

Prediction of replication time zones at single nucleotide resolution in the human genome

Feng Gao, Chun-Ting Zhang*

Department of Physics, Tianjin University, Tianjin 300072, China

Received 6 March 2008; revised 3 June 2008; accepted 4 June 2008

Available online 12 June 2008

Edited by Takashi Gojobori

Abstract The human genome is structured at multiple levels: it is organized into a series of replication time zones, and meanwhile it is composed of isochores. Accumulating evidence suggests a match between these two genome features. Based on newly developed software GC-Profile, we obtained a complete coverage of the human genome by 3198 isochores with boundaries at single nucleotide resolution. Interestingly, the experimentally confirmed replication timing sites in the regions of 1p36.1, 6p21.32, 17q11.2 and 22q12.1 nearly all coincide with the determined isochore boundaries. The precise boundaries of the 3198 isochores are available via the website: <http://tubic.tju.edu.cn/isomap/>.

© 2008 Federation of European Biochemical Societies. Published by Elsevier B.V. All rights reserved.

Keywords: Human genome; Isochore; Replication timing; Segmentation; G + C content

1. Introduction

Since the pioneer work of Bernardi and co-workers, it has been well established that the mammalian genomes are composed of large sequence segments of fairly homogeneous G + C content that was revealed by the analytical ultracentrifugation of bulk DNA in the mid 1970s [1]. The long DNA segments (on average >300 kb) of fairly homogeneous G + C content lately were given the name ‘isochores’ [2]. Since then, the issues of isochores in mammalian and other eukaryotic genomes have attracted wide attention [3–17].

One of challenging problems in isochore studies is to determine how many isochores there are in the human genome, and the solution to this problem largely depends on the definition of isochores. Bernardi and co-workers recently obtained a complete coverage of the human genome by about 3200 isochores [10,11], and they also found the rules according to which the 3200 isochores of the human genome are assembled in high resolution chromosome bands (850 bands) [12]. Being different from Bernardi and co-workers [10,12], here we present an alternative solution. A large body of evidence has shown that the mammalian chromosome is organized into a number of large replication time zones [18,19], ranging from 100 kb to 2 Mb in length roughly [20]. The replication time zones (RT-

zones) themselves are not the minimum replication units; rather, most of them are composed of several tandem clustered replicons [21]. It was reported that the RT-zones with higher G + C content replicate early, whereas the RT-zones with lower G + C content replicate late [6,7]. Many experiments revealed that the replication timing switch sites are consistent with the transitions of G + C content along the chromosomes investigated [6,7]. Very recent experiment shows that two adjacent replication timing sites constitute an RT-zone, which is exactly an isochore [7]. In the current work, isochores in the human genome were determined such that their boundaries coincide with replication timing sites that have been experimentally confirmed. Therefore, other isochore boundaries, in addition to those coincide with confirmed replication timing sites, are possible to have the same role, thereby boosting the number of potential replication timing sites for further experimental validation.

In our previous papers [16,17], only 56 isochores with length longer than 3000 kb were studied, which cover about 21% of the whole human genome, whereas in the present study, 3198 isochores with length longer than 30 kb were identified, which cover nearly 100% (except gaps) of the whole human genome.

2. Materials and methods

The whole human genome sequences were downloaded from <http://hgdownload.cse.ucsc.edu/downloads.html#human> released in March 2006. Two independent methods were used to identify boundaries of isochores. The first method is called the cumulative GC-profile (z' curve) method [15]. For a given genome or chromosome, there is a unique cumulative GC-profile or z' curve corresponding to it. The z' curve or the cumulative GC-profile is used interchangeably in this paper. Note that the essence of cumulative GC-profile is to intuitively display the variations of the G + C content along a genome or chromosome. It is not the G + C content itself. Rather, the derivative of z' curve with respect to the base position n is negatively proportional to the G + C content at the given position, i.e., $G+C \propto -dz'/dn$. Therefore, the average slope of the z' curve within a region reflects the average G + C content of the sequence within this region. If the z' curve in a region is an approximately straight line, the G + C content keeps approximately constant within this region. A jump (drop) in the z' curve indicates a decrease (increase) of the G + C content. A turning point in the z' curve indicates a switch site, at which the G + C content undergoes an abrupt change from a relatively GC-poor (GC-rich) region to a relatively GC-rich (GC-poor) region. The point at which the derivative of the z' curve is not continuous is called a turning point or segmentation point. However, a more convenient method to find the coordinates is to use a newly developed segmentation algorithm [22], which is implemented using a computer program, called GC-Profile [23]. The cumulative GC-profile, the distribution of G + C content, the isochores and their boundary coordinates for each of the human

*Corresponding author. Fax: +86 22 2740 2697.

E-mail addresses: ctzhang@tju.edu.cn (C.-T. Zhang), ren_zhang@yahoo.com (C.-T. Zhang).

chromosomes as well as the features of isochores are available online by visiting <http://tubic.tju.edu.cn/isomap/>. Consequently, the total number of isochores is 3198 for the whole human genome (hg18).

3. Results and discussion

3.1. Comparisons with experimental evidence

In what follows, we will show evidence that the boundaries of some isochores obtained here are in accordance with the known replication timing sites confirmed by experiments. Please refer to the related materials (such as Supplementary Figures 1–4) from <http://tubic.tju.edu.cn/isomap/supplementary> while reading the following text. The first evidence we show here concerns the replication timing switch site in the human major histocompatibility complex (MHC) sequence. One of the replication timing sites was found experimentally within this sequence [6]. Readers are suggested to refer to the UCSC Genome Browser in the region of chr6: 31,647,700–33,232,435, or Supplementary Figure 1. The replication timing switch region is between the gene *NOTCH4* (chr6: 32,270,599–32,299,822) in MHC class III (isochore H3) and the gene *HLA-DRA* (chr6: 32,515,625–32,520,799) in MHC class II (isochore L2). To be more accurate, the switch region should be located within the interval 32,280–32,320 kb (please refer to Fig. 6 or Fig. 8 in Ref. [6]), while the segmentation point obtained by our software GC-Profile, 32,300,166 bp within the 6p21.32 band, is situated exactly within the above interval. Therefore, the boundary between the predicted isochore 6-30 and isochore 6-31 is precisely consistent with the replication timing site confirmed experimentally [6].

The second evidence is related to the transitions in the replication timing in a 340 kb region of the human chromosome R-band 1p36.1 [24]. Please refer to the UCSC Genome Browser in the region of chr1:17,212,621–19,418,116 or Supplementary Figure 2. According to the result by GC-Profile, the precise coordinate of the segmentation point within this region is 19,288,551. Previous experimental evidence showed that there is a switch region (also called replication fork barrier by these authors) between the gene *ALDH4A1* (chr1:19,070,513–19,101,659) and the gene *RBAF600* (chr1: 19,417,172–19,450,633) [24]. Obviously, the predicted replication timing site, i.e., 19,288,551, is exactly situated within the above interval (chr1: 19,101,659–19,417,172).

The third evidence we present here is regarding the replication timing region on the human chromosome 17q11.2, in which the *NFI* gene resides. Schmegner and co-workers found that (i) a transition from a GC-poor isochore to a GC-rich one in the *NFI* region occurs within 5 kb; (ii) at the isochore transition the replication fork is stalled in the mid-S phase of cell cycle, which can be visualized by fiber-FISH techniques as a Y-shaped structure [25]. In other words, the boundary of the two isochores is also the replication timing site. It was found that the boundary between the two isochores is sharp and is located exactly at the 14 kb intergenic region between the gene *NFI* and gene *RAB11FIP4*. Please refer to the figure of UCSC Genome Browser (hg18) in the region of chr17: 26,457,053–26,954,701, or Supplementary Figure 3. The intergenic region is between the gene *NFI* (chr17: 26,446,121–26,728,820) and the gene *RAB11FIP4* (chr17: 26,742,768–26,889,352), with 14 kb in length, while the segmentation point obtained by the software GC-Profile, i.e., 26,735,738, is exactly within this

region, which is just the transition site at which the G + C content varies from 37% to 51% between the isochore 17-74 (26,457,053–26,735,738) and the isochore 17-75 (26,735,739–26,954,701). The two isochores identified here, i.e., the isochore 17-74 and isochore 17-75, are exactly the two isochores studied by Schmegner et al. [25].

The fourth evidence we show here concerns the replication timing of the *MNI/PITPNB* gene region on chromosome 22 [7]. The analysis of the G + C content for this fragment of DNA sequence showed that there are four distinct sub-regions with different G + C content. The proximal sub-region has a G + C content of 50.1%, and the G + C content of the second, third and fourth sub-regions are 39.5%, 53.0% and 39.5%, respectively. The four sub-regions are termed A-, B- C- and D-isochores, respectively. Perhaps, the most surprising finding of the work is that the four isochores match the replication timing zones with such a degree that the authors called the match “perfect” [7]. The authors found that the A- and C-isochores replicate early, whereas the B- and D-isochores replicate late during the S phase of the cell cycle. There are three sharp boundaries (transitions) of the four isochores, i.e., the boundaries between isochores A/B, B/C and C/D, respectively. The analysis of the G + C content (using the window-based method) showed that the transition occurs within a region of few kb. Interestingly, the four replication zones or isochores the authors studied [7] are basically the four isochores obtained in this work. Please refer to the figure of UCSC Genome Browser (hg18) in the region of chr22: 26,129,874–26,866,656, or Supplementary Figure 4. This region consists of four distinct isochores, i.e., the isochores 22-32, 22-33, 22-34 and 22-35, respectively, separated by three sharp isochore transitions. The four isochores (RT-zones) are arranged in the following order: the GC-rich proximal isochore, the GC-poor isochore, the short GC-rich isochore and the distal isochore, respectively. Obviously, the isochore 22-32 corresponds to the isochore A; the isochore 22-33 corresponds to the isochore B; the isochore 22-34 corresponds to the isochore C and the isochore 22-35 corresponds to the isochore D. The first switch region is between the gene *MNI* (chr22: 26,474,266–26,527,486) and the gene *PITPNB* (chr22: 26,577,658–26,645,255), while the segmentation point obtained by the software GC-Profile, i.e., the predicted replication timing site, 26,577,053, is within this intergenic region. The second switch region is near the proximal part of *KIAA1043* (RP3-477H23.1-001) gene (chr22: 26,707,254–26,889,455), while the segmentation point obtained by GC-Profile, 26,703,751, is precisely located here. The third switch region is within the transition between the isochores C and D. The corresponding segmentation point is found to be 26,757,747 by GC-Profile, which is exactly the replication timing site between the isochores C and D confirmed by experiment with isochore C replicating early and isochore D late [7].

Replication timing of the human X-inactivation center (XIC) region on chromosome X was studied 7 years ago before the completion of the Human Genome Project [26]. The authors found two regions where the replication timing changes from the early to late period during the S phase of the cell cycle. The first region they found is located near a large inverted duplication segment proximal to the XIC, and the second is near the XIST locus. However, the predicted positions of replication timing have ~500 kb displacements with those reported by these authors. One of possible explanations for

the inconsistency is that the complete match between replication time zones and isochores is invalid for the XIC region since other factors, in addition to compositional difference, also influence the replication timing of chromosomes. It is also possible that the different versions of the human genome used between the two studies cause the inconsistency.

The above comparison results are summarized in Table 1. In summary, nearly all the known replicating timing sites are in accordance with the boundaries of isochores concerned, suggesting that the number and locations of isochores proposed in this paper are reasonable.

3.2. Relationships between the isochores and chromosome bands

The obtained isochores can be displayed in the UCSC Genome Browser as a custom track [27], and a series of tracks aligned with the genomic sequence, such as chromosome bands, can also be shown together. The chromosome band track represents the approximate location of bands seen on Giemsa-stained chromosomes, and a total of 862 sub-bands were provided at the UCSC Genome Browser (hg18). The isochore maps at the UCSC Genome Browser for each of the human chromosomes can be viewed by clicking on the corresponding link at Supplementary Table S1 by visiting <http://tubic.tju.edu.cn/isomap/supplementary>. There are 3198

isochores obtained by GC-Profile while only 862 chromosome bands are provided at the UCSC Genome Browser. In general, one band, especially for the gene-rich R band, corresponds to multiple isochores. For example, the R band 21q22.3 (chr21: 41,400,001–46,944,323), which is rich in G + C content, CpG islands and genes, covers 20 isochores from 21-19 (chr21: 41,349,486–41,651,246) to 21-38 (chr21: 46,704,406–46,944,323). The one-to-one correspondence between isochores and chromosome bands with a fairly good match can also be found for a number of isochores. For example, isochores X-17 (chrX: 25,307,501–29,896,793), X-18 (chrX: 29,896,794–31,489,771) and X-19 (chrX: 31,489,772–36,628,265) correspond to the bands Xp21.3 (chrX: 24,900,001–29,400,000), Xp21.2 (chrX: 29,400,001–31,500,000) and Xp21.1 (chrX: 31,500,001–37,500,000), respectively. In some cases, multiple bands correspond to one isochore, especially for the bands in the AT-rich region. Take the longest isochore on human chromosome X as an example. The size of this isochore X-60 (chrX: 77,492,741–99,276,124) is 21.78 Mb. The G + C content of this isochore is 36.36%, much lower than the average G + C content (38.98%) of the residing contig. The entire region in which six bands are involved is mainly made up of G bands, and the total size of G bands in this region is 17.41 Mb. Genes and CpG islands are very rare in this region, as expected from the fact that the

Table 1
Comparison between the coordinates of segmentation points obtained by GC-Profile and the boundaries of RT-zones confirmed by experiments

Number	Chromosome	RT-switch region	Segmentation point	Adjoining isochores	Adjacent genes or markers
1	1	19,101,659–19,417,172	19,288,551	1-40 and 1-41	<i>ALDH4A1</i> and <i>RBAF600</i>
2	6	32,299,822–32,515,625	32,300,166	6-30 and 6-31	<i>NOTCH4</i> and <i>HLA-DRA</i>
3	17	26,728,820–26,742,768	26,735,738	17-74 and 17-75	<i>NF1</i> and <i>RAB11FIP4</i>
4	22	~26,577,658	26,577,053	22-32 and 22-33	<i>MN1</i> and <i>PITPNB</i>
5	22	~26,707,254	26,703,751	22-33 and 22-34	<i>PITPNB</i> and <i>KIAA1043</i>
6	22	~26,757,000	26,757,747	22-34 and 22-35	<i>KIAA1043</i>
7	X	71,850,607–72,272,647	71,507,552	X-56 and X-57	<i>PHKA1</i> and <i>DXS227</i>
8	X	~72,963,000	73,440,225	X-57 and X-58	XIST locus (STS Marker SWXD66)

Table 2
Predictions of replication timing sites for each of chromosomes 1–22, X and Y for further experimental validation

Chromosome	Coordinates	Segmentation point	Adjoining isochores	Band	Adjacent gene in UCSC Genome Browser
1	174,139,182	P199	1-200 and 1-201	1q25.1	<i>RFWD2</i>
2	72,967,307	P101	2-101 and 2-102	2p13.2	<i>SPR</i>
3	75,442,188	P94	3-94 and 3-95	3p12.3	<i>CNTN3</i>
4	99,560,557	P62	4-63 and 4-64	4q23	<i>RAP1GDS1</i>
5	17,705,074	P28	5-28 and 5-29	5p15.1	<i>BASP1</i>
6	140,026,729	P115	6-116 and 6-117	6q24.1	<i>CITED2</i>
7	130,836,857	P119	7-120 and 7-121	7q32.3	<i>PODXL</i>
8	50,008,057	P52	8-53 and 8-54	8q11.21	<i>SNAI2</i>
9	106,576,964	P86	9-87 and 9-88	9q31.1	<i>NIPSNAP3B</i>
10	82,483,369	P92	10-93 and 10-94	10q23.1	<i>SH2D4B</i>
11	75,215,289	P116	11-117 and 11-118	11q13.5	<i>UVRAG</i>
12	115,198,321	P104	12-105 and 12-106	12q24.21	<i>MED13L</i>
13	52,950,362	P15	13-15 and 13-16	13q21.1	<i>OLFM4</i>
14	64,950,457	P38	14-38 and 14-39	14q23.3	<i>FUT8</i>
15	75,574,602	P98	15-98 and 15-99	15q24.3	<i>HMG20A</i>
16	8,447,169	P15	16-15 and 16-16	16p13.2	<i>C16orf68</i>
17	2,894,926	P14	17-14 and 17-15	17p13.3	<i>GARNL4</i>
18	33,052,794	P33	18-34 and 18-35	18q12.2	<i>KIAA1328</i>
19	9,755,583	P17	19-17 and 19-18	19p13.2	<i>LOC162993</i>
20	45,722,781	P52	20-54 and 20-55	20q13.12	<i>SULF2</i>
21	14,419,186	P3	21-4 and 21-5	21q11.2	<i>LIP1</i>
22	38,458,793	P57	22-57 and 22-58	22q13.1	<i>ENTHD1</i>
X	3,868,101	P6	X-6 and X-7	Xp22.33	<i>PRKX</i>
Y	6,816,628	P8	Y-8 and Y-9	Yp11.2	<i>AMELY</i>

density of genes or CpG islands is related to the levels of G + C content of the investigated regions.

3.3. Predictions of replication timing sites for further experimental validation

The distribution of isochores can be displayed intuitively with the cumulative GC-profile (z' curve), which is a discrete function of the nucleotide position n in a genome or chromosome. For the basic characteristics of the cumulative GC-profile, please refer to Section 2 for more details. It is possible that other isochore boundaries, in addition to those coincide with experimentally confirmed ones, are also replication timing sites. Based on the features of the replication timing switch regions confirmed by experiments, we predicted 24 most likely replication timing sites (one site for each chromosome) for future experimental validation. These sites are all boundaries of the isochores, in which the z' curves are approximately straight lines and undergo an abrupt change from a relatively GC-poor (GC-rich) region to a relatively GC-rich (GC-poor) region in each boundary. Take the prediction of replication timing site on chromosome 2 as an example. The segmentation point obtained by GC-Profile is 72,967,307, at which the G + C content of the isochore 2-101 (37%) jumps to that of the isochore 2-102 (50%). There is an *RAB11* family interacting protein coding gene *RAB11FIP5* around the predicted boundary of RT-zones, similar to the switch region on chr17 confirmed by experiment, which is between the gene *NFI* and the gene *RAB11FIP4* in *RAB11* family. In that case, the G + C content of the switch region on chromosome 17 changes from 37% (the isochore 17-74) to 51% (the isochore 17-75). More predictions of replication timing sites for each of chromosomes 1–22, X and Y for further experimental validation can be found in Table 2. These validations are important in that if most of them are proved to be true, the present study will represent a new step towards understanding the complicated replication mechanisms of the human and other mammalian genomes.

Acknowledgements: We would like to thank Drs. Ren Zhang and Wen-Xin Zheng for invaluable assistance. The present work was supported in part by NNSF of China (Grant Nos. 90408028 to C.T. Zhang and 10747150 to F. Gao).

References

- [1] Macaya, G., Thiery, J.P. and Bernardi, G. (1976) An approach to the organization of eukaryotic genomes at a macromolecular level. *J. Mol. Biol.* 108, 237–254.
- [2] Saccone, S., De Sario, A., Wiegant, J., Raap, A.K., Della Valle, G. and Bernardi, G. (1993) Correlations between isochores and chromosomal bands in the human genome. *Proc. Natl. Acad. Sci. USA* 90, 11929–11933.
- [3] Bernardi, G. (1995) The human genome: organization and evolutionary history. *Annu. Rev. Genet.* 29, 445–476.
- [4] Bernardi, G. (2001) Misunderstandings about isochores. Part 1. *Gene* 276, 3–13.
- [5] Fukagawa, T., Sugaya, K., Matsumoto, K., Okumura, K., Ando, A., Inoko, H. and Ikemura, T. (1995) A boundary of long-range G + C% mosaic domains in the human MHC locus: pseudoautosomal boundary-like sequence exists near the boundary. *Genomics* 25, 184–191.
- [6] Tenzen, T. et al. (1997) Precise switching of DNA replication timing in the GC content transition area in the human major histocompatibility complex. *Mol. Cell. Biol.* 17, 4043–4050.
- [7] Schmegner, C., Hameister, H., Vogel, W. and Assum, G. (2007) Isochores and replication time zones: a perfect match. *Cytogenet. Genome Res.* 116, 167–172.
- [8] Cohen, N., Dagan, T., Stone, L. and Graur, D. (2005) GC composition of the human genome: in search of isochores. *Mol. Biol. Evol.* 22, 1260–1272.
- [9] Li, W., Bernaola-Galvan, P., Carpena, P. and Oliver, J.L. (2003) Isochores merit the prefix 'iso'. *Comput. Biol. Chem.* 27, 5–10.
- [10] Costantini, M., Clay, O., Auletta, F. and Bernardi, G. (2006) An isochore map of human chromosomes. *Genome Res.* 16, 536–541.
- [11] Costantini, M. and Bernardi, G. (2008) Replication timing, chromosomal bands, and isochores. *Proc. Natl. Acad. Sci. USA* 105, 3433–3437.
- [12] Costantini, M., Clay, O., Federico, C., Saccone, S., Auletta, F. and Bernardi, G. (2007) Human chromosomal bands: nested structure, high-definition map and molecular basis. *Chromosoma* 116, 29–40.
- [13] Oliver, J.L., Bernaola-Galvan, P., Carpena, P. and Roman-Roldan, R. (2001) Isochore chromosome maps of eukaryotic genomes. *Gene* 276, 47–56.
- [14] Li, W., Bernaola-Galvan, P., Haghghi, F. and Grosse, I. (2002) Applications of recursive segmentation to the analysis of DNA sequences. *Comput. Chem.* 26, 491–510.
- [15] Zhang, C.T. and Zhang, R. (2004) Isochore structures in the mouse genome. *Genomics* 83, 384–394.
- [16] Zhang, C.T. and Zhang, R. (2003) An isochore map of the human genome based on the Z curve method. *Gene* 317, 127–135.
- [17] Zheng, W.X. and Zhang, C.T. (2008) Biological implications of isochore boundaries in the human genome. *J. Biomol. Struct. Dynam.* 25, 327–336.
- [18] Selig, S., Okumura, K., Ward, D.C. and Cedar, H. (1992) Delineation of DNA replication time zones by fluorescence in situ hybridization. *EMBO J.* 11, 1217–1225.
- [19] Zink, D. (2006) The temporal program of DNA replication: new insights into old questions. *Chromosoma* 115, 273–287.
- [20] White, E.J. et al. (2004) DNA replication-timing analysis of human chromosome 22 at high resolution and different developmental states. *Proc. Natl. Acad. Sci. USA* 101, 17771–17776.
- [21] Berezney, R., Dubey, D.D. and Huberman, J.A. (2000) Heterogeneity of eukaryotic replicons, replicon clusters, and replication foci. *Chromosoma* 108, 471–484.
- [22] Zhang, C.T., Gao, F. and Zhang, R. (2005) Segmentation algorithm for DNA sequences. *Phys. Rev. E* 72, 041917.
- [23] Gao, F. and Zhang, C.T. (2006) GC-Profile: a web-based tool for visualizing and analyzing the variation of GC content in genomic sequences. *Nucleic Acids Res.* 34, W686–W691.
- [24] Brylawski, B.P., Cohen, S.M., Horne, H., Irani, N., Cordeiro-Stone, M. and Kaufman, D.G. (2004) Transitions in replication timing in a 340 kb region of human chromosomal R-Band 1p36.1. *J. Cell. Biochem.* 92, 755–769.
- [25] Schmegner, C., Berger, A., Vogel, W., Hameister, H. and Assum, G. (2005) An isochore transition zone in the NF1 gene region is a conserved landmark of chromosome structure and function. *Genomics* 86, 439–445.
- [26] Watanabe, Y., Tenzen, T., Nagasaka, Y., Inoko, H. and Ikemura, T. (2000) Replication timing of the human X-inactivation center (XIC) region: correlation with chromosome bands. *Gene* 252, 163–172.
- [27] Karolchik, D. et al. (2003) The UCSC Genome Browser Database. *Nucleic Acids Res.* 31, 51–54.