, 2003.
13 **20**(12): p. 1986-96.
14 23. Antonov, I., et al., *Identification of the nature of reading frame transitions*
15 *observed in prokaryotic genomes.* Nucleic Acids Res, 2013. **41**(13): p. 6514-30.
16 24. Nagano, T., Y. Kikuchi, and Y. Kamio, *High expression of the second lysine*
17 *decarboxylase gene, ldc, in Escherichia coli WC196 due to the recognition of the stop codon*
18 *(TAG), at a position which corresponds to the 33th amino acid residue of sigma(38), as a*
19 *serine residue by the amber suppressor, supD.* Bioscience Biotechnology and Biochemistry,
20 2000. **64**(9): p. 2012-2017.
21 25. Kuriki, Y., *Temperature-Sensitive Amber Suppression of Ompf'-'Lacz Fused*
22 *Gene-Expression in a Supe Mutant of Escherichia-Coli K12.* Fems Microbiology Letters, 1993.
23 **107**(1): p. 71-76.
24 26. Johnston, H.M. and J.R. Roth, *UGA suppressor that maps within a cluster of*
25 *ribosomal protein genes.* J Bacteriol, 1980. **144**(1): p. 300-5.
26 27. Prather, N.E., B.H. Mims, and E.J. Murgola, *supG and supL in Escherichia coli code*
27 *for mutant lysine tRNAs+.* Nucleic Acids Res, 1983. **11**(23): p. 8283-6.
28 28. Chan, T.S. and A. Garen, *Amino acid substitutions resulting from suppression of*
29 *nonsense mutations. V. Tryptophan insertion by the Su9 gene, a suppressor of the UGA*
30 *nonsense triplet.* J Mol Biol, 1970. **49**(1): p. 231-4.
31 29. Seligmann, H., *Undetected antisense tRNAs in mitochondrial genomes?* Biol Direct,
32 2010. **5**: p. 39.
33 30. Seligmann, H., *Avoidance of antisense, antiterminator tRNA anticodons in*
34 *vertebrate mitochondria.* Biosystems, 2010. **101**(1): p. 42-50.
35 31. Seligmann, H., *Pathogenic mutations in antisense mitochondrial tRNAs.* J Theor
36 Biol, 2011. **269**(1): p. 287-96.
37 32. Seligmann, H., *Overlapping genetic codes for overlapping frameshifted genes in*
38 *Testudines, and Lepidochelys olivacea as special case.* Computational Biology and Chemistry,
39 2012. **41**: p. 18-34.
40 33. Seligmann, H., *An overlapping genetic code for frameshifted overlapping genes in*
41 *Drosophila mitochondria: Antisense antitermination tRNAs UAR insert serine.* Journal of
42 Theoretical Biology, 2012. **298**: p. 51-76.

34. Seligmann, H., *Two genetic codes, one genome: Frameshifted primate mitochondrial genes code for additional proteins in presence of antisense antitermination tRNAs.* Biosystems, 2011. **105**(3): p. 271-285.

35. Faure, E., et al., *Probable presence of an ubiquitous cryptic mitochondrial gene on the antisense strand of the cytochrome oxidase I gene.* Biol Direct, 2011. **6**: p. 56.

36. Dabrowski, M., Z. Bukowy-Bieryllo, and E. Zietkiewicz, *Translational readthrough potential of natural termination codons in eucaryotes--The impact of RNA sequence.* RNA Biol, 2015. **12**(9): p. 950-8.

37. Schueren, F. and S. Thoms, *Functional Translational Readthrough: A Systems Biology Perspective.* PLoS Genet, 2016. **12**(8): p. e1006196.

38. Le Roy, F., et al., *A newly discovered function for RNase L in regulating translation termination.* Nat Struct Mol Biol, 2005. **12**(6): p. 505-12.

39. Findley, G.L., A.M. Findley, and S.P. McGlynn, *Symmetry characteristics of the genetic code.* Proc Natl Acad Sci U S A, 1982. **79**(22): p. 7061-5.

40. Frappat, L., P. Sorba, and A. Sciarrino, *Symmetry and codon usage correlations in the genetic code.* Physics Letters A, 1999. **259**(5): p. 339-348.

41. Koch, A.J. and J. Lehmann, *About a symmetry of the genetic code.* Journal of Theoretical Biology, 1997. **189**(2): p. 171-174.

42. Lenstra, R., *Evolution of the genetic code through progressive symmetry breaking.* J Theor Biol, 2014. **347**: p. 95-108.

43. Hornos, J.E.M., Y.M.M. Hornos, and M. Forger, *Symmetry and symmetry breaking: An algebraic approach to the genetic code.* International Journal of Modern Physics B, 1999. **13**(23): p. 2795-2885.

44. Antoneli, F. and M. Forger, *Symmetry breaking in the genetic code: Finite groups.* Mathematical and Computer Modelling, 2011. **53**(7-8): p. 1469-1488.

45. Das, G. and R.H.D. Lyngdoh, *Configuration of wobble base pairs having pyrimidines as anticodon wobble bases: significance for codon degeneracy.* Journal of Biomolecular Structure & Dynamics, 2014. **32**(9): p. 1500-1520.

46. Bizinoto, M.C., et al., *Codon pairs of the HIV-1 vif gene correlate with CD4+T cell count.* Bmc Infectious Diseases, 2013. **13**.

47. Wu, X.M., et al., *Computational identification of rare codons of Escherichia coli based on codon pairs preference.* Bmc Bioinformatics, 2010. **11**.

48. Tats, A., T. Tenson, and M. Remm, *Preferred and avoided codon pairs in three domains of life.* Bmc Genomics, 2008. **9**.

49. Boycheva, S., G. Chkodrov, and I. Ivanov, *Codon pairs in the genome of Escherichia coli.* Bioinformatics, 2003. **19**(8): p. 987-998.

50. Boycheva, S.S. and I.G. Ivanov, *Missing codon pairs in the genome of Escherichia coli.* Biotechnology & Biotechnological Equipment, 2002. **16**(1): p. 142-144.

51. Willie, E. and J. Majewski, *Evidence for codon bias selection at the pre-mRNA level in eukaryotes.* Trends Genet, 2004. **20**(11): p. 534-8.

52. Coghlan, A. and K.H. Wolfe, *Relationship of codon bias to mRNA concentration and protein length in Saccharomyces cerevisiae.* Yeast, 2000. **16**(12): p. 1131-45.

53. Goetz, R.M. and A. Fuglsang, *Correlation of codon bias measures with mRNA levels: analysis of transcriptome data from Escherichia coli.* Biochem Biophys Res Commun, 2005. **327**(1): p. 4-7.

54. Roymondal, U., S. Das, and S. Sahoo, *Predicting gene expression level from relative codon usage bias: an application to Escherichia coli genome.* DNA Res, 2009. **16**(1): p. 13-30.

55. Herbeck, J.T., D.P. Wall, and J.J. Wernegreen, *Gene expression level influences amino acid usage, but not codon usage, in the tsetse fly endosymbiont Wigglesworthia.* Microbiology, 2003. **149**(Pt 9): p. 2585-96.

56. Li, H. and L. Luo, *The relation between codon usage, base correlation and gene expression level in Escherichia coli and yeast.* J Theor Biol, 1996. **181**(2): p. 111-24.

57. Shen, X., S. Chen, and G. Li, *Role for gene sequence, codon bias and mRNA folding energy in modulating structural symmetry of proteins.* Conf Proc IEEE Eng Med Biol Soc, 2013. **2013**: p. 596-9.

58. Pop, C., et al., *Causal signals between codon bias, mRNA structure, and the efficiency of translation and elongation.* Mol Syst Biol, 2014. **10**: p. 770.

59. Martinez-Perez, F., et al., *Influence of codon usage bias on FGLamide-allatostatin mRNA secondary structure.* Peptides, 2011. **32**(3): p. 509-17.

60. Carlini, D.B., Y. Chen, and W. Stephan, *The relationship between third-codon position nucleotide content, codon bias, mRNA secondary structure and gene expression in the drosophilid alcohol dehydrogenase genes Adh and Adhr.* Genetics, 2001. **159**(2): p. 623-33.

61. Subramanian, A. and R.R. Sarkar, *Comparison of codon usage bias across Leishmania and Trypanosomatids to understand mRNA secondary structure, relative protein abundance and pathway functions.* Genomics, 2015. **106**(4): p. 232-41.

62. Griswold, K.E., et al., *Effects of codon usage versus putative 5'-mRNA structure on the expression of Fusarium solani cutinase in the Escherichia coli cytoplasm.* Protein Expr Purif, 2003. **27**(1): p. 134-42.

63. Gambari, R., C. Nastruzzi, and R. Barbieri, *Codon usage and secondary structure of the rabbit alpha-globin mRNA: a hypothesis.* Biomed Biochim Acta, 1990. **49**(2-3): p. S88-93.

64. Zama, M., *Codon usage and secondary structure of mRNA.* Nucleic Acids Symp Ser, 1990(22): p. 93-4.

65. Klumpp, S., J. Dong, and T. Hwa, *On ribosome load, codon bias and protein abundance.* PLoS One, 2012. **7**(11): p. e48542.

66. Zhou, J.H., et al., *The effects of the synonymous codon usage and tRNA abundance on protein folding of the 3C protease of foot-and-mouth disease virus.* Infect Genet Evol, 2013. **16**: p. 270-4.

67. McHardy, A.C., et al., *Comparing expression level-dependent features in codon usage with protein abundance: an analysis of 'predictive proteomics'.* Proteomics, 2004. **4**(1): p. 46-58.

68. Mukhopadhyay, P., S. Basak, and T.C. Ghosh, *Synonymous codon usage in different protein secondary structural classes of human genes: implication for increased non-randomness of GC3 rich genes towards protein stability.* J Biosci, 2007. **32**(5): p. 947-63.

69. Stenoien, H.K. and W. Stephan, *Global mRNA stability is not associated with levels of gene expression in Drosophila melanogaster but shows a negative correlation with codon bias.* J Mol Evol, 2005. **61**(3): p. 306-14.

70. Mishima, Y. and Y. Tomari, *Codon Usage and 3' UTR Length Determine Maternal mRNA Stability in Zebrafish.* Mol Cell, 2016. **61**(6): p. 874-85.

71. Trifonov, E.N., *Evolution of the Genetic Code and the Earliest Proteins.* Origins of Life and Evolution of Biospheres, 2009. **39**(3-4): p. 184-184.

72. Koonin, E.V. and A.S. Novozhilov, *Origin and Evolution of the Genetic Code: The Universal Enigma.* Iubmb Life, 2009. **61**(2): p. 99-111.

73. Archetti, M. and M. Di Giulio, *The evolution of the genetic code took place in an anaerobic environment.* Journal of Theoretical Biology, 2007. **245**(1): p. 169-174.

74. Wiltschi, B. and N. Budisa, *Natural history and experimental evolution of the genetic code.* Applied Microbiology and Biotechnology, 2007. **74**(4): p. 739-753.

75. Travers, A., *The evolution of the genetic code revisited.* Origins of Life and Evolution of the Biosphere, 2006. **36**(5-6): p. 549-555.

76. Knight, R.D. and L.F. Landweber, *The early evolution of the genetic code.* Cell, 2000. **101**(6): p. 569-572.

77. Jimenez-Montano, M.A., *Protein evolution drives the evolution of the genetic code and vice versa.* Biosystems, 1999. **54**(1-2): p. 47-64.

78. Davis, B.K., *Evolution of the genetic code.* Progress in Biophysics & Molecular Biology, 1999. **72**(2): p. 157-243.

79. JimenezSanchez, A., *On the origin and evolution of the genetic code.* Journal of Molecular Evolution, 1995. **41**(6): p. 712-716.

80. Beland, P. and T.F.H. Allen, *The Origin and Evolution of the Genetic-Code.* Journal of Theoretical Biology, 1994. **170**(4): p. 359-365.

81. Baumann, U. and J. Oro, *3 Stages in the Evolution of the Genetic-Code.* Biosystems, 1993. **29**(2-3): p. 133-141.

82. Osawa, S., et al., *Recent-Evidence for Evolution of the Genetic-Code.* Microbiological Reviews, 1992. **56**(1): p. 229-264.

83. Arques, D.G., J.P. Fallot, and C.J. Michel, *An evolutionary analytical model of a complementary circular code simulating the protein coding genes, the 5' and 3' regions.* Bull Math Biol, 1998. **60**(1): p. 163-94.

84. Arques, D.G. and C.J. Michel, *A complementary circular code in the protein coding genes.* J Theor Biol, 1996. **182**(1): p. 45-58.

85. El Soufi, K. and C.J. Michel, *Circular code motifs in genomes of eukaryotes.* J Theor Biol, 2016. **408**: p. 198-212.

86. Michel, C.J., *Circular code motifs in transfer RNAs.* Comput Biol Chem, 2013. **45**: p. 17-29.

87. Michel, C.J., *Circular code motifs in transfer and 16S ribosomal RNAs: a possible translation code in genes.* Comput Biol Chem, 2012. **37**: p. 24-37.

88. Michel, C.J., *An extended genetic scale of reading frame coding.* J Theor Biol, 2015. **365**: p. 164-74.

89. Michel, C.J., *A genetic scale of reading frame coding.* J Theor Biol, 2014. **355**: p. 83-94.

90. Ahmed, A., G. Frey, and C.J. Michel, *Frameshift signals in genes associated with the circular code.* In Silico Biol, 2007. **7**(2): p. 155-68.

91. Ahmed, A., G. Frey, and C.J. Michel, *Essential molecular functions associated with the circular code evolution.* J Theor Biol, 2010. **264**(2): p. 613-22.

92. Michel, C.J., *The maximal C(3) self-complementary trinucleotide circular code X in genes of bacteria, eukaryotes, plasmids and viruses.* J Theor Biol, 2015. **380**: p. 156-77.

93. Michel, C.J., M. Pellegrini, and G. Pirillo, *Maximal dinucleotide and trinucleotide circular codes.* J Theor Biol, 2016. **389**: p. 40-6.

94. Russell, R.D. and A.T. Beckenbach, *Recoding of Translation in Turtle Mitochondrial Genomes: Programmed Frameshift Mutations and Evidence of a Modified Genetic Code.* Journal of Molecular Evolution, 2008. **67**(6): p. 682-695.

95. Masuda, I., M. Matsuzaki, and K. Kita, *Extensive frameshift at all AGG and CCC codons in the mitochondrial cytochrome c oxidase subunit 1 gene of Perkinsus marinus (Alveolata; Dinoflagellata).* Nucleic Acids Research, 2010. **38**(18): p. 6186-6194.

96. Haen, K.M., W. Pett, and D.V. Lavrov, *Eight new mtDNA sequences of glass sponges reveal an extensive usage of+1 frameshifting in mitochondrial translation.* Gene, 2014. **535**(2): p. 336-344.

97. Seligmann, H., *Overlapping genetic codes for overlapping frameshifted genes in Testudines, and Lepidochelys olivacea as special case.* Comput Biol Chem, 2012. **41**: p. 18-34.

98. Seligmann, H., *An overlapping genetic code for frameshifted overlapping genes in Drosophila mitochondria: antisense antitermination tRNAs UAR insert serine.* J Theor Biol, 2012. **298**: p. 51-76.

99. Wenthzel, A.M., M. Stancek, and L.A. Isaksson, *Growth phase dependent stop codon readthrough and shift of translation reading frame in Escherichia coli.* FEBS Lett, 1998. **421**(3): p. 237-42.

100. Namy, O., et al., *Identification of stop codon readthrough genes in Saccharomyces cerevisiae.* Nucleic Acids Research, 2003. **31**(9): p. 2289-2296.

101. Loughran, G., et al., *Evidence of efficient stop codon readthrough in four mammalian genes.* Nucleic Acids Research, 2014. **42**(14): p. 8928-8938.

102. Stiebler, A.C., et al., *Ribosomal Readthrough at a Short UGA Stop Codon Context Triggers Dual Localization of Metabolic Enzymes in Fungi and Animals.* Plos Genetics, 2014. **10**(10).

103. Jungreis, I., et al., *Evidence of abundant stop codon readthrough in Drosophila and other metazoa.* Genome Research, 2011. **21**(12): p. 2096-2113.

104. Howard, M.T., et al., *Readthrough of dystrophin stop codon mutations induced by aminoglycosides.* Annals of Neurology, 2004. **55**(3): p. 422-426.

105. Dunn, J.G., et al., *Ribosome profiling reveals pervasive and regulated stop codon readthrough in Drosophila melanogaster.* Elife, 2013. **2**.

106. Steneberg, P. and C. Samakovlis, *A novel stop codon readthrough mechanism produces functional Headcase protein in Drosophila trachea.* Embo Reports, 2001. **2**(7): p. 593-597.

107. Williams, I., et al., *Genome-wide prediction of stop codon readthrough during translation in the yeast Saccharomyces cerevisiae.* Nucleic Acids Research, 2004. **32**(22): p. 6605-6616.

1	108.	Chen, J., et al., *Dynamic pathways of-1 translational frameshifting.* Nature, 2014.
2	**512**(7514): p. 328-+.
3	109.	Dinman, J.D., *Mechanisms and implications of programmed translational*
4	*frameshifting.* Wiley Interdisciplinary Reviews-Rna, 2012. **3**(5): p. 661-673.
5	110.	Smekalova, Z. and T. Ruml, *Programmed translational frameshifting - Translation*
6	*of alternative products.* Chemicke Listy, 2006. **100**(12): p. 1068-1074.
7	111.	Farabaugh, P.J., *Programmed translational frameshifting.* Microbiological Reviews,
8	1996. **60**(1): p. 103-&.
9	112.	Morgens, D.W. and A.R. Cavalcanti, *An alternative look at code evolution: using*
10	*non-canonical codes to evaluate adaptive and historic models for the origin of the genetic*
11	*code.* J Mol Evol, 2013. **76**(1-2): p. 71-80.
12	113.	Santos, M.A., et al., *Driving change: the evolution of alternative genetic codes.*
13	Trends Genet, 2004. **20**(2): p. 95-102.
14	114.	Ardell, D.H. and G. Sella, *On the evolution of redundancy in genetic codes.* J Mol
15	Evol, 2001. **53**(4-5): p. 269-81.
16	115.	Maeshiro, T. and M. Kimura, *The role of robustness and changeability on the*
17	*origin and evolution of genetic codes.* Proc Natl Acad Sci U S A, 1998. **95**(9): p. 5088-93.
18	116.	Wang, X.W., X.; Chen, G.; Zhang, J.; Liu, Y.; Yang C. , *The shiftability of protein*
19	*coding genes: the genetic code was optimized for frameshift tolerating.* PeerJ PrePrints 2015.
20	**3** p. e806v1.
21	117.	Wang, X., et al., *Why are frameshift homologs widespread within and across*
22	*species?* bioRxiv, 2016.
23

```
                 *         20        *        40
HV1J3  : ----------ATGAGAGTGAAGGGGATCAGGAAGAA--TTA :    29
SIVCZ  : ----------ATGAAAGTAATGGAGAAGAAGAAGAG--AGA :    29
SIVGB  : ATGTCTACAGGAAACGTGTACCAGGAACTAATAAGAAGATAC :    42


                 *         60        *        80
HV1J3  : TCAGCACTTGTGGAGATGGGGCACGATGCTCCTTGGGATATT :    71
SIVCZ  : CTGGAACAGCTTATCCATAATTACAATCATAACAATCATTTT :    71
SIVGB  : CTGGTAGTGGTGAAGAAGCTATACGAAGGTAAGTATGAAGTG :    84


               *         100       *         120
HV1J3  : GATGATCTGTAGTGCTGCAGAACAATTGTGGGTCACAGTC-- :   111
SIVCZ  : GCTAACCCCATGTTTGACCTCTGAGTTATGGGTAACAGTA-- :   111
SIVGB  : TCCAGGTCTTTTTCTTATACTATGTTTA-GCCTACTAGTAGG :   125


             *         140       *         160
HV1J3  : TATTATGGGGTACCTGTGTGGAAAGAAGCAGCCACCACTCTA :   153
SIVCZ  : TATTATGGAGTACCTGTTTGGCATGATGCTGACCCGGTACTC :   153
SIVGB  : TATTATAGGAAAACAATATGTGACAGT-CTTCTATGGAGTAC :   166


           *         180       *         200        *
HV1J3  : TTTTGTGCATCAGATGCTAAAGCATAT---------GATACA :   186
SIVCZ  : TTTTGTGCCTCAGACGCTAAGGCACAT---------AGTACA :   186
SIVGB  : CAGTATGGAA-GGAAGCTAAAACACATTTGATTTGTGCTACA :   207


             220       *         240        *
HV1J3  : GAGGTACATAATGTTTGGGCCACACATGCCTGTGTACCCACA :   228
SIVCZ  : GAGGCTCATAATATTTGGGCCACACAGGCATGTGTACCTACA :   228
SIVGB  : GATAATTCAAGTCTCTGGGTAACCACTAATTGCATACCTTCA :   249


             260       *         280        *
HV1J3  : GACCCCAACCCACAAGAAGTAGTATTGGAAAATGTGACAGAA :   270
SIVCZ  : GATCCCAGTCCTCAGGAAGTATTTCTTCCAAATGTAATAGAA :   270
SIVGB  : TTGCCAGATTATGATGAGGTAGAAATTCCTGATATAAAGGAA :   291


             300       *         320        *
HV1J3  : AAATTTAA------CATGTGGAAAAATAACATGGTAGAACAG :   306
SIVCZ  : TCATTTAA------CATGTGGAAAAATAATATGGTGGACCAA :   306
SIVGB  : AATTTTACAGGACTTATAAGGGAAAATCAGATAGTTTATCAA :   333
```

Fig 1 (A). Alignment of coding sequences of HIV/SIV GP120

1

```
              *         20            *         40
HV1J3 : --------------MRVKGIRKNYQHLWRWGTMLLGILMICSA : 29
SIVCZ : --------------MKVMEKKKRDWNSLSIITIITIILLTPCL : 29
SIVGB : MSTGNVYQELIRRYLVVVKKLYEGKYEVSRSFSYTMFSLLVGI : 43
                      6 V   k k      s    t   t il6

              *         60            *         80
HV1J3 : AEQLWVTVYYGVPVWKEAATTLFCASDAKAYDTEVHNVWATHA : 72
SIVCZ : TSELWVTVYYGVFVWHDADPVLFCASDAKAHSTEAFNIWATQA : 72
SIVGB : IGKQYVTVFYGVPVWKEAKIHLICATDNSS-------LWVTTN : 79
           l5VTV5YGVPVWkeA t LfCA3Daka   te hn6WaT a

            *         100           *         120
HV1J3 : CVPTDPNPQEVVLENVTEKFN--MWKNNMVEQMHEDIISLWDQ : 113
SIVCZ : CVPTDPSPQEVFLPNVIESFN--MWKNNMVDQMHEDIISLWDQ : 113
SIVGB : CIPSLPDYDEVEIPDIKENFTGLIRENQIVYQAWHAMGSMLDT : 122
          C6P3dP pqEV 6p16 E Fn   6wkNn6V Qmhed6iS6wDq

         *         140           *         160          *
HV1J3 : SLKECVALTPLCVTLNCIDWGNDTSPNATNTTSSGGEKMEKGE : 156
SIVCZ : SLKECVELTPLCVTLQCSKANFSQAKNLTNQTSS-----PPLE : 151
SIVGB : ILKPCVKINEYCVKMQCQETENVSATTAKPITTPTTTSTVASS : 165
          sLKPCVK6tPlCVt6qC     n  a natn T3s         e

          180           *         200           *
HV1J3 : MKNCSFNITTSIRDKVQKEHALFY------KHDVVPINNSTKD : 193
SIVCZ : MKNCSFNVTTELRDKKKQVYSLFY------VEDVVNLG----- : 183
SIVGB : TEIYLDVDKNNTEEKVERNHVCRYNITGLCRDSKEEIVTNFEG : 208
          mkncsfn tt  rdKv   h lfY        dvv  6

          220           *         240           *         2
HV1J3 : NIKNDNSTRYRLISCNTSVITQACPKISFEPIPIHYCAPAGFA : 236
SIVCZ : ---NENNT-YRIINCNTTAITQACPKTSFEPIPIHYCAPAGFA : 222
SIVGB : DDVKCENNTCYMNHCNESVNTEDCQKG-LLIRCILGCVPPGYV : 250
            n nnt yr6i CNt3viT2aCpK sfepipIhyCaPaG5a

         60            *         280           *         300
HV1J3 : IIKCNDKKFNGTGPCTNVSTVQCTHGIKPVVSTQLLLNGSLAE : 279
SIVCZ : ILKCNDKDFSGKGKCTNVSTVHCTHGIKPVVTTQLLINGSLAE : 265
SIVGB : MLRYN-EKLNNNKLCSNISAVQCTQHLVATVSSFFGFNGTMHK : 292
          664cNdkkfng g C3N6StVqCThg6kpvV33qll NG36ae

              *         320           *         340
HV1J3 : EEVVIRSENFTDNAK-------TIIVQLKEPVVINCTRPSKTT : 315
SIVCZ : GNITVRVENKSKNTD-------VWIVQLVEAVSLNCHRPGNNT : 301
SIVGB : EGELIPIDDKYRGPEEFHQRKFVYKVPGKYGLKIECHRKGNRS : 335
          e    6r e1k  n          v iVqlke 6 6nChRpgn 3
```

Fig 1 (B). Alignment of protein sequences of HIV/SIV GP120

1

**2A** *Frameshift Orthologs*

Species 1 → (Speciation) → Species 2

Gene A → Gene a

Gene A: 1 2 3

Gene a: 1 2' 3

*Frameshifting*

**2B** *Frameshift Paralog*

Species 1

Gene A → (Gene Duplication) → Gene B

Gene A: 1 2 3

Gene B: 1' 2' 3'

*Frameshifting*

**Fig 2**

**2C**

# UCSC Genome Browser on Human Dec. 2013 (GRCh38/hg38) Assembly

**Fig 2**

## 2D VEGFAA and its frameshifts



```
                        *             20             *
vegfaa   : MNLVVYLIQLFLAALLHLSAVKAAHIPKEGGRSKNDVI : 38
vegfaa-1 : MTWLFI*YSYFSRLSSICLL*RLPTYPKKGERAKMM*F : 35
vegfaa-1-r : MTWLFIWYSYFSRLSSICLLKRLPTYPKKGERAKMMWF : 38
vegfaa-2 : -MLGCLFDTVISRGSPPSVCCKGCPHTQRRGREQK*CD : 36
vegfaa-2-r : -MLGCLFDTVISRGSPPSVCCKGCPHTQRRGREQKWCD : 37
                        sr s              4          4

             40            *             60             *
vegfaa   : PFMDVYKKSACKTRELLVDIIQEYPDEIEHTYIPSCVV : 76
vegfaa-1 : PSWMCIKRVRARPESCW*TSSRSIPMRSSTRTSRPVWF : 72
vegfaa-1-r : PSWMCIKRVRARPESCWSTSSRSIPMRSSTRTSRPVWF : 76
vegfaa-2 : SLHGCV*KECVQDPRAAGRHHPGVSR*DRAHVHPVLCG : 72
vegfaa-2-r : SLHGCVKKECVQDPRAAGRHHPGVSRWDRAHVHPVLCG : 75
             c k4

             80            *            100            *
vegfaa   : LMRCAGCCNDEALECVPTETRNVTMEVLRVKQRVSQHN : 114
vegfaa-1 : SCAVQDAVMMRRSNASRQRHETSLWRCGSSNAYRSII : 110
vegfaa-1-r : SCAVQDAVMMRRSNASRQRHETSLWRCGSSNAYRSII : 114
vegfaa-2 : SHALCRML***GARMRPDRDTKRHYGGAAGQATRIAA* : 106
vegfaa-2-r : SHALCRMLKWWGARMRPDRDTKRHYGGAAGQATRIAAK : 113
             s a                  r

             120           *            140            *
vegfaa   : FQLSFTEHTKCECRPKAEVKAKKENHCEPCSERRKRLY : 152
vegfaa-1 : FS*VSQNTPSVNAGQRQKSKQRKKTTVSLAQREGSACM : 147
vegfaa-1-r : FSWVSQNTPSVNAGQRQKSKQRKKTTVSLAQREGSACM : 152
vegfaa-2 : FSAEFHRTHQV*MQAKGRSQSKERKPL*ALLREKEALV : 142
vegfaa-2-r : FSAEFHRTHQVWMQAKGRSQSKERKPLWALLREKEALV : 151
             Fs     t  v     4  s 4           re a

             160           *            180
vegfaa   : VQDPLTCKCSCKFTQMQCKSRQLELNERTCRCEKPR- : 188
vegfaa-1 : CRTPSPVNAPANSHK-CNASPDNLS*TRELADVKSQD : 182
vegfaa-1-r : CRTPSPVNAPANSHK-CNASPDNLSKTRELADVKSQD : 188
vegfaa-2 : CAGPPHL*MLLQIHTNAMQVQTT*VKRKNLQM*KAKM : 176
vegfaa-2-r : CAGPPHLKMLLQIHTNAMQVQTTWVKRKNLQMWKAKM : 188
             c  P        h              4 l   K
```
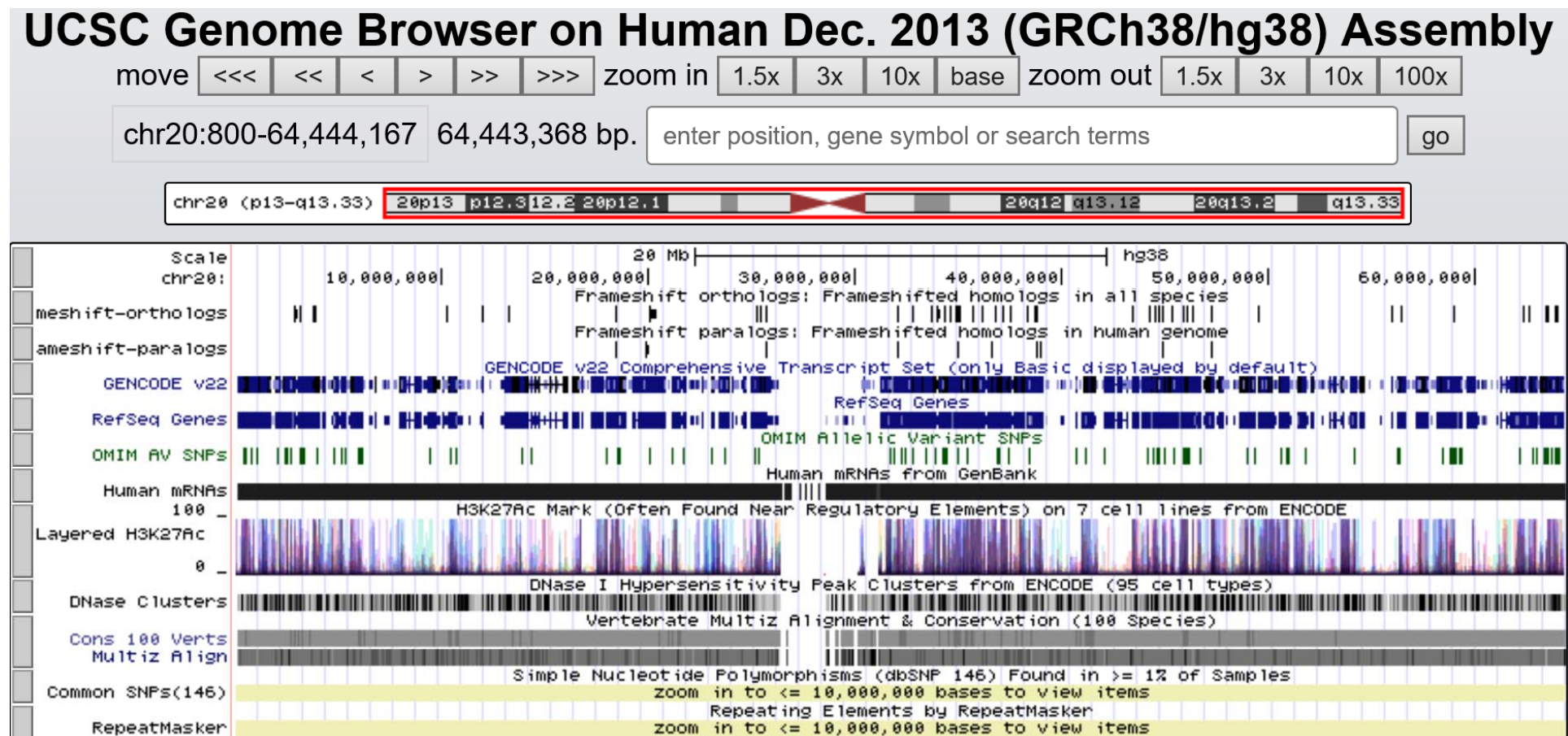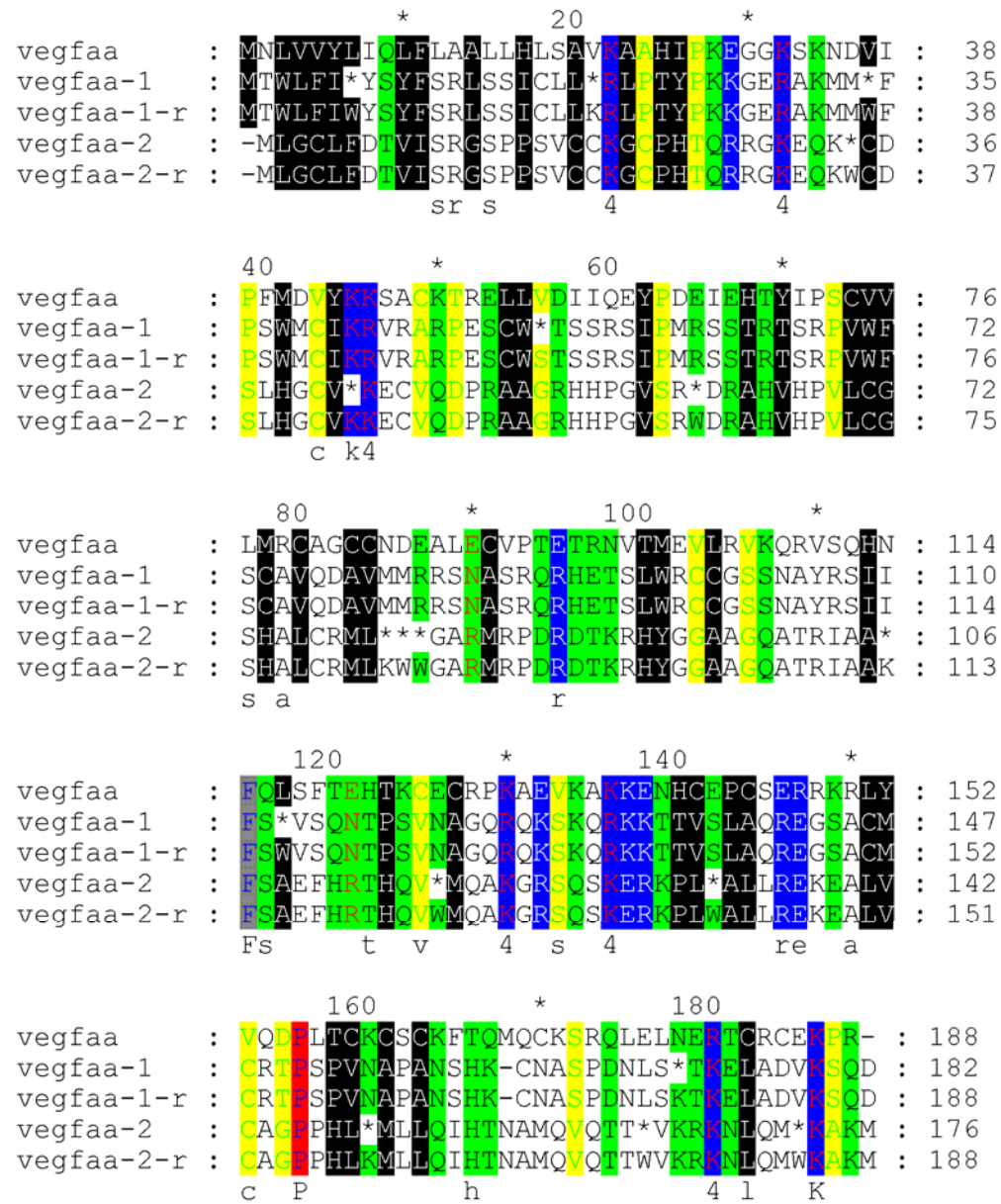
Fig 2

**2E**

```
Score = 63.5 bits (153),   Expect = 3e-09, Method: Compositional matrix adjust.
Identities = 38/66 (58%), Positives = 42/66 (64%), Gaps = 0/66 (0%)

Query  942    RQADRENRPHPGRGYPHG**AHPPVPLRQSAPAGALHADTICRAAGRKSLSLRPAGPQGG   1001
              RQ D E+ PHP       G  A PP+PLRQ AP G  HADT  RA GR+ L LRPAG   G
Sbjct  763    RQTDGEDGPHPDGCDTDGRRADPPLPLRQPAPPGVPHADTAWRATGREPLPLRPAGTANG   822


Query  1002   KTGMAA   1007
              +TGMAA
Sbjct  823    ETGMAA   828
```

*Query translation  is non-readthrough*

```
Score = 64.7 bits (156),   Expect = 1e-09, Method: Compositional matrix adjust.
Identities = 38/66 (58%), Positives = 42/66 (64%), Gaps = 0/66 (0%)

Query  942    RQADRENRPHPGRGYPHGWWAHPPVPLRQSAPAGALHADTICRAAGRKSLSLRPAGPQGG   1001
              RQ D E+ PHP       G  A PP+PLRQ AP G  HADT  RA GR+ L LRPAG   G
Sbjct  767    RQTDGEDGPHPDGCDTDGRRADPPLPLRQPAPPGVPHADTAWRATGREPLPLRPAGTANG   826


Query  1002   KTGMAA   1007
              +TGMAA
Sbjct  827    ETGMAA   832
```

*Query translation is readthrough*

**Fig 2**