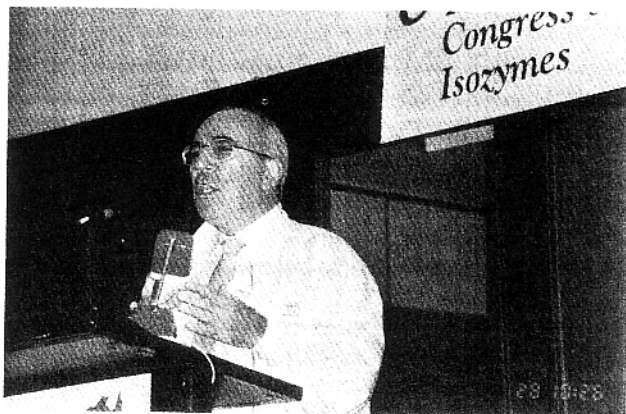


THE DISTRIBUTION OF GENES IN THE HUMAN GENOME

Giorgio Bernardi*

* Laboratoire de Genetique Moleculaire, Institut Jacques Monod, 2 Place Jussieu, 75005 Paris, France



The term *genome* was coined three quarters of a century ago by a German botanist, Winkler¹ to designate the haploid chromosome set. While current textbooks of Molecular Biology do not yet go beyond the purely operational definition of the eukaryotic genome as the sum total of genes and of intergenic sequences, a number of molecular biologists have been thinking for some time that the genome is more than the sum of its parts. This implies the existence of structural and functional interactions between the minority of coding sequences and the majority of non-coding sequences. This general, rather vague concept has been changed into a precise one by the discovery, in our laboratory, of genome properties. These properties, which have mainly been defined by investigations on the nuclear genome of vertebrates (see Refs. 2-5) will be briefly outlined here. They comprise the isochore organization, the compositional patterns of DNA fragments (or molecules) and of coding sequences, the compositional correlations between coding and non-coding sequences and, above all, the gene distribution and its associated functional properties.

The mammalian genomes are mosaics of *isochores* (see Fig. 1), namely of long (>300 Kb) DNA segments that are homogeneous in base composition and range from 30 to 60% GC.^{6,7} This is an extremely wide range, almost as wide as that covered by all bacterial DNAs (25-72% GC). In the human genome, isochores can be assigned to two GC-poor families (L1 and L2) representing 2/3 of the genome, and to three GC-rich families (H1, H2 and H3) forming the remaining 1/3 (Fig. 2).

The *compositional distributions* of large (>100 Kb) genome fragments, such as those forming routine DNA preparations, of exons (and particularly of their third codon positions) and of introns represent *compositional patterns*.^{2,10} These correspond to *genome phenotypes*,³ in that they differ characteristically not only between cold- and warm-blooded vertebrates,

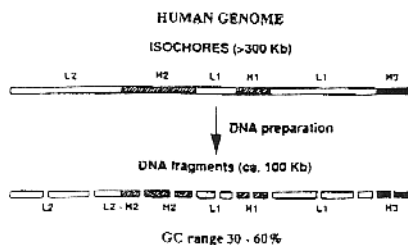


Figure 1. Scheme of the isochore organization of the human genome. This genome which is a typical mammalian genome, is a mosaic of large (>300 Kb) DNA segments, the isochores. These are compositionally homogeneous (above a size of 3 Kb) and can be subdivided into a small number of families, GC-poor (L1 and L2), GC-rich (H1) and (H2), and very GC-rich (H3). The GC-range of the isochores from the human genome is 30-60% (from Ref. 5).

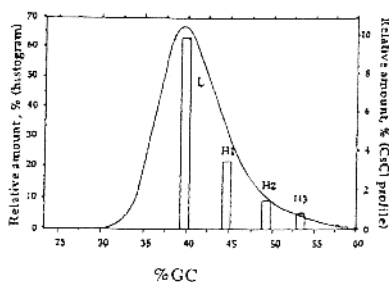


Figure 2. Histogram of the isochore families from the human genome. The relative amounts of major DNA components derived from isochore families L (i.e., L1 + L2), H1, H2, H3 (see Ref. 7) are superimposed on the CsCl profile of human DNA (from Ref. 9).

but also between mammals and birds and even between murids and most other mammals (see Figs. 3 and 4).

*Compositional correlations*² exist between exons (and their codon positions) and isochores (Fig. 5), as well as between exons and introns.¹¹ These correlations concern, therefore, coding and non-coding sequences and are not trivial since coding sequences only make up about 3% of the genome, whereas non-coding sequences correspond to 97% of the genome. The compositional correlations represent a *genomic code*.^{4,5} It should be noted that a *universal correlation* holds among GC levels of codon positions (third positions against first and/or second positions). Both the genomic code and the universal correlations are apparently due to compositional constraints working in the same direction (towards GC or AT), although to different extents on coding and non-coding sequences as well as on different codon positions.

The compositional correlations between GC3 (the GC level of third codon positions) and isochore GC have a practical interest in that they allowed us to position the coding sequence histogram of Fig. 4 relative to the CsCl profile of Fig. 3 and to assess the *gene distribution* in the human genome.^{5,9} In fact, if one divides the relative number of genes per histogram bar by the corresponding relative amount of DNA, one can see that gene concentration is low and constant in GC-poor isochores, increases with increasing GC in isochore families H1 and

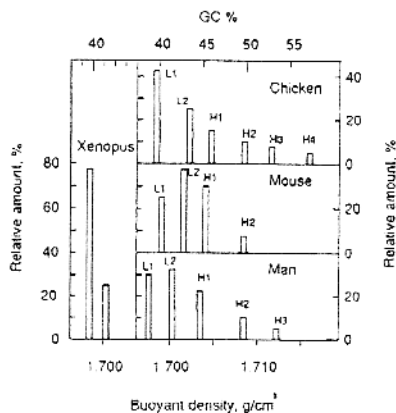


Figure 3. Compositional patterns of vertebrate genomes. Histograms showing the relative amounts, modal buoyant densities and GC levels of the major DNA components from *Xenopus*, chicken, mouse and man, as estimated after fractionation of DNA by preparative density gradient in the presence of a sequence-specific DNA ligand (Ag^+ or BAMD; BAMD is bis (acetato mercuri methyl) dioxane). The major DNA components are the families of large DNA fragments (see Fig. 1) derived from different isochore families. Satellite and minor DNA components (such as rDNA) are not shown in these histograms (from Ref. 5).

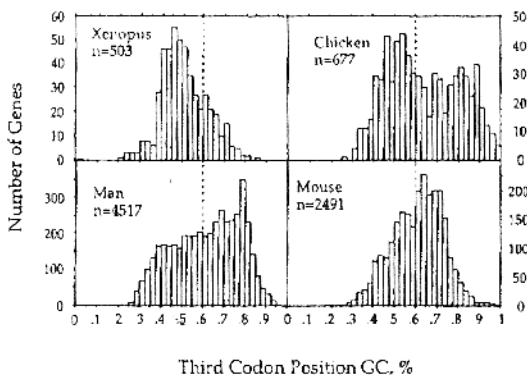


Figure 4. Compositional distribution of third codon positions from vertebrate genes. The number of genes taken into account is indicated. A 2.5% GC window was used. The broken line at 60% GC is shown to provide a reference (from Ref. 5).

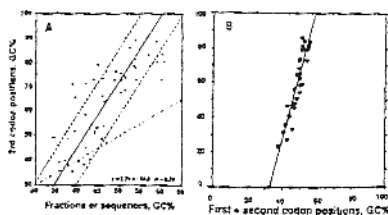


Figure 5. A) GC levels of third codon positions from human genes are plotted against the GC levels of DNA fractions (dots) or extended sequences (circles) in which the genes are located. The correlation coefficient and slope are indicated. The dash-and-point line is the diagonal line (slope=1). GC levels of third codon positions would fall on this line if they were identical to GC levels of surrounding DNA. The broken lines indicate a $\pm 5\%$ GC range around the slope (from Ref. 9). B) Plot of GC levels of third codon positions of genes from prokaryotic and eukaryotic genomes are plotted against GC levels of first + second positions. All values are averaged per genome (or per genome compartment, in the case of compositionally compartmentalized genomes) (from Ref. 12).

H2, and reaches a maximum in isochore family H3, which exhibits at least a 20-fold higher gene concentration compared to GC-poor isochores (Fig. 6).

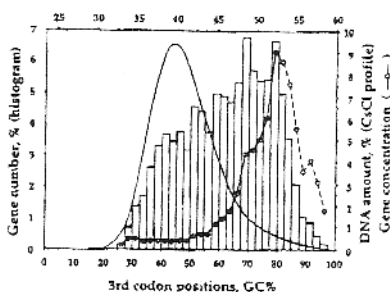


Figure 6. Profile of gene concentration in the human genome as obtained by dividing the relative amounts of genes in each 2.5% GC interval of the histogram by the corresponding relative amounts of DNA deduced from the CsCl profile. The apparent decrease in gene concentration for very high GC values (broken line) is due to the presence of rDNA in that region. The last concentration values are uncertain because they correspond to very low amounts of DNA (from Ref. 5).

The H3 isochore family has been called the *human genome core*,⁵ because it corresponds to the functionally most significant part of the human genome. Indeed, the H3 isochore family is not only endowed with the highest gene (and CpG island) concentration, but also with an open chromatin structure (as witnessed by the accessibility to DNases, as well as by the scarcity of histone H1, the acetylation of histones H3 and H4 and wider nucleosome spacing),¹³ with the highest transcription and recombination levels and with the earliest replication timing. The genes of the genome core have the highest GC3 levels relative to their flanking sequences, have the shortest exons and introns,¹⁴ exhibit an extreme codon usage, and encode proteins characterized by amino acid frequencies differing from those of proteins encoded by GC-poor isochores.¹²

The human genome core is located in T(olomeric)-bands,¹⁵ which are essentially formed by GC-rich isochores (mainly of the H2 and H3 families). In contrast, R'-bands, namely the R(everse) bands exclusive of T-bands, comprise both GC-rich isochores (of the H1 family) and GC-poor isochores. Finally, G(iemsa) bands are formed almost exclusively by GC-poor isochores (see Fig. 7).⁸ The difference in GC level between G-bands and T-bands is about 15%. About 20% of genes are present in G-bands and about 80% in R-bands (60% of them in T-bands). The location of a majority of genes in T-bands is of interest in view of the association of telomeres with the nuclear matrix and envelope.¹⁶

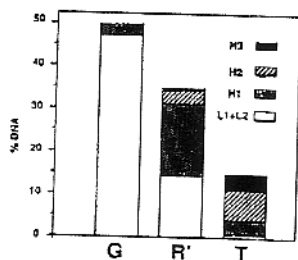


Figure 7. A scheme of the relative amounts of isochore families L1 + L2, H1, H2 and H3 in G-bands, R'-bands and T-bands; R'-bands are R-bands exclusive of T-bands (from Ref. 8).

It should be stressed that the gene distribution reported for the human genome seems to have been conserved in evolution, genes showing their highest concentration in the GC-richest isochores of all vertebrates.⁵

In the case of homologous mammalian genes, it has been possible to show that third codon position synonymous substitutions exhibit frequencies and compositions which strongly suggest natural selection.^{17,18} Under these circumstances, the compositional changes in non-coding sequences, which are correlated with those occurring in third codon positions, suggest that non-coding sequences are not junk DNA, but must fulfill some functional role.

As already mentioned, the compositional pattern of the human genome, which is typical of the genomes of most mammals and similar to the genomes of birds, is strikingly different from the compositional patterns of cold-blooded vertebrates which exhibit a much lower degree of heterogeneity and are characterized by metaphase chromosomes which do exhibit R-banding. These different genome phenotypes of warm- versus cold-blooded vertebrates are due to compositional changes. While the gene-poor, GC-poor isochores have undergone little or no compositional change in vertebrates genomes, the gene-rich, GC-rich isochores are those which underwent compositional changes in evolution.

References

1. H. Winkler, *Verbreitung und Ursache der Parthenogenesis im Pflanzen und Tierreich*, Fischer, Jena, 1920.
2. G. Bernardi, B. Olofsson, J. Filipski, M. Zerial, J. Salinas, G. Cuny, M. Meunier-Rotival, and F. Rodier, "The mosaic genome of warm-blooded vertebrates", *Science* **228** (1985) 953-958.
3. G. Bernardi, and G. Bernardi, "Compositional constraints and genome evolution", *J. Mol. Evol.* **24** (1986) 1-11.

4. G. Bernardi, "Le génome des vertébrés: organisation, fonction, évolution", *Biofutur* 94 (1990) 43-46.
5. G. Bernardi, "The human genome organization and its evolutionary history: a review", *Gene* 135 (1993) 57-66.
6. J.P. Thiery, G. Macaya, and G. Bernardi, "An analysis of eukaryotic genomes by density gradient centrifugation", *J. Mol. Biol.* 108 (1976) 219-235.
7. G. Macaya, J.P. Thiery, and G. Bernardi, "An approach to the organization of eukaryotic genomes at a macromolecular level", *J. Mol. Biol.* 108 (1976) 237-254.
8. S. Saccone, A. De Sario, J. Wiegant, A.K. Rap, G. Della Valle, and G. Bernardi, "Correlations between isochores and chromosomal bands in the human genome", *Proc. Natl. Acad. Sci. USA* 90 (1993) 11929-11933.
9. D. Mouchiroud, G. D'Onofrio, B. Aïssani, G. Macaya, C. Gautier, G. Bernardi, "The distribution of genes in the human genome", *Gene* 100 (1991) 181-187.
10. D. Mouchiroud, G. Fichant, and G. Bernardi "Compositional compartmentalization and gene composition in the genome of vertebrates", *J. Mol. Evol.* 26 (1987) 198-204.
11. B. Aïssani, G. D'Onofrio, D. Mouchiroud, K. Gardiner, C. Gautier and G. Bernardi, "The compositional properties of human genes", *J. Mol. Evol.* 32 (1991) 497-503.
12. G. D'Onofrio, D. Mouchiroud, B. Aïssani, C. Gautier, and G. Bernardi, "Correlations between the compositional properties of human genes, codon usage and amino acid composition of proteins", *J. Mol. Evol.* 32 (1991) 504-510.
13. J. Tazi, and A. Bird, "Alternative chromatin structure at CpG islands", *Cell* 60 (1991) 909-920.
14. L. Duret, D. Mouchiroud, and C. Gautier, "Statistical analysis of vertebrate sequences reveals that long genes are scarce in GC-rich isochores", *J. Mol. Evol.* 40 (1995) 308-317.
15. S. Saccone, A. De Sario, G. Della Valle, and G. Bernardi, "The highest gene concentrations in the human genome are in T-bands of metaphase chromosomes", *Proc. Natl. Acad. Sci. USA* 89 (1992) 4913-4917.
16. T. de Lange, "Human telomeres are attached to the nuclear matrix", *EMBO J.* 11 (1992) 717-724.
17. S. Cacciò, P. Perani, S. Saccone, F. Kadi, and G. Bernardi, "Single-copy sequence homology among the GC-richest isochores of the genomes from warm-blooded vertebrates", *J. Mol. Evol.* 39 (1995) 331-339.
18. S. Zoubak, G. D'Onofrio, S. Cacciò, G. Bernardi, and G. Bernardi, "Specific compositional patterns of synonymous positions in homologous mammalian genes", *J. Mol. Evol.* 40 (1995) 293-307.