Codon catalog usage and the genome hypothesis

R.Grantham, C.Gautier, M.Gouy, R.Mercier and A.Pavé

Equipe Evolution Moléculaire, Laboratoire de Biométrie, Université Lyon I, 69622 Villeurbanne Cedex, France

ABSTRACT

Frequencies for each of the 61 amino acid codons have been determined in every published mRNA sequence of 50 or more codons. The frequencies are shown for each kind of genome and for each individual gene. A surprising consistency of choices exists among genes of the same or similar genomes. Thus each genome, or kind of genome, appears to possess a "system" for choosing between codons. Frameshift genes, however, have widely different choice strategies from normal genes. Our work indicates that the main factors distinguishing between mRNA sequences relate to choices among degenerate bases. These systematic third base choices can therefore be used to establish a new kind of genetic distance, which reflects differences in coding strategy. The choice patterns we find seem compatible with the idea that the genome and not the individual gene is the unit of selection. Each gene in a genome tends to conform to its species' usage of the codon catalog; this is our genome hypothesis.

INTRODUCTION

The genetic code provides options among codons for all 20 amino acids of protein except methionine and tryptophan, which have single codons (see Figure 1 of reference 1). Choices among synonymous codons do not affect the nature of the protein produced but they may relate to expression of a gene. Indeed, mRNA expressivity (rate of synthesis of a protein and the total amount of it made) may be under the control of degenerate base use in the mRNA. This is still largely an untested hypothesis. It is, however, an important one since choices among third bases have often been considered "neutral", that is, of no influence on fitness of the phenotype. Messenger RNA expressivity may well

turn out to be an adaptive phenomenon.

In order to better understand codon usage, we have studied all published mRNA sequences of more than about 50 codons. We here report three analyses on these sequences. Table 1 first shows number, name, symbol, length in codons and reference for each mRNA. Then Figures 1, 2 and 3 reveal the analyses. In Figure 1 mRNA are combined according to genome type (RNA or DNA virus, animal, etc.). Frequencies per thousand for each of the 61 amino acid codons appear in Figures 1 and 2 (initiator and terminator codons are excluded). Absolute frequencies are had by multiplying by number of codons (Table 1). Figure 2 presents frequencies in each individual gene (unless it is the only example of a genome type, in which case it is in Figure 1). This allows comparison of codon use not only among genes in the same or different genomes, but also between individual genes and the genome types of Figure 1. Finally, messengers with similar patterns of use of the code's degeneracy are seen as neighboring points in Figure 3. The spacings in Figure 3 are projections into two dimensions of distances calculated by correspondence analysis (2, 3) on codon third base frequencies. Other work has shown that the codon frequencies themselves give correspondence analysis distances that depend mostly on systematic choices between degenerate bases (4, 5).

DISCUSSION

It is not evident that distances between mRNA based on codon frequencies would depend mainly on degenerate base use. The other two codon positions could dominate in establishing these distances if proteins were relatively more different among themselves than are their mRNA. Elsewhere (4, 5) we have found that distances between proteins, determined by correspondence analysis on their amino acid frequencies, do not agree with distances between their mRNA, calculated from frequencies for the 61 codons in each messenger. However, frequencies of duet coded amino acids can particularly affect distances between mRNA (5). In order to eliminate all possible influence of protein composition, therefore, distances in Figure 3 were determined solely from quartet codons (1). Each of the eight quartet sets of codons has a complete choice of bases in codon position III. Figure 3 results from correspondence analysis on these quartet third base frequencies in each mRNA.

The structure in Figure 3 can be considered at several levels. First, animals lie to the right and viruses to the left. Next, there is grouping by genome or genome type: All papova virus (SV40 and BKV) genes lie together

(no. 51-56); immunoglobins (IG) fall in a nearby zone (no. 73-78). All mRNA of mammals except immunoglobins (MAM-IG) are grouped to the upper right (no. 79-90). Curiously, however, frameshift genes E of ♦x174 and G4 (no.13 and 23) are near this mammal area. Other genes of these phage are mainly to the upper left of the figure. Finally, most IG have different coding strategies from other mammalian genes (compare, for example, MUSBLP and MUSK2, no. 73 and 80).

What do these groupings and distances mean? We can partly answer this question although the exact biological significance is unknown. The horizontal axis of Figure 3 roughly corresponds to GC content of the third bases (5). The rightmost mRNA often use C and G as third bases. Thus, with MUSBLP 90.5% of all quartet codons have C or G in position III, while the value is only 2.8% in the solitary yeast mitochondrial mRNA (no. 64) and 37.0% in MUSK2. This also explains positioning of the above E genes near mammals: their mRNA show more GC in quartet position III (68.9 and 64.7%, respectively) than do other phage mRNA. Likewise, vertical contrast exists between use of A and U. For example, quartet third bases in MUSK2 contain 42.6% A and 20.4% U while those in FDV5 (no. 33) have 4.1% A and 53.1% U. Consequently, in base choices for quartet codons papova virus mRNA resemble IG mRNA but differ greatly from other mammalian genes, for both GC content and use of A (1). Lastly, correspondence analysis groupings as in Figure 3 are very stable, whether the starting data are frequencies for all 61 codons, the 32 quartet codons, or for third bases of all codons or, as here, third bases of quartet codons. Some rationalization of this stable structure is found in the distribution of the isoacceptor tRNA for the codons (5-7). Coordination of codon usage with tRNA gene amplification and expression may eventually aid in understanding the genome hypothesis.

Our general goal in this work is to identify and understand the biological information in nucleic acid sequences. Our hypothesis, however, is so far mostly only descriptive. We observe that mRNA of the same genome are clustered by correspondence analysis on codon frequencies and that their proteins are not similarly grouped by analysis on the amino acid frequencies (4). Indeed, correspondence analysis on the proteins does not suggest classical systematics as does that on mRNA. In this very limited sense, messengers better reflect evolution than proteins do (4, 5). We do not yet know why.

Table 1     mRNA Portfolio

| No. | Species and gene | Symbol | No.codons | Ref. |
|---|---|---|---|---|
| 1 | Phage MS2   gene A | MS2A | 392 | 8 |
| 2 | Phage MS2   coat | MS2C | 129 | 9 |
| 3 | Phage MS2   replicase | MS2R | 544 | 10 |
| 4 | Phage Q beta   coat | QBVC | 79 | 11 |
| 5 | Tobacco mosaic virus   gene A | TMVA | 55 | 12 |
| 6 | Tobacco mosaic virus   coat | TMVC | 158 | 13 |
| 7 | Tobacco mosaic virus   30K protein | TMV30K | 105 | 13 |
| 8 | Turnip yellow mosaic virus   coat | TYMC | 188 | 14 |
| 9 | Phage ΦX 174   gene A | FIXA | 511 | 15,16 |
| 10 | Phage ΦX 174   gene B | FIXB | 119 | 15,16 |
| 11 | Phage ΦX 174   gene C | FIXC | 85 | 15,16 |
| 12 | Phage ΦX 174   gene D | FIXD | 151 | 15,16 |
| 13 | Phage ΦX 174   gene E | FIXE | 90 | 15,16 |
| 14 | Phage ΦX 174   gene F | FIXF | 422 | 15,16 |
| 15 | Phage ΦX 174   gene G | FIXG | 174 | 15,16 |
| 16 | Phage ΦX 174   gene H | FIXH | 326 | 15,16 |
| 17 | Phage ΦX 174   gene J | FIXJ | 37 | 15,16 |
| 18 | Phage ΦX 174   gene K | FIXK | 55 | 17 |
| 19 | Phage G4   gene A | VG4A | 553 | 18 |
| 20 | Phage G4   gene B | VG4B | 119 | 18 |
| 21 | Phage G4   gene C | VG4C | 83 | 18 |
| 22 | Phage G4   gene D | VG4D | 151 | 18 |
| 23 | Phage G4   gene E | VG4E | 95 | 18 |
| 24 | Phage G4   gene F | VG4F | 426 | 18 |
| 25 | Phage G4   gene G | VG4G | 176 | 18 |
| 26 | Phage G4   gene H | VG4H | 336 | 18 |
| 27 | Phage G4   gene J | VG4J | 24 | 18 |
| 28 | Phage G4   gene K | VG4K | 55 | 17 |
| 29 | Phage fd   gene 1 | FDV1 | 347 | 19 |
| 30 | Phage fd   gene 2 | FDV2 | 409 | 19 |
| 31 | Phage fd   gene 3 | FDV3 | 423 | 19 |
| 32 | Phage fd   gene 4 | FDV4 | 425 | 19 |
| 33 | Phage fd   gene 5 | FDV5 | 86 | 19 |
| 34 | Phage fd   gene 6 | FDV6 | 111 | 19 |
| 35 | Phage fd   gene 7 | FDV7 | 32 | 19 |
| 36 | Phage fd   gene 8 | FDV8 | 72 | 19 |
| 37 | Phage fd   gene 10 | FDV10 | 110 | 19 |
| 38 | Phage M13   gene 1 | M131 | 215 | 20 |
| 39 | Phage M13   gene 3 | M133 | 423 | 20 |
| 40 | Phage M13   gene 4 | M134 | 59 | 20 |
| 41 | Phage M13   gene 6 | M136 | 111 | 20 |
| 42 | Phage M13   gene 7 | M137 | 32 | 21 |
| 43 | Phage M13   gene 9 | M139 | 31 | 21 |
| 44 | Phage T7   gene 1 | T7VG1 | 56 | 22 |
| 45 | Phage λ   gene CI | LAMCI | 236 | 23 |
| 46 | Phage λ   gene CII | LAMCII | 96 | 24 |
| 47 | Phage λ   gene cro | LAMCRO | 65 | 24 |
| 48 | Phage λ   gene O | LAMO | 298 | 25 |
| 49 | Phage 434   gene CII | 434CII | 96 | 26 |
| 50 | Phage 434   gene cro | 434CRO | 70 | 26 |
| 51 | Simian virus 40   gene T | S40GT | 626 | 27,28 |

| No. | Species and gene | Symbol | No. codons | Ref. |
|-----|------------------|--------|------------|------|
| 52 | Simian virus 40   gene t | S40PT | 173 | 27,28 |
| 53 | Simian virus 40   gene VP1 | S40VP1 | 361 | 27,28 |
| 54 | Simian virus 40   gene VP2 | S40VP2 | 351 | 27,28 |
| 55 | Simian virus 40   gene VP3 | S40VP3 | 233 | 27,28 |
| 56 | Virus BK   gene t | BKVPT | 171 | 29 |
| 57 | Hepatitis B virus   surface antigen | HBVSA | 225 | 30,31 |
| 58 | Escherichia coli   lac 1 | ECOLAC | 359 | 32 |
| 59 | Escherichia coli   ribosomal protein L11 | ECOL11 | 141 | 33 |
| 60 | Escherichia coli   ribosomal protein L1 | ECOL1 | 233 | 33 |
| 61 | Escherichia coli   ribosomal protein L10 | ECOL10 | 164 | 33 |
| 62 | Escherichia coli   ribosomal protein L7/L12 | ECO712 | 120 | 33 |
| 63 | Salmonella paratyphi   ampr gene | SAPAMP | 285 | 34 |
| 64 | Saccharomyces   subunit 9 mitochondrial ATPase | SACMT9 | 75 | 35 |
| 65 | Saccharomyces   iso-1 cytochrome C | SACCC1 | 108 | 36 |
| 66 | Psammechinus miliaris   H1 histone | PSMH1 | 85 | 37 |
| 67 | Psammechinus miliaris   H2A histone | PSMH2A | 123 | 37 |
| 68 | Psammechinus miliaris   H2B histone | PSMH2B | 102 | 37 |
| 69 | Psammechinus miliaris   H3 histone | PSMH3 | 135 | 37 |
| 70 | Strongylocentrotus purpuratus   H2A histone | SPUH2A | 123 | 38 |
| 71 | Strongylocentrotus purpuratus   H3 histone | SPUH3 | 102 | 38 |
| 72 | Chicken   ovalbumin | GALOVA | 385 | 39 |
| 73 | Mouse   K2 immunoglobulin | MUSK2 | 117 | 40 |
| 74 | Mouse   immunoglobulin light chain MOPC 21 | MUSLC | 107 | 41 |
| 75 | Mouse   immunoglobulin λ1-type light chain | MUSL1 | 173 | 42 |
| 76 | Mouse   VλII immunoglobulin | MUSVL2 | 165 | 43 |
| 77 | Mouse   MOPC-41 light chain immunoglobulin | MUSVL41 | 129 | 44 |
| 78 | Mouse   immunoglobulin γ-1 constant heavy chain | MUSPH21 | 153 | 45 |
| 79 | Mouse   hemoglobin beta | MUSBGL | 146 | 46 |
| 80 | Mouse   beta lipotropin | MUSBLP | 48 | 47 |
| 81 | Rat   growth hormone | RATGH | 215 | 48 |
| 82 | Rat   prolactin hormone | RATPLH | 132 | 49 |
| 83 | Rat   preproinsulin | RATPPI | 107 | 50 |
| 84 | Rabbit   hemoglobin alpha | RABAGL | 141 | 51 |
| 85 | Rabbit   hemoglobin · beta | RABBGL | 146 | 52 |
| 86 | Bovine   corticotropin beta lipotropin | BOVCBL | 264 | 53 |
| 87 | Human   alpha chorionic gonadotropin | HUMACG | 115 | 54 |
| 88 | Human   hemoglobin beta | HUMBGL | 132 | 55 |
| 89 | Human   pregrowth hormone | HUMPGH | 216 | 56 |
| 90 | Human   chorionic sommatomammotropin | HUMCSL | 168 | 57 |

.

| Codon | SS RNA VIRUS | | | SS DNA VIRUS | DS DNA VIRUS | | | |
|---|---|---|---|---|---|---|---|---|
| | ALL mo 1-8 | PHAGE 1-4 | PLANT 5-8 | VIRUS 9-43 | ALL 44-57 | PHAGE 45-50 | PAPOVA 51-56 | HBVSA 57 |
| Arg CGA | 8 | 6 | 11 | 5 | 8 | 16 | 0 | 0 |
| CGC | 10 | 15 | 6 | 15 | 14 | 27 | 0 | 4 |
| CGG | 5 | 6 | 5 | 3 | 1 | 0 | 1 | 4 |
| CGU | 20 | 25 | 14 | 21 | 9 | 15 | 2 | 9 |
| AGA | 12 | 4 | 20 | 5 | 15 | 6 | 28 | 4 |
| AGG | 6 | 3 | 9 | 2 | 15 | 9 | 24 | 0 |
| Leu CUA | 11 | 11 | 10 | 4 | 10 | 5 | 12 | 31 |
| CUC | 17 | 20 | 15 | 15 | 9 | 12 | 3 | 31 |
| CUG | 8 | 11 | 6 | 18 | 17 | 19 | 13 | 31 |
| CUU | 13 | 14 | 12 | 30 | 23 | 25 | 22 | 18 |
| UUA | 16 | 17 | 16 | 24 | 20 | 13 | 28 | 18 |
| UUG | 12 | 6 | 18 | 17 | 16 | 12 | 21 | 18 |
| Ser UCA | 19 | 15 | 24 | 18 | 11 | 10 | 9 | 27 |
| UCC | 12 | 13 | 10 | 13 | 7 | 4 | 7 | 22 |
| UCG | 19 | 21 | 17 | 8 | 4 | 5 | 0 | 13 |
| UCU | 20 | 22 | 19 | 36 | 16 | 14 | 18 | 18 |
| AGC | 13 | 14 | 13 | 4 | 10 | 13 | 7 | 4 |
| AGU | 11 | 4 | 17 | 8 | 14 | 8 | 20 | 22 |
| Thr ACA | 17 | 7 | 27 | 11 | 16 | 17 | 14 | 27 |
| ACC | 26 | 23 | 30 | 12 | 20 | 22 | 18 | 22 |
| ACG | 9 | 9 | 9 | 8 | 4 | 6 | 0 | 13 |
| ACU | 32 | 32 | 31 | 25 | 18 | 13 | 25 | 13 |
| Pro CCA | 10 | 8 | 11 | 7 | 13 | 10 | 13 | 40 |
| CCC | 16 | 10 | 22 | 5 | 8 | 3 | 10 | 31 |
| CCG | 11 | 14 | 8 | 10 | 8 | 14 | 0 | 9 |
| CCU | 14 | 18 | 10 | 17 | 17 | 5 | 29 | 22 |
| Ala GCA | 23 | 29 | 18 | 10 | 22 | 31 | 14 | 4 |
| GCC | 19 | 18 | 19 | 12 | 11 | 15 | 8 | 4 |
| GCG | 17 | 20 | 14 | 15 | 11 | 22 | 0 | 4 |
| GCU | 22 | 24 | 21 | 35 | 40 | 42 | 41 | 13 |
| Gly GGA | 17 | 12 | 21 | 8 | 18 | 13 | 23 | 27 |
| GGC | 11 | 16 | 5 | 22 | 11 | 12 | 11 | 9 |
| GGG | 10 | 16 | 5 | 5 | 11 | 10 | 12 | 18 |
| GGU | 22 | 34 | 10 | 37 | 14 | 16 | 12 | 9 |
| Val GUA | 16 | 19 | 12 | 13 | 11 | 7 | 15 | 9 |
| GUC | 27 | 23 | 31 | 14 | 4 | 6 | 2 | 4 |
| GUG | 16 | 14 | 17 | 8 | 12 | 9 | 15 | 18 |
| GUU | 30 | 36 | 24 | 37 | 24 | 26 | 23 | 18 |
| Lys AAA | 28 | 27 | 29 | 44 | 45 | 48 | 47 | 13 |
| AAG | 28 | 28 | 28 | 19 | 30 | 39 | 24 | 0 |
| Asn AAC | 34 | 40 | 29 | 18 | 24 | 35 | 13 | 9 |
| AAU | 28 | 24 | 32 | 27 | 18 | 11 | 26 | 18 |
| Gln CAA | 17 | 17 | 16 | 26 | 26 | 24 | 29 | 13 |
| CAG | 21 | 29 | 13 | 21 | 21 | 22 | 20 | 18 |
| His CAC | 3 | 5 | 1 | 4 | 4 | 0 | 7 | 4 |
| CAU | 3 | 3 | 3 | 9 | 9 | 8 | 12 | 0 |
| Glu GAA | 18 | 11 | 25 | 23 | 36 | 38 | 39 | 0 |
| GAG | 24 | 26 | 22 | 18 | 32 | 41 | 26 | 9 |
| Asp GAC | 25 | 15 | 34 | 23 | 24 | 26 | 24 | 4 |
| GAU | 25 | 19 | 30 | 25 | 28 | 22 | 38 | 9 |
| Tyr UAC | 21 | 26 | 16 | 10 | 12 | 10 | 15 | 13 |
| UAU | 6 | 5 | 8 | 26 | 15 | 11 | 21 | 13 |
| Cys UGC | 9 | 8 | 10 | 5 | 11 | 3 | 18 | 22 |
| UGU | 3 | 4 | 2 | 8 | 9 | 3 | 10 | 40 |
| Phe UUC | 18 | 20 | 16 | 24 | 16 | 22 | 5 | 40 |
| UUU | 18 | 11 | 25 | 26 | 25 | 16 | 35 | 31 |
| Ile AUA | 15 | 10 | 19 | 8 | 9 | 5 | 12 | 18 |
| AUC | 25 | 25 | 24 | 15 | 17 | 29 | 0 | 31 |
| AUU | 11 | 12 | 11 | 33 | 26 | 26 | 26 | 22 |
| Met AUG | 13 | 12 | 14 | 16 | 29 | 32 | 27 | 22 |
| Trp UGG | 12 | 16 | 7 | 13 | 23 | 17 | 24 | 58 |

Fig. 1(i)

DS DNA

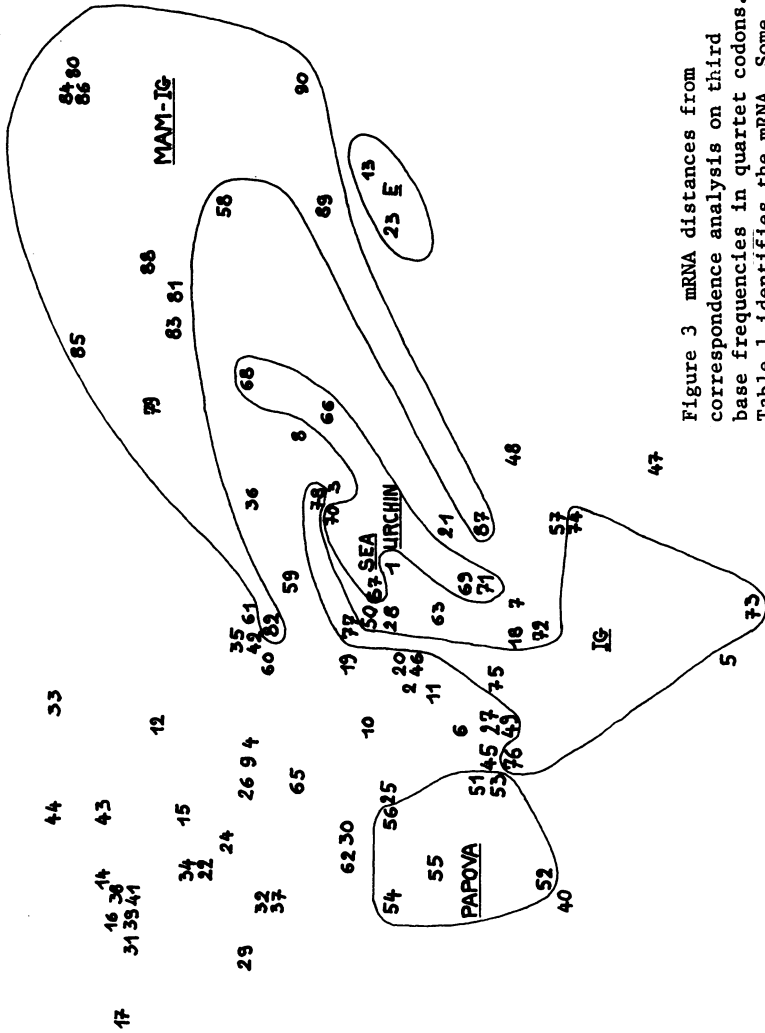| Codon | | BACT 58-63 | MITOCH 64 | YEAST 65 | ALL ANI 65-90 | ANI-MAM 65-72 | ALL MAM 73-90 | IG 73-78 | MAM-IG 79-90 |
|---|---|---|---|---|---|---|---|---|---|
| Arg | CGA | 3 | 0 | 0 | 4 | 8 | 2 | 7 | 0 |
| | CGC | 20 | 0 | 0 | 14 | 26 | 9 | 4 | 12 |
| | CGG | 2 | 0 | 0 | 6 | 5 | 6 | 3 | 8 |
| | CGU | 20 | 0 | 0 | 10 | 21 | 7 | 1 | 9 |
| | AGA | 2 | 13 | 28 | 8 | 14 | 6 | 12 | 3 |
| | AGG | 0 | 0 | 0 | 11 | 12 | 11 | 9 | 11 |
| Leu | CUA | 3 | 13 | 9 | 9 | 9 | 8 | 11 | 7 |
| | CUC | 3 | 0 | 0 | 25 | 27 | 24 | 26 | 24 |
| | CUG | 59 | 0 | 0 | 42 | 21 | 51 | 16 | 68 |
| | CUU | 7 | 0 | 9 | 10 | 16 | 7 | 7 | 7 |
| | UUA | 6 | 147 | 9 | 3 | 1 | 4 | 11 | 1 |
| | UUG | 6 | 0 | 46 | 8 | 10 | 7 | 10 | 6 |
| Ser | UCA | 3 | 67 | 9 | 10 | 9 | 11 | 24 | 5 |
| | UCC | 14 | 0 | 0 | 17 | 11 | 20 | 17 | 22 |
| | UCG | 4 | 0 | 9 | 3 | 2 | 4 | 1 | 5 |
| | UCU | 18 | 0 | 19 | 17 | 17 | 17 | 29 | 11 |
| | AGC | 9 | 0 | 0 | 20 | 19 | 21 | 27 | 18 |
| | AGU | 6 | 0 | 0 | 15 | 11 | 17 | 33 | 9 |
| Thr | ACA | 4 | 27 | 19 | 13 | 13 | 13 | 26 | 7 |
| | ACC | 21 | 0 | 28 | 26 | 28 | 25 | 30 | 23 |
| | ACG | 6 | 0 | 0 | 8 | 10 | 7 | 4 | 9 |
| | ACU | 24 | 0 | 28 | 18 | 13 | 19 | 36 | 11 |
| Pro | CCA | 6 | 13 | 28 | 12 | 13 | 12 | 20 | 7 |
| | CCC | 5 | 0 | 0 | ·17 | 12 | 18 | 13 | 21 |
| | CCG | 26 | 0 | 0 | 6 | 2 | 8 | 5 | 9 |
| | CCU | 1 | 13 | 9 | 12 | 11 | 13 | 13 | 12 |
| Ala | GCA | 40 | 53 | 0 | 16 | 29 | 11 | 20 | 6 |
| | GCC | 18 | 13 | 37 | 38 | 49 | 34 | 22 | 40 |
| | GCG | 28 | 0 | 0 | 6 | 6 | 5 | 1 | 8 |
| | GCU | 70 | 67 | 28 | 25 | 31 | 23 | 21 | 24 |
| Gly | GGA | 4 | 27 | 0 | 14 | 23 | 10 | 20 | 5 |
| | GGC | 31 | 0 | 19 | 26 | 21 | 28 | 16 | 34 |
| | GGG | 5 | 0 | 19 | 11 | 13 | 11 | 9 | 11 |
| | GGU | 31 | 107 | 74 | 19 | 19 | 18 | 25 | 15 |
| Val | GUA | 34 | 67 | 0 | 4 | 8 | 3 | 5 | 1 |
| | GUC | 9 | 0 | 0 | 19 | 25 | 16 | 23 | 13 |
| | GUG | 14 | 0 | 19 | 31 | 22 | 35 | 17 | 44 |
| | GUU | 42 | 13 | 9 | 9 | 12 | 7 | 7 | 7 |
| Lys | AAA | 65 | 27 | 56 | 19 | 29 | 15 | 17 | 14 |
| | AAG | 11 | 0 | 93 | 54 | 74 | 45 | 24 | 56 |
| Asn | AAC | 21 | 13 | 46 | 27 | 21 | 29 | 29 | 29 |
| | AAU | 7 | 13 | 19 | 9 | 6 | 10 | 13 | 9 |
| Gln | CAA | 14 | 13 | 19 | 12 | 17 | 10 | 11 | 9 |
| | CAG | 23 | 0 | 0 | 31 | 29 | 32 | 30 | 32 |
| His | CAC | 6 | 0 | 19 | 19 | 9 | 22 | 8 | 29 |
| | CAU | 4 | 0 | 19 | 10 | 9 | 10 | 11 | 10 |
| Glu | GAA | 57 | 13 | 46 | 22 | 23 | 22 | 21 | 23 |
| | GAG | 14 | 0 | 19 | 35 | 33 | 36 | 23 | 43 |
| Asp | GAC | 30 | 13 | 28 | 23 | 11 | 27 | 22 | 30 |
| | GAU | 19 | 0 | 9 | 16 | 11 | 18 | 23 | 16 |
| Tyr | UAC | 10 | 0 | 28 | 21 | 23 | 20 | 20 | 20 |
| | UAU | 4 | 13 | 19 | 12 | 3 | 15 | 17 | 14 |
| Cys | UGC | 4 | 0 | 9 | 9 | 1 | 12 | 6 | 16 |
| | UGU | 2 | 13 | 19 | 9 | 2 | 11 | 18 | 8 |
| Phe | UUC | 12 | 80 | 19 | 29 | 18 | 33 | 30 | 35 |
| | UUU | 12 | 13 | 19 | 14 | 9 | 16 | 13 | 18 |
| Ile | AUA | 3 | 0 | 0 | 4 | 3 | 4 | 7 | 3 |
| | AUC | 28 | 27 | 19 | 25 | 40 | 20 | 21 | 19 |
| | AUU | 18 | 93 | 19 | 11 | 10 | 11 | 18 | 8 |
| Met | AUG | 29 | 27 | 19 | 16 | 16 | 16 | 14 | 16 |
| Trp | UGG | 3 | 0 | 9 | 12 | 1 | 16 | 25 | 12 |

Fig 1 (ii)

Fig. 2 (i)

Fig. 2 (ii)

Fig. 2 (iii)

Fig. 2 (iv)

Figure 3 mRNA distances from correspondence analysis on third base frequencies in quartet codons. Some interesting genome types have been encircled (here "by eye"; elsewhere groupings are made by automatic classification (4, 5)).

REFERENCES

1. Grantham, R. (1978) FEBS lett. 95, 1-11.
2. Benzécri, J.P. (1973) In L'Analyse des données 2. l'analyse des correspondances, Dunod, Paris.
3. Hill, M.O. (1974) Appl. Statist. 23, 340-354.
4. Grantham, R. and Gautier, C. (1980) Naturwissenschaften, in press.
5. Grantham, R., Gautier, C. and Gouy, M. in preparation.
6. Garel, J.P. (1974) J. Theor. Biol. 43, 211-225.
7. Osterman, L.A. (1979) Biochimie 61, 323-342.
8. Fiers, W., Contreras, R., Duerinck, F., Haegeman, G., Merregaert, J., Min Jou, W., Raeymakers, A., Volckaert, G., Ysebaert, M., Van de Kerckhove, J., Nolf, F. and Van Montagu, M. (1975) Nature 256, 273-278.
9. Min Jou, W., Haegeman, W., Ysebaert, M. and Fiers, W. (1972) 247, 82-88.
10. Fiers, W., Contreras, R., Duerinck, F., Haegeman, G., Iserentant, D., Merregaert, J., Min Jou, W., Molemans, F., Raeymaekers, A., Van den Berghe, A., Volckaert, G. and Ysebaert, M. (1976) Nature 260, 500-507.
11. Escarmis, C., Sastry, P.A. and Billeter, M.A. (1978) J. Biol. Chem. 253, 8390-8396.
12. Jonard, G., Richards, K., Mohier, E. and Gerlinger, P. (1978) Eur. J. Biochem. 4, 521-529.
13. Guilley, H., Jonard, G., Kukla, B. and Richards, K.E. (1979) Nucleic Acids Res. 6, 1287-1308.
14. Guilley, H. and Briand, J. (1978) Cell 15, 113-122.
15. Sanger, F., Air, G., Barrell, B., Brown, N., Coulson, A., Fiddes, J., Hutchinson III, C., Slocombe, P. and Smith, M. (1977) Nature 265, 687-695.
16. Fiddes, J.C. (1977) Sci. American 237, 55-67.
17. Shaw, D.C., Walker, J.E., Northrop, F.D., Barrell, B.G., Godson, G.N. and Fiddes, J.C. (1978) Nature 272, 510-515.
18. Godson, G., Barrell, B., Staden, R. and Fiddes, J. (1978) Nature 276, 236-247.
19. Sugimoto, K., Sugusaki, H., Okamoto, T. and Takanami, M. (1978) Nucleic Acids Res. 5, 4495-4510.
20. Van Wezenbeek, P. and Schoenmakers, J.G.G. (1979) Nucleic Acids Res. 6, 2799-2818.
21. Huselbos, T. and Schoenmakers, J.G.G. (1978) Nucleic Acids Res. 5, 4677-4698.
22. Mc Connell, D.J. (1979) Nucleic Acids Res. 6, 3491-3503.
23. Sauer, R. (1978) Nature 276, 301-302.
24. Schwarz, E., Scherer, G., Hobom, G. and Kössel, H. (1978) Nature 272, 410-414.
25. Scherer, G. (1978) Nucleic Acids Res. 5, 3141-3156.
26. Grosschedl, R. and Schwarz, E. (1979) Nucleic Acis Res. 6, 867-881.
27. Fiers, W., Contreras, R., Haegeman, G., Rogiers, R., Van de Voorde, A., Van Heuverswyn, H., Van Herreweghe, J., Volckaert, G. and Ysebaert, M. (1978) Nature 273, 113-120.
28. Reddy, V.B., Thimmappaya, B., Dhar, R., Subramanian, K.N., Zain, B.S., Pan, J., Ghosh, P.K., Celma, M.L. and Weissman, S.M. (1978) Science 200, 494-502.
29. Dhar, R., Seif, I. and Khoury, G. (1979) Proc. Nat. Acad. Sci. USA 76, 565-569.
30. Valenzuela, P., Gray, P., Quiroga, M., Zaldivar, J., Goodman, H.M. and Rutter, W.J. (1979) Nature 280, 815-819.

31. Charnay, P., Mandart, E., Hampe, A., Fitoussi, F., Tiollais, P. and Galibert , F. (1979) Nucleic Acids Res. 7, 335-346.
32. Farabaugh, P. (1978) Nature 274, 765-769.
33. Post, L.E., Strycharz, G.D., Nomura, M., Lewis, H. and Dennis, P.P. (1979) Proc. Nat. Acad. Sci. USA 76, 1697-1701.
34. Sutcliffe, J. (1979) Proc. Nat. Acad. Sci. USA 75, 3737-3741.
35. Hensgens, L.A.M., Grivell, L.A., Borst, P. and Bos, J.L. (1979) Proc. Nat. Acad. Sci. USA 76, 1663-1667.
36. Smith, M., Leung, D.W., Gillam, S., Astell, C.R., Montgomery, D.L. and Hall, B.D. (1979) Cell 16, 753-761.
37. Schaffner, W., Keenz, G., Daetwyler, H., Telford, J., Smith, H. and Birnstiel, M. (1978) Cell 14, 655-671.
38. Sures, I., Lowry, J. and Kedes, L.H. (1978) Cell 15, 1033-1044.
39. Mc Reynolds, L., O'Malley, B.W., Nisbet, A.D., Fothergill, J.E., Givol, D., Fields, S., Robertson, M. and Brownlee, G.G. (1978) Nature 273, 723-728.
40. Seidman, J., Leder, A., Nau, M., Norman, B. and Leder, P. (1978) Science 202, 11-17.
41. Hamlyn, P.H., Brownlee, G.G., Cheng, C.C., Gait, M.J. and Milstein, C. (1978) Cell 15, 1067-1075.
42. Bernard, O., Hozumi, N. and Tonegawa, S. (1978) Cell 15, 1133-1144
43. Tonegawa, S., Maxam, A.M., Tizard,  , Bernard, O. and Gilbert, W. (1978) Proc. Nat. Acad. Sci. USA 75, 1485-1489.
44. Seidman, J.G., Max, E.E. and Leder, P. (1979) Nature 280, 370-375.
45. Rogers, J., Clarke, P. and Salser, W. (1979) Nucleic Acids Res. 6, 3305-3321.
46. Konkel, D.A., Tilghman, S.M. and Leder, P. (1978) Cell 15, 1125-1132.
47. Roberts, J.L., Seeburg, P.H., Shine, J., Herbert, E., Baxter, J.D. and Goodman, H.M. (1979) Proc. Nat. Acad. Sci. USA 76, 2153-2157.
48. Seeburg, P.H., Shine, J., Martial, J.A., Baxter, J.D. and Goodman, H.M. (1977) Nature 270, 486-494.
49. Gubbins, E.J., Maurer, R.A., Hartley, J.L. and Donelson, J.E. (1979) Nucleic Acids Res. 6, 915-930.
50. Ullrich, A., Shine, J., Chirgwin, J., Pictet, R., Tischer, E., Rutter, W.J. and Goodman, H.M. (1977) Science 196, 1313-1319.
51. Heindell, H.C., Liu, A., Paddock, G.V., Studnika, G.M. and Salser, W.A. (1978) Cell 15, 43-54.
52. Efstratiadis A., Kafatos, F.C. and Maniatis, T. (1977) Cell 10, 571-585.
53. Nakanishi, S., Inoue, A., Kita, T., Nakumura, M., Chang, A.C.Y., Cohen, S.N. and Numa, S. (1979) Nature 278, 423-427.
54. Fiddes, J.C. and Goodman, H.M. (1979) Nature 281, 351-356.
55. Marotta, C.A., Wilson, J.T., Forget, B.G. and Weissman, S.M. (1977) J. Biol. Chem. 252, 5040-5050.
56. Martial, J.A., Hallewell, R.A., Baxter, J.D. and Goodman, H.M. (1979) Science 205, 602-607.
57. Shine, J., Seeburg, P., Martial, J., Baxter, J. and Goodman, H.M. '1977) Nature 270, 494-499.