

Relationship Between G + C in Silent Sites of Codons and Amino Acid Composition of Human Proteins

David W. Collins and Thomas H. Jukes

Space Sciences Laboratory, University of California, Berkeley, 6701 San Pablo Avenue, Oakland, CA 94608, USA

Summary. We have investigated the relationship between the G + C content of silent (synonymous) sites in codons and the amino acid composition of encoded proteins for approximately 1,600 human genes. There are positive correlations between silent site G + C and the proportions of codons for Arg, Pro, Ala, Trp, His, Gln, and Leu and negative ones for Tyr, Phe, Asn, Ile, Lys, Asp, Thr, and Glu. The median proteins coded by groups of genes that differ in silent-site G + C content also differ in amino acid composition, as do some proteins coded by homologous genes. The pattern of compositional change can be largely explained by directional mutation pressure, the genetic code, and differences in the frequencies of accepted amino acid substitutions; the shifts in protein composition are likely to be selectively neutral.

Key words: G + C content — Silent sites — GC pressure — Directional mutation pressure — Human genome — Codon usage — Amino acid composition — Neutral theory

The G + C content of DNA varies widely among organisms, especially microorganisms. Among bacterial species the mean G + C content ranges from approximately 25% to 75%, but the amount of intragenomic heterogeneity is small (Rolfe and Meselson 1959; Sueoka et al. 1959). Vertebrate genomes,

however, are characterized by a relatively large intragenomic variability in base composition (Sueoka 1959) and this variability is reflected in sequence data (e.g., Ikemura and Aota 1988; Ikemura et al. 1990). It is likely that the human and other vertebrate genomes are mosaics of fairly large regions of similar base composition ("isochores") belonging to a few distinct classes (Bernardi et al. 1985). These may be associated with chromosomal banding. GC-rich genes are most often associated with chromosomal R-(Giemsa pale, quinacrine dull) bands and AT-rich genes with G-(Giemsa dark, quinacrine bright) bands (Ikemura and Wada 1991). Human genes with high G + C in third codon (mostly silent) positions are typically surrounded by GC-rich sequences, whereas those with low third position G + C are embedded in relatively AT-rich sequences (Aota and Ikemura 1986).

Sueoka (1962, 1988, 1992) has proposed that variations in G + C content may be due to directional mutational pressure. According to this theory, the G + C content of DNA is determined by the effective base conversion rates u (G·C to A·T) and v (A·T to G·C). If these rates are unequal, a directional mutation pressure will result. For example, if v is greater than u , GC pressure will result; the G + C content at equilibrium is $v/(u + v)$. These pressures could result from biases in the process of DNA replication and repair (reviewed in Filipitski 1990). GC pressure was demonstrated experimentally by Cox and Yanofsky (1967), who found that the genomic G + C content of *E. coli* increased

following mutations in the Treffers mutator gene (*mutT*). A mutation in *mutT* results in a high frequency of A·T to C·G transversions. The normal *mutT* protein degrades 8-oxo-7,8-dihydro-2'-dGTP, which may otherwise be incorporated opposite A during DNA replication (Maki and Sekiguchi 1992). The human genome may be subject to a range of mutational biases which could be responsible for its mosaic structure (Filipski 1990; Sueoka 1992).

Changes in genomic G + C content are accompanied by small but significant shifts in the amino acid composition of proteins. Sueoka (1961) examined bacterial species ranging in G + C content from 35% to 72% and showed that mean DNA base composition was correlated with the amino acid composition of total bacterial protein. The protein from more GC-rich organisms was found to have more Pro, Ala, Gly, Arg, and less Lys, Asn, Tyr, Phe, Ile. These effects have also been observed in comparisons of homologous gene sequences from bacteria and mitochondria (Jukes and Bhushan 1986; Muto and Osawa 1987; Filipski 1990) and viruses (Karlin et al. 1990). Typically, large differences in silent-site G + C content are accompanied by smaller differences in the G + C content of replacement sites (mostly codon first and second positions) (e.g., Jukes and Bhushan 1986). For samples of human genes, positive correlations have been found between the G + C content of codon third positions and the G + C content of codon first and second positions (Sueoka 1988, 1992; Aissani et al. 1991; D'Onofrio et al. 1991). In the present work, we attempt to characterize the influence of this effect on the composition of human proteins.

Methods and Results

Sample of Human Genes. We used codon usage tables to calculate the silent-site G + C content and amino acid composition coded by human genes. Wada et al. (1991) have compiled codon usage from the GenBank genetic sequence data bank (Bilofsky and Burks 1988). We obtained a magnetic tape containing codon usage for 2,681 human gene sequences extracted from release 69.0 of GenBank (Wada et al. 1991; T. Ikemura personal communication).

The same gene may be represented by more than one entry and LOCUS name in GenBank. As an example, an identical 963-bp mRNA sequence is present in six different GenBank files, variously defined as endonexin II (GenBank LOCUS name HUMENN), lipocortin-V (HUMLC5), blood coagulation inhibitor (HUMBIC), placental anticoagulant protein PAP (HUMATC), placental anticoagulant PP4 (HUMPAP4), and vascular anticoagulant (HUMVAC) (ACCESSION numbers m18366, d00172, j03745, m21731, m19384, x12454).

The following procedure was used to filter out data from redundant and nearly identical sequence submissions. A database of codon tables for 2,681 GenBank entries was sorted according to the following criteria: (1) total number of codons, (2) the Euclidean distance of the counts of the four nucleotides from

the origin, (3) the number of occurrences of the nucleotide A, and (4) the number of occurrences of the nucleotide C. Identical and nearly identical sequences necessarily have similar codon usage; the effect of the sorting procedure is to bring these together. The amino acid counts specified by each database entry were then compared to the those of the preceding entry in the sorted list, using a simple distance measure (Collins et al. 1992):

$$\text{Distance} = \left[(L_1 - L_2)^2 + \sum_{i=1}^{20} (N_{i1} - N_{i2})^2 \right]^{1/2}$$

where N_{i1} and N_{i2} represent the counts of synonymous codons for each of the 20 amino acids and L_1 and L_2 represent the total number of codons. Additionally, the terms in the sum corresponding to Cys and Trp were multiplied by five and those corresponding to Gly, Pro, Phe, Leu, and Tyr were multiplied by two. If the calculated distance was less than 10, the corresponding record was deleted from the sample. The choice of 10 was arbitrary; 937 records were deleted on this basis. This procedure also removed closely related sequences that do not represent redundant data-bank submissions. For example, HUMCNPG, green cone photoreceptor pigment, was retained but HUMCNPR, the closely related red cone photoreceptor having 96% amino acid sequence identity (Nathans et al. 1986), was excluded. Entries having fewer than 100 codons were also deleted; there were 137 of these. The filtered sample represents 1,607 human gene sequences containing 762,710 codons. These are listed in order of increasing silent-site G + C content and identified by GenBank LOCUS name in the Appendix. A floppy disk containing the codon usage data is available upon request.

Silent Sites, Replacement Sites, GC- and AT-Coded Amino Acids

We define codon silent sites as those nucleotide positions that may potentially undergo substitution without changing the meaning of the codon. In the "universal" genetic code, these are the first-position nucleotides of Arg (CGR, AGR) and Leu (CTR, TTR) codons¹ and the third-position nucleotides of all codons except Met (ATG) and Trp (TGG). The other codon positions are replacement sites. These distinctions help to differentiate between changes in base composition which influence amino acid composition and those which only produce changes in synonymous codon usage. For example, a correlation between the G + C content of codon first and third positions is ambiguous: it may indicate either changes in amino acid coding or changes in synonymous codon usage, since both first and third positions are a mix of silent and replacement sites. There are 61 codons for amino acids. Of these, most (53/61) codon first positions are replacement sites, and most (59/61) third positions are silent. All second positions are replacement

¹ R = A or G, Y = C or T(U), N = any base. Silent sites are underlined.

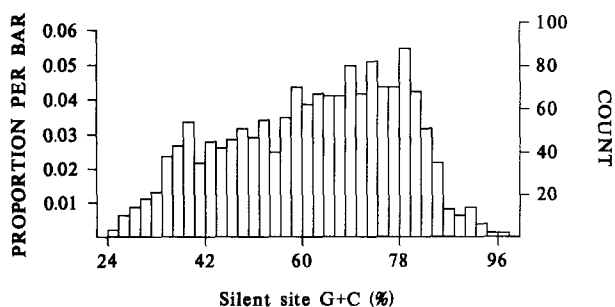


Fig. 1. G + C content of codon silent sites, 1,607 human genes. Silent sites are defined in the text. The sample was culled from a larger collection of 2,681 coding sequences from release 69.0 of the GenBank database.

sites. Because the genetic code provides an alternative codon containing either a G·C or A·T base pair for all amino acids having two or more codons, the G + C content of silent sites is free to vary from 0 to 100%, independently of amino acid composition.

GC-coded and AT-coded amino acids are defined as those with replacement sites occupied exclusively by G/C and A/T, respectively. The GC-coded amino acids are Ala (GCN), Arg (CGN, AGR), Gly (GGN), and Pro (CCN) and the AT-coded amino acids are Asn (AA \bar{Y}), Ile (AT \bar{Y} , ATA), Lys (AAR), Phe (TT \bar{Y}), and Tyr (TAY). Note that Met (ATG) is not AT-coded because the third position is a replacement site occupied by "G." Similarly, Arg (CGN, AGR) is GC-coded because the first positions of AGR codons are silent sites.

Correlation Between Silent-Site G + C Content and Amino Acid Composition

For each gene in the sample we used the corresponding codon table to calculate the G + C content of silent and replacement sites and the amino acid composition of the protein. This was done with the aid of microcomputers and spreadsheet software.

Figure 1 is a histogram representing the silent G + C contents of the sample of genes. The shape of this broad, possibly multimodal distribution is familiar from previous studies of third-position (mostly silent) G + C content in human genes (Karlin et al. 1990; Bernardi and Bernardi 1991). The silent-site G + C content of the sample ranges from 24.3% (GenBank LOCUS HUMACADM, acyl-CoA dehydrogenase) to 96.6% (HUMMIFA, migration inhibitory factor MIF). This wide variation in synonymous codon usage is typical of vertebrate genomes (Ikemura and Aota 1988; Ikemura et al. 1990; Wada et al. 1991).

The G + C content of silent vs. replacement sites for the sample of genes is plotted in Fig. 2. Although

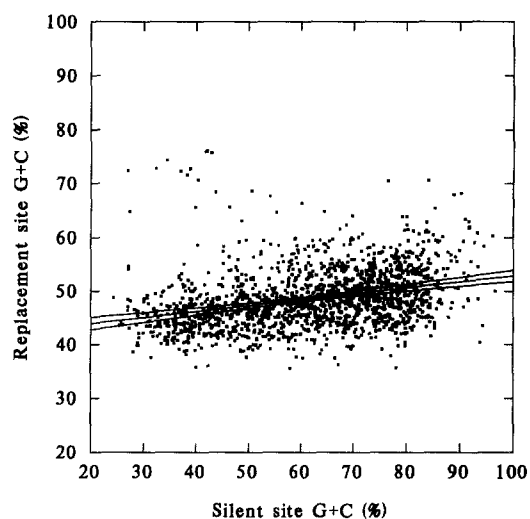


Fig. 2. G + C content of replacement vs silent sites in codons, 1,607 human genes. A regression line is plotted with 99% confidence intervals ($r = 0.31$, $P < 0.001$). No outlying points have been excluded.

the points are widely scattered, there is a significant positive association ($r = 0.31$, $P < 0.001$). Figures 3 and 4 show the relationships between silent G + C and the fractions of GC-coded (Ala + Arg + Gly + Pro) and AT-coded amino acids (Asn + Ile + Lys + Phe + Tyr) for the sample of genes. The correlations are significantly positive and negative, respectively.

Table 1 summarizes the correlations between silent G + C and the fraction of codons for each of the 20 amino acids. The levels of the individual AT-coded amino acids are all negatively correlated with silent-site G + C. The levels of the individual GC-coded amino acids are positively and significantly correlated with silent G + C with the exception of Gly. The sample contains several genes (e.g., collagens, proline-rich proteins) that are extreme outliers with respect to content of glycine and proline codons. These are mostly low in silent G + C and exert a large influence on the correlations. When the 20 most proline-rich proteins are removed from the sample, the correlation between Gly and silent G + C is significantly positive ($P < 0.001$).

Among amino acids with "mixed" coding (i.e., replacement sites occupied by both G·C and A·T), Trp, His, Gln, and Leu are positively correlated with silent G + C; Asp, Thr, and Glu are negatively correlated. The correlations of Met, Cys, Val, and Ser are not significantly different from zero.

Average Protein Composition

The preceding results suggest that the average composition of human proteins may vary, depending on the pattern of synonymous codon usage of the

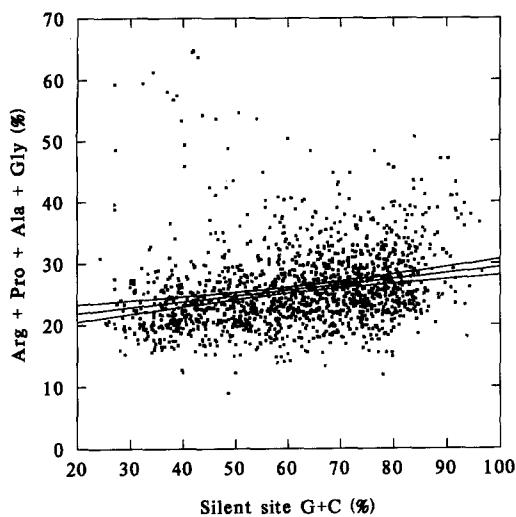


Fig. 3. Percentage of codons for GC-coded amino acids (Arg, Pro, Ala, Gly) vs silent-site G + C content, 1,607 human genes. The regression line is shown with a 99% confidence interval ($R = 0.22$, $P < 0.001$).

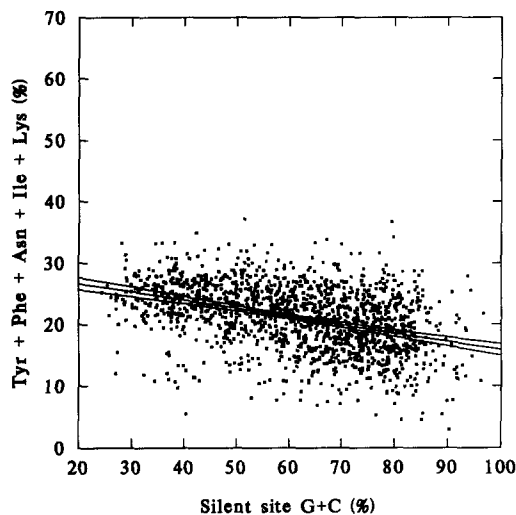


Fig. 4. Percentage of codons for AT-coded amino acids (Tyr, Phe, Asn, Ile, Lys) vs silent-site G + C content, 1,607 human genes ($R = 0.41$, $P < 0.001$).

genes. In order to quantify these differences, we divided the sample of genes into three groups (low, medium, high) on the basis of silent site G + C content. The low group consisted of 473 genes having from 24% to 52% silent G + C; the medium group consisted of 578 genes having from 52% to 71% silent G + C; and the high group consisted of 556 genes having 71–97% silent G + C. The divisions were chosen to place, as nearly as possible, equal numbers of codons into three groups; the decision to use three groups was arbitrary. We then compared the average proteins coded by these groups.

The composite codon usage of the three groups is

presented in Table 1. Increases in the proportion of the GC-coded amino acids Ala and Arg (but not Gly and Pro) and decreases in AT-coded amino acids are present across the three groups (from low to high). The largest changes are for Lys and Arg; the ratio of Lys to Arg codons is 1.47:1 for the low group and decreases to 0.81:1 for the high group. There is a stepwise increase in the ratio of GC- to AT-coded amino acids across the three groups. For amino acids with mixed coding, the proportions are relatively constant except for Leu, which shows a stepwise increase (Table 1).

The distributions of amino acid content of samples of proteins are often skewed to the right; consequently the median rather than mean composition may provide a better representation. The median protein compositions coded by the three groups of genes are shown in Figs. 5 and 6. Regular increases in the fractions of GC-coded amino acids (Arg, Pro, Ala, Gly) and decreases in AT-coded amino acids (Tyr, Phe, Asn, Ile, Lys) are apparent (Fig. 5). Among GC-coded amino acids, the increases in Arg and Ala are greater than for Pro and Gly. For AT-coded amino acids the decreases in Lys, Ile, and Asn are larger than for Tyr and Phe. For amino acids with mixed coding, there are small but stepwise decreases in Thr and Glu, a small increase in Gln, and a larger increase in Leu.

Homologous Genes—Steroid Hormone Receptors, Ras Proto-Oncogenes

There are relatively few opportunities to compare homologous human genes differing greatly in base composition, but some examples are presented in Table 2. The amino acid composition and silent-site G + C content coded by three *ras* genes are compared, as are those of three members of the steroid hormone receptor superfamily. Silent G + C is again more subject to increase than is G + C in replacement sites. Even in these small samples, the trends in amino acid composition resemble those in Table 1. The ratio of GC- to AT-coded amino acids increases with silent site G + C.

In the case of the three *ras* proteins, which are very similar in sequence, most of the change in amino acid composition is concentrated in a 25-residue C-terminal variable region which is “not required for any of the known biochemical functions of the protein” (de Vos et al. 1988). In *K-ras* (30% silent G + C) this region contains 13 AT-coded amino acids (12 Lys + 1 Ile) and one GC-coded amino acid (Gly). In *H-ras* (82% silent G + C) there are only 4 AT-coded amino acids (3 Lys, 1 Asn) and 6 GC-coded amino acids (3 Pro, 2 Gly, 1 Arg). Thus

Table 1. Average composition of protein coded by 1,607 human genes^a

| Group Silent G + C No. genes | No. of codons | | | 1,607 human genes protein composition | | |
|--|----------------------|-------------------------|-----------------------|--|-----------|---------------------------------|
| | Low 24–52% 473 | Medium 52–71% 578 | High 71–97% 556 | Avg. (%) | SD (%) | (<i>r</i>) vs silent G + C |
| <i>(a) AT-coded amino acids</i> | | | | | | |
| Tyr (TAY) | 7,999 | 7,619 | 7,123 | 3.0 | 1.3 | –0.11*** |
| Phe (TTY) | 9,784 | 9,625 | 9,356 | 3.8 | 1.6 | –0.06* |
| Asn (AAY) | 11,272 | 9,904 | 8,501 | 3.8 | 1.6 | –0.33*** |
| Ile (ATY/A) | 13,121 | 11,300 | 10,178 | 4.5 | 1.9 | –0.29*** |
| Lys (AAR) | 16,986 | 14,436 | 12,092 | 6.0 | 2.9 | –0.29*** |
| <i>(b) Amino acids with mixed coding</i> | | | | | | |
| Trp (TGG) | 3,000 | 3,447 | 3,674 | 1.3 | 1.0 | 0.11*** |
| Met (ATG) | 5,161 | 4,929 | 5,045 | 2.0 | 1.1 | –.03 |
| Cys (TGY) | 5,978 | 5,757 | 6,212 | 2.4 | 2.1 | .04 |
| His (CAY) | 5,838 | 6,235 | 6,181 | 2.4 | 1.2 | .06* |
| Gln (CAR) | 10,879 | 11,252 | 11,798 | 4.4 | 1.9 | .06* |
| Asp (GAY) | 12,977 | 12,862 | 12,718 | 5.0 | 1.9 | –0.09*** |
| Thr (ACN) | 14,496 | 14,242 | 13,490 | 5.5 | 1.8 | –0.09*** |
| Val (GTN) | 15,868 | 15,760 | 16,307 | 6.3 | 1.9 | .02 |
| Glu (GAR) | 17,506 | 17,555 | 17,186 | 6.8 | 3.1 | –0.07** |
| Ser (TCN/AGY) | 18,672 | 19,152 | 18,219 | 7.3 | 2.5 | .00 |
| Leu (CTN/TTR) | 21,714 | 23,962 | 26,202 | 9.7 | 2.8 | 0.29*** |
| <i>(c) GC-coded amino acids</i> | | | | | | |
| Arg (CGN/AGR) | 11,534 | 13,640 | 14,963 | 5.5 | 2.4 | 0.23*** |
| Pro (CCN) | 15,156 | 15,150 | 15,513 | 5.8 | 3.1 | 0.06* |
| Ala (GCN) | 16,288 | 17,651 | 19,805 | 7.3 | 2.7 | 0.25*** |
| Gly (GGN) | 18,883 | 18,855 | 18,488 | 7.3 | 3.4 | 0.02 |
| Total | 253,112 | 253,333 | 253,051 | 100.0 | | |
| <i>(c)/(a)</i> | 1.05 | 1.23 | 1.46 | 1.23 | | |

^a The first three columns list the total number of codons for each amino acid coded by the “low,” “medium,” and “high” groups of genes. The low, medium, and high groups consist of genes with silent G + C contents ranging from 24 to 52% (473 genes), 52 to 71% (578 genes), and 71 to 97% (556 genes), respectively. The sequence files are identified by GenBank LOCUS name in

the Appendix. Initiation codons are not included in “Met” and stop codons have been omitted. The fourth column lists the mean protein composition coded by the sample of 1,607 genes. The last column lists the correlation coefficient (*r*) for silent-site G + C content (%) and amino acid abundance (%) in the sample of genes (**P* < 0.05, ***P* < .01, ****P* < 0.001)

the nonfunctional region is more subject to change, as would be expected from the neutral theory.

The effect of GC pressure on amino acid composition should be related to the fraction of sites in the protein which is free to accept amino acid substitutions. If an amino acid sequence is perfectly conserved, only the pattern of synonymous codon usage may “respond” to GC pressure. For example, the human genome contains two genes for calmodulin having divergent codon usage (26% and 68% G + C in silent sites), but in spite of this disparity, these sequences code identical proteins (Fischer et al. 1988).

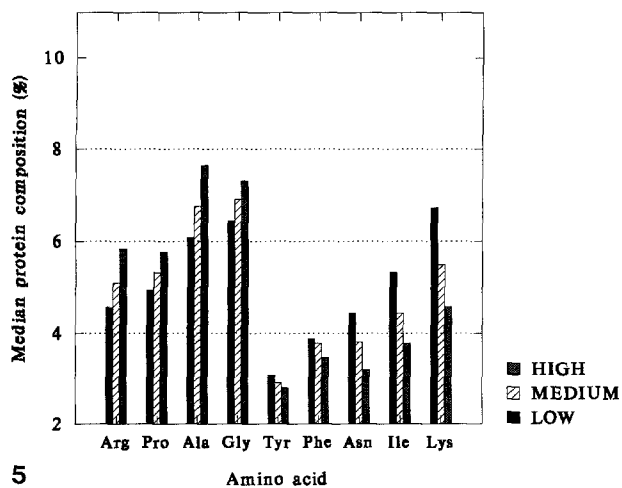
Discussion

Amino Acid Mutability

In general, we have found that GC-coded amino acids increase with silent G + C, and AT-coded

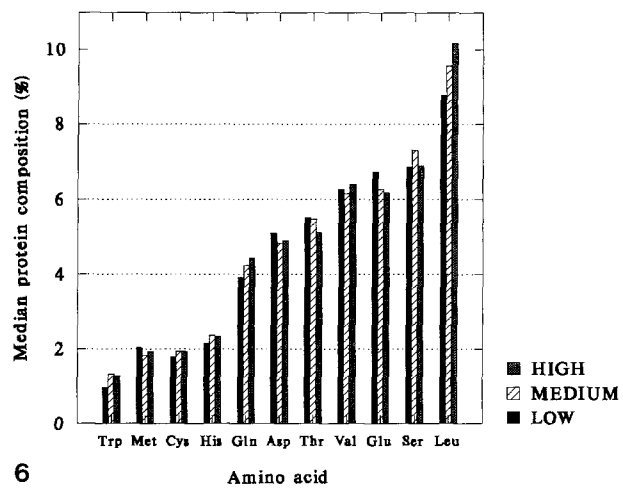
amino acids decrease. However, the “response” of the individual amino acids differs substantially in magnitude (Table 1, Figs. 5, 6). Among GC-coded amino acids, the increases in Pro and Gly with silent G + C (Table 1, Fig. 5) are smaller than those of Ala and Arg. For AT-coded amino acids, the decreases in Tyr and Phe are small relative to Asn, Ile, Lys (Table 1, Fig. 6). Similar results were obtained by Karlin et al. (1990), who compared the amino acid usage of homologous proteins in related viruses differing greatly in genomic G + C content. Several amino acids with “mixed” coding are also correlated with silent-site G + C content.

The pattern of compositional change may be explained by considering differences in rates of amino acid substitution in addition to coding. For example, amino acid exchanges involving Pro (56) and Gly (49) residues are fixed less frequently during evolution than those involving Arg (65) and Ala



5

Amino acid



6

Amino acid

Figs. 5, 6. Median protein composition coded by genes with low, medium, and high silent G + C content. Content of GC-coded (Arg, Pro, Ala, Gly) and AT-coded (Tyr, Phe, Asn, Ile, Lys) amino acids (Fig. 5). Content of amino acids with "mixed" coding (Trp, Met, Cys, His, Gln, Asp, Thr, Val, Glu, Ser, Leu) (Fig. 6).

Table 2. Amino acid composition and silent G + C content of two groups of dispersed homologous human genes^a

| Gene | <i>ras</i> oncogenes | | | Steroid hormone receptors | | |
|-------------------|----------------------|---------------|---------------|---------------------------|------|------|
| | K- <i>ras</i> | N- <i>ras</i> | H- <i>ras</i> | HAPRA | VDR | EAR3 |
| Silent G + C (%) | 29.6 | 44.1 | 81.5 | 50.5 | 76.1 | 85.4 |
| Replace G + C (%) | 42.2 | 43.9 | 45.9 | 45.2 | 45.4 | 53.5 |
| Chromosome | 12p12.1 | 1p22 | 11p15.5 | 3 | ? | 5 |
| Total codons | 188 | 188 | 188 | 447 | 426 | 424 |
| | (No. codons) | | | (No. codons) | | |
| <i>AT-coded</i> | | | | | | |
| Tyr (TAY) | 9 | 9 | 9 | 10 | 8 | 12 |
| Phe (TTY) | 6 | 6 | 5 | 16 | 18 | 14 |
| Asn (AAY) | 4 | 6 | 5 | 13 | 13 | 16 |
| Ile (ATY/A) | 15 | 11 | 11 | 25 | 24 | 26 |
| Lys (AAR) | 16 | 12 | 11 | 32 | 24 | 14 |
| (a) Total | 50 | 44 | 41 | 96 | 87 | 82 |
| <i>GC-coded</i> | | | | | | |
| Arg (CGN,AGR) | 12 | 10 | 12 | 23 | 29 | 26 |
| Pro (CCN) | 5 | 5 | 6 | 28 | 22 | 31 |
| Ala (GCN) | 9 | 10 | 11 | 21 | 20 | 36 |
| Gly (GGN) | 11 | 14 | 13 | 20 | 20 | 37 |
| (b) Total | 37 | 39 | 42 | 92 | 91 | 130 |
| (b)/(a) | 0.74 | 0.89 | 1.02 | 0.96 | 1.05 | 1.59 |

^a K-*ras* = c-Ki-*ras* oncogene (GenBank accession number K03209, Hirai et al. 1985), N-*ras* = N-*ras* proto-oncogene (X00642, Brown et al. 1984), H-*ras* = c-Ha-*ras*1 proto-oncogene (J00277, Sekiya et al. 1984), HAPRA = hepatocellular carcinoma hormone receptor (Y00291, de The et al. 1987), VDR = vitamin D receptor (J03258, Baker et al. 1988), EAR3 = v-erbA-related hormone receptor (X12795, Miyajima et al. 1988). Sequences were obtained from the GenBank database (Bilofsky and Burks 1988)

(100). (The numbers represent the relative mutabilities of the amino acids given in Dayhoff 1978.) Accordingly, the fractions of Pro and Gly may be expected to evolve at a lower rate in response to mutational biases relative to other GC-coded amino acids. Similarly, mutations involving Phe (41) and Tyr (41) are accepted less often than are those involving the other AT-coded amino acids: Asn (134), Ile (96), and Lys (56) (Dayhoff 1978). The large in-

crease and decrease in Arg and Lys, respectively (Table 1, Fig. 5), may be linked; Lys and Arg are readily exchanged through single point mutations from Arg (AGR) to Lys (AAR).

The Case of Leucine

A similar argument may explain the increase in Leu (CTN, TTR) (Table 1, Fig. 6) with silent G + C

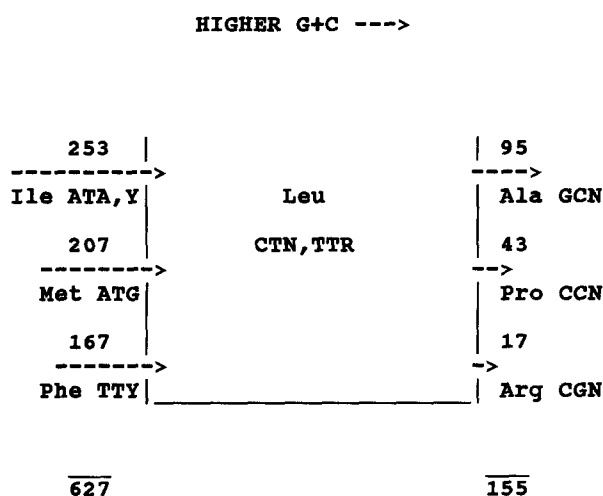


Fig. 7. Schematic representation of a possible mechanism for increase in the fraction of leucine codons under GC pressure. A mutational bias favoring G + C (GC pressure) would preferentially suggest mutations in the directions of the arrows. New leucine codons would be generated from codons with lower G + C (Ile, Met, Phe), but leucine codons would also be lost through mutation to GC-coded amino acids (Ala, Pro, Arg). However, the latter mutations are fixed less often than the former, causing a net increase in the number of Leu codons. The numbers above each arrow indicate the relative frequency of accepted point mutations between leucine and the indicated amino acid (from Dayhoff 1978).

despite relatively AT-rich mixed coding. If the underlying pattern of nucleotide substitution is biased in the direction of G + C (GC pressure), one might expect new leucine codons to be generated through point mutations of Ile (ATA, ATY), Phe (TTY), and Met (ATG), but simultaneously eliminated via mutation of other Leu codons to Pro (CCN), Arg (CGN), and Ala (GCN, two-base change). These pathways are not selectively equivalent, however, and therefore may not balance. Exchanges between Leu and Ile/Phe/Met (lower GC) are fixed about four times as frequently as those involving Ala/Pro/Arg (higher GC) (Dayhoff 1978). Consequently, the fraction of leucine codons may be predicted to increase under GC pressure, rather than remain constant. This scheme is diagrammed in Fig. 7.

Compositional Change and Neutral Evolution

There are two broad viewpoints on heterogeneity in the G + C content of DNA. These differ fundamentally with regard to the role of natural selection. The first viewpoint holds that compositional variation is maintained by positive Darwinian selection for G + C content per se or for particular amino acids in proteins (Bernardi and Bernardi 1986a,b). In contrast, the second viewpoint attributes compositional variation to directional mutation biases, rather than to selective advantage (Sueoka 1988); the composi-

tional shifts are viewed as selectively neutral phenomena. The neutral theory of molecular evolution postulates that "nucleotide substitutions inherently take place in DNA as a result of point mutations followed by random genetic drift" (Jukes and Kimura 1984). The substitution rate is highest in the absence of selective constraints—for example, in pseudogenes. Lower rates indicate the presence of constraints imposed by negative selection, which rejects deleterious mutations.

Bernardi and Bernardi (1986a,b) have proposed that GC-rich regions of chromosomes in warm-blooded vertebrates have evolved to accommodate higher body temperatures because G·C pairs have a higher melting point than A·T pairs. They offer selectionist explanations for the observed regional variations, attributing them to "compositional constraints." Higher GC levels in mRNAs are proposed to "increase their base-pairing and stability." Additionally, increases in G + C content of coding sequences "lead to thermodynamically more stable proteins" (Bernardi and Bernardi, 1986b). Specifically, alanine and arginine "are most frequently acquired in thermophiles: and produce an increase in heat stability, while serine and lysine, which diminish stability, are correspondingly lost."

Sueoka (1992) has criticized the proposal that large-scale variations in base composition are maintained by natural selection because point mutations away from the hypothetical optimal G + C content would have to be considered deleterious on the basis of an infinitesimal effect on the heat stability of DNA. Furthermore, some thermophilic organisms are high in A + T (Filipski 1990).

Alanine and arginine may accumulate under GC pressure via neutral mutations because they are GC-coded, rather than via positive selection operating at the protein level. Filipski compared sequences from (1) *Thermus thermophilus* (high GC, thermophilic) and *Saccharomyces cerevisiae* (high AT, mesophilic) and (2) from *Streptomyces limosus*, *Streptomyces hygroscopicus* (high GC, mesophilic), and *Dictyoglomus thermophilium* (high AT, thermophilic). Alanine content was found to correlate with genomic G + C rather than with the thermophilicity of the organisms. We have found that serine, listed with lysine as decreasing thermal stability, is not negatively correlated with G + C in silent sites (Table 1, Fig. 6). The decrease in lysine that accompanies increasing GC may proceed by selectively neutral mutations from Lys (AAR) to Arg (AGR) suggested by GC pressure, rather than from selective advantage.

According to the neutral theory of molecular evolution, nucleotide changes take place and are fixed by genetic drift. If mutations in the direction AT → GC predominate, then the GC content of

DNA will increase, although this trend may be opposed by negative selection to some extent. The neutral theory predicts that this increase will take place most readily in silent positions of codons and this is the case. There is a much wider variation in G + C content of codon silent sites than replacement sites (Fig. 2). Replacement substitutions are fixed less frequently than silent substitutions (Li et al. 1985), and there are varying degrees of constraint acting against amino acid exchanges. However, many replacements are neutral, when they do not interfere appreciably with protein function. This occurs in many proteins—for example, in globins (Perutz 1983). In the case of three members of the *ras* family of proto-oncogenes, the variation in amino acid composition is concentrated in a non-essential region.

Conclusion

We have determined the actual changes in coded amino acid content that accompany increases in the silent G + C content of human genes and have observed correlations between silent-site G + C and the levels of the majority of the 20 amino acids. Amino acids with codons high in G + C were higher in genes with high content of G + C in silent sites, and the opposite was true for amino acids with codons high in A + T. The largest changes are for amino acids that may undergo conservative replacements, e.g., Arg/Lys and Ile/Leu. Less mutable amino acids such as Phe show smaller increases. Codons with mixed G + C content in replacement sites were in some cases positively and in other cases negatively correlated with G + C content of silent sites. Positive correlation is conspicuous in the case of leucine, and we suggest that this results from the lower selective cost of mutations to leucine from isoleucine, methionine, and phenylalanine relative to mutations of leucine to amino acids with codons high in G + C.

The strength of these associations, both positive and negative, can be accounted for by the genetic code, amino acid mutability, and the neutral theory.

Acknowledgments. We thank Dr. Toshimichi Ikemura for kindly providing us with a magnetic tape containing codon usage from GenBank and Dr. Noboru Sueoka and an anonymous referee for comments. This work was supported by NIH grant R01 HG00312-03.

References

- Aïssani B et al (1991) The compositional properties of human genes. *J Mol Evol* 32:493–503
- Aota S, Ikemura T (1986) Diversity of G + C content at the third position of codons in vertebrate genes and its cause. *Nucleic Acids Res* 14:6345–6355
- Baker AR et al (1988) Cloning and expression of full-length cDNA encoding human vitamin D receptor. *Proc Natl Acad Sci USA* 85:3294–3298
- Bernardi G et al (1985) The mosaic genome of vertebrates. *Science* 228:953–958
- Bernardi G, Bernardi G (1986a) Compositional constraints and genome evolution. *J Mol Evol* 24:1–11
- Bernardi G, Bernardi G (1986b) The human genome and its evolutionary context. *Cold Spring Harbor Symp Quant Biol* 51:479–487
- Bernardi G, Bernardi G (1991) Compositional properties of nuclear genes from cold-blooded vertebrates. *J Mol Evol* 33:57–67
- Bilofsky HS, Burks C (1988) The GenBank genetic sequence data bank. *Nucleic Acids Res* 16:1861–1864
- Brown R et al (1984) Mechanism of activation of an N-ras gene in the human fibrosarcoma cell line HT1080. *EMBO J* 3:1321–1326
- Collins DW et al (1992) Numerical classification of coding sequences. *Nucleic Acids Res* 20(6):1405–1410
- Cox EC, Yanofsky C (1967) Altered base ratios in the DNA of an *Escherichia coli* mutator strain. *Proc Natl Acad Sci USA* 58:1895–1902
- Dayhoff MO (1978) Atlas of protein sequence and structure, vol 5, suppl 3, National Biomedical Research Foundation, Silver Spring, MD
- de The H et al (1987) A novel steroid thyroid hormone receptor-related gene inappropriately expressed in human hepatocellular carcinoma. *Nature* 330:667–670
- de Vos AM et al (1988) Three-dimensional structure of an oncogene protein: catalytic domain of human c-H-ras p21. *Science* 239:888–893
- D'Onofrio G et al (1991) Correlation between the compositional properties of human genes, codon usage, and amino acid composition of proteins. *J Mol Evol* 32:504–510
- Filipinski J (1990) Evolution of DNA sequence. Contribution of mutation bias and selection to the origin of chromosomal compartments. *Adv Mutagenesis Res* 2:1–54
- Fischer R et al (1988) Multiple divergent mRNAs code for a single human calmodulin. *J Biol Chem* 263:17055–17062
- Hirai H et al (1985) Activation of the c-K-ras oncogene in a human pancreas carcinoma. *Biochem Biophys Res Commun* 127:168–174
- Ikemura T, Aota S (1988) Global variation in G + C content along vertebrate genome DNA. Possible correlation with chromosome band structures. *J Mol Biol* 203:1–13
- Ikemura T et al (1990) Giant G + C% mosaic structures of the human genome found by arrangement of GenBank human DNA sequences according to genetic positions. *Genomics* 8:207–216
- Ikemura T, Wada K (1991) Evident diversity of codon usage patterns of human genes with respect to chromosome banding patterns and chromosome numbers; relation between nucleotide sequence data and cytogenetic data. *Nucleic Acids Res* 19:4333–4339
- Jukes TH, Kimura M (1984) Evolutionary constraints and the neutral theory. *J Mol Evol* 21:90–92
- Jukes TH, Bhushan V (1986) Silent nucleotide substitutions and G + C content of some mitochondrial and bacterial genes. *J Mol Evol* 24:39–44
- Karlin S et al (1990) Contrasts in codon usage of latent versus productive genes of Epstein-Barr virus: data and hypotheses. *J Virol* 64(9):4264–4273
- Li W-H et al (1985) A new method for estimating synonymous and nonsynonymous rates of nucleotide substitution considering the relative likelihood of nucleotide and codon changes. *Mol Biol Evol* 2(2):150–174

- Maki H, Sekiguchi M (1992) MutT protein specifically hydrolyses a potent mutagenic substrate for DNA synthesis. *Nature* 355:273-275
- Miyajima N et al (1988) Identification of two novel members of *erbA* superfamily by molecular cloning: the gene products of the two are highly related to each other. *Nucleic Acids Res* 16:11057-11074
- Muto A, Osawa S (1987) The guanine and cytosine content of genomic DNA and bacterial evolution. *Proc Natl Acad Sci USA* 84:166-169
- Nathans J et al (1986) Molecular genetics of human color vision: the genes encoding blue, green and red pigments. *Nature* 322:193-202
- Perutz M (1983) Species adaptation in a protein molecule. *Mol Biol Evol* 1:1-28
- Rolfe R, Meselson M (1959) The relative homogeneity of microbial DNA. *Proc Natl Acad Sci USA* 45:1039-1043
- Sekiya T et al (1984) Molecular cloning and the total nucleotide sequence of the human *c-Ha-ras-1* gene activated in a melanoma from a Japanese patient. *Proc Natl Acad Sci USA* 81:5384-5388
- Sueoka N et al (1959) Heterogeneity in deoxyribonucleic acids. II. Dependency of the density of deoxyribonucleic acids on guanine-cytosine content. *Nature* 183:1429-1431
- Sueoka N (1961) Correlation between base composition of deoxyribonucleic acid and amino acid composition of protein. *Proc Natl Acad Sci USA* 47:1141-1149
- Sueoka N (1962) On the genetic basis of variation and heterogeneity of DNA base composition. *Proc Natl Acad Sci USA* 48:582-592
- Sueoka N (1988) Directional mutation pressure and neutral molecular evolution. *Proc Natl Acad Sci* 85:2633-2657
- Sueoka N (1992) Directional mutation pressure, selective constraints and genetic equilibria. *J Mol Evol* 34:95-114
- Wada K et al (1991) Codon usage tabulated from the GenBank genetic sequence data. *Nucleic Acids Res* 19 (Suppl):1981-1986
- 34 LAC109*
34 PYHBASB
35 FUMA
35 TKFER
35 ANTLF3
35 EIF2A
35 ASP
35 LCA
35 CBF
35 IFNRG
35 DGIGLY
35 CCC5
35 ETFA
35 RP19
35 GLYCQ
35 ANTC36
36 NKG2B
36 ZFY
36 NKG2A
36 PLA2A
36 TFIIS
36 UG2BA
36 GAPA
36 C9M
36 LCAR
36 DCKATPB
36 PHSR2
36 HSF2
36 FISP
36 VIPHM6*
36 GST
36 RELH2
36 FNRB
36 ARAA5*
36 NKG2C
36 SCF
37 HELB#1
37 CS1PA
37 CREB
37 SCFA2
37 HMGCOA
37 KGF
37 HGFHL
37 PRPD
37 PHSR1
37 CREBA#2
37 BFDNA
37 FBRG#1
37 DCREB
37 FBRG#2
37 PTH2
37 TPX1A
37 TRA1
37 PRPC4B
37 HPRTA1A
38 POLYAB
38 ROSA
38 APOH
38 GAGMR
38 SNEXIN
38 PCCR
38 ILIRA
38 AFP
38 VLA2A
38 CBP
38 P70S6KB
38 CNR
38 ALBGC#2
38 ALBAF1
38 CERP
38 NUCLEO
38 COR2M
38 BRAF
38 POLB
38 G25KA
38 FABP
38 UDPGTA
38 PRPE
38 NKG2D
38 HFSP
38 PC1Q1
38 PP2AB
38 DUG
38 DBLTP
38 HNRNPA
38 PKCBA
38 MRNAEN
39 CALLA
39 ELF2
39 PLAST
39 LSZA
39 NF1AA
39 ARF2A
39 C7A
39 PRPF
39 HELA#2
39 PAM12
39 UDPGT
39 FVIIIC
39 EBUR13*
39 FASANT
39 VTNR
39 CAM1V
39 P70S6KA
39 PPKKA
39 FBRA
39 ALPHLA
39 ADH5CHI
39 ADH4C1
39 B1LYM
39 FOL5*
39 EB2CR2
39 CFTRM
40 IDE#2
40 LAMP1B
40 OATC
40 ZFX
40 DAFB
40 CLGNA
40 MRA
40 SNRAA
40 OMGPA
40 DAFA#2
40 ANTN
40 TS11
40 ASNS
40 GHR
40 IL2S3*
40 EGFRE
40 U1C
40 GCRB
40 PRPH2
40 PAP4
40 GCRA
40 QBPCA
41 RALA
41 CSAE
41 ANX3
41 NMTDC
41 POLDNAA
41 ADH5C3
41 TGFBB
41 COOTAA
41 ELI#1
41 TFRR
41 LDHBR
41 VCAM1
41 MHCC6A
41 ALDC13*
42 INTA6R
42 AMD
42 RAL
42 BCKDHA
42 LKHA4
42 ELASF
42 SIALO
42 STEAA
42 EL20#1*
42 EL20#2*
42 KUP
42 TSH2*
42 MDRI
42 PDHBET
42 GABAAA1
42 ECPG
42 SGII
42 LGTPA
42 IFNG
42 HRGA
42 STROM2
42 CYCLINA
43 KINI10#1*
43 ADH6
43 ELIA
43 LIPCR
43 HMG1
43 RIREM1
43 STROMR
43 HF10
43 I156KD
43 EL20#3*
43 CD2A
43 EDNMRN
43 IL5R
43 METPOA
43 HPBA
43 OTC10*
43 FVA
43 PUMP1
43 AACTA1
43 ELF1B
43 CAPPTA
43 EMP4I
43 CN2
43 ATPSY
44 HEXB13*
44 ANTC2
44 DMDR
44 CAIR
44 BPPTK
44 COL1A42*
44 ALBP
44 RPS12
44 CARMC

Received February 27, 1992/Revised June 3, 1992

Appendix 1. Coding sequence data

| | | |
|------------|-------------|------------|
| 24 ACADM | 30 RASK25* | 32 TR29 |
| 25 BFXIII | 30 MALENAD | 33 RBS |
| 26 CAM | 30 PMMPP1 | 33 CNRA1 |
| 26 HELB#3 | 30 MUT13* | 33 PROSA |
| 27 COL3AI | 30 KAD | 33 UCHL3A |
| 27 EVI2B3P | 31 AMY12* | 33 SNONR |
| 27 CYES1 | 31 PROS30 | 33 SNOAR |
| 27 DLDH | 31 MCP | 33 BCAT |
| 27 RNPB1A | 31 CHEB | 33 33DPTP |
| 27 RNPA2A | 31 AMYA110* | 33 TR211 |
| 27 RASRPB | 31 SEM | 33 MRL3R |
| 27 CA1XIA | 31 IL6GP | 33 PROP2AA |
| 28 FAPAPC | 31 CD46Q | 33 RAB2 |
| 28 HPLK | 31 PTYPH | 33 CNRAB |
| 28 IREBPA | 31 CHEBG4* | 34 MSCA |
| 29 ANTARNP | 31 EVI22 | 34 SNRNP |
| 29 6ORO | 31 IL76* | 34 LFA3R |
| 29 HAAB | 32 CDC2 | 34 SPROTR |
| 29 MCH | 32 GAPB | 34 TOPII |
| 29 SBLA | 32 NPM | 34 CIX |
| 29 GLYCA | 32 TPPIIA | 34 PRPC |
| 29 PMPCA | 32 MITF1 | 34 FIX |
| 29 P68A | 32 MACT | 34 LDHX |
| 29 IFNRA | 32 SRTR2A | 34 DBLPRO |
| 29 KRASM | 32 C1A2 | 34 CALCBP |

| | | | | | |
|-------------|------------|-------------|-------------|--------------|-------------|
| 44 TCRGA | 47 PLNHR | 50 PSG3A | 53 XCGD | 55 KER673* | 58 QRE |
| 44 XYPFLA | 47 GABAAS | 50 ZNF7 | 53 SAPR | 56 ACK110 | 58 CRPGA |
| 44 RPS6 | 47 TYRM | 50 PZPHEP | 53 CLMF35 | 56 APOD | 58 AHSB |
| 44 TRGC64 | 48 MAP4 | 50 PAPB | 53 CFCGRI6* | 56 CYCR | 58 BGPAA |
| 44 TGLIP | 48 HAP2A | 50 MAOAA | 53 MEA | 56 HPRG5* | 58 RIBIIR |
| 44 7AH | 48 IGB7 | 50 B2M2* | 53 ACYLHYD | 56 PSPB | 58 CGM7 |
| 44 RCYP3 | 48 CDC25HS | 50 ET | 53 PDGFARA | 56 IFNA01 | 58 MHHSPHO |
| 44 CMYBLA | 48 RPL17 | 50 P18 | 53 GPPSBA | 56 LSP1Q1 | 58 FIBAA |
| 44 RASAB | 48 PRG | 50 PSGB1A | 53 IL2RB | 56 HISAB | 58 LIC |
| 44 MDR3 | 48 IREBP | 50 SPDMAT | 53 RPS4X | 56 IFNAN | 58 PFKM23 |
| 44 UMPS | 48 RNPC2A | 50 SPP | 53 SPTA01 | 56 GTLPA | 58 HXMA7* |
| 44 LDHA7* | 48 HIS2AZ | 50 C8AS | 53 PSBGAA | 56 THRAA | 58 HUGBR1 |
| 44 TCRGAD | 48 LHHCGR | 50 RASAA | 53 CEASG5 | 56 CX43 | 58 RNP7008* |
| 44 THD | 48 TYRA | 50 ADH229* | 53 BMP3A | 56 IDSX | 58 5AR |
| 44 ZFX1 | 48 ATPFIB | 50 TUMP | 53 SLK | 56 NCAX | 58 PHIDYIN |
| 44 GC | 48 C1RS | 50 ARGCAA | 53 LPLR | 56 CD1A | 58 CRYBA6* |
| 44 HPROT | 48 PKCAMD | 50 THYMA | 53 MLCAB | 56 TRHA | 58 BGAL |
| 44 KUPMR | 48 PAI2A | 51 HAPRA | 53 NAKATP2* | 56 MYOHP | 58 UROD |
| 44 CYPHLP | 48 C1S | 51 HK2A | 53 CA2 | 56 CD1A6* | 58 IFNAII |
| 44 RPS6A | 48 ATPSYB | 51 AIXIII | 53 CNTFG | 56 CKMT | 58 INTAZ |
| 44 ELI#2 | 48 ALIPOA | 51 H1T | 53 GPBPS | 56 GDHL | 58 YAVREB |
| 44 ARG8* | 48 CTAP3 | 51 5NUASE | 53 ATCT2* | 56 TF | 59 MHCAG1 |
| 45 HGFR | 48 7B2 | 51 RPS25 | 53 PSG6A | 56 GFB#3 | 59 HLD0BR |
| 45 CYP3N | 48 HNRNA | 51 TCRT3E | 54 PSG12 | 56 GBP1 | 59 LYB2 |
| 45 SRICPA | 48 CATHL | 51 PROAF | 54 LAMBB | 57 CYB5 | 59 ARSBX |
| 45 TCRGC | 48 CD44E | 51 ARP450 | 54 CDW44A | 57 TROPSR | 59 IFNATD |
| 45 ODCA2 | 48 U2AR | 51 IL1P | 54 A2M | 57 GABAAB1 | 59 P40CYT |
| 45 CCG1 | 48 SF2P33 | 51 PPARP1 | 54 GPPSBD | 57 RPS14 | 59 MHDRA |
| 45 GLUCG2 | 48 CYPB | 51 THYMAA | 54 A1CKII | 57 AFPA4 | 59 CMYCQ |
| 45 CR1 | 48 TGFB1B | 51 PDHA12* | 54 PSGA | 57 RPZH21 | 59 BAT3A |
| 45 LAMBP | 48 TOP19* | 51 CD1C6* | 54 TGFB2A | 57 FSH3* | 59 PCC |
| 45 SOD1 | 49 PLG24* | 51 CATR | 54 RAG1 | 57 IFNF | 59 CRIPTO |
| 45 IL7AA | 49 CCBL | 51 CSPG1A | 54 CD38 | 57 CERBA | 59 CLMF40 |
| 45 ACKII | 49 TEF1 | 51 P65 | 54 P47 | 57 AMYB19* | 59 ASAM |
| 45 H33G4* | 49 CDR34 | 51 RPS4Y | 54 NLF1 | 57 PAP4A | 59 TSHX |
| 45 TCRDR | 49 PPEPB | 51 HK1A | 54 TCAYE | 57 HUGBR2 | 59 CARAA |
| 45 GAP43A | 49 TFIIDA | 51 RASAC | 54 HPA2B | 57 BMP2A | 59 ASPAT |
| 45 MCR | 49 LBP | 51 VINC | 54 UPCP12* | 57 BGALRP | 59 ORF |
| 45 HSC70 | 49 RETAA | 51 CYCAA | 54 CHROMB | 57 TK14 | 59 USFMR |
| 46 ADH1CB | 49 TCR3G6* | 52 SYTA | 54 M6PR | 57 DOCKP | 59 CGM1A |
| 46 KITCR | 49 TFIID | 52 PIP | 54 ETSR | 57 RASFAB | 59 CD284#2* |
| 46 RODSA | 49 TBG | 52 RAFR | 54 HP1G5* | 57 A20 | 59 LCTHA |
| 46 BPGM3* | 49 BRANK2 | 52 CD1C | 54 BDNFC | 57 HSP90B | 59 ANTNC |
| 46 PGK2G | 49 CYL | 52 PSBGA14* | 54 ATPAR | 57 MNSOD | 59 SSARO |
| 46 OSTRO | 49 MDMCSF | 52 CD1B6* | 54 C4A2 | 57 APA4R | 59 GLYPL |
| 46 LDL100 | 49 PRLR | 52 FGF53* | 54 IFNAH | 57 BGPAB#1 | 59 IGFIR |
| 46 HMG14 | 49 ARNTA | 52 EAP | 54 NKSFP35 | 57 ALFUC | 59 FCGRA |
| 46 MAC2 | 49 GMP140 | 52 FCRHA | 54 AGPRO | 58 TRO | 59 ALDOB9* |
| 46 BCGF | 49 CYP19 | 52 RPL31 | 54 FDX | 58 BGPI | 59 LACTA |
| 46 COL8A1 | 49 LYN | 52 MAOB | 54 ACHRM2 | 58 HER3A | 59 IFNB2R |
| 46 TCRGR | 49 REPA | 52 PHH | 54 SYB1A5* | 58 MHDRARN | 59 GLI |
| 46 RPS24A#2 | 49 C8BS | 52 RHOB | 55 FSHRE | 58 SAA | 59 PKCL |
| 46 FXI | 49 CCCR | 52 CAP | 55 ALAS1R | 58 PALC | 59 ANT2X |
| 47 TC0BI | 49 ELAM9* | 52 PLASTA | 55 CD59A | 58 PALFAP | 59 ACROS |
| 47 ASF#1 | 49 FCERI | 52 FGF5A#2 | 55 BDNF | 58 PROP2AB#3 | 59 MBPA |
| 47 TDTA | 49 A2TPI | 52 C4BAA | 55 PPROA | 58 HCPB | 59 GLUSYN |
| 47 IGHDO | 49 LAMB | 52 IL2RA | 55 IFNAIP | 58 BLAST1 | 59 HXBP1 |
| 47 DONT11* | 50 PORAC | 52 ITI2 | 55 MLC3F | 58 EMBPA | 59 STSB |
| 47 CYPAX | 50 YB1A | 52 OCT1A | 55 TROP3R | 58 HLRPR | 59 PTHL3* |
| 47 INTB6A | 50 LAMB33* | 52 TCAR | 55 SPTCS | 58 PTPAAA | 60 IGGK |
| 47 MONAP | 50 CALBR | 52 TFPB | 55 SAA1A | 58 IFNAB | 60 BFR |
| 47 NFKB34 | 50 GPPSBC | 52 SCAR | 55 F13A15* | 58 BCTHA | 60 OSF1 |
| 47 MGPA | 50 FMO1 | 52 GSTC | 55 SLIPG | 58 FCGRB | 60 KSAMAA |
| 47 EF1AR | 50 RASAD | 52 MIC2A | 55 NMOR | 58 GIF | 60 PPARP2 |
| 47 EP2AA | 50 GRP78 | 52 GLYCA4* | 55 BAT2B4* | 58 PPPB1A | 60 PLP6* |
| 47 ISG2* | 50 BADPTA | 53 IGFBPS | 55 TM2CEA | 58 PTKJAK1 | 60 MHRD5* |
| 47 LYAM9* | 50 AMPD1 | 53 CAMPR2 | 55 FERC | 58 IFNB1 | 60 PP15 |
| 47 AGALAR | 50 PNU4* | 53 RPA70KD | 55 CHYMASE | 58 AGG | 60 ETS1A |

| | | | | | |
|-------------|-------------|------------|-------------|--------------|--------------|
| 60 MRcox4* | 62 CD43#2 | 64 D1D0 | 66 QM | 68 IL4 | 70 CNPB5* |
| 60 CGM1B | 62 DHPR | 64 PINCAM | 66 MHBA123 | 68 KALX | 70 FBP |
| 60 ERCC3A | 62 APOC2G | 64 ANDREC | 66 KER654* | 68 CYP178* | 70 MYC3L |
| 60 CEAF | 62 OAS08* | 64 CSK2B | 66 MH3C2 | 68 SOMI | 70 C4AA2* |
| 60 CA1V | 62 HDC | 64 MYCTR | 66 PKM2L | 68 ATCT4 | 70 ENKB4* |
| 60 1433 | 62 CDW40 | 64 IRF1 | 66 MYOD1R | 68 REGB | 70 HOXB |
| 60 GFIAB4* | 62 ETMAGA | 64 UKPM | 66 TAU1 | 69 ANTCd | 70 CBG |
| 60 FCREC | 62 GAST | 64 ET3 | 67 CD53 | 69 TKR | 70 TGFB3A |
| 60 RNPAB | 62 PKCB2A | 64 NGFBA2 | 67 MYF5 | 69 CYPBA | 70 FOS |
| 60 BGPAB#2 | 63 ERG2 | 64 MLC2 | 67 BMP2B | 69 GIPX6* | 70 COX4AA |
| 60 BHA14* | 63 CGPRA | 65 CTHG | 67 RCC1B | 69 MLC3NM | 70 TPO15* |
| 60 GP34M | 63 FERG2* | 65 MGDMT | 67 UBIQAA | 69 FNRA | 70 G3PDC |
| 60 TS1 | 63 PPARP0 | 65 ALASR | 67 GPIIIA | 69 TMPKMR | 70 A1GLY2 |
| 60 FCRII | 63 OAS07* | 65 EGFRN | 67 MYHC | 69 NM23H2S | 71 TPOB |
| 61 MHDQADR | 63 CRFBP | 65 NKG5 | 67 ZP3 | 69 RNP70K | 71 LOX15A |
| 61 THRR | 63 FMLP | 65 RHPAA | 67 HPS12 | 69 HLADRBA | 71 GFRIL |
| 61 VPF | 63 MBPC | 65 ALAS2R | 67 A1ATB | 69 CRcMUT | 71 CYP45C |
| 61 MBPZ | 63 PRPS2 | 65 ENOA | 67 TAUa | 69 UBA52P | 71 SYB2A5* |
| 61 IGGFCRA | 63 PBGDR2 | 65 GF1A | 67 TRL | 69 IBSUB | 71 PLAT |
| 61 GRFCIG | 63 TUBAK | 65 ETS2A | 67 CAM3X1* | 69 INTB5A | 71 CRYGX2* |
| 61 HLADQA | 63 GABAR | 65 RETSA#2 | 67 ASM | 69 TRKR | 71 TGFA |
| 61 MBP17K | 63 NT3A | 65 HCF2 | 67 A1ACM | 69 CSF1M3* | 71 GLCB |
| 61 PGP95 | 63 GYPCAA | 65 RPL32 | 67 C1INHb | 69 GLUTRA | 71 INSR |
| 61 HBP | 63 KBLOOD | 65 RIBIR | 67 CRPR | 69 TBBM40 | 71 PDGA7* |
| 61 G6PA | 63 HODB3* | 65 GHG | 67 MYCL2A | 69 CD19W07* | 71 FGF2H |
| 61 DBI | 63 RPS17 | 65 GLI3A | 67 TCII | 69 IGLV | 71 P45SC9* |
| 61 FLAP5* | 63 HELAGT#2 | 65 GCB | 67 ANT1 | 69 PSAG | 71 HIS4 |
| 61 HSDI | 63 ILRA | 65 CTSE | 67 AUTAN64 | 69 HMGIA | 71 PEC12L |
| 61 GALT | 63 PBGDR | 65 P42SA | 67 IGFIB | 69 PTCAA | 71 FGF3H |
| 61 IIP | 63 ANPCR | 65 PTHL | 67 FABPLA | 69 INTB7A | 71 ANTLA |
| 61 UNG | 63 OGCB | 65 EGR2A | 67 3OCTR | 69 MHDRBA | 71 CRYGA2* |
| 61 LCT17#2* | 63 PROZII | 65 GALTA | 67 ACTSG7* | 69 ADRBR | 71 ANK |
| 61 PEMP | 63 MBPB | 65 GFIAB5* | 67 I6REC | 69 B61 | 71 AGPIA |
| 61 IL1RAA | 63 3B5H5E | 65 AP2AA | 67 MYLCA#1 | 69 FBPB | 71 ERP |
| 61 ATP | 63 HBEGF | 65 NID | 67 A1GP06* | 69 TCBYZ | 71 TIMPR |
| 61 THYP | 63 THROMR | 65 TPI | 67 CALCR4* | 69 SCYLP | 71 PROPERD |
| 61 IL1C | 63 BN51 | 65 PROT2 | 67 HEMOB | 69 BETGLA | 71 CP45IV |
| 61 CSFM | 63 FAB | 65 SECPA | 67 PAR | 69 CTSB | 71 VIM |
| 61 IMP | 63 EGFRS | 65 CSDF1 | 67 FLA1A | 69 TCBXA | 71 INsRA |
| 61 CAIII7* | 63 FIGRE | 65 MAX#1 | 67 PPTRH03* | 69 ROSSAA | 71 ACHRB |
| 61 KRT10A | 63 HEM1 | 65 ANFA | 67 ALD | 69 A1MICR | 71 MLC1SA |
| 61 CPT | 63 FIGRD | 65 PTHL4* | 67 MHSXA | 69 RNAGLA | 71 TNFR |
| 61 TCAXB | 64 RISDAD | 66 PRL7* | 67 RFPA | 69 IL10 | 71 TCBYY |
| 61 FGF4H | 64 MXA | 66 CYPIIE | 67 CBPE | 69 HLASBA#1 | 71 KER2A |
| 61 MUCAB | 64 TGLH | 66 A1AR2* | 68 HOX329 | 69 UBA52C | 71 MUC18A |
| 61 IRP | 64 GHVA#2 | 66 PRPOA | 68 LAMP1A | 69 RNP7011#2 | 71 AP2 |
| 61 PLPDM | 64 FOSB#1 | 66 PRP | 68 REN10* | 69 IGLBV | 71 CRE |
| 61 AT3X6* | 64 MAS | 66 ACTCA4* | 68 CYP2BA | 69 SRAA | 71 LT |
| 61 COXCA | 64 FOLLI2* | 66 PC2A | 68 ALAD | 70 MYOL1 | 71 NAGA |
| 61 ALDCG | 64 HBGF3* | 66 DNFA | 68 INHBA | 70 CALRTR | 71 TNFRII |
| 61 SMCK | 64 SHBGA | 66 GFIBP | 68 RENT3* | 70 KTRAN | 72 INTLEU8 |
| 62 NOXF | 64 ERG11 | 66 HBB#2 | 68 IGFBP1A | 70 P5ONFKB | 72 TPAR |
| 62 MXB | 64 ASPX | 66 HBB#3 | 68 CALCR5* | 70 ICOA | 72 INCP3* |
| 62 PGRR | 64 PROZI | 66 FAH | 68 GFII | 70 SYN | 72 TCF1A |
| 62 HRSR | 64 HELAGT#1 | 66 ENOG | 68 HCR | 70 RNP7011#1 | 72 CAD |
| 62 P53C | 64 ARX | 66 XRCC1 | 68 UBI13 | 70 MHDRBU | 72 NFH4* |
| 62 T519#2 | 64 MHCSE | 66 CA6 | 68 BPIAA | 70 LAPA | 72 ALDHIR |
| 62 MLN | 64 IGLAM2 | 66 FAPS | 68 TRBP | 70 MAC1A | 72 VACB |
| 62 GRP5E | 64 HBB#4 | 66 PTHL2* | 68 A1ATZ | 70 ERMCF | 72 PDGFA6* |
| 62 AR | 64 CRYABA | 66 MYLA1 | 68 ACHRA | 70 VTNSP | 72 HA44G |
| 62 TNsCN | 64 GSTMUA | 66 HEXKIN | 68 IL4R | 70 PHK | 72 SYN1E13#1 |
| 62 EGFAA | 64 MAX#2 | 66 FOSB#2 | 68 LMGP | 70 IGFBP4 | 72 ASGPR2 |
| 62 HBB222* | 64 SGLT1 | 66 APOAII | 68 ET2A | 70 SPRO | 72 PPE |
| 62 TM30R | 64 RPOLAA | 66 DCDK | 68 PP11A | 70 ABL | 72 HISAC |
| 62 AFH | 64 GCB1 | 66 NFM | 68 PLC | 70 INHA | 72 CLG4Q13* |
| 62 WRSAA | 64 NK4 | 66 GPIBAA | 68 HIS3PRM | 70 VIL2 | 72 GOAQ10#2* |
| 62 ALR | 64 P42LA | 66 CYP2BB | 68 FGL2 | 70 KERC15 | 72 VWFR1 |
| 62 PLAX | 64 ARB | 66 GBR#2 | 68 PCAR | 70 ANTCd9 | 72 FGFAA |

| | | | | | |
|-------------|--------------|-------------|------------|-------------|-------------|
| 72 POVRA | 74 ERBT1 | 75 GH#1 | 77 CSP40 | 79 CMP8* | 81 FESFPS |
| 72 ELA308* | 74 BMYB | 75 MRP14 | 77 IGHAE2* | 79 ENDOA2#1 | 81 GCSFR |
| 72 HBGFA | 74 GCSFRD | 75 IGHAF | 77 C45AII | 79 GTUB#1 | 81 TIMP2 |
| 72 ALRM | 74 ASA | 76 PROTP | 77 MHB27D | 79 CERA | 81 P58GTA |
| 72 FCER | 74 STROL3 | 76 BNPA | 77 ELS2 | 79 CP21OHC | 81 HNF1 |
| 72 TCF1B | 74 OCT2A#1 | 76 TPMYOC | 78 RBP | 79 LYL1B | 81 P971 |
| 72 INT07* | 74 TGFBC | 76 GATA | 78 C8G | 79 OXYGR | 81 KER7E9* |
| 72 EPP10* | 74 GMCSFRB | 76 GCSFR4 | 78 CYPIIF | 79 ALPI1 | 81 ATPA23* |
| 72 SAPD1 | 74 GNAS6#3* | 76 COLIP | 78 GFIBPA | 79 CHRМ | 81 PNMT |
| 72 CMOS | 74 GSA2R | 76 DNASEI | 78 CYP2DG | 79 FFI2A | 81 PCD |
| 72 JGEBFR | 74 FMSCPO | 76 NFLG | 78 LEC14K | 79 ACTAR | 81 UBILP |
| 72 KERE9* | 74 CHRAA | 76 THRA1A | 78 PGP | 79 ALPP | 81 BMYH7 |
| 72 ICAMA1M | 74 ELA3A | 76 LCKAA | 78 CPIIA3A | 79 MHDNDRW | 81 TGASE |
| 72 ADAM2 | 74 ATP1A2 | 76 MHCSA#2 | 78 MHCCB2 | 79 CRYGQ2* | 81 CYPDB1 |
| 72 ASL | 74 RALBPC | 76 HLDQWB | 78 COF | 79 HIS10G | 81 CYP2D6 |
| 73 VILLR | 74 C3 | 76 LOX5A | 78 PSPBA | 80 P3A | 81 TUBBM |
| 73 GNBPB3 | 74 TROPIA | 76 GAPJR | 78 LOX5 | 80 ENDOA2#5 | 81 TROPI |
| 73 SA#1 | 74 ARAF1R | 76 PGSR | 78 RDS | 80 FFI2B | 81 GAA |
| 73 CYP27 | 74 RAP2 | 76 CYPB3* | 78 CENPB | 80 AHCY2 | 81 IL11 |
| 73 GLBA | 74 CNTFR | 76 TNS | 78 ASLA | 80 FVII#2 | 81 ASGPR1 |
| 73 PDEAA | 74 YUBG1 | 76 PLPSPC#1 | 78 LAMAR | 80 GFAP | 81 ARF1BA |
| 73 THB | 74 TCPTK | 76 GLYSA | 78 SCL | 80 CRF | 81 ACHRM4 |
| 73 CD8B | 74 XEH | 76 G19P1A | 78 CYIIA4A | 80 KERE9 | 81 GGT |
| 73 PF4A | 74 GSA1R | 76 VDR | 78 IL2AB | 80 LIPBSA | 81 INSPR |
| 73 HAAG | 74 PLA2A2* | 76 GHV | 78 TNFAA | 80 AMIPEP | 81 NADPHO |
| 73 MYCRT | 74 NAGB | 76 SERDHY | 78 SRF | 80 TNCS | 82 H2B2H2#1 |
| 73 C1R | 74 CRBP2* | 76 HPSNA | 78 SISPDG | 80 ZNFBPAA | 82 BCR |
| 73 ERCC1 | 74 LDLR18* | 76 HISH4 | 78 PEPD | 80 RPS11 | 82 CANPR |
| 73 CD14R | 74 THBP | 76 H2B2H2#2 | 78 4COLA | 80 MLC2A | 82 ADRA2RA |
| 73 TCF1C | 74 PAIR | 76 LCTHB | 78 MAL | 80 IFP | 82 INTBE4 |
| 73 TFLS | 74 CS3 | 76 PLPSPC#2 | 78 NKB | 80 THY1A | 82 DBHRB |
| 73 ALDH13* | 74 GRFP5#2* | 76 TRK2H | 78 PP14A | 80 A2PIG6* | 82 GATA3M |
| 73 IRGT | 74 PLA | 76 ELAP2B | 78 MHCW1C | 80 MYP | 82 LHB |
| 73 H4 | 74 ALDA | 76 PANMU | 78 TK | 80 GLUTRN | 82 DES |
| 73 PGEX11* | 75 CNPG6* | 77 PYGM20* | 78 CNPDEA | 80 RNAPII | 82 GATA3R |
| 73 PROF | 75 GRFP5#1* | 77 FGR | 78 RYR | 80 CFXII3* | 82 PEROXP |
| 73 GADD45 | 75 ABPA | 77 NORTR | 78 PKBR | 80 KER19 | 82 TNC2* |
| 73 ACP5 | 75 HCKA | 77 KERIA8* | 78 MGI3* | 80 ALPHA | 82 ECK |
| 73 TIR | 75 FIBUA | 77 SPARC | 78 FKBP13 | 80 ATPGG | 82 KER8 |
| 73 HOXA | 75 ARF5A | 77 GH#2 | 78 EDG | 80 GLAA | 82 PRC7* |
| 73 AK1 | 75 PDGFRA | 77 SPTB#2 | 78 ACHE | 80 GIP2A8* | 82 BCL2B |
| 73 CGBS08 | 75 ICAM4* | 77 APOJ | 79 CYP345 | 80 SCL7* | 82 18D |
| 73 TBP1 | 75 IGLR141 | 77 PEPC9* | 79 IRBPG | 80 PDECGF | 82 BCL2A |
| 73 KERUHS | 75 KER18 | 77 PEPCA9* | 79 IRBPG4* | 80 AICEB | 82 GAA01* |
| 73 HISAA | 75 RARG | 77 RAPIGAP | 79 MH6 | 80 HKATPC | 82 PRCM |
| 73 HER2A | 75 MHDPL | 77 TNTSA | 79 MHA3 | 80 CKMA | 82 CTF1 |
| 73 IL2RBC | 75 CETP7* | 77 GROB5 | 79 MHCGE2* | 80 TACEA | 82 MZF1 |
| 73 RETPO | 75 MYCM | 77 MH | 79 LCAT | 80 GGTX | 82 TCXAAA |
| 73 PKLR | 75 GLUT5 | 77 SPARC10* | 79 PRF1A | 80 MDP4 | 83 LAP |
| 73 AFL2* | 75 ADXR#2 | 77 ETR103 | 79 CYP4A7* | 80 MHGM | 83 ALIFA |
| 73 CYC1 | 75 927A | 77 OCS3* | 79 SCAD | 80 ALPPB | 83 INT2 |
| 73 HLAATE | 75 ADREDB#1* | 77 ASFA | 79 GCSFG | 80 LAMC | 83 APRTA |
| 73 PKA | 75 CSFGM | 77 KER18* | 79 PRPHOS1 | 80 BCL3AA | 83 INTB4R |
| 73 SB2BR | 75 PLAPL | 77 SPTB#1 | 79 P2A | 80 IFNIN3 | 83 DKERB |
| 73 LAPI1* | 75 ADREDB#2* | 77 KDEL | 79 GP | 80 ISK | 83 RAR |
| 73 ANGG5* | 75 HBL0D | 77 NEKAR5* | 79 FGFR4 | 80 MHCACA#1 | 83 RASR2* |
| 73 LIGAA | 75 HLA1EA | 77 CLI | 79 INV2 | 80 E12A | 83 SHIIC |
| 74 APOC3B | 75 BCTHB | 77 ANTP | 79 TRNB | 80 CD37 | 83 GSTPI |
| 74 MHDRODQ | 75 GCSFR3 | 77 ERYTH | 79 PEP9* | 80 CD7 | 83 CYS3A3* |
| 74 POVRB | 75 U1RNPA | 77 APHOL | 79 NMYCA#1 | 80 TFAA | 83 CGB |
| 74 LSP1A | 75 PPR | 77 MHTRP | 79 CEAPX | 80 ACHRG8* | 83 RETREC |
| 74 GPIIBA | 75 GROG5 | 77 PIM1 | 79 TACEB | 80 C5AAR | 83 UDPG |
| 74 PLCA | 75 NLK | 77 BMP1A | 79 FX8* | 81 GCSF | 83 18U |
| 74 GNAS6#2* | 75 UMOD | 77 HLADZA | 79 ARP1 | 81 FVII#1 | 83 UDPCNA#2 |
| 74 MHDRO3* | 75 PDGFA7* | 77 LDLRRL | 79 HISH2B | 81 EDHB17 | 83 GKNASE |
| 74 CINHP | 75 RHOC9 | 77 ARF6A | 79 HLA11E | 81 MGPHB | 83 A2MGRAP |
| 74 FIBUB | 75 CCK3* | 77 TROPA | 79 MHCA1A | 81 DRD2A | 83 NGFR |
| 74 HOM4 | 75 FIBUC | 77 IMPH | 79 TRKPOA | 81 ATPK14* | 83 JUNCAA |

| | | | | | |
|-------------|-------------|------------|------------|-----------|------------|
| 83 CYTOK | 84 IGFBP5A | 85 PGPIX#1 | 86 IHRP | 89 ISGA2* | 92 JUNDR |
| 83 PFKLA | 84 GAAA | 85 BGN3* | 87 SPERSYN | 89 SKIR | 92 OTNPI |
| 83 TGFB | 84 PLAKO | 85 MAG | 87 MYOD | 90 HBA4#2 | 92 OTCB |
| 83 FKMKA | 84 THX | 85 GXA | 87 PCHSUCA | 90 APOE4 | 92 D4DOP |
| 83 TRY | 84 LORAA | 85 HPGI | 87 INHBB2* | 90 GPIB | 92 ADRA2R |
| 83 HD5DR | 84 INT1G | 85 MIS | 87 FGFR3 | 91 ADRA2C | 93 EAR2 |
| 84 ETR101 | 84 APOAIB | 85 EAR3 | 88 GLYPIC | 91 GA733A | 93 HST |
| 84 PROSYN#1 | 84 EF2AB | 85 HIP3K | 88 IGFBP5 | 91 HSP70D | 93 VPNP |
| 84 PCA | 84 ARB3A | 85 HPBS | 88 BCL1 | 91 MHHSP | 94 HBA1 |
| 84 COMTA | 84 THM | 86 ELFT | 88 NCBLCA | 91 CPGISL | 94 GTPBRPA |
| 84 MYELA | 84 G6PDG13* | 86 HTH1R | 88 POMC9* | 91 G0S2A | 94 IDB |
| 84 OPS | 84 RACPC | 86 ACTGA | 88 SAACT | 92 CKB | 95 ADRB1 |
| 84 SRC11* | 84 R2IMP | 86 MYCPOB | 88 TBB5 | 92 GLTH1 | 96 SODEC |
| 84 HOX14 | 84 PTBMR | 86 LEUELA | 89 GDF1#2 | 92 RHOB6 | 97 MIFA |
| 84 ELFTL | 84 THR | 86 CNP | 89 HSP27 | | |
| 84 SPYRAT | 84 RACB | 86 CTRP | | | |
| 84 PKCGA | 85 PGAMM2* | 86 TGFB | | | |
| 84 APOA4B | 85 TAPA1 | 86 GDF1#1 | | | |
| 84 G6PDA | 85 JUNA | 86 BIGFII | | | |
| 84 TRPY1B | 85 GA16 | 86 CATD5* | | | |

^a 1,607 coding sequences are identified by GenBank LOCUS name and are listed in order of increasing silent site G + C (“%G + C”). The first three characters of each LOCUS name (HUM) have been omitted.