# BREAKTHROUGHS AND VIEWS

## The Isochores in Human Chromosomes 21 and 22

Oliver Clay and Giorgio Bernardi[1]

*Laboratory of Molecular Evolution, Stazione Zoologica Anton Dohrn, Villa Comunale, 80121 Naples, Italy*

In a recent *Biochemical and Biophysical Research Communications* article entitled "No Isochores in the Human Chromosomes 21 and 22?," Häring and Kypr (1) have questioned whether isochores and homogeneity of GC level exist in the contiguously sequenced human chromosomes 21 and 22. Their communication is structured in two parts. In the first part, they assume, without justification, that a sequence can be "homogeneous" only if randomization no longer lowers the variance in GC level among its segments. Since they find "the variations to be higher everywhere compared to the randomized sequences," they incorrectly conclude that "the (G + C) content is certainly not homogeneous on the isochore scale in the two human chromosomes." Unfortunately, both the title and the abstract mention only the results from this first part of their analysis.

The second part of their paper, summarized in its Fig. 4, tells quite a different story. Indeed, this figure shows "[r]egions of homogenous (G + C) content longer than 300 kbp" along human chromosomes 21 and 22, at the isochore scale, in clear contradiction to what is concluded in the abstract, and in answer to the question posed in their title. The authors make no attempt to reconcile these more positive results with the contradictory assertions in the first part of their paper. Similarly, the definitions of "homogeneous" used in the two parts of their paper are in contradiction to one another, and the authors make no attempt to explain this. The communication in question is, unfortunately, flawed in both parts, and isochores continue to exist on human chromosomes 21 and 22.

The flaw in the first part of the paper is a simple one, namely an unreasonably strict definition of homogeneity. The DNA sequences within individual isochores show variations of GC level that are small compared to those in the entire human genome, but much greater than those expected from a "random" sequence in which nucleotides are independent and identically distributed. This result, shown two decades ago (2) by density gradient ultracentrifugation, can be explained by the presence of large-scale correlations in the sequences (3, 4, and references therein). Randomizing such sequences simply breaks up the correlations and, as a consequence, decimates the variance. The authors' first definition of "homogeneous" therefore excludes not only human and *E. coli* DNA, as they indeed observe, but it also excludes any natural DNA except for repetitive satellites, and, more generally, any binary sequences in which long-range correlations are present. It is a pity that the authors do not openly re-evaluate their initially adopted definition of "homogeneous" in the light of their telling results. At first they still leave open the possibility that "the (G + C) content variation as defined above is not the proper quantity to identify the isochores," but unfortunately they do not mention it again under Discussion, let alone in the abstract, so that the reader is misled.

The flaws in the second part of the paper are more subtle. The authors, realizing that the randomization criterion cannot lead to isochores, now explore less stringent variants of their criterion, for example in which the isochores "had to have a minimum length of 300 kbp inside which the (G + C) content fluctuated within 2% of the (A + C + G + T) content." In their analyses, fluctuation thresholds are allowed to vary from 1 to 3%, and fluctuations are assessed for iterated extensions of 2.5 kb to 250 kb. For this class of criteria, Häring and Kypr consistently obtain only relatively short homogeneous regions in human, except where they use 250 kb extensions and an extrinsically imposed lower bound of 3 Mb on the regions.

A benchmark test for isochore prediction methods is the long, ≈7 Mb, GC poor and gene-poor isochore in chromosome 21 (5), of which the two halves have almost identical GC distributions for fragments in the range 100 bp to 100 kb. With 250 kb extensions and a lower bound of 3 Mb, the authors recover this isochore, and observe no partitioning of *E. coli* into isochores, but (by construction) miss any isochores shorter than 3

[1] To whom correspondence should be addressed. Fax: +39 081 245 5807. E-mail: bernardi@alpha.szn.it.

Mb in human. On the other hand, when they use shorter extensions and do not impose the 3 Mb bound, they obtain only much shorter "isochores" (<2 Mb) within this isochore, and as many as 8 "isochores" within the relatively homogeneous genome of *E. coli,* in which one would expect only one or, at most, two such regions.

It is well known that isochores in the human genome have different lengths (see Ref. (6) for examples and a discussion). The inability of Häring and Kypr's methods to recognize both short and long isochores using a single criterion leaves the reader with only one remaining option: that none of the measures of GC content variation used are "the proper quantity to identify the isochores."

Sequences containing long-range correlations, such as DNA, exhibit fluctuations that can be an order of magnitude higher than for randomized sequences consisting of independent, identically distributed nucleotides. It would be beyond the scope of this comment to discuss possible modifications to the authors' tests required by the presence of long-range correlations (4) or by the more pronounced fluctuations in GC rich isochores compared to GC poor isochores (2), or to discuss existing, successful methods of detecting isochores at the DNA sequence level. Such methods, which employ recursive segmentation (7), will be discussed in a forthcoming special issue of *Gene.* Their existence renders questionable the final conclusion of the paper, namely that "[i]n any case, the present analysis demonstrates that the isochores should be defined in unambiguous molecular terms to be useful in up-to-date genome analysis." This conclusion, which is echoed in the abstract, is ambiguous (as regards the nature of "molecular"), inaccurate (the only demonstration given by the analysis is that the authors' methods or criteria do not succeed) and not up-to-date (other methods are not mentioned).

## REFERENCES

1. Häring, D., and Kypr, J. (2001) *Biochem. Biophys. Res. Commun.* **280,** 567–573.
2. Cuny, G., Soriano, P., Macaya, G., and Bernardi, G. (1981) *Eur. J. Biochem.* **115,** 227–233.
3. Li, W., Stolovitzky, G., Bernaola-Galván, P., and Oliver, J. L. (1998) *Genome Res.* **8,** 916–928.
4. Beran, J. (1994) Statistics for Long-Memory Processes, Chapman & Hall/CRC, Boca Raton, FL.
5. Hattori, M., Fujiyama, A., Taylor, T. D., *et al.* (2000) *Nature* **405,** 311–319.
6. Bernardi, G. (1995) *Annu. Rev. Genet.* **29,** 445–476.
7. Bernaola-Galván, P., Román-Roldán, R., and Oliver, J. L. (1996) *Phys. Rev. E* **53,** 5181–5189.