

RESEARCH

Open Access



Microsatellite density landscapes illustrate short tandem repeats aggregation in the complete reference human genome

Yun Xia^{1†}, Douyue Li^{1†}, Tingyi Chen^{1†}, Saichao Pan^{1†}, Hanrou Huang^{1†}, Wenxiang Zhang^{1†}, Yulin Liang^{1†}, Yongzhuo Fu¹, Zhuli Peng¹, Hongxi Zhang¹, Liang Zhang¹, Shan Peng¹, Ruixue Shi¹, Xingxin He¹, Siqian Zhou¹, Weili Jiao¹, Xiangyan Zhao¹, Xiaolong Wu¹, Lan Zhou¹, Jingyu Zhou¹, Qingjian Ouyang¹, You Tian¹, Xiaoping Jiang¹, Yi Zhou¹, Shiying Tang¹, Junxiong Shen¹, Kazusato Ohshima² and Zhongyang Tan^{1*}

Abstract

Background Microsatellites are increasingly realized to have biological significance in human genome and health in past decades, the assembled complete reference sequence of human genome T2T-CHM13 brought great help for a comprehensive study of short tandem repeats in the human genome.

Results Microsatellites density landscapes of all 24 chromosomes were built here for the first complete reference sequence of human genome T2T-CHM13. These landscapes showed that short tandem repeats (STRs) are prone to aggregate characteristically to form a large number of STRs density peaks. We classified 8,823 High Microsatellites Density Peaks (HMDPs), 35,257 Middle Microsatellites Density Peaks (MMDPs) and 199,649 Low Microsatellites Density Peaks (LMDPs) on the 24 chromosomes; and also classified the motif types of every microsatellites density peak. These STRs density aggregation peaks are mainly composing of a single motif, and AT is the most dominant motif, followed by AATGG and CCATT motifs. And 514 genomic regions were characterized by microsatellite density feature in the full T2T-CHM13 genome.

Conclusions These landscape maps exhibited that microsatellites aggregate in many genomic positions to form a large number of microsatellite density peaks with composing of mainly single motif type in the complete reference genome, indicating that the local microsatellites density varies enormously along the every chromosome of T2T-CHM13.

Keywords Human genome, Microsatellite density, STRs aggregation, Landscape

Background

Tandem repeat biology is causing a revolution in genetics in the past two decades as many studies of understanding the evolution and biological functions of tandem repeats in the genomes of thousand species and also human diseases [1–4]. Indeed, more than half of the human genome is constituted of repetitive sequences making them the most challenging genomic regions to study, and repetitive sequences are generally classified into interspersed repeats and tandem repeats [5]. The Short Tandem

[†]Yun Xia, Douyue Li, Tingyi Chen, Saichao Pan, Hanrou Huang, Wenxiang Zhang and Yulin Liang co-First author, contributed equally to this work.

*Correspondence:

Zhongyang Tan
zhongyangtan@yeah.net

¹ Bioinformatic Center, College of Biology, Hunan University, Lushan Road (S), Yuelu District, Changsha 410082, China

² Faculty of Agriculture, Saga University, Saga 840-8502, Japan



Repeats (STRs), also called microsatellites, happen ubiquitously in human genome scattered in coding and non-coding regions, which commonly occur with repeat units of 1–6 base pairs and own the highest mutational rate in genome [1, 3, 6]. Microsatellites are increasingly realized to have biological significance in human genome and health, and are reported to involve in more than 30 disorders and several cancers [7–10], regulate gene expression in healthy genomes [11, 12], and also be related to genetic plasticity and missing heritability [13–20]. Slipped-strand mispairing was suggested as a major mechanism for tandem repeat occurrence [21, 22], and we formerly presented a folded slippage model for short tandem repeats occurring mechanism, predicting that micro-disturbing in the process of replication may provide a lot of chances for the template chain folded to produce short tandem repeats, which are possibly selected and fixated for different biological significance in the long-history evolution [23].

Although short tandem repeats have been studied in human genome for several decades, many sequenced columns of human genome so far contain a large number of unsequenced gaps, and these unsequenced gaps are often composed of short tandem repeats [3, 5]. To date, studies about short tandem repeats on complete, gap-free human genomes have been very lacking. The recent assembly of T2T-CHM13 reference removes the gap filled regions of autosomes and Chromosome X [5, 24–28], and over 50% gaps in Chromosome Y [29], therefore it represents a truly complete human genome, this complete reference sequence of human genome T2T-CHM13 will certainly bring great helps for comprehensive study short tandem repeats in human genome. And so far, there is still no unified threshold for microsatellites studies, the repeats in animal genomes are generally considered to be longer than 12 bp [30, 31], and 3% of human genomes was reported consist of microsatellites under that standard [32]. And microsatellites of human sequences comprised of GA/TC/GC/AT bases are investigated using seq-requester microsatellite [24]. However, we applied a threshold of 6, 3, 3, 3, 3, 3 iterations for mono- to hexanucleotide repeat motifs for analyzing microsatellites in the complete reference human genome T2T-CHM13, the threshold was widely used to study microsatellites in the genome sequences of viruses, mitochondrial and chloroplast [2, 33–35]. And we have first applied the threshold to investigate microsatellites in the Y-DNA of the human reference genome (GRCh38, NC_000024.10) at 1 Kbp resolution by Differential Calculator of Microsatellite version 2.0 (DCM 2.0) method, revealing an exact distributional feature of STRs in every local bins of 1 Kbp sequence of the chromosome Y [36]. Herein, we built 24 Microsatellite landscape maps in all 24 chromosomes

of the first complete human genome T2T-CHM13 (CHM13), including the positions and motif types of all classified High Microsatellite Density Peaks (HMDPs), Middle Microsatellite Density Peaks (MMDPs), and Low Microsatellite Density Peaks (LMDPs).

Materials and methods

Genome sequences

The complete 24 chromosome sequences of the human reference genome T2T-CHM13v2.0 were collected from GenBank, and the accession No. of the complete 24 chromosome sequences were listed in Table S2. All chromosome sequences of human reference genome GRCh38.p14 were also obtained from GenBank (Table S2).

STRs identification

The Imperfect Microsatellite Extractor (IMEx 2.1) [37] was applied to identify STRs in all the sequences of T2T-CHM13v2.0 and GRCh38.p14 genomes. The threshold for extracting perfect tandem repeats was set at iterations of 6, 3, 3, 3, 3, 3 for mono-, di-, tri-, tetra-, penta- and hexa- repeat motif respectively.

Calculate local microsatellite density at 1 kb resolution

The local microsatellite density of 24 complete chromosome sequences of T2T-CHM13 were counted by the new promoting program of Differential Calculator of Microsatellites version 3.0 (<https://github.com/zhongyangtan/DCM.git>), which calculate the local STRs density by divided the every chromosome sequence into a large number of differential-units (bins) with size of 1 kb, so microsatellite position-related Differential-unit_{1kb} (D_1) Relative Density (pD_1RD) was calculated for every 1 kb sequence of the 24 chromosome sequences of T2T-CHM13, and also for GRCh38.p14, which was proven to be a better resolution to calculate local STRs density of human genomic sequence [36]. The calculation formula for pD_1RD is:

$$pD_1RD_i = \frac{m_i}{1kb} \times 1000$$

In this formula, pD_1RD_i is the pD_1RD in the i -th 1 kb bins of the genome, m_i is the size (bp) of microsatellite in the i -th 1 kb bins.

Identifying Microsatellite Density Peaks

After the microsatellite density in every 1 kb bins, i.e., pD_1RD , of the genome was calculated. The adjacent bins with similar density ranges were merged, and each merged density bin is referred to as a peak. When the pD_1RD values of each bin in adjacent bins is greater than or equal to 90 and less than 150 ($150 > pD_1RD \geq 90$), those adjacent bins was identify as a Low Microsatellite Density

Peak (LMDP); the pD_1RD values of each bin in adjacent bins is greater than or equal to 150 and less than 300 ($300 > pD_1RD \geq 150$), those adjacent bins was identified as a Middle Microsatellite Density Peak (MMDP); and the pD_1RD values of each bin in adjacent bins is greater than or equal to 300 ($pD_1RD \geq 300$), those adjacent bins was identify as a High Microsatellite Density Peak (HMDP).

Building the STRs density landscape maps

A series of pD_1RD values were obtained by T2T-CHM13 chromosome sequences by the DCM v3.0. Maps displayed the series of pD_1RD values from telomere of p-arm to telomere of q-arm of every chromosome by using the UCSC Genome Browser's bedGraph format. Then, 24 landscape maps of STRs density for T2T-CHM13 were built in the viewer of UCSC Genome Browser window, in which a full chromosomal view of local STRs density is shown for every bin of 1 kb sequence of the chromosome from telomere of p-arm to telomere of q-arm. 24 landscape maps of STRs density for GRCh38.p14 were also built.

Sorting Microsatellite Density Peaks

High Microsatellite Density Peaks (HMDPs) ($pD_1RD \geq 300$), Middle Microsatellite Density Peaks (MMDPs) ($300 > pD_1RD \geq 150$) and Low Microsatellite Density Peaks (LMDPs) ($150 > pD_1RD \geq 90$) were sorted by the new developed program Microsatellite Density Peaks Sorter version 1.0 (MDPS.v1.0, <https://github.com/zhongyangtan/MDPS.git>). And the MDPS.v1.0 also further sorted HMDPs, MMDPs and LMDPs into different Main Motif Type (MMT) and sub-Main Motif Type (sub-MMT). (Table S5, S6).

Mapping Microsatellite Density Peaks

The sorted data of exact position, corresponding motif type and peak name of HMDPs, MMDPs and LMDPs in every chromosome were transferred by the UCSC Genome Browser's bed format and displayed under the landscapes with three tracks of HMDPs, MMDPs and LMDPs in the viewer of UCSC Genome Browser window.

Comparison of HMDPs between T2T-CHM13.v2.0 and GRCh38.p14

The similar HMDPs between T2T-CHM13.v2.0 and GRCh38.p14, were firstly determined allelic position by MUMmer v3.23 [38], then comparing motif type by new developed program Genome HMDPs Comparator version 1.0 (GHC.v1.0, <https://github.com/zhongyangtan/GHC.git>), finally, aligned by ClustalX v2.1 [39].

The division of genomic regions

The full genome of T2T-CHM13 was divided into different Genome Regions according to local microsatellites density features (Table S9). The Genomic Regions with the extremum pD_1RD value ≥ 300 between bins were classified as High Variable Microsatellite Density Region (HVMD-R), those with ($300 >$ the extremum pD_1RD value ≥ 150) between bins were classified as Middle Variable Microsatellite Density Region (MVMD-R), and those with ($150 >$ the extremum pD_1RD value ≥ 90) between bins were classified as Low Variable Microsatellite Density Region (LVMD-R) subclass. The Genomic Regions with the extremum pD_1RD value < 90 between bins were classified as relatively Even microsatellite density (E-) Region class. Genomic Regions clustering with single dominant motif type of peaks were classified as Peak Cluster (PC-) Region class and Telomere Repeat (T-) Region class (located in the telomere regions).

Result

Microsatellite Density landscapes

The recently released complete assemblies of T2T-CHM13 including all 22 autosomes, chromosome X and Y, comprise of 3,117,275,501 bp of DNA; and we obtained 22,198,470 microsatellites with total size of 183,133,984 bp and microsatellites relative density (RD) is 58.75 in the full genome (Table 1, Table S2 & S3). The widely used method of relative density for analyzing microsatellites was proved limited for analyzing STRs in big sequence like human genome; therefore, to explore the exact distributing feature, the microsatellites landscapes at 1 Kbp resolution, were formerly comprehensively surveyed in the Y-DNA of reference human GRCh38 by the differential calculator of microsatellites [36]. Herein, STRs of the complete reference human genome T2T-CHM13 including all chromosomes were investigated by the differential calculate method, in which the values of position related relative density at 1 kilo-base resolution (pD_1RD) were calculated in every bin unit of 1 kb DNA sequence, and the pD_1RD value was suggested to be possibly a better way to estimate the local microsatellite density variation in human genome [36]. The bedgraph format tool of UCSC Genome Browser was used to map the exact STRs distribution features with the pD_1RD data of the full T2T-CHM13 genome, a series of panoramic landscapes of microsatellites density were obtained to precisely exhibit the local microsatellites density in every bin of 1 kb sequence for the 22 autosome, chromosome X and Y of T2T-CHM13 in viewers of Genome browser window (Figs. 1 and 2, Fig. S1, Fig. S2, Table S1). The 24 landscapes revealed that the relative density in every 1 kb sequence (the pD_1RD value) varies



Fig. 1 The landscape of STR Density, HMDP (High Microsatellite Density Peak), MMDP (Middle Microsatellite Density Peak) and LMDP (Low Microsatellite Density Peak) in T2T-CHM13 chromosome 18 with smallest HMDPs. **A.** Landscape of STR Density is displayed on the first track, and the microsatellite density values are shown on the ordinate. **B.** HMDPs are displayed in pack mode on the second track, including its Motif Type, location (shown in colored vertical line) and name. The HMDP name consists of "Genome name", "chr No.", "HMDP No.", "HMDP size" and "sub-Main Motif Type (sub-MMT)". **C-D.** MMDPs and LMDPs are displayed in dense mode on the third and last track. Chromosome 18 with the lowest number of HMDP, is shown here, and other chromosomes of the T2T-CHM13 are shown in Figure S2

enormously in different site along the 24 chromosomes even from 0 to 970, and the STRs is prone to accumulate to a large numbers of different extent microsatellite density peaks; these density peaks are found genome-wide, and were classified into 8,823 high microsatellites density peaks (HMDPs) ($pD_1RD \geq 300$), 35,257 middle microsatellites density peaks (MMDPs) ($300 > pD_1RD \geq 150$) and 199, 649 low microsatellites density peaks (LMDPs) ($150 > pD_1RD \geq 90$) (Table 2, Table S4.01-S4.06); and the bed format tool of UCSC Genome Browser was used to display the exact position and motif types of these HMDPs, MMDPs and LMDPs (Fig. 2, Fig. S2). Similarly, the microsatellites landscapes were comprehensively surveyed in the reference human genome GRCh38 as comparison (Table 2, Table S14 & S15, Fig. S12).

High Microsatellite Density Peaks (HMDPs)

The local relative density values ($pD_1RD \geq 300$) of those HMDPs are approximately six times or more of the average relative density ($RD = 58.75$) of the full T2T-CHM13 genome, the significant statistical bias of microsatellites relative density in local genome region hints their

importance to human genome (Fig. 3A, Fig. S2.01). It was found that all these HMDPs occur genome widely in all the 24 chromosomes of the T2T-CHM13, especially pervade in the two arms of all chromosomes at different intervals, except absent at most centromeric region and some pericentromeric region. The largest quantity of HMDPs was found in ChrY with 2,697 HMDPs, secondly in Chr13, the minimum number of HMDPs in Chr19 with 127 HMDPs identified, and only 305 HMDPs were found on chr1 with longest sequence; therefore, the distribution of HMDPs is not related to chromosome size directly. Furthermore, 89.26% of the HMDPs are identified in intergenic, 10.74% in intron region, but no HMDPs in exon region (Fig. 3B-a, Table S10.01).

Main Motif Types (MMTs) of HMDPs

Analysis of the composition of the motifs in every HMDP demonstrated that same type of motif is prone to accumulate in the same HMDP, and therefore, we classified these HMDPs by the main motif composition in each HMDP; for example, when the AT motif account for the main motif composition in a HMDP, it was named

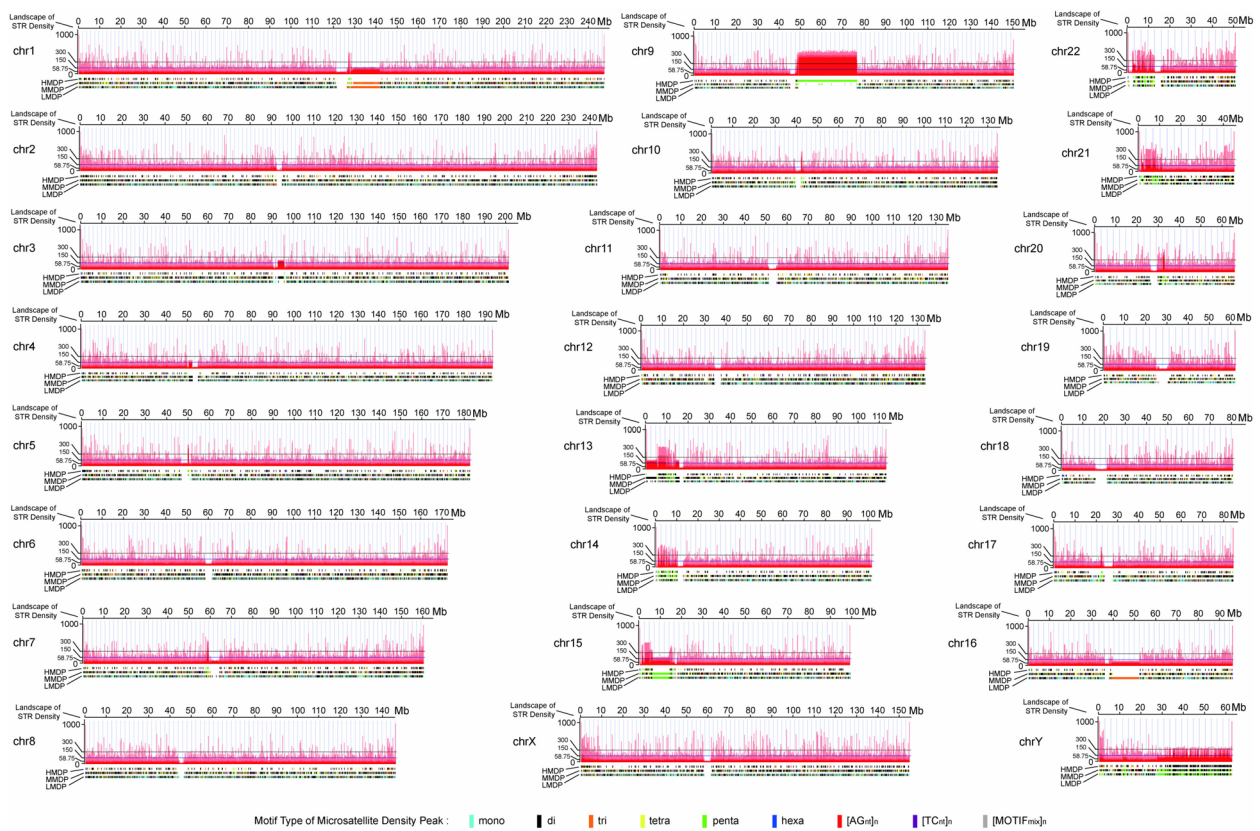


Fig. 2 Landscape map of STR Density of 24 chromosomes of T2T-CHM13. Landscape of each chromosome, first track displayed landscape of STR Density, second track displayed HMDP, third track displayed MMDP, last track displayed LMDP

Table 1 Statistic of microsatellite density landscape information on the human genome T2T-CHM13. The assembly information and microsatellite statistics of T2T-CHM13v2.0 and GRCh38.p14

Genome	Assembled size (bp)	Sequenced size (bp)	Gap number	Gap size (bp)	Microsatellite number	Microsatellite size (bp)	Microsatellite Relative Density (bp/Kbp)
CHM13	3,117,275,501	3,117,275,501	0	0	22,198,470	183,133,984	58.75
GRCh38	3,088,269,832	2,937,639,396	783	150,630,436	19,672,293	151,222,113	51.48

(AT)_n Main Motif Type (MMT) HMDP (Table S4.). And the MMT HMDPs can also be further classified into three sub-types: the high-percentage Main Motif Type (hMMT) HMDPs (main motif ≥ 66.7%), the middle-percentage Main Motif Type (mMMT) HMDPs (66.7% > main motif ≥ 50%), and the low-percentage Main Motif Type (lMMT) HMDPs (50% > main motif ≥ 33.4%) (Table 3, Table S4-S8). All 8,823 HMDPs be classified into Main Motif Type (MMT) class and Mix Motif Type class. Though the microsatellites generally comprising tandem repeats with motifs of 1–6 bp makes that there are possible total 964 motif types (Table S16), we actually identified only 87 MMTs from 8819 HMDPs belong to

MMT class, and these 87 MMTs were further classified into 8 subclasses: mono-, di-, tri-, tetra-, penta-, hexa-, AG- and TC-MMT subclass. The di-MMT subclass includes 6944 HMDPs is the most abundant HMDPs subclass, secondly is the penta-MMT subclass. Notably, 86.1% of the total 8823 HMDPs are the hMMT HMDPs (Fig. 3B-b, Table S10.01), suggesting that high percentage of same repeat motif is prone to accumulated into same high microsatellites density peak (HMDP) in the complete reference human genome T2T-CHM13. Moreover, many MMT HMDPs usually appears in pairs with motif reverse complementary in the complete genome like the example of 386 (AATGG)_n MMT HMDPs reversely

Table 2 Statistic of microsatellite density landscape information on the human genome T2T-CHM13. The list of microsatellites density peaks in T2T-CHM13v2.0 and GRCh38.p14

Chr No	HMDPs		MMDPs		LMDPs		Chr No	HMDPs		MMDPs		LMDPs	
	CHM13	GRCh38	CHM13	GRCh38	CHM13	GRCh38		CHM13	GRCh38	CHM13	GRCh38	CHM13	GRCh38
01	305	199	2,703	1,968	16,891	14,086	13	587	100	1,682	825	7,098	5,976
02	357	294	2,015	1,972	14,802	14,395	14	342	83	1,205	728	6,005	5,405
03	231	180	1,475	1,404	11,798	11,624	15	507	67	2,638	637	7,772	4,968
04	234	230	1,955	1,579	11,889	11,421	16	130	105	1,000	989	7,688	5,509
05	258	198	1,455	1,372	10,580	10,572	17	135	89	880	786	5,771	5,681
06	200	158	1,379	1,365	10,172	10,128	18	125	91	643	645	4,429	4,420
07	249	197	1,581	1,506	10,002	9,978	19	127	105	1,046	956	5,102	4,869
08	188	156	1,291	1,313	8,481	8,352	20	184	75	806	650	3,691	3,753
09	237	121	1,058	997	7,239	7,199	21	449	87	1,088	505	3,009	2,392
10	208	167	1,328	1,203	8,226	8,143	22	378	99	1,039	462	3,074	2,403
11	191	143	1,125	1,060	7,581	7,706	X	313	276	1,673	1,648	9,465	9,470
12	191	155	1,215	1,203	8,470	8,433	Y	2,697	105	3,977	537	10,414	1,965
							Total	8,823	3,480	36,257	26,310	199,649	178,848

complementary pairing with 311 (CCATT)_n MMT HMDPs (Table 3).

[AT]_n^h motifs type are the most abundant HMDPs

Among the all hMMT HMDPs, 4828 HMDPs are [AT]_n^h motif type, representing 63.5% of the di-MMT subclass HMDPs and also 54.7% of all 8823 HMDPs, therefore, it is the most abundant hMMT HMDPs (Table 3, Fig. 3C, Table S10.02). These [AT]_n^h HMDPs were classified into two distributional patterns by average distance (AD): sparse-distributional (AD: 1,434,206 bp) and dense-distributional (AD: 4,476 bp) patterns, there are 1886 [AT]_n^h HMDPs scatter sparsely and 2942 [AT]_n^h HMDPs array densely in the full T2T-CHM13 genome (Fig. 3D-a,b,c, Table S10.03, Fig. S3.01). As a whole, the distributional features can be summed up to 3 [AT]_n^h HMDPs chromosomal distribution models, one model is that only [AT]_n^h HMDPs sparse-distribution appear in the two arms of the 19 chromosomes, another is that [AT]_n^h HMDPs dense-distribution occur in the short p-arm and sparse-distribution in the long q-arm as found in the five acrocentric chromosomes, the other is the [AT]_n^h HMDPs distribution in the Chr Y which 35 [AT]_n^h HMDPs with sparse-distribution are found in Chr Y euchromatin region and 2575 [AT]_n^h HMDPs with dense-distributions in Chr Y heterochromatin region. However, the [AT]_n^h HMDPs is abundantly representing 54.7% of all HMDPs, but 21.3% of all MMDPs and only 2.7% of all LMDPs (Fig. 3D-d), suggesting that AT motif tandem repeats is easy to accumulate at high density in the complete reference genome. And (AT)_n repeats aggregation are reported to related to chromosomal structure and rearrangements [40–42].

The [AATGG]_n^h and [CCATT]_n^h penta-hMMT HMDPs account for about half sequence size of all hMMT HMDPs

There are 773 penta-motif HMDPs identified in the full T2T-CHM13 genome, and ranks secondly in all 8 subclasses of the MMT HMDP class; however, the reverse complementary [AATGG]_n^h and [CCATT]_n^h hMMT HMDPs are found 331 and 252 in numbers accounting for 75.4% of all penta-hMMT HMDPs (Table 3, Table S4 & S5). Though the large majority of the two types are small in size that is lower than 5 bins (5 kbp); the HMDPs with size that is larger than 5 bin, are mainly [AATGG]_n^h and [CCATT]_n^h hMMT HMDPs (Fig. 3C-b & c, Table S10.02); therefore, total size of [AATGG]_n^h and [CCATT]_n^h hMMT HMDPs are 3789 and 3362 bins respectively, and combination of two hMMT HMDPs accounting for about 43.55% of the size of all hMMT HMDPs, implied their importance to human genome structure. Analysis of occurred locations of the [AATGG]_n^h and [CCATT]_n^h hMMT HMDPs reveals 4 distributional models about the two hMMT HMDPs in the T2T-CHM13 genome: (a) no [AATGG]_n^h and [CCATT]_n^h hMMT HMDPs distribution occur in Chr 6, 8, 11, 12, 18, 19 and X, (b) both the two hMMT HMDPs appear in short arm of the five acrocentric chromosomes, (c) pericentric distribution of the two hMMT HMDPs in Chr 1, 2, 3, 4, 5, 7, 9, 10, 17 and Y, (d) the [AATGG]_n^h and [CCATT]_n^h hMMT HMDPs are found on the main arms of the chromosomes (on arm distribution) that are majorly in Chr 1, 2, 16, 20 and 22) (Fig. 3E-a, Table S10.04, Fig. S3.02). Remarkably, Chr 9 is clustered a large segment of [AATGG]_n^h and [CCATT]_n^h hMMT HMDPs at the pericentromeric region of q-arm (Fig. 3E-b), this region well corresponds

Table 3 Statistic of microsatellite density landscape information on the human genome T2T-CHM1.3. The motif types of HMDP in T2T-CHM1.3v2.0

Main Motif Type (MMT) class										
di-MMT subclass	tri-MMT subclass		tetra-MMT subclass		penta-MMT subclass		hexa-MMT subclass		mono-MMT subclass	
	MMT	num ^c	MMT	num	MMT	num	MMT	num	MMT	num
(AT) _n 5063	[AT] _n ^h 4,828	(GGT) _n 64	(CTT) _n 215	(GGCT) _n 6	(AATGG) _n 386	[AATGG] _n ^h 331	(AAAAT) _n 4	(CCCTAA) _n 26	(G) _n 8	
	[AT] _n ^m 188	(ACC) _n 56	(AAAG) _n 130	(AGCC) _n 3	[AATGG] _n ^m 50	[AATGG] _n ^m 50	(ATTT) _n 2	(TTAGGG) _n 25	(C) _n 6	
	[AT] _n ^l 47	(AGG) _n 39	(ATCC) _n 69	(ACTC) _n 3	[AATGG] _n ^l 5	[AATGG] _n ^l 5	(ATAT) _n 2	(CCCTCT) _n 8	(T) _n 1	
(CT) _n 905	[CT] _n ^h 748	(CCT) _n 37	(GGAT) _n 59	(GAGT) _n 3	(CCATT) _n 311	[CCATT] _n ^h 252	(AATAT) _n 1	(AGAGGG) _n 4	3	15
	[CT] _n ^m 95	(CTT) _n 25	(AAGG) _n 54	(ACCC) _n 2	[CCATT] _n ^m 56	[CCATT] _n ^m 56	(CCTC) _n 2	(ACCATC) _n 1		
	[CT] _n ^l 62	(AAG) _n 23	(CCTT) _n 35	(GGGT) _n 1	[CCATT] _n ^l 3	[CCATT] _n ^l 3	(AGAGG) _n 1	(GATGGT) _n 1		
(AG) _n 101	[AG] _n ^h 77	(ATC) _n 12	(AGGG) _n 28	(ACCT) _n 1		(AAAAG) _n 1	(AGATAT) _n 1	(AGATAT) _n 1	AG-MMT subclass	
	[AG] _n ^m 15	(GAT) _n 12	(CCCT) _n 27	(AGGT) _n 1	MMT	num	(CTTTT) _n 1	(GAGGGT) _n 1	MMT	num
	[AG] _n ^l 9	(CCG) _n 3	(AGAT) _n 16	(AATG) _n 3	(ACTCC) _n 10	(AGCCC) _n 2	(AGCC) _n 2	(ATAGGT) _n 1	(AG) _n ^h	33
(GT) _n 443	[GT] _n ^h 414	(CGG) _n 2	(ATCT) _n 10	(ATGT) _n 3	(GGAGT) _n 8	(ATCCT) _n 2	(ATCCT) _n 2	(GGAGGT) _n 1	1	33
	[GT] _n ^m 15	(ACT) _n 3	(CTGT) _n 1	(CTGT) _n 1	(AAAGG) _n 8	(GGGAT) _n 2	(GGGAT) _n 2	(CCAATC) _n 1		
	[GT] _n ^l 14	(AGT) _n 1	21	670	(CCTTT) _n 5	(ACACC) _n 1	(ACACC) _n 1	(ATCCTC) _n 1	TC-MMT subclass	
(AC) _n 431	[AC] _n ^h 394	(AGC) _n 2	(AGGG) _n 2		(AGGG) _n 8	(ATAGT) _n 1	(ATAGT) _n 1	(GGATGT) _n 1	MMT	num
	[AC] _n ^m 25	(ATT) _n 1			(CCCTC) _n 4	(GCCCT) _n 1	(GCCCT) _n 1	(CCCCCT) _n 1	(TC) _n ^h	30
	[AC] _n ^l 12	(GTT) _n 1			(AAGGG) _n 5	(CTGGT) _n 1	(CTGGT) _n 1	14	73	30
(CG) _n 1	[CG] _n ^h 1	15	281		(CCCTT) _n 3	(GGTTT) _n 1	(GGTTT) _n 1	26	773	
	[CG] _n ^m 0									
	[CG] _n ^l 0									
6 ^d	6,944 ^e									
MMTs/HMDPs in MMT class: 87/8819 (mono-: 3/15, di-: 6/6944, tri-: 15/281, tetra-: 21/670, penta-: 26/773, hexa-: 14/73, AG-: 1/33, TC-: 1/30)										
Mix Motif Type (Mix) class										
		Motif type	num							
		Mix	4							HMDPs in Mix class: 4

^{a-c} Main Motif Type (MMT), sub-Main Motif Type (sub-MMT), and HMDP number (num); ^dMMT number in each subclass, ^eHMDP number in each subclass
 The brace represents the sub-MMT corresponding to MMT, Reverse complementary MMT in dotted box
 Total HMDPs: 8823 (HMDPs in MMT class: 8819 (h-: 7600, m-: 749, l-: 407, AG-: 33, TC-: 30); HMDPs in Mix class: 4)

to the heterochromatic regions consisting of Classical Human Satellite III [43, 44].

Comparison of HMDPs between T2T-CHM13 and GRCh38

The microsatellite density landscapes of 24 chromosomes of the reference human genome GRCh38.p14 were also built here, As the reference genome GRCh38 contained gaps and collapsed tandem repeats (Table 1, Table S14), the total sequence assembly of GRCh38.p14 is 2,937,639,396 bp, with unsequenced gap size of 150,630,436 bp. And only 3,480 HMDPs, 26,310 MMDPs and 178,848 LMDPs were identified in the microsatellite density landscapes of GRCh38.p14 genome (Table 2, Table S14); so HMDPs in GRCh38 genome are shown to be greatly different from that of T2T-CHM13, but MMDPs and LMDPs in GRCh38 genome are relatively close to those in T2T-CHM13 (Table 2). A detailed comparison of HMDPs in microsatellite density landscapes between full T2T-CHM13 and GRCh38 full genome, revealed that 1233 and 1290 are high and low-middle similar corresponding to those HMDPs of GRCh38 respectively (Fig. 3F-b,c,d, Table S15, Fig. S12), above 2,000 HMDPs correspond to middle/low microsatellite density peaks (M/LMDPs) and no peaks region of GRCh38, more than half of HMDPs of T2T-CHM13 are corresponding to un-sequenced gaps region of GRCh38 (Fig. 3F-d, Table S15); suggesting that HMDPs alleles variety possibly contribute a lot to the individual diversity of human genome in spite of that the no-sequenced gap region may influence the comparing result.

Genomic regions divided by local microsatellites density features

Because the landscapes graphed by the differential calculating microsatellites density method, the STRs density distribution feature can be visualized in every 1 kb sequence local genomic region; we discovered from the landscapes that the full genome of T2T-CHM13 may be divided into 514 different microsatellites density related Genomic Regions (Table 4, Fig. 4A, Fig. S4, Table S9). These Genomic Regions comprise of 4 Genomic Region class: Variable microsatellite density (V-) Region class,

Even microsatellite density (E-) Region class, Peaks Cluster (PC-) Region class and Telomere repeat (T-) Region class. The V- Region class was further divided into High Variable Microsatellite Density Region (HVMD-R), Middle Variable Microsatellite Density Region (MVMD-R) and Low Variable Microsatellite Density Region (LVMD-R) subclass, the 154 Regions of this class principally correspond to gene coding regions and also some satellite distribution region (Table 4). The E- Region class includes relatively Even and Average Microsatellite Density Region (EAMD-R) subclass and relatively Even and extreme Low Microsatellite Density Region (EeLMD-R) subclass, 67 Regions of this class mainly match α satellites and other satellites. The PC- Region class are subdivided into di-, tri-, tetra-, and penta- subclasses, Regions of this class showed that special motif type of Microsatellite Density Peaks (MDPs) usually gather in proximity to form peaks cluster. And the hexa-nucleotide motif microsatellite cluster Region class was discovered only locating in the telomeric region of T2T-CHM13, the 24 (TTAGGG)_n-Rs were found in the telomere of q-arm of all chromosomes, as well as the reverse complementary 23 (CCCTAA)_n-Rs were found in the telomere of p-arm of all chromosomes except the H(AATGG)_n-R in the p-arm of Chr13 (For simplicity, the Region name abbreviation are only shown here and below, Region name nomenclature please refer to Table S9).

The broad high variable microsatellite density region

The microsatellite density landscape maps display that the density of microsatellites on the main arms of most chromosomes has a huge rhythmic change, showing that HMDPs, MMDPs and LMDPs appear alternately, such regions are summarized as broad High Variable Microsatellite Density Regions (b-HVMD-Rs) with every genomic region size more than 10,000,000 bp (Fig. 4A, Table S9). Total 42 b-HVMD-Rs were classified in the HVMD-R subclass, all chromosomes contain two b-HVMD-Rs in main central part of two arms except that the five acrocentric chromosomes include only one b-HVMD-R in each of their q-arm and ChrY include one in p-arm (Fig. 4B-a, Table S9). Comparing all b-HVMD-Rs with

(See figure on next page.)

Fig. 3 Summary of HMDPs in T2T-CHM13. **A.** An [AT]_n^h HMDP in the T2T-CHM13. **B.** (a) Ratio of HMDP corresponding to intron and intergenic regions. (b) Ratio of hMMT, mMMT and lMMT HMDPs and other HMDPs. **C.** Statistic of HMDP number and bin, (a) The percentage of peak number of mono-, di-, tri-, tetra-, penta- and hexa- hMMT HMDPs in all hMMT HMDPs; (b, c) The number and total bin number of hMMT HMDP in each Peak sequence size (bins) group. **D.** [AT]_n^h HMDP distribution, (a) The 3 [AT]_n^h HMDP distribution models; (b) The distribution map of interval distance of [AT]_n^h Sparse distribution (S-distribution) and Dense distribution (D-distribution) HMDPs, the ordinate indicated the logarithm base 10 of the interval distance (bp), and the values of the interval distance are listed in Table S10.03; (c) The number of [AT]_n^h Sparse distribution and Dense distribution HMDPs; (d) The percentage of [AT]_n^h HMDPs, MMDPs and LMDPs to all HMDPs, MMDPs and LMDPs. **E.** [AATGG]_n^h and [CCATT]_n^h HMDP distribution, (a) The 4 distributional models of [CCATT]_n^h and [AATGG]_n^h HMDPs; (b) A big cluster of [AATGG]_n^h and [CCATT]_n^h motif type HMDPs at the p-arm of Chr9. **F.** Comparison of landscapes T2T-CHM13.v2.0 and GRCh38.p14, (a) Comparison in Chr 1; (b) identity of a pair low similar HMDPs; (c) identity of a pair high similar HMDPs; (d) corresponding HMDPs of the T2T-CHM13 corresponding to MDPs of GRCh38.p14

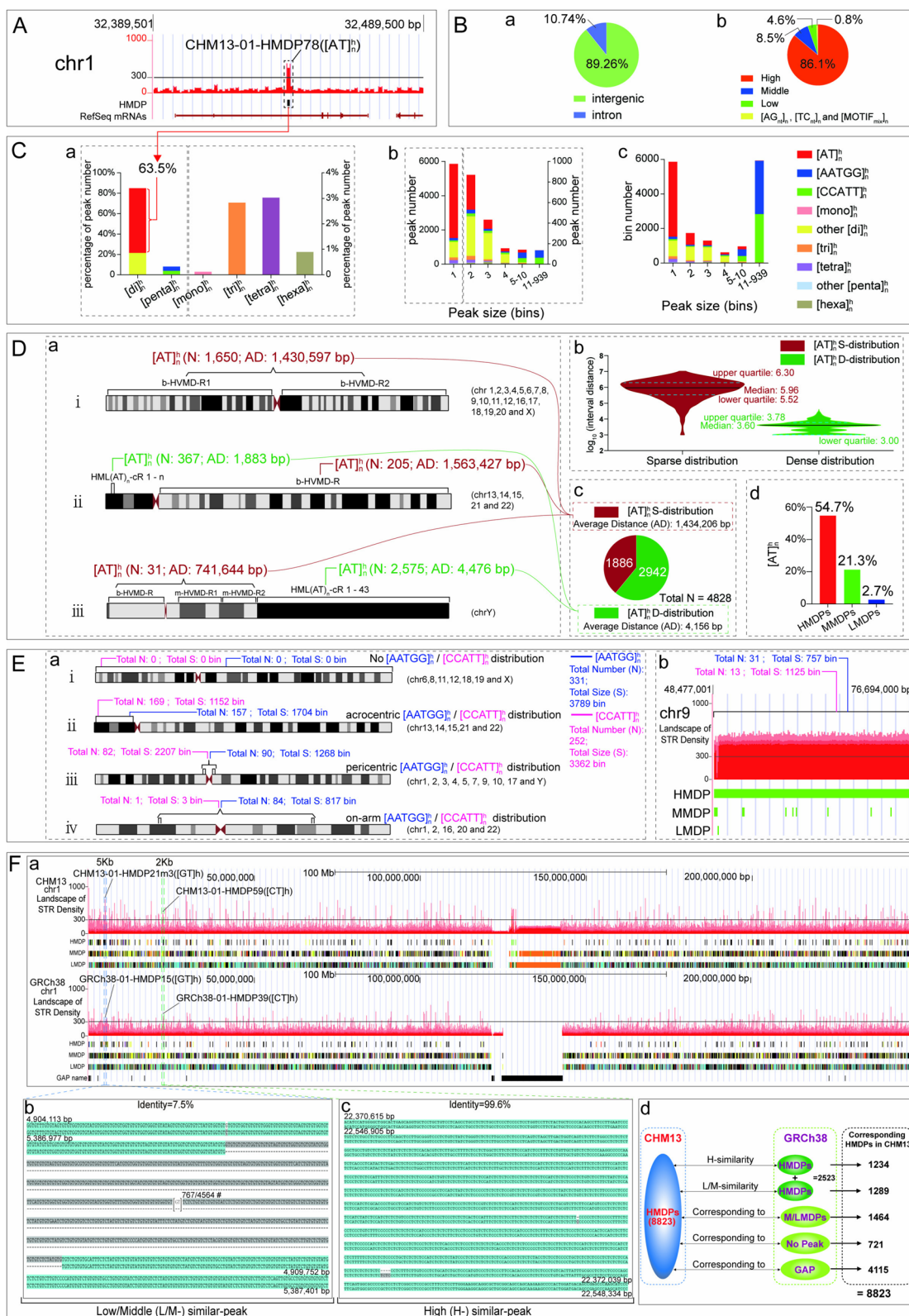


Fig. 3 (See legend on previous page.)

CAT/Liftoff Genes, RefSeq mRNAs and CenSat Annotation tracks, demonstrated that the b-HVMD-Rs well corresponding to broad gene coding regions in every chromosome (Table 4, Table S9). The sum size of the 42 b-HVMD-Rs is 2,850,947,000 bp, representing 91.46% of the full sequence size of T2T-CHM13 genome; but only 4118 of the 8823 HMDPs located in the intergenic and intron region of all b-HVMD-Rs (Fig. 4B-c, Table S10 & S11).

Microsatellite density Peaks Clustering (PC-) regions

In the 246 Genomic Regions of PC- Region class (Table 4, Fig. 4C, Fig. S6-S9); we clarified 70 Regions in the di-subclass, these Regions were classified into 4 Region types as HML(CT[#])_n-cR, HML(AT)_n-cR, M(AT)_n-cR and ML(AT)_n-cR; the 5 HML(CT[#])_n-cR are found in the near central part of short arm of every the five acrocentric chromosomes and correspond to the five rRNA genes coding regions; the 53 HML(AT)_n-cR occur in the five short arm of the five acrocentric chromosomes and the heterochromatin region of ChrY; and the 3 M(AT)_n-cR and 9 ML(AT)_n-cR appear mainly in the five short arm of the five acrocentric chromosomes and centromeric regions of Chr 3 and 4. Secondly, 13 Regions in the tri-subclass were clarified into 4 Region types as ML(ATC)_n-cR, L(ATC)_n-cR, ML(GAT)_n-cR and L(GAT)_n-cR, which are mainly presenting at the pericentromeric regions of chr1, 2, 7, 10, 16, and 17. Then, only 5 Regions were identified in the tetra-subclass as 4 Region types as ML(AAGG)_n-cR, L(AAGG)_n-cR, HML(AGCC)_n-cR and HML(GGCT)_n-cR, occurring at Chr2, 3 and 22. At last, we sorted out 158 Regions in the penta-subclass as 10 Region types; therein 8 Region types included 156 Regions containing the reverse complementary motif of (AATGG)_n or (CCATT)_n, were found at almost all chromosomes except chr6, 8, 11, 12, 16, 19, and X, and this is consistent with the observation of aforesaid no [AATGG]_n^h and [CCATT]_n^h motif type HMDPs distribution; the other two Region types include only two Regions also containing another two reverse complementary motif of (ACTCC)_n and (GGAGT)_n, are only found

at the short arm of the acrocentric chromosome 14 and 15. The comparing with CenSat Annotation tracks in UCSC genome browser reveals that di-, tri-, tetra-, and penta- subclasses mainly correspond to diverse classical human satellite sequences region (Table 4, Table S9).

HML(AT)_n-cR and ML(AATGG)_n-cR alternatively connecting to form the Y heterochromatin

Previously we have constructed the landscapes of the incomplete Y-DNA sequence of human genome GRCh38, which includes many unsequenced gaps especially a big gap of more than 30,000,000 bp in the heterochromatin [36]. Herein, we made the landscape of Y-DNA again but for the first complete sequence of human genome T2T-CHM13 [29], illustrating the local microsatellite relative density feature of the Y heterochromatin (Fig. 4D), in which Genomic Region HML(AT)_n-cR clustering of (AT)_n HMDPs, MMDPs and LMDPs alternatively link to Region ML(AATGG)_n-cR clustering (AATGG)_n MMDPs and LMDPs; therefore, the Y heterochromatin contains 43 HML(AT)_n-cRs and 36 ML(AATGG)_n-cRs alternatively connected at intervals (Fig. 4D), and only 2 ML(CCATT)_n-cRs and 5 L(AATGG)_n-cRs replace ML(AATGG)_n-cRs to separate HML(AT)_n-cRs. The above-mentioned densely distributing 2575 [AT]_n^h HMDPs scatter among the 43 HML(AT)_n-cRs and jaggedly separated by (AATGG)_n MMDPs and LMDPs. Comparing the landscape track with CenSat Annotation track in UCSC genome browser showed that HML(AT)_n-cRs and ML(AATGG)_n-cRs well correspond to alternating pattern of classical human satellite 1B (hsat1B) and classical human satellite 3 (hsat3) in the heterochromatin Yq12 region [29] (Fig. 4D, Fig. S10).

Similar local microsatellites density features in short arms of the five acrocentric chromosomes

The landscapes showed that local microsatellites density features in the five short arm of the five acrocentric chromosomes is more complicate than other autosomes. A palisade arranging HML(CT[#])_n cluster Region situates in central or near central part of the short arm of these

(See figure on next page.)

Fig. 4 Summary of Genomic Regions characterized by microsatellite density in T2T-CHM13. **A.** Chromosome 1 of the T2T-CHM13 characterized into 14 Genomic Regions. **B.** broad High Variable Microsatellite Density Region (b-HVMD-R), (a) example of a b-HVMD-R (chromosome 11, including RefSeq mRNAs and CAT / Liftoff Genes Track). (b) HMDPs corresponding to gene intron and intergenic region; (c) the number of HMDPs in b-HVMD-Rs and other regions. **C.** Feature maps of Peaks Cluster (PC-) Regions (a: di-, b: tri-, c: tetra- and d: penta-subclass). **D.** HML(AT)_n-cRs and ML(AATGG)_n-cRs alternatively connecting in the Chr Y heterochromatin. **E.** Similar local microsatellites density features in the short arms of acrocentric chromosomes (HML(CT[#])_n-cR in center). **F.** Microsatellite density feature of centromeres. The ordinate shown average pD₁RD value, the large font number: the average value of pD₁RD of each centromere (colors showing different value ranges), 58.75: the average pD₁RD of full genome. **G.** The (CT[#])_n HMDPs groups separate rRNA genes in the five HML(CT[#])_n-cR, (a) Partial enlarged view of HML(CT[#])_n-cR in Chr15; (b) Sequence size of Duplication Segment Units (DSUs), rRNA-Groups (rRNA-Gs) and HMDP-Groups (HMDP-Gs) in HML(CT[#])_n-cRs; (c) Comparison of identity of DSUs, rRNA-Gs and HMDP-Gs



Fig. 4 (See legend on previous page.)

Table 4 Statistic of microsatellite density landscape information on the human genome T2T-CHM13. Genomic regions characterized by microsatellite density in T2T-CHM13v2.0

Region class	Region subclass	Region type	Region num	Corresponding annotation	Located chr	Region class	Region subclass	Region type	Region num	Corresponding annotation	Located chr	
V- ^a (154)	HVMD-R (78)	b-HVMD-R	42	Broad GCR; (ct, censat, etc.) ^e	all chr	PC- ^c (246)	tetra- (5)	ML(AAGG) _n -cR	1	censat, ct	22	
		m-HVMD-R	7	Medium GCR; (censat, ct, etc.)	13,14,16,20,21,Y				L(AAGG) _n -cR	1	censat	3
		s-HVMD-R	29	Short GCR; (ct, censat, etc.)	2,9,12-15,21,22,Y				HML(AGCC) _n -cR	2	censat	2
		s-WVMD-R	37	Sparse GCR; (ct, censat, etc.)	1-4,7,10,13-15,17,21,22,Y				HML(GGCT) _n -cR	1	censat	2
E- ^b (67)	LWMD-R (39)	m-LWMD-R	1	bsat	22		penta- (158)	ML(AATGG) _n -cR	44	hsat3	14,15,17,Y	
		s-LWMD-R	38	Sparse GCR; (ct, censat, etc.)	1,3,5-7,10,12-15,17,19,21,22			HML(AATGG) _n -cR	39	hsat3, oSat (mon), hsat2, ct	1-3,7,10,13-15,17,20-22,Y	
		m-EAMD-R	9	oSat (hor, dhor, mon), bsat, ct	1,5-7,10,12,13,16,19			HML(AATGG) _n -cR	10	hsat3, ct	1,9,10,14,15,21,22	
		s-EAMD-R	22	oSat (mon, hor, dhor), ct	3,6,7,14,15,19,21,22,Y			H(AATGG) _n -cR	4	hsat3, censat; (p-TRR) ^f	13,14,21,22	
EeLWMD-R (36)		m-EeLWMD-R	12	oSat (hor, dhor, mon), ct	2,4,8,9,11,14,15,17,18,20,22,X			HML(AATGG) _n -cR	1	hsat3	7	
		s-EeLWMD-R	24	bsat, oSat (hor, mon), censat, ct	3-5,9,10,13-15,21,22,Y			HML(CCATT) _n -cR	44	hsat3, oSat (mon, dhor), etc	1,2,4,7,9,10,13-15,21,22,Y	
		di- (70)	HML(CT ^g) _n -cR	5	rRNA GCR; (ct, censat)	13,14,15,21,22			HML(CCATT) _n -cR	10	hsat3	2,5,7,13,14,15,17,22,Y
			HML(AT) _n -cR	53	hsat18, hsat1A, hsat3, ct, bsat	13,14,15,21,22,Y			ML(CCATT) _n -cR	4	hsat3	15,22,Y
tri- (13)		M(AT) _n -cR	3	hsat1A	3,14,15			M(CTCC) _n -cR	1	hsat3	15	
		ML(AT) _n -cR	9	hsat1A, bsat	3,4,13,14,15,21,22			M(GGAGT) _n -cR	1	hsat3	14	
		L(ATC) _n -cR	5	hsat2, ct	1,2	T- ^d (47)	hexa- (47)	H(CCCTA) _n -R	16	p-TRR; (censat)	1-3,6,7,9-12,16,18-21,X,Y	
		ML(GAT) _n -cR	4	hsat2	2,7,10			H(MCCCTAA) _n -R	7	p-TRR; (censat)	4,5,8,14,15,17,22	
		L(GAT) _n -cR	2	hsat2	1,2,17			H(TTAGGG) _n -R	18	q-TRR ^h ; (censat)	1-7,11-14,17-19,21,22,X,Y	
					7,16			HM(TTAGGG) _n -R	6	q-TRR	8,9,10,15,16,20	
										Total regions	514	

^a Variable microsatellite density (V-) region class; ^b Even microsatellite density (E-) region class; ^c Peaks Cluster (PC-) region class; ^d Telomere repeat (T-) region class; ^e Gene Coding Region (GCR). Outside the parentheses were dominant annotations, inside the parentheses were secondary annotations;

^f Telomere Repeat Region in chromosome p-arm (p-TRR); ^g Telomere Repeat Region in chromosome q-arm (q-TRR). Reverse complementary motif region type in dotted box

chromosomes, a set of HML(CCATT)_n-cR, H(AATGG)_n-cR, M(AT)_n-cR, HML(AATGG)_n-cR array in left side, 2HML(CCATT)_n-cR, one HML(AT)_n-cR and several other Regions lie orderly at the right side of the central HML(CT[#])_n-cR; though these Regions are different in length, they arrange in similar order (Fig. 4E, Fig. S4 & S5). Comparing with CAT/Liftoff Genes, RefSeq mRNAs and CenSat Annotation tracks, revealed that the HML(CT[#])_n-cR well correspond to the rRNA gene cluster Regions. Moreover, there are also other PC- Regions arraying in either side of the HML(CT[#])_n-cR. In a word, the local microsatellites density features arraying in the landscapes are very similar in the five short arms of the acrocentric chromosomes.

Even and Low microsatellites density mainly distribute in the centromere

Although the centromeres are known as comprised of large arrays of tandem repeated alpha satellite [27] and most Genomic Regions of the T2T-CHM13 genome exhibit large variations local STRs density (ie. the pD₁RD value) in the microsatellite density landscape maps, the centromeric regions showed that the feature of local relative microsatellite density is relatively even and low (Fig. 4F, Fig. S10); the average pD₁RD of microsatellite in all centromeric regions is 42 and far lower than average relative microsatellite density with value of 58.75 in the full genome, the lowest average pD₁RD is only 12.3 in centromere of Chr2, the average pD₁RD value of centromeres in 11 chromosomes (chr18, 11, 2, 15, 17, X, 8, 22, 20, 9, 14) are lower than 30, the average pD₁RD value of centromeres in 6 chromosomes (chr19, 1, 6, 12, 13, 21) are between 30 to 42 (the average pD₁RD of all centromeres), and 7 centromeres in chromosomes (chr 4, 3, 7, 5, 10, 16, Y) with the average pD₁RD value higher than 42. Two M(AT)_n-cRs are embedded in the landscapes of centromere in Chr 3 and 4, and cause the average pD₁RD of centromeres rise to 138.4 and 60.9 respectively.

(CT[#])_n HMDPs groups separate rRNA genes in the five acrocentric short arms

As aforesaid, the five short arm of the acrocentric chromosomes contain five HML(CT[#])_n-cRs in their central part, series HMDPs group with (CT)_n as main motif were found clustering palisadingly and followed by several MMDPs and LMDPs, forming a tandem Duplication Segment Unit (DSU), these DSUs connects tandemly composing these HML(CT[#])_n-cRs. The size of these DSUs are between 43,000–47,000 bp considered as rDNA array [24], and there are 76 DSUs found in the short arm of Chr13, 16 DSUs in Chr14, 50 DSUs in Chr15, 56 DSUs in Chr21, 21 DSUs in Chr22, and total 219 DSUs. Further comparison show that every HML(CT[#])_n-cR comprise of

the rRNA genes Groups region and the [CT[#]]_n HMDPs Groups region, remarkably, the [CT[#]]_n HMDPs Groups region well one-by-one separate every the rRNA genes Group in every DSU of the five HML(CT)_n-cR Regions (Fig. 4G-a, Fig. S5). Alignment results showed that the identities among the 219 DSUs, 219 rRNA genes Groups and 219 HMDPs Groups are higher than 0.82, 0.99 and 0.68 respectively (Fig. 4G-b & c, Fig. S5), suggesting the [CT[#]]_n HMDPs Groups region may functionally separate the rRNA genes region but be more variable than rRNA genes region.

Discussion

This work built the first comprehensive genome wide microsatellite density landscape maps of the first complete human reference genome T2T-CHM13 [24, 29]. These landscape maps exhibited that microsatellites aggregate in many genomic positions to form a large number of microsatellite density peaks, and these peaks array along the every human chromosomes arranged like notes of a beautiful piece of music; therefore, the 24 microsatellite density landscapes together look like to form a symphony of human life (Fig. 2). Though transcriptional and epigenetic state of human repeat elements were comprehensive analyzed [5], we mainly focus on the short tandem repeats of human genome, and our works revealed that the microsatellite density landscapes are compatible with human chromosome structure at a high extent, suggesting these landscapes possibly contribute to the further exploration the relationship between STRs aggregation and human genome structure, variation, evolution and so on.

The high microsatellites density peak (HMDP) aggregated predominantly by a single motif, is the most statistically significant short tandem repeats aggregation phenomenon. We identified 8823 HMDPs in T2T-CHM13, much more than the number of 3480 HMDPs in GRCh38; the HMDP number was substantially increased in T2T-CHM13, mainly corresponding to unassembled regions in GRCh38, especially the heterochromatin of Chr Y and the five short arms of the acrocentric Chromosomes. Almost all the HMDPs are single motif dominated, particularly, 86.1% of the 8823 HMDPs are high percentage dominated by a single motif. And there should be total 964 Main Motif Type of HMDPs ranging from mono- to hexa- main motif types, actually, we only identified 87 single motif dominated Main Motif Types (MMTs); these statistical significances are strongly implied the importance of aggregation of short tandem repeats on gene expression, regulation, epigenetics, genetic architecture and evolution of human genome. What is more, the (AT)_n single motif dominated Main Motif Types (MMTs) account for more than half of the total number of HMDPs, and the

reverse complementary motif (AATGG)_n and (CCATT)_n dominated two Main Motif Types (MMTs) account for approximately half of the total bins of HMDPs; suggesting that the three MMT HMDPs are worthy to be further explored for their roles in gene function and chromosome structure of human genome.

In addition, we classified 514 characterized by microsatellite density in the complete reference genome, the Variable microsatellite density (V-) region class almost overlap the gene coding region, the Even microsatellite density (E-) region class mainly cover centromere, pericentromere and p-arm of acrocentric chromosomes, especially, the m-EAMD-R and m-EeLMD-R well correspond to centromeres (Table S9.01); the 246 PC-regions generally locate in pericentromere, p-arm of acrocentric chromosomes and heterochromatin of Chr Y; and the reverse complementary repeat motif (CCCTAA)_n and (TTAGGG)_n telomere repeat regions almost dominate in p-arms and q-arms respectively. Meanwhile, similar genomic regions structural characters were observed in the five acrocentric short arms by the view of microsatellite density features, and alternative microsatellite density regions structure observed in heterochromatin Chr Y [29]. all these indicate that local regional microsatellite density variation may be related to human genome structure. Furthermore, most (microsatellite density) peaks cluster (PC-) regions correspond to different Classical human satellites (hsat); but most the even microsatellite density regions especially in centromere, where microsatellite density are limited to very low and even, correspond to different α satellite higher-order repeats (α Sat hor) (Table 4); suggesting that although microsatellites (STRs) are highly related to satellites (long tandem repeats), the mechanism of their occurrence are likely completely different.

The 219 near identical Duplication Segment Units (DSUs) of approximately 45-kbp that contain segments of encoding 45S rRNAs, are embedded in the short arms of the 5 acrocentric chromosomes [5]; we identified that the HMDPs also well separate every 45S rRNA gene of about 13,300 bp. All the HMDPs were classified here in the intergenic and intron regions, implying that the HMDPs possibly own the basic role of separating genes and exons in the genome and may be very worthwhile to further explore. The first complete reference T2T-CHM13 provided us the opportunity to study short tandem repeat sequences in human genome comprehensively, our results of microsatellite density landscape maps revealed that short tandem repeats are tend to aggregate characteristically throughout the full genome of T2T-CHM13, it may be very helpful to deepen the exploring of the mysterious roles of tandem repeats to human genome's structure, evolution, regulation, variation and also human diseases.

Conclusions

These landscape maps exhibited that microsatellites aggregate in many genomic positions to form a large number of microsatellite density peaks with composing of mainly single motif type in the complete reference genome, indicating that the local microsatellites density varies enormously along the every chromosome of T2T-CHM13.

Abbreviations

STRs	Short Tandem Repeats
HMDPs	High Microsatellites Density Peaks
MMDPs	Middle Microsatellites Density Peaks
LMDPs	Low Microsatellites Density Peaks
MDPs	Microsatellite Density Peaks
DCM 2.0	Differential Calculator of Microsatellites version 2.0
IMEx 2.1	Imperfect Microsatellite Extractor version 2.1
pD ₁ RD	Position-related D ₁ -Relative Density
MDPS.v1.0	Microsatellite Density Peaks Sorter version 1.0
MMT	Main Motif Type
sub-MMT	Sub-Main Motif Type
GHC.v1.0	Genome HMDPs Comparator version 1.0
RD	Relative Density
hMMT	High-percentage Main Motif Type
mMMT	Middle-percentage Main Motif Type
IMMT	Low-percentage Main Motif Type
AD	Average Distance
V-region class	Variable microsatellite density region class
region class	Even microsatellite density region class
PC-region class	Peaks Cluster region class
T-region class	Telomere repeat region class
HVMD-R	High Variable Microsatellite Density Region
MVMD-R	Middle Variable Microsatellite Density Region
LVMD-R	Low Variable Microsatellite Density Region
EAMD-R	Even and Average Microsatellite Density Region
EeLMD-R	Even and extreme Low Microsatellite Density Region
DSUs	Duplication Segment Units
b-HVMD-Rs	Broad High Variable Microsatellite Density Regions
hsat	Human satellites
α Sat hor	α Satellite higher-order repeats

Supplementary Information

The online version contains supplementary material available at <https://doi.org/10.1186/s12864-024-10843-9>.

Supplementary Material 1.

Supplementary Material 2.

Acknowledgements

We thank Yongliang Wu (The Hong Kong Polytechnic University), Jingchen Xie, Junyu Zeng and Wenjing Zhou (Hunan University), Yasha He (South Central University), Puxi Tan (Xi'an Jiaotong University) for helping in data analysis.

Authors' contributions

Z.T. designed and directed this study. Z.T., K.O., Y.X., D.L., T.C., Sc.P., H.H., W.Z. and Y.L. prepared the manuscript. T.C., D.L., Sc.P., H.H., Y.F., W.Z. and Z.T. developed the statistical method and programs. Y.X., D.L., T.C., Sc.P., H.H., W.Z., Y.L., Z.P., H.Z., Li.Z., S.P., R.S., X.H., S.Z., W.J., X.Z., X.W., La.Z., J.Z., Q.O., Y.T., X.J., Y.Z., S.T., J.S. and Z.T. performed the data analysis, K.O. and Z.T. edited this manuscript. All authors read and approved the final version of the manuscript.

Funding

No fund.

Availability of data and materials

The Web version of microsatellite density landscape maps of 24 chromosomes of T2T-CHM13v2.0 can be accessed in (<http://genome.ucsc.edu/s/zhongyangtan/CHM13v2.0>).

The Web version of microsatellite density landscape maps of 24 chromosomes of GRCh38.p14 can be accessed in (<http://genome.ucsc.edu/s/zhongyangtan/GRCh38.p14>).

All data generated or analysed during this study are included in this published article and its supplementary information files.

Declarations**Ethics approval and consent to participate**

Not applicable.

Consent for publication

Not applicable.

Competing interests

The authors declare no competing interests.

Received: 18 May 2024 Accepted: 26 September 2024

Published online: 14 October 2024

References

- Ellegren H. Microsatellites: simple sequences with complex evolution. *Nat Rev Genet.* 2004;5(6):435–45.
- Zhao X, Tian Y, Yang R, Feng H, Ouyang Q, Tian Y, et al. Coevolution between simple sequence repeats (SSRs) and virus genome size. *BMC Genomics.* 2012;13: 435.
- Hannan AJ. Tandem repeats mediating genetic plasticity in health and disease. *Nat Rev Genet.* 2018;19(5):286–98.
- Hannan AJ. Repeat DNA expands our understanding of autism spectrum disorder. *Nature.* 2021;589(7841):200–2.
- Hoyt SJ, Storer JM. From telomere to telomere: the transcriptional and epigenetic state of human repeat elements. *Science.* 2022;376(6588):eabk3112.
- Hartl DL. Molecular melodies in high and low C. *Nat Rev Genet.* 2000;1(2):145–9.
- Kim TM, Laird PW, Park PJ. The landscape of microsatellite instability in colorectal and endometrial cancer genomes. *Cell.* 2013;155(4):858–68.
- Hause RJ, Pritchard CC, Shendure J, Salipante SJ. Classification and characterization of microsatellite instability across 18 cancer types. *Nat Med.* 2016;22(11):1342–50.
- Priestley P, Baber J, Lolkema MP, Steeghs N, de Bruijn E, Shale C, et al. Pan-cancer whole-genome analyses of metastatic solid tumours. *Nature.* 2019;575(7781):210–6.
- van Wietmarschen N, Sridharan S, Nathan WJ, Tubbs A, Chan EM, Callen E, et al. Repeat expansions confer WRN dependence in microsatellite-unstable cancers. *Nature.* 2020;586(7828):292–8.
- Gymrek M, Willems T, Guilmatre A, Zeng H, Markus B, Georgiev S, et al. Abundant contribution of short tandem repeats to gene expression variation in humans. *Nat Genet.* 2016;48(1):22–9.
- Quilez J, Guilmatre A, Garg P, Highnam G, Gymrek M, Erlich Y, et al. Polymorphic tandem repeats within gene promoters act as modifiers of gene expression and DNA methylation in humans. *Nucleic Acids Res.* 2016;44(8):3750–62.
- Verstrepen KJ, Jansen A, Lewitter F, Fink GR. Intragenic tandem repeats generate functional variability. *Nat Genet.* 2005;37(9):986–90.
- Fondon JW 3rd, Hammock EA, Hannan AJ, King DG. Simple sequence repeats: genetic modulators of brain function and behavior. *Trends Neurosci.* 2008;31(7):328–34.
- Hannan AJ. Tandem repeat polymorphisms: modulators of disease susceptibility and candidates for “missing heritability.” *Trends Genet.* 2010;26(2):59–65.
- Nasrallah MP, Cho G, Simonet JC, Putt ME, Kitamura K, Golden JA. Differential effects of a polyalanine tract expansion in *Arx* on neural development and gene expression. *Hum Mol Genet.* 2012;21(5):1090–8.
- Willems T, Gymrek M, Highnam G, Genomes Project C, Mittelman D, Erlich Y. The landscape of human STR variation. *Genome Res.* 2014;24(11):1894–904.
- Willems T, Zielinski D, Yuan J, Gordon A, Gymrek M, Erlich Y. Genome-wide profiling of heritable and de novo STR variations. *Nat Methods.* 2017;14(6):590–2.
- Mallick S, Li H, Lipson M, Mathieson I, Gymrek M, Racimo F, et al. The simons genome diversity project: 300 genomes from 142 diverse populations. *Nature.* 2016;538(7624):201–6.
- Gymrek M, Willems T, Reich D, Erlich Y. Interpreting short tandem repeat variations in humans using mutational constraint. *Nat Genet.* 2017;49(10):1495–501.
- Levinson G, Gutman GA. High frequencies of short frameshifts in poly-CA/TG tandem repeats borne by bacteriophage M13 in *Escherichia coli* K-12. *Nucleic Acids Res.* 1987;15(13):5323–38.
- Schlötterer C, Tautz D. Slippage synthesis of simple sequence DNA. *Nucleic Acids Res.* 1992;20(2):211–5.
- Zhang H, Li D, Zhao X, Pan S, Wu X, Peng S, et al. Relatively semi-conservative replication and a folded slippage model for short tandem repeats. *BMC Genomics.* 2020;21(1):563.
- Nurk S, Koren S, Rhie A, Rautiainen M, Bizkadez AV, Mikheenko A, et al. The complete sequence of a human genome. *Science.* 2022;376(6588):44–53.
- Gershman A, Sauria MEG, Guitart X, Vollger MR, Hook PW, Hoyt SJ, et al. Epigenetic patterns in a complete human genome. *Science.* 2022;376(6588):eabj5089.
- Vollger MR, Guitart X, Dishuck PC, Mercuri L, Harvey WT, Gershman A, et al. Segmental duplications and their variation in a complete human genome. *Science.* 2022;376(6588):eabj6965.
- Altomose N, Logsdon GA, Bizkadez AV, Sidhwani P, Langley SA, Caldas GV, et al. Complete genomic and epigenetic maps of human centromeres. *Science.* 2022;376(6588):eab14178.
- Aganezov S, Yan SM, Soto DC, Kirsche M, Zarate S, Avdeyev P, et al. A complete reference genome improves analysis of human genetic variation. *Science.* 2022;376(6588):eab13533.
- Rhie A, Nurk S, Cechova M, Hoyt SJ, Taylor DJ, Altomose N, et al. The complete sequence of a human Y chromosome. *Nature.* 2023;621:344–54.
- Lei Y, Zhou Y, Price M, Song Z. Genome-wide characterization of microsatellite DNA in fishes: survey and analysis of their abundance and frequency in genome-specific regions. *BMC Genomics.* 2021;22(1):421.
- Qi WH, Yan CC, Li WJ, Jiang XM, Li GZ, Zhang XY, et al. Distinct patterns of simple sequence repeats and GC distribution in intragenic and intergenic regions of primate genomes. *Aging.* 2016;8(11):2635–54.
- Subramanian S, Mishra RK, Singh L. Genome-wide analysis of microsatellite repeats in humans: their abundance and density in specific genomic regions. *Genome Biol.* 2003;4(2): R13.
- de Freitas KEJ, Busanello C, Viana VE, Pegoraro C, de Carvalho VF, da Maia LC, et al. An empirical analysis of mtSSRs: could microsatellite distribution patterns explain the evolution of mitogenomes in plants? *Funct Integr Genomics.* 2022;22(1):35–53.
- Chen M, Tan Z, Zeng G, Peng J. Comprehensive analysis of simple sequence repeats in pre-miRNAs. *Mol Biol Evol.* 2010;27(10):2227–32.
- Sahu BP, Majee P, Singh RR, Sahoo N, Nayak D. Genome-wide identification and characterization of microsatellite markers within the *Avipoxviruses*. *3 Biotech.* 2022;12(5):113.
- Li D, Pan S, Zhang H, Fu Y, Peng Z, Zhang L, et al. A comprehensive microsatellite landscape of human Y-DNA at kilobase resolution. *BMC Genomics.* 2021;22(1):76.
- Mudunuri SB, Nagarajaram HA. IMEX: imperfect microsatellite extractor. *Bioinformatics.* 2007;23(10):1181–7.
- Kurtz S, Phillippy A, Delcher AL, Smoot M, Shumway M, Antonescu C, et al. Versatile and open software for comparing large genomes. *Genome Biol.* 2004;5(2): R12.
- Larkin MA, Blackshields G, Brown NP, Chenna R, McGettigan PA, McWilliam H, et al. Clustal W and Clustal X version 2.0. *Bioinformatics.* 2007;23(21):2947–8.
- Helmrich A, Stout-Weider K, Hermann K, Schrock E, Heiden T. Common fragile sites are conserved features of human and mouse chromosomes and relate to large active genes. *Genome Res.* 2006;16(10):1222–30.
- Irony-Tur Sinai M, Salamon A, Stanleigh N, Goldberg T, Weiss A, Wang YH, et al. AT-dinucleotide rich sequences drive fragile site formation. *Nucleic Acids Res.* 2019;47(18):9685–95.

42. Inagaki H, Ohye T, Kogo H, Yamada K, Kowa H, Shaikh TH, et al. Palindromic AT-rich repeat in the NF1 gene is hypervariable in humans and evolutionarily conserved in primates. *Hum Mutat.* 2005;26(4):332–42.
43. Ramesh KH, Verma RS. Breakpoints in alpha, beta, and satellite III DNA sequences of chromosome 9 result in a variety of pericentric inversions. *J Med Genet.* 1996;33(5):395–8.
44. Starke H, Seidel J, Henn W, Reichardt S, Volleth M, Stumm M, et al. Homologous sequences at human chromosome 9 bands p12 and q13–21.1 are involved in different patterns of pericentric rearrangements. *Eur J Hum Genet.* 2002;10(12):790–800.

Publisher's Note

Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.