

The complete sequence of a human Y chromosome

<https://doi.org/10.1038/s41586-023-06457-y>

Received: 2 December 2022

Accepted: 19 July 2023

Published online: 23 August 2023

 Check for updates

Arang Rhie^{1,53}, Sergey Nurk^{1,51,53}, Monika Cechova^{2,3,53}, Savannah J. Hoyt^{4,53}, Dylan J. Taylor^{5,53}, Nicolas Altemose⁶, Paul W. Hook⁷, Sergey Koren¹, Mikko Rautiainen¹, Ivan A. Alexandrov^{8,9,52}, Jamie Allen¹⁰, Mobin Asri¹¹, Andrey V. Bzikadze¹², Nae-Chyun Chen¹³, Chen-Shan Chin^{14,15}, Mark Diekhans¹¹, Paul Flicek^{10,16}, Giulio Formenti¹⁷, Arkarachai Fungtammasan¹⁸, Carlos Garcia Giron¹⁰, Erik Garrison¹⁹, Ariel Gershman⁷, Jennifer L. Gerton^{20,21}, Patrick G. S. Grady⁴, Andrea Guarracino^{19,22}, Leanne Haggerty¹⁰, Reza Halabian²³, Nancy F. Hansen^{1,24}, Robert Harris²⁵, Gabrielle A. Hartley⁴, William T. Harvey²⁶, Marina Haukness¹¹, Jakob Heinz⁷, Thibaut Hourlier¹⁰, Robert M. Hubley²⁷, Sarah E. Hunt¹⁰, Stephen Hwang²⁸, Miten Jain²⁹, Rupesh K. Kesharwani³⁰, Alexandra P. Lewis²⁶, Heng Li^{31,32}, Glennis A. Logsdon²⁶, Julian K. Lucas^{3,11}, Wojciech Makalowski²³, Christopher Markovic³³, Fergal J. Martin¹⁰, Ann M. Mc Cartney¹, Rajiv C. McCoy⁵, Jennifer McDaniel³⁴, Brandy M. McNulty^{3,11}, Paul Medvedev^{35,36,37}, Alla Mikheenko^{9,38}, Katherine M. Munson²⁶, Terence D. Murphy³⁹, Hugh E. Olsen^{3,11}, Nathan D. Olson³⁴, Luis F. Paulin³⁰, David Porubsky²⁶, Tamara Potapova²⁰, Fedor Ryabov⁴⁰, Steven L. Salzberg⁴¹, Michael E. G. Sauria⁵, Fritz J. Sedlazeck^{30,42}, Kishwar Shafin⁴³, Valery A. Shepelev⁴⁴, Alaina Shumate⁷, Jessica M. Storer²⁷, Likhitha Surapaneni¹⁰, Angela M. Taravella Oill⁴⁵, Françoise Thibaud-Nissen³⁹, Winston Timp⁷, Marta Tomaszewicz^{25,46}, Mitchell R. Vollger²⁶, Brian P. Walenz¹, Allison C. Watwood²⁵, Matthias H. Weissensteiner²⁵, Aaron M. Wenger⁴⁷, Melissa A. Wilson⁴⁵, Samantha Zarate¹³, Yiming Zhu³⁰, Justin M. Zook³⁴, Evan E. Eichler^{26,48}, Rachel J. O'Neill^{4,49,50}, Michael C. Schatz^{5,13}, Karen H. Miga^{3,11}, Kateryna D. Makova²⁵ & Adam M. Phillippy[✉]

The human Y chromosome has been notoriously difficult to sequence and assemble because of its complex repeat structure that includes long palindromes, tandem repeats and segmental duplications^{1–3}. As a result, more than half of the Y chromosome is missing from the GRCh38 reference sequence and it remains the last human chromosome to be finished^{4,5}. Here, the Telomere-to-Telomere (T2T) consortium presents the complete 62,460,029-base-pair sequence of a human Y chromosome from the HG002 genome (T2T-Y) that corrects multiple errors in GRCh38-Y and adds over 30 million base pairs of sequence to the reference, showing the complete ampliconic structures of gene families *TSPY*, *DAZ* and *RBMY*; 41 additional protein-coding genes, mostly from the *TSPY* family; and an alternating pattern of human satellite 1 and 3 blocks in the heterochromatic Yq12 region. We have combined T2T-Y with a previous assembly of the CHM13 genome⁴ and mapped available population variation, clinical variants and functional genomics data to produce a complete and comprehensive reference sequence for all 24 human chromosomes.

The human Y chromosome plays critical roles in fertility and hosts genes important for spermatogenesis, as well as *SRY*, the mammalian sex-determining locus⁶. However, in the human reference genome, GRCh38, the Y chromosome remains the most incomplete chromosome with over 50% of bases represented by gaps. These multimegabase gaps have persisted for decades and represent sequences flanking the endogenous model centromere, parts of the ampliconic regions and large heterochromatic regions. The architecture of the Y chromosome, specifically the presence of large, tandemly arrayed and inverted repeats (IRs) (palindromes)¹, makes assembly difficult and hinders the study of rearrangements, inversions, duplications and deletions in several critical regions such as Azoospermia factor-a (AZFa), AZFb

and AZFc (azoospermia factors that are linked to clinical phenotypes, including infertility)⁷.

Following the first complete assemblies of chromosomes X⁸ and 8 (ref. 9), the Telomere-to-Telomere (T2T) consortium successfully assembled all chromosomes of the CHM13 cell line⁴. This first complete human genome assembly was enabled by innovative technological improvements in the generation of Pacific Biosciences (PacBio) high-fidelity reads (HiFi)¹⁰ and Oxford Nanopore ultralong reads (ONT)¹¹, the development of better assembly algorithms for utilization of HiFi reads and generation of assembly graphs¹², the use of ONT reads for better graph resolution¹³, new methods for validating and polishing^{14–18} and a coordinated curation effort to finish the assembly.

A list of affiliations appears at the end of the paper.

Table 1 | Comparison of GRCh38-Y and T2T-Y

		GRCh38-Y	T2T-Y	%Δ
Assembly	Total bases	57,264,655	62,460,029	+9.1
	Assigned bases	57,227,415	62,460,029	+9.1
	Unlocalized bases	37,240	0	
	No. of gaps	56	0	
	No. of N bases	30,812,366	0	
Annotation	No. of genes	589	693	+17.7
	Protein coding	66	106	+60.6
	No. of additional genes	6	110	
	Protein coding	1	41	
	No. of transcripts	681	883	+29.7
	Protein coding	372	488	+31.2
	No. of additional transcripts	4	206	
	Protein coding	4	120	
Ampliconic gene copy numbers	BPY2	4 (3, 0)	4 (3, 0)	0
	CDY	26 (4, 0)	26 (4, 0)	0
	DAZ	4 (4, 0)	4 (4, 0)	0
	HSFY	8 (2, 0)	8 (2, 0)	0
	PRY	8 (2, 0)	8 (2, 0)	0
	RBMY	32 (6, 4)	34 (6, 4)	+3.3
	TSPY	25 (7, 0)	66 (45, 0)	+164.0
	VCY	2 (2, 0)	2 (2, 0)	0
	XKRY	8 (0, 2)	8 (0, 2)	0
Haplogroup	Haplogroup	R-L20 (R1b1a2a1a2b1a1)	J-L816 (J1a2b3a1)	
	Ancestry	European	Ashkenazi Jewish	
Repetitive bases	SINE	2,625,350	4,385,917	+67.1
	Retroposon	18,506	18,500	-0.0
	LINE	6,378,323	6,456,888	+1.2
	LTR	4,604,368	4,613,537	+0.2
	DNA/Rolling-circle	2,626,425	4,387,030	+67.0
	Satellite	1,578,773	14,522,636	+819.9
	Simple repeat	1,124,311	21,568,381	+1,818.4
	Other	705,062	972,612	+37.9
	All repeat classes	17,501,283	53,004,524	+202.9
	% Repetitive (non-N)	66.3	84.9	+28.1
Accessible with short reads		13,785,359	14,363,623	+4.2

Annotation statistics for GRCh38-Y are taken from the RefSeq (v.110) annotation, and T2T-Y statistics from a lifted and curated combination of RefSeq (v.110) and GENCODE (v.35) annotations. Numbers of additional genes/transcripts are those found exclusively in one assembly compared with the other. Ampliconic gene copy numbers are shown as X(Y,Z) where X is the total number of annotated genes, Y protein-coding genes and Z transcribed pseudogenes. %Δ is percentage change from GRCh38-Y to T2T-Y. Blank spaces indicate not applicable.

Having been derived from a complete hydatidiform mole, CHM13 has a 46,XX karyotype but is almost entirely homozygous. This simplified assembly of its genome but prevented assembly of a Y chromosome.

In parallel, with the goal of including broader genomic diversity across populations¹⁹, the Human Pangenome Reference Consortium (HPRC) has evaluated various methods for the generation of high-quality diploid genome assemblies²⁰ using a well-characterized human genome, HG002, which has previously been assembled²¹ and is commonly used for benchmarking by the Genome in a Bottle (GIAB) consortium²². Using this rich set of data, and integrating the lessons learned from assembling CHM13, we successfully reconstructed the complete sequence of the HG002 Y chromosome, hereafter referred to as T2T-Y.

Here we analyse the composition of the newly assembled pseudo-autosomal regions (PARs), ampliconic and palindromic sequences, centromeric satellites and q-arm heterochromatin of a complete Y chromosome. We have annotated T2T-Y and combined it with the previous

T2T-CHM13 assembly to form a new, complete reference for all human chromosomes, T2T-CHM13+Y. To enable the use of this new reference sequence we have lifted over available variation datasets from ClinVar²³, GWAS²⁴, dbSNP²⁵ and gnomAD²⁶. In addition, we have recalled variants from the 1000 Genomes Project (1KGP)²⁷ and Simons Genome Diversity Panel (SGDP)²⁸ data, as well as epigenetic profiles from ENCODE data²⁹. These experiments demonstrate improved mappability and variant calling for XY individuals when using T2T-Y as a reference.

Assembly and validation of T2T-Y

Assembly of the HG002 Y chromosome followed the strategy used for the T2T-CHM13 genome⁴ (Supplementary Table 1 and Supplementary Fig. 1). We used PacBio HiFi reads (60× haploid genome coverage) and ONT ultralong reads (90× in reads longer than 100 kb) generated from HG002. An assembly string graph was first constructed for the whole

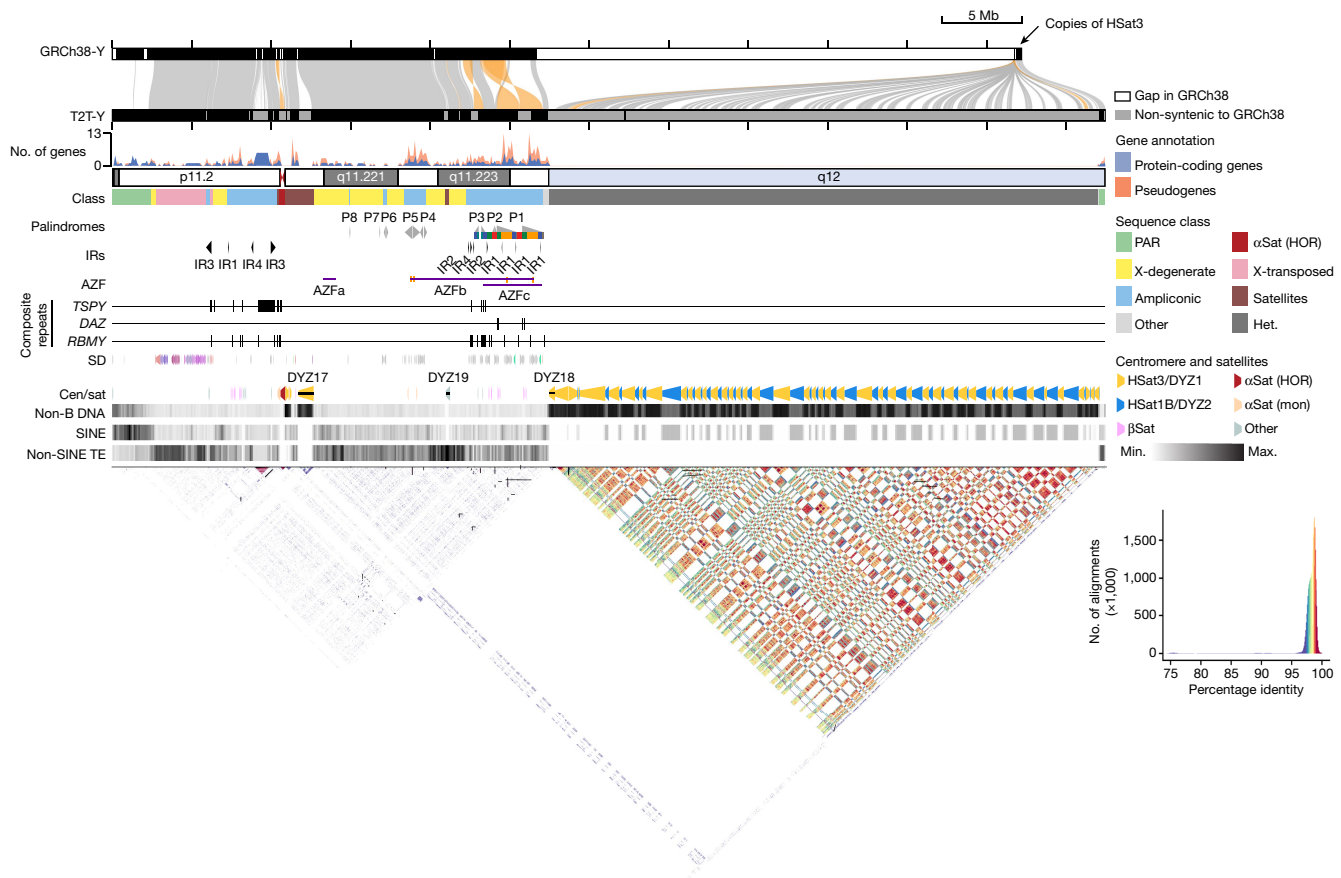


Fig. 1 | Structure of a complete Y chromosome. From top to bottom, alignment of GRCh38-Y and T2T-Y. Regions with sequence identity over 95% are connected and coloured by alignment direction (grey, forward; orange, reverse). Gene density plot shows enriched protein-coding genes in ampliconic sequences. Sequence class, palindromes, IRs and AZFa–AZFc are annotated. Composite repeat arrays are named after the contained ampliconic genes. Segmental duplications (SDs) are coloured by duplication types defined in DupMasker⁷⁴. Centromere (cen) and satellite (sat) annotations highlight the alternating

HSat1 and HSat3 pattern comprising Yq12. Non-B DNA track shows that regions forming alternative sequence structures are enriched in centromeric and satellite repeats. SINE, including *AluY*, are highly enriched in PAR1. All other non-SINE TEs are found only in the euchromatin. All repeats within T2T-Y are visualized by StainedGlass⁷⁵, with similar repeats coloured by percentage identity in the style of an alignment dotplot. Het., heterochromatic; mon, monomeric.

HG002 genome using PacBio HiFi reads. The ChrX and ChrY string graph components shared connections to one another at the PARs, but to no other chromosomes in the genome and could be independently analysed (Extended Data Fig. 1a). The remaining tangles in these XY subgraphs were resolved using ONT reads (Extended Data Fig. 1b). ChrX and ChrY chromosomal walks were identified using haplotype-specific k-mers from parental Illumina reads (Extended Data Fig. 1c), and a consensus sequence was computed for each. Pseudoautosomal region 1 (PAR1) was enriched for GA-microsatellites, which reduced HiFi coverage in this region and led to a more fragmented graph (due to a known HiFi sequencing bias¹²). These gaps were manually patched using a de novo assembly of trio-binned parental ONT reads¹⁴.

The ChrY draft assembly was further polished and validated using sequencing reads from Illumina (66× haploid genome coverage), HiFi (84×) and ONT (250×). During four rounds of polishing, 1,520 small and ten large (over 50 base) errors were detected and corrected (Extended Data Fig. 2a). Conservatively filtered long-read alignments identified two potential assembly issues remaining in the satellite (HSat) arrays around positions 40 and 59.1 Mb, and Strand-seq^{30,31} identified one inversion error within palindromic sequence P5 around position 18.8 Mb (Extended Data Fig. 2b,c, Supplementary Table 2 and Supplementary Figs. 2–4). The validation signal at the two HSat positions was ambiguous, and the P5 inversion appears as a true recurrent inversion³², so these regions were noted but left uncorrected in this release. The remaining

sequences showed no signs of collapse or false duplication, with even HiFi coverage (mean $39.3 \times \pm$ s.d. 12.5 on ChrXY) except for regions associated with known sequencing biases¹⁷, all of which had supporting ONT coverage (reads over 25 kb, mean $78.1 \times \pm$ s.d. 13.6 on ChrXY). The base error is estimated at less than one error per 10 Mb (Phred Q73.8, Supplementary Table 3). Mapped HiFi and ONT reads from the paternal HG003 genome are also consistent with the HG002 T2T-Y assembly, suggesting that no large, structural variants were introduced during cell line immortalization and culture (Supplementary Fig. 5).

The resulting T2T-Y assembly is 62,460,029 bases in length with no gaps or model sequences, showing the previously uncharacterized 30 Mb (approximately) of sequence within the heterochromatic region of the q-arm (Table 1). In comparison, ChrY in the human reference genome (GRCh38-Y) consists of two sequences, with the longer sequence totalling 57.2 Mb (NC_000024.10) and for which 53.8% (30.8 Mb) of the bases are unresolved gaps. The shorter GRCh38-Y sequence (NT_187395.1) is 37.2 kb in length, not placed in the primary Y assembly and has been omitted from most previous genomic studies. The PAR1 (2.77 Mb) and PAR2 (329.5 kb) sequences in GRCh38-Y are duplicated from ChrX rather than assembled de novo, and the centromere is represented by a 227 kb model sequence. Direct sequence comparison between T2T-Y and GRCh38-Y yields an average sequence identity of around 99.8% in the alignable regions, but with multiple structural differences, including an incorrectly oriented centromere model for GRCh38-Y (Fig. 1 and

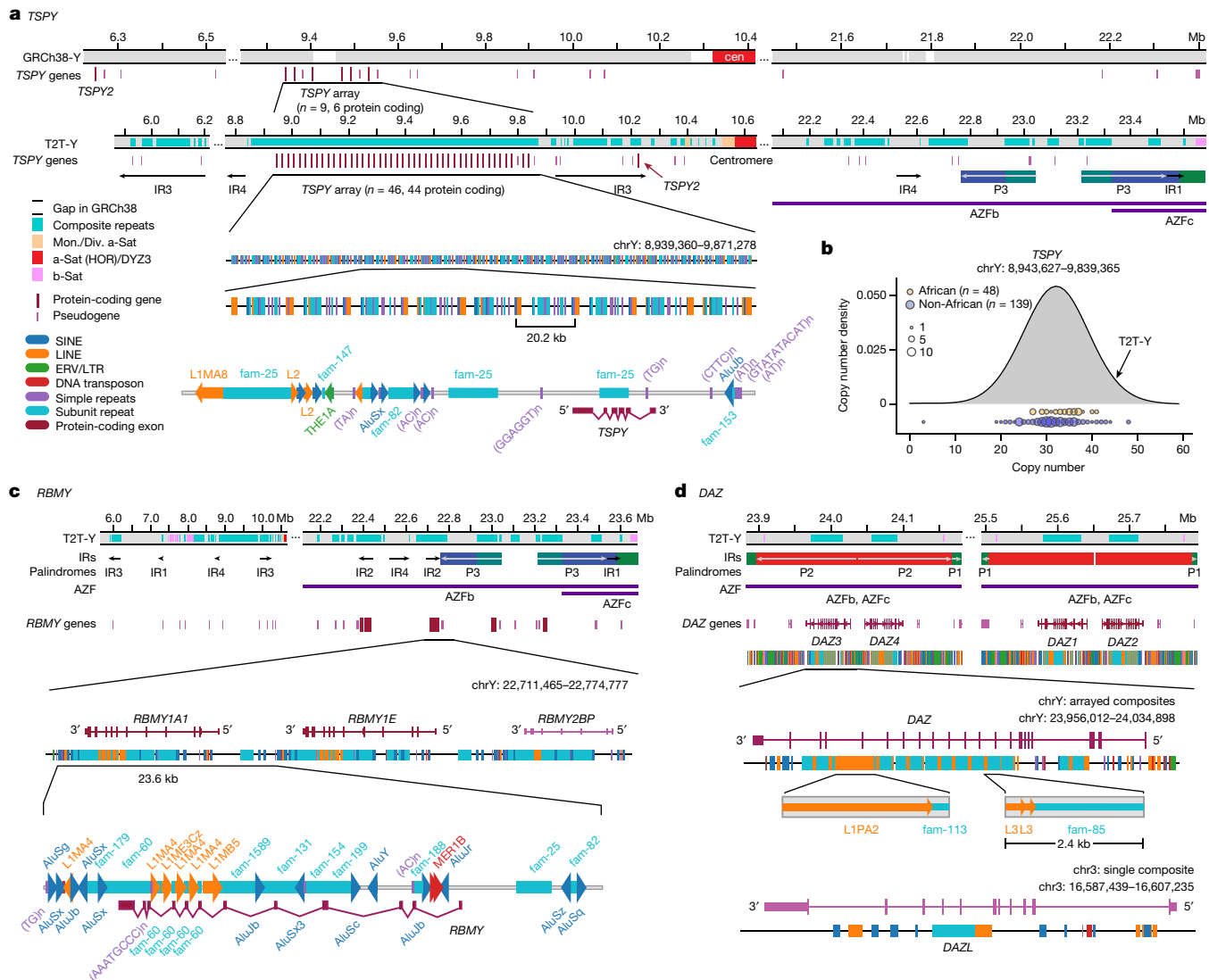


Fig. 2 | Ampliconic genes forming composite repeats. **a**, T2T-Y has 44 *TSPY* protein-coding genes, organized in a single continuous array and a single *TSPY2* copy, compared with GRCh38-Y which has a gap in the *TSPY* array. T2T-Y shows a more regularized array and recovers additional *TSPY* pseudogenes not present in GRCh38-Y. **b**, Copy number differences of *TSPY* protein-coding copies found

in the SGDP. **c**, Repeat composition of the *RBMY* gene family. **d**, Repeat composition of the *DAZ* gene family, with one extra copy annotated on Chr3 that is missing L1PA2. Whereas *TSPY* and *RBMY* genes are found within repeat composites forming arrays, *DAZ*-associated composites are embedded within the introns of the gene.

Extended Data Fig. 3). We identified the Y-chromosome haplogroup of HG002 as J-L816 (J1) and that of GRCh38 as R-L20 (R1b). These haplogroups are most commonly found among Ashkenazi Jews³³ and Europeans³⁴, respectively, consistent with the established ancestry of these genomes. T2T-Y was combined with the T2T-CHM13v1.1 assembly to create a new Y-bearing reference, T2T-CHM13v2.0, referred to here as T2T-CHM13+Y.

Comprehensive annotation of the Y Gene annotation

We annotated T2T-CHM13+Y by mapping RefSeq (v.110) and GENCODE (v.35) annotations from GRCh38 and performed hand-curation of the ampliconic gene arrays (Fig. 1 and Supplementary Tables 4 and 5). NCBI RefSeq and EBI Ensembl generated additional de novo annotations using HG002 full-length cDNA sequencing (Iso-seq) transcriptomes from B lymphocyte and induced pluripotent stem cell lines, as well as tissue-specific expression data from other publicly available sources (Supplementary Table 1 and Supplementary Figs. 6 and 7).

Our annotation of T2T-Y totals 693 genes and 883 transcripts, of which 106 genes (488 transcripts) are predicted to be protein coding (Table 1 and Supplementary Table 4). In addition to containing all genes annotated in GRCh38-Y, T2T-Y contains an additional 110 genes, among which 41 are predicted to be protein coding. The majority of these protein-coding genes (38 of 41) are additional copies of *TSPY*, one of the nine ampliconic gene families, filling the corresponding gap in GRCh38-Y (Table 1). The annotated ampliconic gene copies in T2T-Y were largely concordant with copy numbers estimated from Illumina reads and droplet digital PCR (ddPCR)³⁵, confirming the accurate copy number representation of the ampliconic genes in T2T-Y (Supplementary Tables 6–9). RNA sequencing data confirmed expression of the annotated ampliconic genes in testis³⁶. Only six genes differed in their annotation between GRCh38-Y and T2T-Y, due to presumed Y haplogroup differences (Supplementary Table 10).

Repeat annotation

We generated comprehensive repeat annotations, incorporating repeat models previously updated with CHM13 (ref. 37), as well as 29 previously

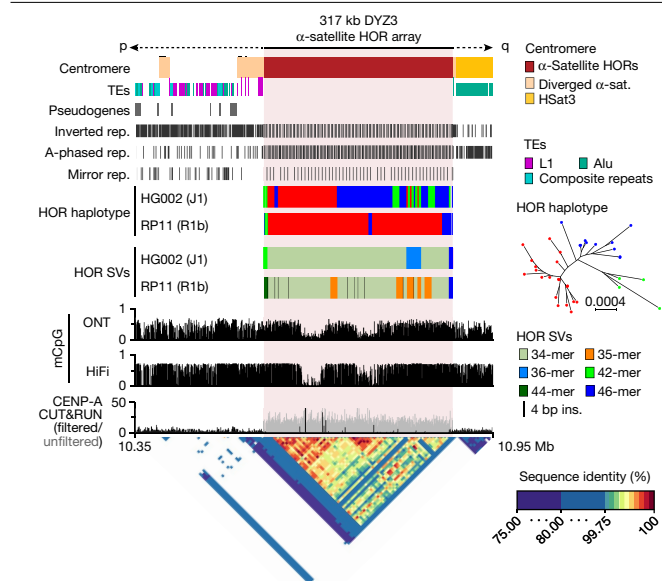


Fig. 3 | Structure of the T2T-Y centromere. No TEs were found within the DY33 array, whereas L1s (upstream) and *Alus* (downstream) were found within the diverged alpha-satellites (drawn taller than the other TEs). A periodic non-B DNA motif pattern is shown within the HOR array. The HG002-Y (T2T-Y) HOR haplotypes and SVs show a different long-range structure and organization compared with a previously assembled centromere from RP11-Y⁴⁵. Three major HOR haplotypes were identified in HG002-Y based on their phylogenetic distance (red, blue and green). RP11-Y has no 36-mer variants but does have a number of 35-mers containing internal duplications. Histograms show the fraction of methylated CpG sites called by both ONT and HiFi, with two hypomethylated CDRs supported by CENP-A binding signal from CUT&RUN⁷⁶. A StainedGlass dotplot illustrates high similarity within the HOR array (99.5–100%).

unknown repeats identified in T2T-Y (Extended Data Fig. 4a and Supplementary Table 11). The newly added sequences increased the percentage of identifiable repeats on the Y chromosome from 66.3 to 84.9%, or 17.5 Mb of non-N bases in GRCh38-Y compared with 53 Mb of bases in T2T-Y (Table 1, Supplementary Tables 12 and 13 and Supplementary Fig. 8). Whereas short-interspersed nuclear elements (SINEs)—specifically *Alus*—are found embedded as part of the human satellite 1 (HSat1) units across most of the q-arm, other transposable elements (TEs: long-interspersed nuclear elements (LINEs), long-terminal repeats (LTRs), SINE-VNTR-*Alus* (SVAs), DNA transposons and Rolling circles) are completely absent (Fig. 1). Moreover, TE distribution biases typify different subregions of ChrY because *Alus* are enriched in the PAR1 region whereas other TEs (particularly L1s) are more abundant in the X-transposed region (XTR)¹ (Extended Data Fig. 4b,c and Supplementary Table 14). The DY33 region is annotated by RepeatMasker entirely as LTRs (Extended Data Fig. 4c), but further analyses indicate that this is a satellite array spanning 265 kb whose 125-base monomeric consensus is derived from an expanded portion of a LTR12B sequence³⁸. Repeat discovery and annotation of T2T-Y also allowed for improved annotation of ChrX in both HG002 and CHM13, particularly in the PAR regions, adding about 33 kb of satellite annotations per ChrX (Supplementary Table 15).

In addition, we searched for TE-driven transductions mediated by L1s and SVAs. We detected six potential 3' L1 transductions within T2T-Y yet no SVA-driven DNA transductions (Supplementary Table 16). Despite a genome-wide investigation of both T2T-CHM13+Y and GRCh38, we were not able to locate any potential donor elements, which confirms a previous analysis that found no evidence for DNA transduction between the Y and the remainder of the genome³⁹. The transduction rate in T2T-Y was also much lower (0.096 per 1 Mb) than that observed in the CHM13

autosomes (average 6.9 per 1 Mb) and ChrX (10.19 per 1 Mb)³⁷ (Supplementary Note 1).

In T2T-Y we identified a total of 825,526 repetitive sequence motifs capable of forming alternative DNA structures (non-B DNA) compared with only 138,640 in GRCh38-Y (Extended Data Fig. 5, Supplementary Table 17 and Supplementary Note 2). This nearly sixfold increase is largely attributed to our use of new and improved experimental and computational methodology because non-B DNA motifs, which might form structures during sequencing, are notoriously difficult to sequence through⁴⁰. We found a particular enrichment of these motifs at the newly sequenced centromeric region (see below) and heterochromatic region on the Yq arm (Fig. 1).

Amplificonic genes in composite repeats

Composite repeats are a type of segmental duplication that are typically arranged in tandem arrays, probably derived through unequal crossing over that contributed to their increased copy numbers^{1,35}. The *TSPY*, *RBMV* and *DAZ* amplificonic gene families are all associated with composite repeats on the Y chromosome, and the T2T-Y assembly provides an opportunity to analyse the complete structure of these arrays (Fig. 2).

TSPY contains the largest number of protein-coding copies on the Y chromosome and is expressed only in testis. Expression level of this gene is dosage dependent and the copy number is polymorphic between individuals⁴¹. In GRCh38-Y, the *TSPY* array includes a 40 kb gap and a limited number of intact protein-coding copies. Our T2T-Y assembly resolved 45 protein-coding *TSPY* copies, including *TSPY2*, which was found downstream of the *TSPY* array in the distal part of the proximal inverted-repeat IR3, in contrast to GRCh38-Y where it is located upstream, possibly due to translocations between the IR3 pairs. The distal positioning of *TSPY2* in HG002 was confirmed among all other Y haplogroups except R and Q, which match the proximal positioning of GRCh38-Y³². All 44 protein-coding copies in the *TSPY* array are embedded in an array of composite repeat units (roughly 20.2 kb in size, matching previous reports^{1,41}), with one composite unit per gene (Fig. 2a and Supplementary Table 18). Each unit includes five new repeat annotations (fam-*), several retro-elements in the LINE, SINE and LTR classes and simple repeats. This 931 kb array is the largest gene-containing composite repeat array in the T2T-CHM13+Y assembly outside of the ribosomal DNA locus, and the third largest overall (the first being the rDNA arrays followed by an LSAU-BSAT composite array on Chr22 (ref. 37)).

Data from 187 SGDP samples confirmed high *TSPY* sequence conservation, but copy number varied from 10 to 40 (Fig. 2b). Phylogenetic analysis, using protein-coding *TSPY*s from a Sumatran orangutan (*Pongo abelii*) and a Silvery gibbon (*Hylobates moloch*) as outgroups, confirmed that all protein-coding *TSPY* copies (including *TSPY2*) originated from the same branch, which is separated from the majority (all but one) of *TSPY* pseudogenes (Extended Data Fig. 6). This result contradicts earlier findings, which concluded that *TSPY2* originated from a different lineage⁴².

The composite structure of *RBMV* is similar to that of *TSPY* (one composite unit per gene), is comparable in size (with *RBMV* at 23.6 kb) and includes LINEs, SINEs, simple repeats and eight new repeat annotations (Fig. 2c). By contrast, the *DAZ* locus is structured such that the entire repeat array, consisting of 2.4 kb composite units each containing a new repeat annotation and a fragmented L3, falls within one gene annotation (Fig. 2d). Out of the three composite arrays described here, *DAZ* is the only one also found on an autosome (Chr3, *DAZL*)⁴³, although as a single unit and lacking the young LINE1 (L1PA2) insertion of the ChrY *DAZ* copies.

Centromere

Normal human centromeres are enriched for an AT-rich satellite family (roughly 171 base monomer), known as alpha-satellite, typically

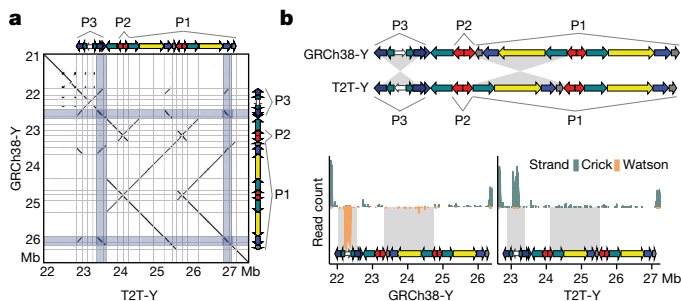


Fig. 4 | Comparison of the palindromic structure of the P1–P3 region. **a**, GRCh38-Y and T2T-Y alignment dotplot and schematics of the palindromes. Frequently recombining IRs in the AZFc region are highlighted in light blue. Deletion of AZFc between IRs is known to cause spermatogenic failure⁴⁹. Self-dotplot of T2T-Y with AZFb and AZFc annotation is available in Supplementary Fig. 15. **b**, Top, schematic of the palindromes. Two inversions were found, one in P3 and one between P1 and P2. Bottom, Strand-seq signal from HG002 confirming the inverted orientation of P3 and P1 in T2T-Y compared with GRCh38-Y.

arranged into higher-order repeat (HOR) structures and surrounded by more diverged alpha- and other satellite classes⁴⁴. Each HOR copy is nearly identical and comprises a tandemly arrayed set of monomers. We annotated 366 kb of alpha-satellite in T2T-Y, spanning 317 kb of the DY23 HOR array. Whereas the individual units within the HOR array are highly similar (99.5–100%), three HOR subtypes were identified from the full-length repeat units based on their monomer structure (red, blue and green HOR haplotypes in Fig. 3, Supplementary Figs. 9–13, Supplementary Tables 19 and 20 and Supplementary Data 1–3). The majority of the T2T-Y centromeric array is composed of 34-mer HORs with a small expansion of a 36-mer, and with longer HOR variants observed in the flanking p-arm (42-mer) and q-arm (46-mer). These variants are structurally different from the RP11 centromere, which is the basis for the GRCh38-Y centromere model and was recently finished by ONT sequencing⁴⁵ (Fig. 3).

Methylated CpG sites called by both HiFi and ONT reads show two adjacent regions of hypomethylation (separated by approximately 100 kb) in the centromeric dip region (CDR) (Fig. 3), which has been reported to coincide with CENP-A binding and is the putative site of kinetochore assembly⁴⁴. In the T2T-Y centromere, the presence of two distinct hypomethylated dips per chromatin fibre was confirmed by inspection of single-molecule ONT reads (Supplementary Fig. 14). A similar pattern of multiple methylation dips within a single centromere was observed in other T2T-CHM13 chromosomes, including Chr11 and Chr20 (ref. 46). In addition, the HORs contained abundant inverted, A-phased and mirror repeat motifs forming a periodic pattern occurring every 5.7 kb (Fig. 3 and Supplementary Table 17). Such non-B DNA motifs—IRs in particular, potentially forming cruciforms—are hypothesized to play a functional role in defining human Y centromeres⁴⁷, and their presence is confirmed here at the sequence level.

Sequence classes and palindromes

We annotated sequence classes on T2T-Y as ampliconic, X-degenerate, X-transposed, pseudoautosomal, heterochromatic and other, in accordance with Skaletsky et al.¹ In addition, we were able to classify a more precise annotation for the satellites (including DY217 and DY219) and the centromere (Fig. 1 and Supplementary Table 21). The X-degenerate and ampliconic regions were estimated to be 8.67 and 10.08 Mb in length, respectively, in concordance with previous findings¹. The T2T-Y ampliconic region contains eight palindromes, with palindromes P4–P8 highly concordant with GRCh38-Y (that is, in terms of arm, spacer length and sequence identity). Arm-to-arm identity of these five T2T-Y palindromes nested within X-degenerate regions

ranged from 99.84 to 99.96% (Supplementary Tables 22 and 23). Palindromes P1–P3 harbour the AZFc region, which contains genes critical for sperm production⁴⁸. We discovered a large polymorphic inversion (over 1.9 Mb) in respect to GRCh38-Y that probably arose from a single, non-allelic homologous recombination event. Using Strand-seq, we were able to locate the breakpoints at two ‘red’ amplicons (naming according to Kuroda-Kawaguchi et al.⁴⁹): one forming the P2 palindrome and the other inside the P1 palindrome (Fig. 4). This rearrangement was previously annotated as the ‘gr/rg’ (green–red/red–green) inversion with variable breakpoints and was confirmed to be present across six Y-chromosome haplogroups out of 44 genealogical branches⁵⁰. Another inversion was detected in P3, which was recently reported as a recurrent variation in human⁵¹ (Extended Data Fig. 7a). Although inversions between amplicons are believed to serve as substrates for subsequent AZFc deletions and duplications that might affect sperm production^{50,52–54}, pinpointing the breakpoints and measuring the frequency of polymorphic inversions proved difficult because of the large size and high identity of the palindromic arms.

Composition of the q-arm heterochromatin

The human Y chromosome contains a large heterochromatic region at the distal end of the q-arm (Yq12), which consists almost entirely of two interspersed satellite sequences classically referred to as DY21 and DY22 (refs. 55–58). The single largest gap in GRCh38-Y is at Yq12, with minimal representation of DY21 and DY22, mostly in unplaced scaffolds. Here we uncovered the detailed structure of the Yq12 region at single-base resolution, characterizing over 20 Mb of DY21 and 14 Mb of DY22 repeats. In T2T-Y, DY21 and DY22 are interspersed in 86 large blocks, with DY21 blocks ranging from 80 to 1,600 kb (median 370 kb) and DY22 blocks ranging from 20 to 1,200 kb (median 230 kb). DY22 blocks appear more abundant at the distal end of Yq12, and this trend is also visible in metaphase chromosome spreads with fluorescence in situ hybridization (FISH) (Fig. 5a,b and Fig. 1, cen/sat track). Yq12 is highly variable in size and sequence structure between individuals^{59–61}, and both number and size of these satellite blocks are expected to vary considerably.

DY21 is composed of a Y-specific subfamily of HSat3 sequences that occurs primarily as nested tandem repeats (of about 3.6 kb) derived from an ancestral tandem repeat of the pentamer CATTC⁶². DY22 is composed of an unrelated satellite family, HSat1B, and comprises a tandem repeat of roughly 2.5 kb made up of three parts: an ancient *AluY* subunit (20% diverged from the *AluY* consensus), an extremely AT-rich region (over 85% AT) and a more genomic copy (GC)-rich region^{58,62,63}. The vast majority of repeat instances were over 98% identical, with slightly higher divergence at the more peripheral satellite blocks (Fig. 5c). A detailed comparison of the sequences within T2T-Y showed recent structural rearrangements, including iterative, tandem duplications as large as 5 Mb spanning multiple blocks of DY21 and DY22 (Extended Data Fig. 8). These structural rearrangement patterns are consistent with evolution by unequal exchange mechanisms. In addition, approximately 15% of Strand-seq libraries showed sister chromatid exchanges within the Yq12 heterochromatic region (Extended Data Fig. 7b).

Whereas HSat3 is present across multiple chromosomes, including the acrocentric short arms, HSat1B is present almost exclusively on the Y and the acrocentric short arms in smaller amounts, with the exception of Chr10 (ref. 64). Whereas HSat1B carries an *AluY*-derived subunit as part of its composite repeat unit, some HSat3 arrays are tightly associated with *Alu* sequences, with blocks of HSat3 intermingled with *Alu* fragments, including *AluY*. Phylogenetic analyses place the ChrY HSat1B *AluY* subunits in a cluster with *AluY* subunits found in HSat1B sequences on the acrocentric chromosomes, with the highly homogenized ChrY copies appearing as a single cluster. Given the topology of this tree, it appears that the HSat1B sequences found on the acrocentric chromosomes were derived from the Y-linked HSat1B, with seeding events leading to local expansion and homogenization

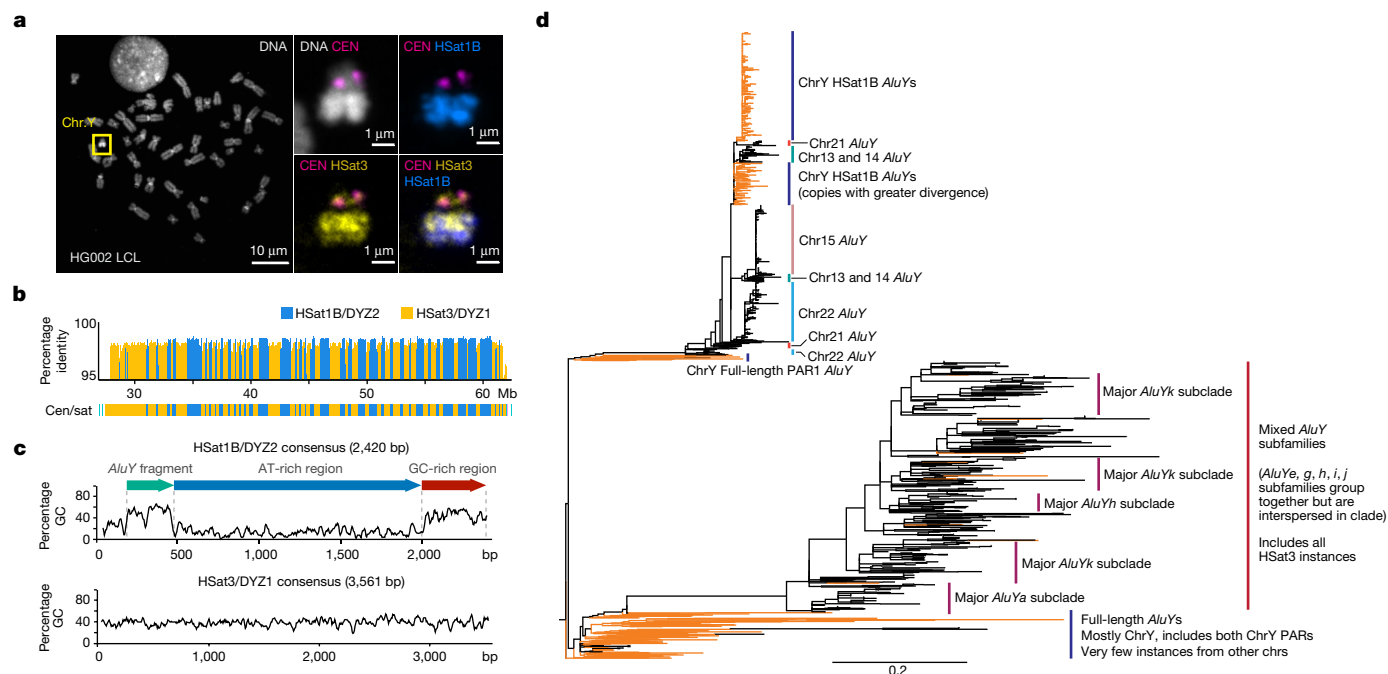


Fig. 5 | Heterochromatic region of Yq12. a, FISH painting of the Y chromosome, centromere/DYZ3 (magenta), HSat1B (blue) and HSat3 (yellow). Top left, overall chromosome labelling by DNA dye (DAPI), with ChrY highlighted in an HG002-derived lymphoblastoid cell line (GM24385). Middle and right, ChrY labelled with FISH probes recognizing centromeric alpha-satellite/DYZ3 (magenta), HSat3/DYZ1 (yellow) and HSat1B/DYZ2 (blue). In concordance with the T2T-Y assembly, HSat3 probes indicate the presence at DYZ17 (close to centromere) as well as a slight enrichment to the proximal part of Yq12 (DYZ1), whereas HSat1B is present only in Yq12 and is more enriched towards the distal part (DYZ2). Maximum-intensity projections are shown in all panels.

The results of this experiment were replicated using two different sets of PCR probes. Fifteen large-field images containing at least 20 spreads were analysed per condition. **b**, Percentage identity of each DYZ2/DYZ1 repeat unit to its consensus sequence. **c**, Percentage GC sequence composition of HSat1B/DYZ2 and HSat3/DYZ1 repeat units and position of an ancient *AluY* fragment in DYZ2. **d**, Phylogenetic tree of *AluY* sequences associated with HSat1B and HSat3, rooted on *AluSc8*. Tree represents subsampling of *AluY* elements, both full length and truncated, including *AluY* sequences found within HSat1B units and associated with HSat3 arrays. Elements located on ChrY are denoted by orange branches. Scale bar, 0.2 substitutions per site on a branch of the same length.

(Fig. 5d, upper branches). The *AluY* fragments found interspersed with HSat3 on the Y chromosome also phylogenetically cluster with *AluY* fragments associated with HSat3 on the acrocentric chromosomes. However, there is no evidence for local homogenization of HSat3-*Alu* fragments; likewise, there is no support for phylogenetic clustering by either subgroup or chromosome. Based on the deep divide between HSat1B and HSat3 clades in the tree for both ChrY and the acrocentric chromosomes, we conclude that the initial seeding events that created these arrays were independent of one another yet were derived from *AluY* elements from PAR1 (Fig. 5d, lower branches).

Improved variant calling for XY samples

We performed short-read alignment and variant calling for 3,202 samples (1,603 XX; 1,599 XY) from the IKG Phase 3, including 1,233 unrelated XY samples averaging at least 30 \times coverage of 150 bp paired-end reads²⁷. This set of 1,233 XY samples spans all 26 geographically diverse IKG populations and 35 distinct Y-chromosome haplogroups (Supplementary Table 24). To more accurately represent the diploid nature of the PARs, we completely hard-masked ChrY in XX samples and ChrY PARs in XY samples, thereby forcing any reads originating from the ChrY PAR to align to the ChrX PAR (Supplementary Tables 25–28 and Extended Data Fig. 9). Diploid genotypes were then called within the PAR for both XX and XY samples⁶⁵ (Extended Data Fig. 10a). Aside from this modification, the alignment and variant-calling pipeline mirrored our previous analysis based on GRCh38-Y⁶⁶.

Across all 1,233 unrelated XY samples we observed improved alignment to T2T-Y, including a higher number of mapped reads (increase of 1.4 million reads on average, s.d. 432,115; Fig. 6a), a higher proportion of

properly paired reads (increase of 1.4% on average, s.d. 1.4%; Fig. 6b) and a lower proportion of mismatched bases (decrease of 0.6% on average, s.d. 0.06%; Fig. 6c) per sample relative to GRCh38-Y (Supplementary Table 29). Within syntenic regions of the two Y chromosome assemblies, the number of variants per sample declined for samples from all Y haplogroups with the exception of haplogroup R (haplogroup of GRCh38-Y), with the greatest reduction observed for samples of haplogroup J1 (haplogroup of T2T-Y; Fig. 6d and Extended Data Fig. 10b,c). Selecting one individual each from the J1, R1b and E1b haplogroups, we compared per-variant read depth and allele balance for both references (Fig. 6e). In all three samples we observed more variants with excessive read depth and abnormal allele balance on GRCh38-Y, corresponding to putative collapsed duplications (Supplementary Table 30 and Fig. 6f,g). We replicated these analyses using an additional 279 samples across 142 populations from the SGDP²⁸ and found similarly improved mappings and variant discovery using T2T-Y (Extended Data Fig. 10d,e and Supplementary Figs. 16–18).

Due to genomic repeats, accuracy of short-read variant calling is heterogeneous across the genome. One approach towards improving reliability is to restrict analysis to ‘accessible’ regions based on various alignment metrics. To this end, we followed published protocols to generate a short-read accessibility mask for T2T-Y based on patterns of normalized read depth, mapping quality and base-calling quality⁶⁷. Our masks show that, although the heterochromatic long arm (Yq12) remains largely inaccessible to short-read analysis, T2T-Y still adds 578 kb of accessible sequence compared with GRCh38-Y (increase of 4.2%; Table 1).

Taken together, these analyses indicate that the complete T2T-Y reference improves short-read alignment and variant calling across

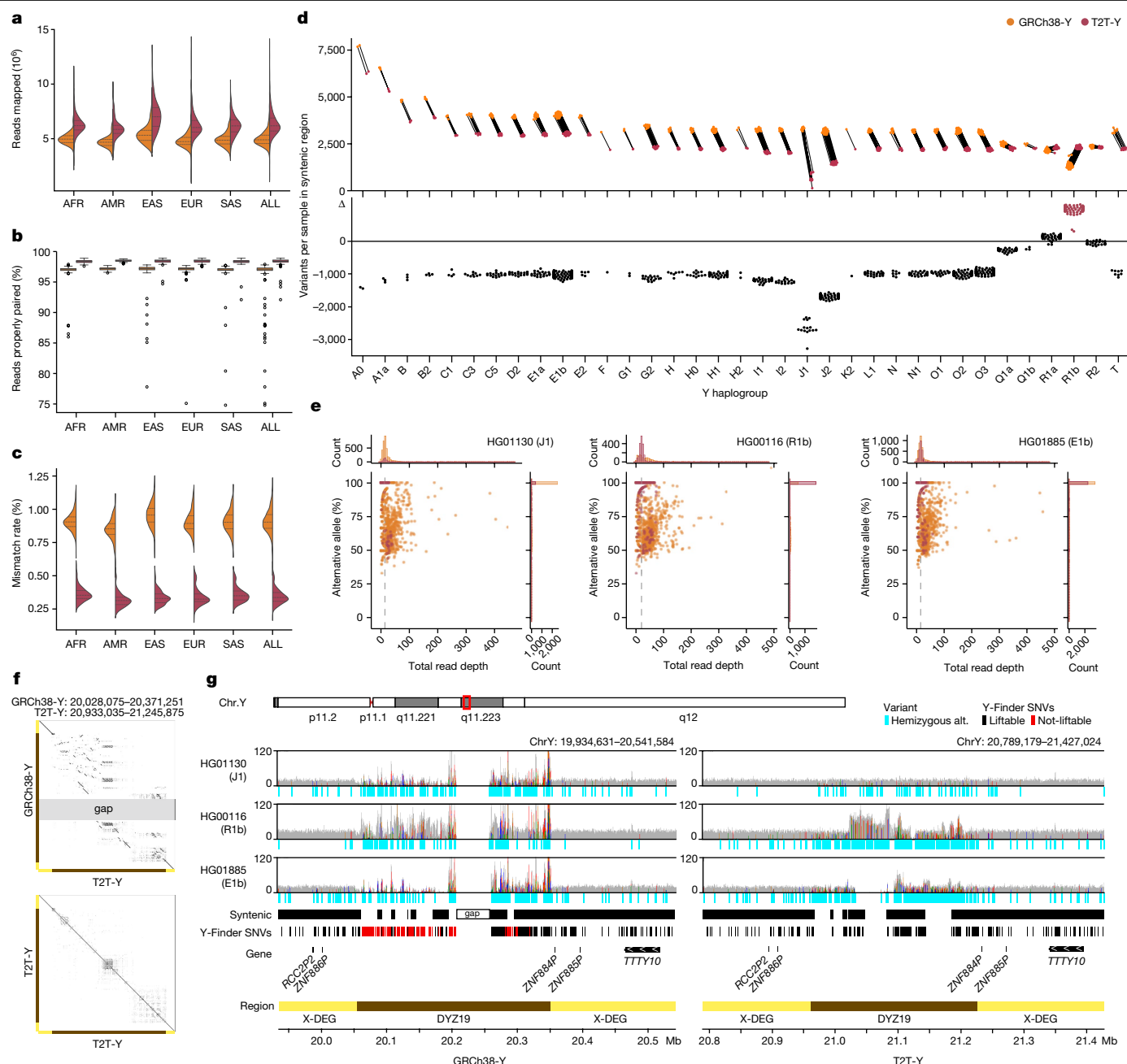


Fig. 6 | Short-read mappability and variant-calling improvements on T2T-Y. In all plots (a–g), GRCh38-Y is orange and T2T-Y is maroon. **a–c**, The complete sequence of T2T-Y improves short-read alignment of the 1KGP dataset by increased number of reads mapped (**a**), higher proportion of reads properly paired (**b**) and lower mismatch rate compared with GRCh38-Y (**c**). Results are grouped by 1KGP super-population code: AFR, African; AMR, American; EAS, East Asian; EUR, European; SAS, South Asian; ALL, all five super-populations. In **b**, bars represent the first, second (median) and third quartiles of the data and whiskers are bound to the 1.5× interquartile range. Data outside of whisker ranges are shown as dots. **d**, Number of called variants within syntenic regions is reduced on T2T-Y for all haplogroups except R1 (haplogroup of GRCh38-Y). **e**, Further investigation of three samples (J1, R1b and E1b) shows a higher number of variants called with excessive read depth and variable alternative allele fractions for GRCh38-Y. Each dot represents a variant, with the percentage

of alternate alleles as a function of total read depth. Dashed line represents the median coverage on T2T-Y, close to the expected one-copy coverage. **f**, Dotplot of the DY19 array between GRCh38-Y and T2T-Y, and self-dotplot of T2T-Y. Large rearrangements are observed, with multiple inversions proximal to the gap in GRCh38-Y with respect to T2T-Y (top) whereas more identical, tandem duplications are visible in T2T-Y (bottom). **g**, Read pile-ups and variants on DY19 for GRCh38-Y (left) and T2T-Y (right) as shown by IGV⁷⁷. Grey histogram shows the mapped read coverage, with coloured lines indicating non-reference bases with over 60% allele frequency. Regardless of haplogroup, the incomplete DY19 array in GRCh38-Y hinders interpretation. Syntenic regions between the two Ys are marked and single-nucleotide variation (SNV) sites used to identify Y haplogroup lineages in Y-Finder are shown below, with variants liftable from GRCh38-Y to T2T-Y in black, not-liftable in red.

populations and corrects errors in GRCh38-Y; however, acknowledging the rich resources available on GRCh38, we also curated a one-to-one whole-genome alignment between each GRCh reference (GRCh37 and GRCh38) and T2T-CHM13+Y to enable lifting annotations in either

direction. The vast majority of genetic variants in ClinVar (13 March 2022 release), dbSNP (build 155) and GWAS Catalog (v.1.0 release) were successfully lifted to T2T-CHM13+Y (99.2/97.8/98.9% overall and 100/95.0/100% for ChrY, respectively; Supplementary Table 31).

Accessibility masks and lifted annotations are provided, along with variant calls as a resource for future studies (Data Availability).

Contamination of genomic databases

Human DNA sequences can sometimes appear as contaminants in the assembled genomes of other species. In microbial studies, the human reference sequence has been used to screen out contaminating human DNA; however, due to the incomplete nature of the current reference, some human fragments are missed and mistakenly annotated as bacterial proteins, leading to thousands of spurious proteins in public databases^{68,69}. For example, a recent analysis of nearly 5,000 human whole-genome datasets found an unexpected linkage between multiple bacterial species and human samples of XY karyotype, including 77,647 100-mers that were significantly enriched in the XY samples⁷⁰. The authors hypothesized that these bacterial genomes were not actually present in the samples, but rather the effect was caused by real human ChrY sequences matching contaminated bacterial genome database entries. We compared XY-enriched 100-mers from the study by Chrisman et al.⁷⁰ with the T2T-Y chromosome and found that, as predicted, more than 95% of them had near-perfect matches to the complete T2T-Y sequence.

We further tested the entire NCBI RefSeq bacterial genome database (release 213, July 2022, totalling 69,122 species with 40,758,769 contig or scaffold accessions) and identified all 64-mers that appeared in both the bacterial database and T2T-Y. We found 4,179 and 5,148 potentially contaminated sequences matching GRCh38-Y and T2T-Y, respectively (Extended Data Fig. 11a). The sequences were relatively short in length (under 1 kb), as is typical of contaminating genomic segments (Extended Data Fig. 11b). The vast majority of contaminated sequences found only with T2T-Y localized to the newly added HSat1B and HSat3 repeats (Extended Data Fig. 11c and Supplementary Table 32). Repeats are common sources of database contamination because their high copy number increases the chance they will be sequenced and assembled. We predict that this human-derived sequence contamination issue includes sequence from all human chromosomes and extends to all sequence databases, including non-microbial genomes.

Discussion

Owing to its highly repetitive structure, the human Y chromosome is the last of the human chromosomes to be completed from telomere to telomere. Here we present T2T-Y, a complete and gapless assembly of the Y chromosome from the HG002 benchmarking genome, along with a full annotation of its gene, repeat and organizational structure. We have combined T2T-Y with the previous T2T-CHM13 assembly to construct a new reference, T2T-CHM13+Y, that is inclusive of all human chromosomes. This assembly, along with all of the annotation resources presented here, is available for use as an alternative reference via NCBI and the UCSC Genome Browser⁷¹ (Data Availability).

Our analysis of the T2T-CHM13+Y reference assembly shows a reduction in false-positive variant calls for XY-bearing samples due to the correction of collapsed, incomplete, misassembled or otherwise inaccurate sequences in GRCh38-Y. Given the history of the GRCh38-Y assembly and its reliance on BAC libraries, we see no feasible means for its completion and suggest T2T-Y as a more suitable ChrY reference going forward. We recommend the use of T2T-CHM13 when mapping reads from XX samples, and ChrY PAR-masked T2T-CHM13+Y when mapping XY samples (Supplementary Note 3).

The completion of ampliconic and otherwise highly repetitive regions of ChrY will also require updates to existing gene annotations that are based on the incomplete GRCh38-Y assembly. How to label and refer to genes within variable-size ampliconic arrays, like *TSPY*, is an open question. Moreover, the highly repetitive sequences pose new challenges in regard to computational tools developed on GRCh38.

One example is the inconsistent methylation pattern observed in the satellite-enriched heterochromatin region, in which both HiFi and ONT are prone to sequencing biases, hindering accurate biological interpretation (Supplementary Note 4 and Supplementary Fig. 19). Lastly, we have noted the improved detection of human contamination in genomic databases using T2T-CHM13+Y and recommend a full contamination audit of public genome databases using this updated human reference. Taken together, these results illustrate the importance of using a complete human reference genome for bioinformatic analyses.

Construction of the T2T-Y assembly challenged the assembly methods developed for the haploid CHM13 genome and spurred the development of new, automated methods for diploid human genome assembly. In particular, the PARs of the HG002 sex chromosomes required phasing akin to heterozygous, diploid haplotypes, and the palindromic and heterochromatic regions of ChrY required expert curation of the initial assembly string graph. Lessons learned from our assembly of T2T-Y informed the development of the Verkko assembler⁷², which automates the integration of HiFi and ONT data for diploid human genome assembly. The companion study of Hallast et al.³² successfully used Verkko to generate 43 near-T2T assemblies from a diverse panel of human Y chromosomes, showing dynamic structural changes within this chromosome over the past 180,000 years of human evolution. Ultimately, as the complete, accurate and gapless assembly of diploid human genomes becomes routine, we expect that 'reference genomes' will become known simply as 'genomes'.

Projects such as the HPRC⁷³ are in the process of generating high-coverage HiFi and ONT sequencing for hundreds of additional human samples, and the assembly of these diverse, complete human genomes, along with similar quality assemblies of the non-human primates, will provide an unparalleled view of human variation and evolution. With the availability of complete, diploid human genome assemblies, association between phenotype and genotype will finally move beyond small variants alone and be made inclusive of all complex, structural genome variation.

Online content

Any methods, additional references, Nature Portfolio reporting summaries, source data, extended data, supplementary information, acknowledgements, peer review information; details of author contributions and competing interests; and statements of data and code availability are available at <https://doi.org/10.1038/s41586-023-06457-y>.

1. Skaletsky, H. et al. The male-specific region of the human Y chromosome is a mosaic of discrete sequence classes. *Nature* **423**, 825–837 (2003).
2. Miga, K. H. et al. Centromere reference models for human chromosomes X and Y satellite arrays. *Genome Res.* **24**, 697–707 (2014).
3. Vollger, M. R. et al. Segmental duplications and their variation in a complete human genome. *Science* **376**, eabj6965 (2022).
4. Nurk, S. et al. The complete sequence of a human genome. *Science* **376**, 44–53 (2022).
5. Schneider, V. A. et al. Evaluation of GRCh38 and de novo haploid genome assemblies demonstrates the enduring quality of the reference assembly. *Genome Res.* **27**, 849–864 (2017).
6. Gustafson, M. L. & Donahoe, P. K. Male sex determination: current concepts of male sexual differentiation. *Annu. Rev. Med.* **45**, 505–524 (1994).
7. Vog, P. H. et al. Human Y chromosome azoospermia factors (AZF) mapped to different subregions in Yq11. *Hum. Mol. Genet.* **5**, 933–943 (1996).
8. Miga, K. H. et al. Telomere-to-telomere assembly of a complete human X chromosome. *Nature* **585**, 79–84 (2020).
9. Logsdon, G. A. et al. The structure, function and evolution of a complete human chromosome 8. *Nature* **593**, 101–107 (2021).
10. Wenger, A. M. et al. Accurate circular consensus long-read sequencing improves variant detection and assembly of a human genome. *Nat. Biotechnol.* **37**, 1155–1162 (2019).
11. Jain, M. et al. Nanopore sequencing and assembly of a human genome with ultra-long reads. *Nat. Biotechnol.* **36**, 338–345 (2018).
12. Nurk, S. et al. HiCanu: accurate assembly of segmental duplications, satellites, and allelic variants from high-fidelity long reads. *Genome Res.* **30**, 1291–1305 (2020).
13. Rautiainen, M. & Marschall, T. GraphAligner: rapid and versatile sequence-to-graph alignment. *Genome Biol.* **21**, 253 (2020).
14. Formenti, G. et al. Merfin: improved variant filtering, assembly evaluation and polishing via k-mer validation. *Nat. Methods* **19**, 696–704 (2022).

15. Kirsche, M. et al. Jasmine and Iris: population-scale structural variant comparison and analysis. *Nat. Methods* **20**, 408–417 (2023).
16. Jain, C., Rhie, A., Hansen, N. F., Koren, S. & Phillippy, A. M. Long-read mapping to repetitive reference sequences using Winnowmap2. *Nat. Methods* **19**, 705–710 (2022).
17. Mc Cartney, A. M. et al. Chasing perfection: validation and polishing strategies for telomere-to-telomere genome assemblies. *Nat. Methods* **19**, 687–695 (2022).
18. Rhie, A., Walenz, B. P., Koren, S. & Phillippy, A. M. Merqury: reference-free quality, completeness, and phasing assessment for genome assemblies. *Genome Biol.* **21**, 245 (2020).
19. Wang, T. et al. The Human Pangenome Project: a global resource to map genomic diversity. *Nature* **604**, 437–446 (2022).
20. Jarvis, E. D. et al. Semi-automated assembly of high-quality diploid human reference genomes. *Nature* **611**, 519–531 (2022).
21. Shumate, A. et al. Assembly and annotation of an Ashkenazi human reference genome. *Genome Biol.* **21**, 129 (2020).
22. Zook, J. M. et al. Extensive sequencing of seven human genomes to characterize benchmark reference materials. *Sci. Data* **3**, 160025 (2016).
23. Landrum, M. J. et al. ClinVar: improvements to accessing data. *Nucleic Acids Res.* **48**, D835–D844 (2020).
24. Welter, D. et al. The NHGRI GWAS Catalog, a curated resource of SNP-trait associations. *Nucleic Acids Res.* **42**, D1001–D1006 (2014).
25. Smigielski, E. M., Sirotkin, K., Ward, M. & Sherry, S. T. dbSNP: a database of single nucleotide polymorphisms. *Nucleic Acids Res.* **28**, 352–355 (2000).
26. Karczewski, K. J. et al. The mutational constraint spectrum quantified from variation in 141,456 humans. *Nature* **581**, 434–443 (2020).
27. Byrsk-Bishop, M. et al. High-coverage whole-genome sequencing of the expanded 1000 Genomes Project cohort including 602 trios. *Cell* **185**, 3426–3440 (2022).
28. Mallick, S. et al. The Simons Genome Diversity Project: 300 genomes from 142 diverse populations. *Nature* **538**, 201–206 (2016).
29. Dunham, I. et al. An integrated encyclopedia of DNA elements in the human genome. *Nature* **489**, 57–74 (2012).
30. Ebert, P. et al. Haplotype-resolved diverse human genomes and integrated analysis of structural variation. *Science* **372**, eabf7117 (2021).
31. Sanders, A. D. et al. Single-cell analysis of structural variations and complex rearrangements with tri-channel processing. *Nat. Biotechnol.* **38**, 343–354 (2020).
32. Hallast, P. et al. Assembly of 43 human Y chromosomes reveals extensive complexity and variation. *Nature* <https://doi.org/10.1038/s41586-023-06425-6> (2023).
33. Hammer, M. F. et al. Extended Y chromosome haplotypes resolve multiple and unique lineages of the Jewish priesthood. *Hum. Genet.* **126**, 707 (2009).
34. Poznik, G. D. et al. Punctuated bursts in human male demography inferred from 1,244 worldwide Y-chromosome sequences. *Nat. Genet.* **48**, 593–599 (2016).
35. Vegesna, R., Tomaszewicz, M., Medvedev, P. & Makova, K. D. Dosage regulation, and variation in gene expression and copy number of human Y chromosome ampliconic genes. *PLoS Genet.* **15**, e1008369 (2019).
36. NCBI RefSeq v10 Browser. *Homo sapiens* isolate NA24385 chromosome Y, alternate assembly T2T-CHM13v2.0. <https://tinyurl.com/bdfudexn> (2022).
37. Hoyt, S. J. et al. From telomere to telomere: the transcriptional and epigenetic state of human repeat elements. *Science* **376**, eabk3112 (2022).
38. Warburton, P. E. et al. Analysis of the largest tandemly repeated DNA families in the human genome. *BMC Genomics* **9**, 533 (2008).
39. Halabian, R. & Makatowski, W. A map of 3' DNA transduction variants mediated by non-LTR retroelements on 3202 human genomes. *Biology* **11**, 1032 (2022).
40. Weissensteiner, M. H. et al. Accurate sequencing of DNA motifs able to form alternative (non-B) structures. *Genome Res.* **33**, 907–922 (2023).
41. Tyler-Smith, C., Taylor, L. & Müller, U. Structure of a hypervariable tandemly repeated DNA sequence on the short arm of the human Y chromosome. *J. Mol. Biol.* **203**, 837–848 (1988).
42. Xue, Y. & Tyler-Smith, C. An exceptional gene: evolution of the TSPY gene family in humans and other great apes. *Genes* **2**, 36–47 (2011).
43. Saxena, R. et al. Four DAZ genes in two clusters found in the AZFc region of the human Y chromosome. *Genomics* **67**, 256–267 (2000).
44. Altemose, N. et al. Complete genomic and epigenetic maps of human centromeres. *Science* **376**, eabl4178 (2022).
45. Jain, M. et al. Linear assembly of a human centromere on the Y chromosome. *Nat. Biotechnol.* **36**, 321–323 (2018).
46. Gershman, A. et al. Epigenetic patterns in a complete human genome. *Science* **376**, eabj5089 (2022).
47. Kasinathan, S. & Henikoff, S. Non-B-form DNA is enriched at centromeres. *Mol. Biol. Evol.* **35**, 949–962 (2018).
48. Nailwal, M. & Chauhan, J. B. Azoospermia factor C subregion of the Y chromosome. *J. Hum. Reprod. Sci.* **10**, 256 (2017).
49. Kuroda-Kawaguchi, T. et al. The AZFc region of the Y chromosome features massive palindromes and uniform recurrent deletions in infertile men. *Nat. Genet.* **29**, 279–286 (2001).
50. Repping, S. et al. A family of human Y chromosomes has dispersed throughout northern Eurasia despite a 1.8-Mb deletion in the azoospermia factor c region. *Genomics* **83**, 1046–1052 (2004).
51. Porubsky, D. et al. Recurrent inversion polymorphisms in humans associate with genetic instability and genomic disorders. *Cell* **185**, 1986–2005 (2022).
52. Teitz, L. S., Pyntikova, T., Skaletsky, H. & Page, D. C. Selection has countered high mutability to preserve the ancestral copy number of Y chromosome amplicons in diverse human lineages. *Am. J. Hum. Genet.* **103**, 261–275 (2018).
53. Jobling, M. A. Copy number variation on the human Y chromosome. *Cytogenet. Genome Res.* **123**, 253–262 (2008).
54. Navarro-Costa, P., Plancha, C. E. & Gonçalves, J. Genetic dissection of the AZF regions of the human Y chromosome: thriller or filler for male (in)fertility? *Biomed Res. Int.* **2010**, e936569 (2010).
55. Evans, H. J., Gosden, J. R., Mitchell, A. R. & Buckland, R. A. Location of human satellite DNAs on the Y chromosome. *Nature* **251**, 346–347 (1974).
56. Schmid, M., Guttenbach, M., Nanda, I., Studer, R. & Eppelen, J. T. Organization of DY22 repetitive DNA on the human Y chromosome. *Genomics* **6**, 212–218 (1990).
57. Manz, E., Alkan, M., Bühler, E. & Schmidtk, J. Arrangement of DY21 and DY22 repeats on the human Y-chromosome: a case with presence of DY21 and absence of DY22. *Mol. Cell. Probes* **6**, 257–259 (1992).
58. Altemose, N. A classical revival: human satellite DNAs enter the genomics era. *Semin. Cell Dev. Biol.* **128**, 2–14 (2022).
59. Gripenberg, U. Size variation and orientation of the human Y chromosome. *Chromosoma* **15**, 618–629 (1964).
60. Mathias, N., Bayés, M. & Tyler-Smith, C. Highly informative compound haplotypes for the human Y chromosome. *Hum. Mol. Genet.* **3**, 115–123 (1994).
61. Altemose, N., Miga, K. H., Maggioni, M. & Willard, H. F. Genomic characterization of large heterochromatic gaps in the human genome assembly. *PLoS Comput. Biol.* **10**, e1003628 (2014).
62. Cooke, H. Repeated sequence specific to human males. *Nature* **262**, 182–186 (1976).
63. Frommer, M., Prosser, J. & Vincent, P. C. Human satellite I sequences include a male specific 2.47 kb tandemly repeated unit containing one Alu family member per repeat. *Nucleic Acids Res.* **12**, 2887–2900 (1984).
64. Babcock, M., Yatsenko, S., Stankiewicz, P., Lupski, J. R. & Morrow, B. E. AT-rich repeats associated with chromosome 22q11.2 rearrangement disorders shape human genome architecture on Yq12. *Genome Res.* **17**, 451–460 (2007).
65. Webster, T. H. et al. Identifying, understanding, and correcting technical artifacts on the sex chromosomes in next-generation sequencing data. *GigaScience* **8**, giz074 (2019).
66. Aganezov, S. et al. A complete reference genome improves analysis of human genetic variation. *Science* **376**, eabl3533 (2022).
67. Bekritsky, M. A., Colombo, C. & Eberle, M. A. Identifying genomic regions with high quality single nucleotide variant calling. *Illumina* https://www.illumina.com/content/illumina-marketing/amr/en_US/science/genomics-research/articles/identifying-genomic-regions-with-high-quality-single-nucleotide-.html (2023).
68. Breitwieser, F. P., Perte, M., Zimin, A. V. & Salzberg, S. L. Human contamination in bacterial genomes has created thousands of spurious proteins. *Genome Res.* **29**, 954–960 (2019).
69. Steinegger, M. & Salzberg, S. L. Terminating contamination: large-scale search identifies more than 2,000,000 contaminated entries in GenBank. *Genome Biol.* **21**, 115 (2020).
70. Chrisman, B. et al. The human “contaminome”: bacterial, viral, and computational contamination in whole genome sequences from 1000 families. *Sci. Rep.* **12**, 9863 (2022).
71. Kent, W. J. et al. The Human Genome Browser at UCSC. *Genome Res.* **12**, 996–1006 (2002).
72. Rautiainen, M. et al. Telomere-to-telomere assembly of diploid chromosomes with Verkko. *Nat. Biotechnol.* <https://doi.org/10.1038/s41587-023-01662-6> (2023).
73. Liao, W.-W. et al. A draft human pangenome reference. *Nature* **617**, 312–324 (2023).
74. Jiang, Z., Hubley, R., Smit, A. & Eichler, E. E. DupMasker: a tool for annotating primate segmental duplications. *Genome Res.* **18**, 1362–1368 (2008).
75. Vollger, M. R., Kerpedjiev, P., Phillippy, A. M. & Eichler, E. E. StainedGlass: interactive visualization of massive tandem repeat structures with identity heatmaps. *Bioinformatics* **38**, 2049–2051 (2022).
76. Skene, P. J. & Henikoff, S. An efficient targeted nuclease strategy for high-resolution mapping of DNA binding sites. *eLife* **6**, e21856 (2017).
77. Robinson, J. T. et al. Integrative genomics viewer. *Nat. Biotechnol.* **29**, 24–26 (2011).

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

© This is a U.S. Government work and not under copyright protection in the US; foreign copyright protection may apply 2023

¹Genome Informatics Section, Computational and Statistical Genomics Branch, National Human Genome Research Institute, National Institutes of Health, Bethesda, MD, USA. ²Faculty of Informatics, Masaryk University, Brno, Czech Republic. ³Department of Biomolecular Engineering, University of California Santa Cruz, Santa Cruz, CA, USA. ⁴Department of Molecular and Cell Biology, University of Connecticut, Storrs, CT, USA. ⁵Department of Biology, Johns Hopkins University, Baltimore, MD, USA. ⁶Department of Molecular and Cell Biology, University of California, Berkeley, CA, USA. ⁷Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD, USA. ⁸Federal Research Center of Biotechnology of the Russian Academy of Sciences, Moscow, Russia. ⁹Center for Algorithmic Biotechnology, Saint Petersburg State University, St Petersburg, Russia. ¹⁰European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge, UK. ¹¹UC Santa Cruz Genomics Institute, University of California Santa Cruz, Santa Cruz, CA, USA. ¹²Graduate Program in Bioinformatics and Systems Biology, University of California, San Diego, CA, USA. ¹³Department of Computer Science, Johns Hopkins University, Baltimore, MD, USA. ¹⁴GeneDX Holdings Corp, Stamford, CT, USA. ¹⁵Foundation of Biological Data Science, Belmont, CA, USA. ¹⁶Department of Genetics, University of Cambridge, Cambridge, UK. ¹⁷The Rockefeller University, New York, NY, USA. ¹⁸DNAnexus, Inc., Mountain View, CA, USA. ¹⁹Department of Genetics, Genomics and Informatics, University of Tennessee Health Science Center, Memphis, TN, USA. ²⁰Stowers Institute for Medical Research, Kansas City, MO, USA. ²¹University of Kansas Medical Center, Kansas City, MO, USA. ²²Genomics Research Centre, Human Technopole, Milan, Italy. ²³Institute of Bioinformatics, Faculty of Medicine, University of Münster, Münster, Germany. ²⁴Cancer Genetics and Comparative Genomics Branch, National Human Genome Research Institute, National Institutes of Health, Bethesda, MD, USA. ²⁵Department of Biology, Pennsylvania State University, University Park, PA, USA. ²⁶Department of Genome Sciences, University of Washington School of Medicine, Seattle, WA, USA. ²⁷Institute for Systems Biology, Seattle, WA, USA. ²⁸XDBio Program,

Johns Hopkins University, Baltimore, MD, USA. ²⁹Department of Bioengineering, Department of Physics, Northeastern University, Boston, MA, USA. ³⁰Human Genome Sequencing Center, Baylor College of Medicine, One Baylor Plaza, Houston, TX, USA. ³¹Department of Data Sciences, Dana-Farber Cancer Institute, Boston, MA, USA. ³²Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA. ³³Genome Technology Access Center at the McDonnell Genome Institute, Washington University, St. Louis, MO, USA. ³⁴Biosystems and Biomaterials Division, National Institute of Standards and Technology, Gaithersburg, MD, USA. ³⁵Department of Computer Science and Engineering, Pennsylvania State University, University Park, PA, USA. ³⁶Department of Biochemistry and Molecular Biology, Pennsylvania State University, University Park, PA, USA. ³⁷Center for Computational Biology and Bioinformatics, Pennsylvania State University, University Park, PA, USA. ³⁸UCL Queen Square Institute of Neurology, UCL, London, UK. ³⁹National Center for Biotechnology Information, National Library of Medicine, National Institutes of Health, Bethesda, MD, USA. ⁴⁰Masters Program in National Research University Higher School of Economics, Moscow, Russia.

⁴¹Departments of Biomedical Engineering, Computer Science, and Biostatistics, Johns Hopkins University, Baltimore, MD, USA. ⁴²Department of Computer Science, Rice University, Houston, TX, USA. ⁴³Google Inc., Mountain View, CA, USA. ⁴⁴Institute of Molecular Genetics, Moscow, Russia. ⁴⁵Center for Evolution and Medicine, School of Life Sciences, Arizona State University, Tempe, AZ, USA. ⁴⁶Department of Biomedical Engineering, Pennsylvania State University, State College, PA, USA. ⁴⁷Pacific Biosciences, Menlo Park, CA, USA. ⁴⁸Investigator, Howard Hughes Medical Institute, University of Washington, Seattle, WA, USA. ⁴⁹Institute for Systems Genomics, University of Connecticut, Storrs, CT, USA. ⁵⁰Department of Genetics and Genome Sciences, UConn Health, Farmington, CT, USA. ⁵¹Present address: Oxford Nanopore Technologies Inc., Oxford, UK. ⁵²Present address: Department of Anatomy and Anthropology and Department of Human Molecular Genetics and Biochemistry, Sackler Faculty of Medicine, Tel Aviv University, Tel Aviv-Yafo, Israel. ⁵³These authors contributed equally: Arang Rhie, Sergey Nurk, Monika Cechova, Savannah J. Hoyt, Dylan J. Taylor. ⁵²e-mail: adam.phillippy@nih.gov

Methods

This section provides a brief summary of the methods. Refer to the Supplementary Methods for further details.

HG002 cell line

HG002 cell lines (GM24385, GM26105 and GM27730) were purchased from the Coriell Institute and used for sequencing and generation of data used in this study. HG002 DNA is available as a reference material from the National Institute of Standards and Technology (NIST), and the associated cell lines have previously been consented for both research use and commercial redistribution. More details can be found at <https://www.nist.gov/programs-projects/genome-bottle> and <https://www.coriell.org/1/NIGMS/Collections/NIST-Reference-Materials>. The authenticity of the cell line and DNA was subsequently confirmed by comparison of assembly-based variant calls with the HG002 GIAB truth set and karyotyping. The cell lines were not tested for mycoplasma contamination. The final product of this study (T2T-Y assembly) was tested for contamination, but none was identified except the Epstein–Barr virus (EBV) used for immortalization, which was found as an external chromosomal component.

Sequencing

Seventeen PacBio HiFi WGS runs were generated from GM24385 using the SMRTbell Express Template Prep Kit 2.0 on the Sequel II platform, after size selection for 15–25 kb fragments. All ONT WGS runs are from the study of Jarvis et al.²⁰, having been generated using protocols from Shafin et al.⁷⁸ and Jain et al.¹¹.

RNA was extracted from three cell lines to generate Iso-seq reads: EBV-immortalized lymphoblastoid cell line (GM24385), induced pluripotent stem cell of the EBV-immortalized lymphoblastoid cell line (GM26105) and induced pluripotent stem cell derived directly from peripheral blood mononuclear cells (GM27730). Iso-seq data were generated on the Sequel II platform and processed using Lima and IsoSeq3.

Specific runs used either in assembly or validation, along with their accessions, are available in Supplementary Table 1, and DNA preparation and library generation information is available in Supplementary Methods.

Assembly and validation

An assembly string graph was first constructed using PacBio HiFi reads (roughly 60×) and processed using custom pruning procedures. Due to high sequence similarity within PAR1 and PAR2, the HG002 ChrX and ChrY string graph components shared connections to one another in the PARs but to no other chromosomes in the genome. The remaining tangles in these sex-chromosome subgraphs were resolved using ONT reads longer than 100 kb (around 90×). A semiautomated repeat-resolution strategy utilized GraphAligner¹³ to map ultralong ONT data to the HiFi assembly graph and identify the correct traversals. To resolve the PAR regions, ChrX and ChrY chromosomal walks were identified using homopolymer compressed trio-binned k-mers from parental Illumina reads⁷⁹, and a consensus sequence was computed for each. Remaining coverage gaps caused by HiFi sequencing biases were patched using a de novo Flye assembly of trio-binned paternal ONT reads^{14,80}. For new projects we now recommend the automated Verkko pipeline⁷², which is able to replicate the semimanual T2T-Y assembly presented here.

For polishing, ChrXY draft assemblies were combined with T2T-CHM13v1.1 autosomes to prevent mapping biases caused by the incompletely resolved autosomes of HG002 (that is, T2T-CHM13+XY). For further polishing and validation we used 66× Illumina, 84× HiFi and 250× ONT (being haploid, the effective coverage on X and Y is half of those depths). Small corrections were identified using DeepVariant^{81,82} and filtered with Merfin¹⁴. Large errors were identified with Sniffles⁸³ and cuteSV⁸⁴ and through comparison with the HPRC-HG002v1

assembly²⁰. All large errors were patched using marker-assisted HiFi and ONT reads. Assembly issues were identified using repeat-aware, long-read alignments from Winnowmap2 (ref. 16) (filtered with globally unique markers¹⁷) and VerityMap⁸⁵ (guided by locally unique markers). Coverage summaries were obtained using scripts from T2T-CHM13 assembly evaluation¹⁷. Putative collapses and inversion errors were identified using Strand-seq data. Raw sequencing reads from 65 Strand-seq libraries^{30,31} were aligned to both GRCh38 and T2T-CHM13+XY with the Burrows–Wheeler algorithm⁸⁶, then processed with breakpointR^{30,87} to identify inversion errors. Recurrent inversions were identified by comparison with results from Porubsky et al.⁵¹. To further confirm integrity of ChrY in the HG002 cell line, we aligned publicly available GIAB²² HiFi and ONT reads from the paternal HG003 genome (including from the PacBio Revio platform⁸⁸) and performed the same coverage analysis. Base error rate was measured by Merqury using a hybrid k-mer set from Illumina and HiFi reads^{17,18} (Supplementary Table 3).

Comparison with GRCh38-Y

Y haplogroup identification. The Y-chromosome haplogroup of HG002 and GRCh38 was identified using yhaplo⁸⁹, which builds a tree from phylogenetically informative single-nucleotide polymorphisms (SNPs) that accumulate in the non-recombining portion of the Y. The Y haplogroups of 1KGP samples were identified by Y-Finder⁹⁰ using SNP calls on GRCh38 from Aganezov et al.⁶⁶.

Alignments between GRCh38 and HG002 Y assemblies. Alignments between the GRCh38-Y and T2T-Y assemblies for the purposes of visualization with Saffire were generated with minimap2 (ref. 91). The pairwise mapping format (PAF) output was then processed with rustybam and visualized with Saffire. DupMasker⁷⁴ and dna-brnn⁹² annotations were generated using Rhodonite (<https://doi.org/10.5281/zenodo.6036498>).

A complementary alignment was generated with LASTZ⁹³ after softmasking repeats from WindowMasker⁹⁴. The alignment dotplot and best identity were plotted using R (<https://github.com/arangrhie/T2T-HG002Y/tree/main/alignments/lastz>). Regions along T2T-Y were coloured according to their class.

To visualize three large structural differences of the three ChrY assemblies (GRCh37-Y, GRCh38-Y and T2T-Y) we used the Pangenomics Research Tool Kit⁹⁵ to construct principal bundles representing contiguous and conserved sequences among the pangenome contigs.

Gene annotation

GENCODE v.35 CAT Liftoff annotation. Preliminary gene annotation was performed by mapping GENCODE v.35 (ref. 96) annotations from GRCh38-Y to T2T-Y using a Cactus⁹⁷ alignment with Chimp as an outgroup. Iso-seq reads were aligned and assembled with Stringtie2 (ref. 98), aligned to the assembly with TransMap⁹⁹ and used as input for CAT¹⁰⁰ along with the GENCODE v.35 annotation. GENCODE v.35 Y annotations were mapped with Liftoff¹⁰¹ then intersected with Bedtools¹⁰² to isolate genes mapped by Liftoff to ChrY that were not in the CAT annotation.

De novo RefSeq v.110 and GENCODE v.38 annotation. In the meantime, a de novo RefSeq annotation was performed on both GRCh38 and T2T-CHM13+Y and released (v.110) as previously described for other vertebrate genomes^{103,104}. A total of 82,862 curated RefSeq transcripts, 345,700 complementary DNAs, 8.65 million expressed sequence tags, 9.7 billion RNA sequencing reads and 83 million PacBio Iso-seq and ONT reads from over 30 distinct tissues were retrieved from Sequence Read Archive and were tentatively aligned to the assembly using either Splign¹⁰⁵ or minimap2 (ref. 91).

Simultaneously, an Ensembl gene annotation was performed by a mapping subset of the genes from GENCODE v.38 (ref. 96) using

minimap2 (ref. 91) and MAFFT¹⁰⁶. Transcripts with low coverage or identity below 98% were realigned using Exonerate¹⁰⁷. Genes in potential recent duplications or collapsed paralogues were adjusted accordingly.

RefSeq Liftoff and curated ampliconic gene annotation. Because the additional copies of ampliconic genes hindered comparison with known genes in GRCh38-Y with differing gene IDs and names, we performed one more annotation by mapping GRCh38 RefSeq v.110 annotations with Liftoff to T2T-CHM13+Y. We compared ampliconic gene family annotation results from those of GENCODE CAT/Liftoff and assigned gene names, followed by best gene coverage and identity, including introns. Later, based on discussions with authors of a companion paper³², we adjusted the gene names for three protein-coding annotations based on exon sequence identity (Supplementary Table 5).

Validation of ampliconic protein-coding genes. Copy numbers for each ampliconic gene family in both the GRCh38-Y and T2T-Y assemblies were estimated using an adapted application of AmpliCoNE³⁵. Copy numbers of these gene families were previously estimated for HG002 using Illumina reads from GIAB¹⁰⁸ and ddPCR³⁵. The only notable difference was in regard to *TSPY* copy number, in which we identified 45 intact protein-coding copies. Copy number was slightly higher in the assembly than in estimates derived from Illumina reads and ddPCR (45 versus 40 and 42, respectively). The in silico PCR primer search matched all 44 protein-coding copies in the *TSPY* gene array and *TSPY2*, and two pseudogenes at the 3' end of the *TSPY* array that we were unable to avoid in the ddPCR primer design. We conclude that our AmpliCoNE, ddPCR and in silico PCR estimates agree with ampliconic gene annotations in the T2T-Y assembly (Table 1).

Repeat annotation

Segmental duplications. Segmental duplication annotations were created using the same methods as in Vollger et al. without modification³. In brief, segmental duplications in T2T-CHM13+Y were identified using SEDEF¹⁰⁹ after repeat masking with Tandem Repeats Finder¹¹⁰ and RepeatMasker¹¹¹.

Repeat model discovery and annotation. A three-step repeat annotation was performed to annotate new repeat models on ChrY. First, RepeatMasker was performed on the T2T-Y assembly using the Dfam 3.3 library¹¹², hard-masking previously annotated repeats. Second, RepeatModeler analysis was performed on the remaining unmasked regions to identify new repeat model consensus, which were subjected to extension and filtering and used as a library for a secondary RepeatMasker run. Two methods were primarily used to identify new satellite repeats: ULTRA¹¹³ and NTRprism⁴⁴. Unannotated regions over 5 kb were identified via Bedtools¹¹⁴ by subtracting repeat annotations from the first and last steps above. Regions were manually curated in UCSC Genome Browser to check for any feature overlap (for example, gene annotations). Tandemly repeated sequences were detected and assessed with a combination of ULTRA, NTRprism and the TRF GUI version¹¹⁰ to determine the best monomer consensus for a given satellite model. Lastly, the compilation pipeline laid out in Hoyt et al.³⁷ was followed to avoid potential false positives by simply masking with a combined library of new repeat models and known repeat models (Dfam library). The same three-step, repeat-annotation pipeline was also applied to GRCh38-Y. Repeats were summarized using buildSummary.pl¹¹⁵ at both class and family level (Table 1 and Supplementary Table 12), and at the subfamily level for new repeats (Supplementary Table 11), in both T2T-Y and GRCh38-Y.

Composite repeats. Composite elements were defined and characterized as described in Hoyt et al.³⁷ as repeating units consisting of three or more repeated sequences, including TEs, simple repeats, composite subunits and/or satellites that are found as a tandem array in at least

one location in the genome. BLAT¹¹⁶ was used to locate other composite unit copies across T2T-Y and cross-reference these with their associated gene annotations (CAT/Liftoff). Identification of potentially active, full-length TEs (SINEs, LINEs and retroposons are AluY, L1Hs and SVA_E/F) across T2T-Y and GRCh38-Y was done following the methods of Hoyt et al.³⁷.

Satellite annotation. Centromeric satellite (cen/sat) annotations were generated as in Altemose et al.⁴⁴, with several refinements tailored to include annotations of the entire ChrY. Major satellite types were extracted from the RepeatMasker track, with features merged for the same satellite type within 10 kb of each other. For HSat2 and HSat3, a specialized annotation tool was used (https://github.com/altemose/chm13_hsat)⁴⁴. DYZ19 and HSat1B were annotated using RepeatMasker annotations. Exact boundaries between HSat3 and HSat1B (aka DYZ1 and DYZ2) were manually refined.

Transduction analysis. We utilized the same approach as in Hoyt et al.³⁷ to identify putative DNA transductions mediated by retro-elements. Briefly, L1s and SVAs were identified in T2T-Y to detect target site duplications and 3' transduction signatures using a modified version of TSDfinder¹¹⁷. Then, we removed transductions residing in segmental duplications and masked the transduced sequences using RepeatMasker¹¹¹. To find the potential progenitor of each transduction within T2T-CHM13+Y and GRCh38, offspring sequences were aligned to the corresponding databases of full-length L1s and SVAs using BLAST¹¹⁸.

Non-B DNA motif annotation. To predict sequence motifs with the potential to form alternative DNA structures (non-B DNA) we used nBMST¹¹⁹ for repeat motifs (A-phased, direct, inverted and mirror repeats and STRs) and Z-DNA motifs^{120,121}. G4-motifs were detected with Quadron¹²², which also yields a score that quantifies the stability of a predicted G4 structure based on a machine-learning algorithm trained on empirical datasets. Motifs with a Quadron score of 19 or more were considered stable and thus used throughout our analysis. Non-B motifs were intersected with other existing annotations of T2T-Y (gene annotations, satellite repeats and CpG islands) using Bedtools¹¹⁴. Rideogram¹²³ was used to generate these visualized tracks, as well as the three composite repeat tracks. GraphPad Prism¹²⁴ was used to generate TE composition per sequence class plots.

Data visualization

For Figs. 1 and 3 the alignment of GRCh38-Y and T2T-Y was visualized with Saffire¹²⁵. Segmental duplications are coloured by duplication types defined in DupMasker⁷⁴. IGV⁷⁷ was used to draw ideograms, sequence classes, palindromes, IRs and AZF. Bedtools¹¹⁴ was used to calculate density (across each gene type), base pair coverage (across each repeat class) and average CpG methylation frequency per 100 kb window. Dotplots coloured by identity were generated with StainedGlass⁷⁵.

TSPY gene family analysis

TSPY copy number estimation from SGRP. Copy number of the *TSPY* gene was estimated as in Vollger et al.³ In brief we applied the fastCN pipeline¹²⁶, which uses read depth as a proxy. Short-read sequence data were processed into 36 bp non-overlapping fragments and mapped using mrsFAST¹²⁷ to a T2T-CHM13+Y reference masked with TRF and RepeatMasker. Read depth across the genome was corrected for genomic copy bias, and copy number determined using linear regression on read depth versus known fixed copy number control regions. Finally, integer genotypes for *TSPY* were generated by taking a weighted average of copy number estimates from windows overlapping the locus.

Phylogenetic tree analysis of TSPY genes. All curated protein-coding and pseudogene *TSPY* copies (including introns) from the CAT/

Liftoff and RefSeq/Liftoff annotations were used. For outgroup rooting of the tree, *TSPY* sequences were used from *Hylobates moloch* (NW_022611649.1)¹²⁸ and *P. abelii* (KP141780.1)¹²⁹. Alignment was carried out in MAFFT¹⁰⁶. Phylogenetic analysis was run in RAXML-NG¹³⁰ with 200 bootstrap replicates and rapid bootstrap approximation. Consensus bootstrap values were then mapped to the highest-likelihood phylogeny in Geneious¹³¹ and visualized in FigTree¹³².

Centromere analysis

T2T-Y was processed using the standard alpha-satellite tools as described in Altemose et al.⁴⁴. The S4CYHIL (DYZ3) alpha-satellite HOR was re-examined and redefined for this Article to take into account its polymorphic variants both known from previous literature^{133,134} and shown by the recent complete centromere assembly of RP11 (ref. 45).

The CENP-A CUT&RUN data were aligned to the T2T-CHM13+Y assembly as previously described in Altemose et al.⁴⁴. Alignments were filtered using the single-copy k-mer locus filtering method as described in Hoyt et al.³⁷ through the use of the UCSC Genome Browser tool overlapSelect.

HG002 ONT ultralong data were again basecalled using Guppy v.6.1.2, Remora to obtain CpG methylation data (Supplementary Table 1). Modbams were converted to FASTQ files and aligned with Winnowmap2 (ref. 16). HG002 ONT nanoNOMe data were generated in Gershman et al.⁴⁶ and analysed with nanopolish¹³⁵. The probability of methylation for each CpG site in PacBio HiFi reads was assigned using primrose in SMRT Link v11.1 available at <https://www.pacb.com/support/software-downloads/>. A newer version of primrose is available as jasmine (<https://github.com/PacificBiosciences/jasmine>). Reads were aligned with pbmm2 (<https://github.com/PacificBiosciences/pbmm2>). The percentage of methylated reads at each reference genome position was calculated using pb-CpG-tools (<https://github.com/PacificBiosciences/pb-CpG-tools>). Resulting modbams were reprocessed identically to Remora-called ONT data to collect comparable aggregated native CpG methylation data. The CDR was manually annotated as the area where CpG methylation is lower than the flanking, active, alpha-satellite (Supplementary Fig. 14).

Sequence classes on the Y chromosome

The X-degenerate and ampliconic regions were annotated using either exact boundaries of palindromes or intrachromosomal identity, as defined in Skaletsky et al.¹, with adjusted borders based on gene annotations. T2T-Y was split into 5 kb sliding windows (step size 1 kb) and these sequences were mapped back to T2T-Y using Winnowmap2 (ref. 16). Following exclusion of self-alignments, windows with identity over 50% were considered indicative of ampliconic regions if present consecutively.

For the schematic representations shown in Fig. 4, amplicons from Teitz et al.⁵² were mapped to GRCh38-Y and T2T-Y assemblies with Winnowmap2 (ref. 16) to identify homologous regions. Approximate boundaries of palindrome P4–P8 arms were manually selected using Gepard¹³⁶ and further refined based on alignment of palindromic arms and adjacent flanks against each other (arm1 to the reverse complement of arm2) using global alignment with Stretcher¹³⁷.

For AZFa, sequences between two *HERV15* genes (including genes *USP9Y* and *DDX3Y*) were used to determine AZFa boundaries¹³⁸. The boundaries of AZFb and AZFc were defined by the amplicon units P5/proximal P1 deletion (yel3/yel1) and by the b2/b4 deletion. A self-dotplot of the T2T-Y assembly was used with word size of 100 in Gepard¹³⁶. Breakpoints were identified as illustrated in figure 2 of Navarro-Costa et al.⁵⁴, as shown in Supplementary Fig. 15.

The PAR and X-transposed regions were initially identified using LASTZ⁹³ alignments between HG002 X, HG002 Y and CHM13 X. Exact boundaries were later refined using minimap2 (ref. 91) alignments.

Yq12 heterochromatin region

DYZ1 and DYZ2. DYZ1 and DYZ2 consensus sequences were generated by multiple sequence alignment using kalign¹³⁹ and converted to a profile hidden Markov model with HMMER¹⁴⁰. Dotplots in Extended Data Fig. 8 were produced using dottup in the EMBOSS software package¹³⁷.

Phylogenetic tree analysis of AluY. The AluY tree was rooted on the RepeatMasker/Dfam-derived consensus sequence for AluSc8. Analysis was run on a MAFFT¹⁰⁶-derived alignment using RAXML-NG¹³⁰ with 100 non-parametric bootstrap replications. Note that in the AluY subfamily clade ('Mixed AluY Subfamilies') there are scattered elements across the group even though the majority are represented in the labelled subclades.

Short-read variant calling on T2T-CHM13+Y

Impact of masking PAR and XTR on variant calling. Simulated paired-end sequence reads were generated using NEAT¹⁴¹. Variants from 10 XY and 10 XX European individuals were collected from high-coverage variant calls of 1KGP²⁷ and used for benchmarking. Reads were processed with bbduk¹⁴² and mapped using bwa⁸⁶ to two versions of GRCh38: X and Y both unmasked (default), and sex chromosome aware (SCC aware⁶⁵). Masking was performed on PAR¹⁴³ or on both PAR and XTR^{65,144}. Mapping quality was assessed on ChrX in each 50 kb window, sliding 10 kb using Bedtools¹¹⁴. Variant calling was performed with GATK¹⁴⁵ and compared against the chosen variants used in simulation of reads.

Mappability comparison and variant calling in 1KGP samples. Using the National Human Genome Research Institute (NHGRI) Genomic Data Science Analysis, Visualization, and Informatics Lab-Space (AnVIL)¹⁴⁶, we performed short-read alignment and variant calling for the 3,202 samples in 1KGP²⁷ using the T2T-CHM13+Y assembly as a reference. These samples were sequenced to at least 30× coverage by the New York Genome Center, with alignment and variant calling previously performed on the GRCh38 reference. We largely followed the short-read alignment and variant-calling pipeline used for analysis of T2T-CHM13v1.0 (ref. 66), except that we used SCC references for all XX and XY individuals using XYalign⁶⁵. In the XX-specific reference, the entire Y chromosome is masked whereas in the XY-specific reference only Y-PARs are masked. For all analyses, measures of mappability (reads mapped, reads properly paired, mismatch rate) were assessed with Samtools¹⁴⁷, and variant counts and allele frequencies were assessed with bcftools¹⁴⁷. Variants in syntenic regions between GRCh38-Y and T2T-Y were further subsetted with Bedtools¹¹⁴.

Putative collapsed regions in GRCh38-Y. Three individuals' variant calls and corresponding bam files from the 1KGP dataset were downloaded from AnVIL—one individual each from the J1, R1b and E1b haplogroups (HG01130, HG00116 and HG01885, respectively). Variant calls on ChrY syntenic region were subset using bcftools¹⁴⁷. From the variant call format (VCF) file, allelic read depth (defined as AD field) and reference allele depth (first value in the depth (DP) field) were extracted using a custom script along with each variant's chromosomal position and visualized with R. Coverage tracks of the bam files were collected with IGVtoolkit⁷⁷ and samtools¹⁴⁷. Variants from HG00116 on GRCh38-Y (R1b, therefore least structural variations expected) were further aggregated as an 'excessive variant region' when non-reference alleles were present and merged within 50 kb. Coverage, variant calls and excessive variant regions were manually inspected on GRCh38-Y and marked as a 'putative collapse' if the region (1) had an excessive number of variants called for all three samples, (2) overlapped with a known gap in GRCh38 and (3) did not overlap the palindromic region (where there were substantial rearrangements between GRCh38-Y and T2T-Y).

Article

Mapping and variant calling of SGDP samples. The SGDP includes 279 open-access, high-coverage genomes from 130 diverse populations²⁸. Compared with 1KGP, SGDP includes 118 additional populations with samples sequenced to an average of 43× coverage using a shared PCR-free Illumina library. SGDP samples were aligned and genotyped to T2T-CHM13+Y and GRCh38 on AnVIL¹⁴⁶ following the same pipeline used for analysis of 1KGP samples.

Curated syntenic region and liftover chains. The initial chain file was generated using nf-LO¹⁴⁸ with minimap2 (ref. 91) alignments. Alignments were filtered and converted to PAF using chaintools. Alignments of non-homologous chromosomes were removed. Overlapping alignments in the query sequence was removed with rustybam to create 1:1 alignments. PAF alignments were converted back to chain format.

In addition to minimap2-based, whole-genome alignment we applied a wfmash-based pipeline¹⁴⁹ to validate the chain file. This pipeline starts with a wfmash¹⁴⁹ whole-genome alignment of T2T-CHM13+Y and the masked and filtered GRCh38 assembly, and identifies 1-to-1 homologous regions at least 5 kb long with a nucleotide identity of at least 95%. Similarly, the resulting chain was postprocessed to obtain 1:1 alignments using rustybam and the paf2chain tool. All PAF files with full CIGAR strings were then inspected with Saffire for quality investigation. The minimap2- and wfmash-based chains showed high consistency over the genomes.

Datasets and resources for T2T-CHM13+Y

Lifting over resources from GRCh38 to CHM13+Y. Using the curated chain file, we lifted over dbSNP build 155 (ref. 150), the 13 March 2022 release of Clinvar^{23,151} and GWAS Catalog v.1.0 (refs. 24,152) from the GRCh38 primary assembly to T2T-CHM13v2.0 (T2T-CHM13+Y). Liftover was performed as previously described⁶⁶ using GATK Picard¹⁵³ LiftoverVcf and the alignment chain described above.

ENCODE. Reads were obtained from the ENCODE dataset²⁹ and mapped with Bowtie2 (ref. 154). Alignments were filtered using Samtools to remove unmapped or single-end mapped reads and those with a mapping quality score under 2. PCR duplicates were identified and removed with Picard tools ‘mark duplicates’. Alignments were then filtered for the presence of unique k-mers. Bigwig coverage tracks and enrichment tracks were created using deepTools2 bamCoverage¹⁵⁵.

gnomAD. Genome-wide variant data from the Genome Aggregation Database (gnomAD) release v.3.1.2 was lifted over from GRCh38 to each assembly using CrossMap¹⁵⁶. The chain files used were created from the GRCh38-based HAL file, downloaded from the cactus-minigraph alignment of Liao et al.⁷³ The resulting variant-call formats were annotated with predicted molecular consequence and transcript-specific variant deleteriousness scores from PolyPhen-2 and SIFT using Ensembl Variant Effect Predictor.

Human Y chromosome contamination in bacterial genomes

Screening against the study of Chrisman et al. We used MUMmer¹⁵⁷ to compare 73,691 bacterial 100-mers reported as enriched in human males by Chrisman et al.⁷⁰ with the T2T-Y assembly. We found that, as predicted, more than 95% of the 100-mers had near-perfect matches, defined as an exact match of 50 bp or longer, to the complete T2T-Y sequence. The nucmer programme from MUMmer was run with default options, except to specify -l 50 for an exact match length of 50 or more and -c 50 so that it reported matches as short as 50 bp.

Screening with 64-mers. Meryl¹⁸ was used to compare 64-mers from NCBI RefSeq release 213 (July 2022) with T2T-Y and GRCh38-Y. Each bacterial contig was annotated with the number of matching k-mers in T2T-Y, GRCh38-Y and the number of k-mers in the contig with a match. Each position in the reference chromosomes was annotated

with the multiplicity of the k-mer at that position in the RefSeq contigs, and with the number of contigs containing the k-mer. Hits per query were filtered to retain only contigs with more than 20 k-mer matches or with more than 10% of the contig sequence covered by k-mer matches. The queries at each reference position were combined and accumulated into 10 kb windows and converted to an interval wiggle file for visualization. RefSeq sequence entries with hits were retrieved using seqrequester and categorized using 64-mers built from HSat1B and HSat3 annotations. The first and second words in sequence entry names were extracted to visualize the taxonomic abundance of microbial genomes in a pie chart using Kronatools¹⁵⁸ (Extended Data Fig. 11c).

Reporting summary

Further information on research design is available in the Nature Portfolio Reporting Summary linked to this article.

Data availability

The T2T-CHM13v2.0 (T2T-CHM13+Y) assembly, reference analysis set, complete list of resources—including gene annotation, repeat annotation, epigenetic profiles, variant-calling results from 1KGP and SGDP, gnomAD, ClinVar, GWAS and dbSNP datasets—are available for download at <https://github.com/marbl/CHM13>. The assembly is also available from NCBI and EBI with GenBank accession GCA_009914755.4. Annotation and associated resources are also browsable as ‘hs1’ from the UCSC Genome Browser (http://genome.ucsc.edu/cgi-bin/hgTracks?db=hub_3671779_hsl1), the Ensembl Genome Browser (<https://projects.ensembl.org/hprc/>) (assembly name T2T-CHM13v2.0) and NCBI data-hub (https://www.ncbi.nlm.nih.gov/data-hub/genome/GCF_009914755.1/). Potential assembly issues are listed and can be tracked at <https://github.com/marbl/CHM13-issues>. 1KGP and SGDP short-read alignments and variant calls are available within AnVIL at https://anvil.terra.bio/#workspaces/anvil-datastorage/AnVIL_T2T_CHRY. Original data from the Gerton lab underlying this manuscript can be accessed from the Stowers Original Data Repository at <http://www.stowers.org/research/publications/libpb-2358>. Sequencing data used in this study are listed in Supplementary Table 1.

Code availability

Custom codes developed for data analysis and visualization are available at <https://github.com/arangrhie/T2T-HG002Y>, https://github.com/snurk/sg_sandbox and <https://github.com/schatzlab/t2t-chm13-chry> and are deposited with Zenodo¹⁵⁹. Software and parameters used are stated in the Supplementary Methods with further details.

78. Shafin, K. et al. Nanopore sequencing and the Shasta toolkit enable efficient de novo assembly of eleven human genomes. *Nat. Biotechnol.* **38**, 1044–1053 (2020).
79. Koren, S. et al. De novo assembly of haplotype-resolved genomes with trio binning. *Nat. Biotechnol.* **36**, 1174–1182 (2018).
80. Kolmogorov, M., Yuan, J., Lin, Y. & Pevzner, P. A. Assembly of long, error-prone reads using repeat graphs. *Nat. Biotechnol.* **37**, 540–546 (2019).
81. Poplin, R. et al. A universal SNP and small-indel variant caller using deep neural networks. *Nat. Biotechnol.* **36**, 983–987 (2018).
82. Shafin, K. et al. Haplotype-aware variant calling with PEPPER-Margin-DeepVariant enables high accuracy in nanopore long-reads. *Nat. Methods* **18**, 1322–1332 (2021).
83. Sedlazeck, F. J. et al. Accurate detection of complex structural variations using single molecule sequencing. *Nat. Methods* **15**, 461–468 (2018).
84. Jiang, T. et al. Long-read-based human genomic structural variation detection with cuteSV. *Genome Biol.* **21**, 189 (2020).
85. Bzikadze, A. V., Mikheenko, A. & Pevzner, P. A. Fast and accurate mapping of long reads to complete genome assemblies with VerityMap. *Genome Res.* **32**, 2107–2118 (2022).
86. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics* **25**, 1754–1760 (2009).
87. Porubsky, D. et al. breakpointR: an R/Bioconductor package to localize strand state changes in Strand-seq data. *Bioinformatics* **36**, 1260–1261 (2020).
88. PacBio Revio WGS Dataset. *Homo sapiens* – GIAB trio HG002-4. <https://downloads.pacbcloud.com/public/revio/2022Q4/> (2022).
89. Poznik, D. yhaplo | Identifying Y-chromosome haplogroups. GitHub <https://github.com/23andMe/yhaplo> (2022).

90. Tseng, B. et al. Y-SNP Haplogroup Hierarchy Finder: a web tool for Y-SNP haplogroup assignment. *J. Hum. Genet.* **67**, 487–493 (2022).
91. Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100 (2018).
92. Li, H. Identifying centromeric satellites with dna-brnn. *Bioinformatics* **35**, 4408–4410 (2019).
93. Harris, R. S. *Improved Pairwise Alignment of Genomic DNA* (Pennsylvania State Univ., 2007).
94. Morgulis, A., Gertz, E. M., Schäffer, A. A. & Agarwala, R. WindowMasker: window-based masker for sequenced genomes. *Bioinformatics* **22**, 134–141 (2006).
95. Chin, C.-S. et al. Multiscale analysis of pangenomes enables improved representation of genomic diversity for repetitive and clinically relevant genes. *Nat. Methods* <https://doi.org/10.1038/s41592-023-01914-y> (2023).
96. Frankish, A. et al. GENCODE 2021. *Nucleic Acids Res.* **49**, D916–D923 (2021).
97. Armstrong, J. et al. Progressive Cactus is a multiple-genome aligner for the thousand-genome era. *Nature* **587**, 246–251 (2020).
98. Kovaka, S. et al. Transcriptome assembly from long-read RNA-seq alignments with StringTie2. *Genome Biol.* **20**, 278 (2019).
99. Stanke, M., Diekhans, M., Baertsch, R. & Haussler, D. Using native and syntenically mapped cDNA alignments to improve de novo gene finding. *Bioinformatics* **24**, 637–644 (2008).
100. Fiddes, I. T. et al. Comparative Annotation Toolkit (CAT)—simultaneous clade and personal genome annotation. *Genome Res.* **28**, 1029–1038 (2018).
101. Shumate, A. & Salzberg, S. L. Liftoff: accurate mapping of gene annotations. *Bioinformatics* **37**, 1639–1643 (2021).
102. Dale, R. K., Pedersen, B. S. & Quinlan, A. R. Pybedtools: a flexible Python library for manipulating genomic datasets and annotations. *Bioinformatics* **27**, 3423–3424 (2011).
103. Rhie, A. et al. Towards complete and error-free genome assemblies of all vertebrate species. *Nature* **592**, 737–746 (2021).
104. Pruitt, K. D. et al. RefSeq: an update on mammalian reference sequences. *Nucleic Acids Res.* **42**, D756–D763 (2014).
105. Kapustin, Y., Souvorov, A., Tatusova, T. & Lipman, D. S. Splein: algorithms for computing spliced alignments with identification of paralogs. *Biol. Direct* **3**, 20 (2008).
106. Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Mol. Biol. Evol.* **30**, 772–80 (2013).
107. Slater, G. S. C. & Birney, E. Automated generation of heuristics for biological sequence comparison. *BMC Bioinformatics* **6**, 31 (2005).
108. Zook, J. M. et al. Integrating human sequence data sets provides a resource of benchmark SNP and indel genotype calls. *Nat. Biotechnol.* **32**, 246–251 (2014).
109. Numanagić, I. et al. Fast characterization of segmental duplications in genome assemblies. *Bioinformatics* **34**, i706–i714 (2018).
110. Benson, G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res.* **27**, 573–580 (1999).
111. Arian, F. A. S., Hubley, R. & Green, P. RepeatMasker Open-4.0 2013-2015. <http://www.repeatmasker.org> (2015).
112. Storer, J., Hubley, R., Rosen, J., Wheeler, T. J. & Smit, A. F. The Dfam community resource of transposable element families, sequence models, and genome annotations. *Mob. DNA* **12**, 2 (2021).
113. Olson, D. & Wheeler, T. ULTRA: a model based tool to detect tandem repeats. *ACM BCB* **2018**, 37–46 (2018).
114. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841–842 (2010).
115. Storer, J. M., Hubley, R., Rosen, J. & Smit, A. F. A. Curation guidelines for de novo generated transposable element families. *Curr. Protoc.* **1**, e154 (2021).
116. Kent, W. J. BLAT—the BLAST-like alignment tool. *Genome Res.* **12**, 656–664 (2002).
117. Szak, S. T. et al. Molecular archeology of L1 insertions in the human genome. *Genome Biol.* **3**, research0052.1 (2002).
118. Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990).
119. Cer, R. Z. et al. Searching for non-B DNA-forming motifs using nBMST (non-B DNA motif search tool). *Curr. Protoc. Hum. Genet.* **73**, 18.71–18.72 (2012).
120. Zou, X. et al. Short inverted repeats contribute to localized mutability in human somatic cells. *Nucleic Acids Res.* **45**, 11213–11221 (2017).
121. Svetec Miklenić, M. et al. Size-dependent antirecombinogenic effect of short spacers on palindromic recombination. *DNA Repair* **90**, 102848 (2020).
122. Sahakyan, A. B. et al. Machine learning model for sequence-driven DNA G-quadruplex formation. *Sci. Rep.* **7**, 14535 (2017).
123. Hao, Z. et al. Rlideogram: drawing SVG graphics to visualize and map genome-wide data on the ideograms. *PeerJ Comput. Sci.* **6**, e251 (2020).
124. Dotmatics. GraphPad Prism v9.1.0 for Windows. <https://www.graphpad.com> (16 March 2021).
125. Vollger, M. R. SafFire. *GitHub* <https://github.com/mrvollger/SafFire> (2022).
126. Pendleton, A. L. et al. Comparison of village dog and wolf genomes highlights the role of the neural crest in dog domestication. *BMC Biol.* **16**, 64 (2018).
127. Hach, F. et al. mrsFAST: a cache-oblivious algorithm for short-read mapping. *Nat. Methods* **7**, 576–577 (2010).
128. Escalona, M. et al. Whole-genome sequence and assembly of the Javan gibbon (*Hylobates moloch*). *J. Hered.* **114**, 35–43 (2023).
129. Cortez, D. et al. Origins and functional evolution of Y chromosomes across mammals. *Nature* **508**, 488–493 (2014).
130. Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* **30**, 1312–1313 (2014).
131. Dotmatics. Geneious v2019.2.3. <https://www.geneious.com/> (2019).
132. Rambaut et al. FigTree v1.4.4. <http://tree.bio.ed.ac.uk/software/figtree/> (2018).
133. Tyler-Smith, C. & Brown, W. R. A. Structure of the major block of alphoid satellite DNA on the human Y chromosome. *J. Mol. Biol.* **195**, 457–470 (1987).
134. Shepelev, V. A. et al. Annotation of suprachromosomal families reveals uncommon types of alpha satellite organization in pericentromeric regions of hg38 human genome assembly. *Genomics Data* **5**, 139–146 (2015).
135. Lee, I. et al. Simultaneous profiling of chromatin accessibility and methylation on human cell lines with nanopore sequencing. *Nat. Methods* **17**, 1191–1199 (2020).
136. Krumsiek, J., Arnold, R. & Rattei, T. Gepard: a rapid and sensitive tool for creating dotplots on genome scale. *Bioinformatics* **23**, 1026–1028 (2007).
137. Rice, P., Longden, I. & Bleasby, A. EMBOSS: The European Molecular Biology Open Software Suite. *Trends Genet.* **16**, 276–277 (2000).
138. Sun, C. et al. Deletion of azoospermia factor (AZFa) region of human Y chromosome caused by recombination between HERV15 proviruses. *Hum. Mol. Genet.* **9**, 2291–2296 (2000).
139. Lassmann, T. Kalign 3: multiple sequence alignment of large datasets. *Bioinformatics* **36**, 1928–1929 (2020).
140. Wheeler, T. J. & Eddy, S. R. nhmmer: DNA homology search with profile HMMs. *Bioinformatics* **29**, 2487–2489 (2013).
141. Stephens, Z. D. et al. Simulating next-generation sequencing datasets from empirical mutation and sequencing models. *PLoS ONE* **11**, e0167047 (2016).
142. Bushnell, B. B. BMap: a fast, accurate, splice-aware aligner. *OSTI.gov* <https://www.osti.gov/biblio/1241166> (2017).
143. Aken, B. L. et al. Ensembl 2017. *Nucleic Acids Res.* **45**, D635–D642 (2017).
144. Poznik, G. D. et al. Sequencing Y chromosomes resolves discrepancy in time to common ancestor of males versus females. *Science* **341**, 562–565 (2013).
145. McKenna, A. et al. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).
146. Schatz, M. C. et al. Inverting the model of genomics data sharing with the NHGRI Genomic Data Science Analysis, Visualization, and Informatics Lab-space. *Cell Genomics* **2**, 100085 (2022).
147. Danecek, P. et al. Twelve years of SAMtools and BCFtools. *GigaScience* **10**, giab008 (2021).
148. Talenti, A. & Prendergast, J. nf-LO: a scalable, containerized workflow for genome-to-genome lift over. *Genome Biol. Evol.* **13**, evab183 (2021).
149. Guaracino, A., Mwaniki, N., Marco-Sola, S., & Garrison, E. wfmash: whole-chromosome pairwise alignment using the hierarchical wavefront algorithm. *GitHub* <https://github.com/ekg/wfmash> (2021).
150. Sherry, S. T., Ward, M. & Sirotkin, K. dbSNP—database for single nucleotide polymorphisms and other classes of minor genetic variation. *Genome Res.* **9**, 677–679 (1999).
151. Landrum, M. J. et al. ClinVar: improving access to variant interpretations and supporting evidence. *Nucleic Acids Res.* **46**, D1062–D1067 (2018).
152. Bunioello, A. et al. The NHGRI-EBI GWAS Catalog of published genome-wide association studies, targeted arrays and summary statistics 2019. *Nucleic Acids Res.* **47**, D1005–D1012 (2019).
153. Van der Auwera, G. A. & O'Connor, B. D. *Genomics in the Cloud: Using Docker, GATK, and WDL in Terra* (O'Reilly Media, 2020).
154. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat. Methods* **9**, 357–359 (2012).
155. Ramirez, F. et al. deepTools2: a next generation web server for deep-sequencing data analysis. *Nucleic Acids Res.* **44**, W160–W165 (2016).
156. Zhao, H. et al. CrossMap: a versatile tool for coordinate conversion between genome assemblies. *Bioinformatics* **30**, 1006–1007 (2014).
157. Marçais, G. et al. MUMmer4: a fast and versatile genome alignment system. *PLoS Comput. Biol.* **14**, e1005944 (2018).
158. Ondov, B. D., Bergman, N. H. & Phillippy, A. M. Interactive metagenomic visualization in a Web browser. *BMC Bioinformatics* **12**, 385 (2011).
159. Rhie, A. Repositories for the analysis of T2T-Y and T2T-CHM13v2.0. *Zenodo* <https://doi.org/10.5281/zenodo.8136598> (2023).
160. Falconer, E. et al. DNA template strand sequencing of single-cells maps genomic rearrangements at high resolution. *Nat. Methods* **9**, 1107–1112 (2012).

Acknowledgements We thank P. Hallast, M. C. Loftus, M. K. Konkel, P. Ebert, T. Marschall and C. Lee for coordination and discussions. J.C.-I. Lee for sharing the GRCh38-Y coordinates used in Y-Finder and members of the Telomere-to-Telomere consortium and HPRC for constructive feedback. This work utilized the computational resources of the National Institutes of Health (NIH) HPC Biowulf cluster (<https://hpc.nih.gov>). Computational resources were partially provided by the e-INFRA CZ project (no. 90140), supported by the Ministry of Education, Youth and Sports of the Czech Republic and Computational Biology Core, Institute for Systems Genomics, University of Connecticut. Certain commercial equipment, instruments and materials are identified to specify adequately experimental conditions or reported results. Such identification does not imply recommendation or endorsement by the NIST, nor does it imply that the equipment, instruments or materials identified are necessarily the best available for the purpose. We thank the Intramural Research Program of NHGRI, NIH no. HG200398 (A.R., S.N., S.K., M.R., A.M.M., B.P.W. and A.M.P.); NIH no. GM123312 (S.J.H., P.G.S.G., G.A.H. and R.J.O.); NIH no. GM130691 (P.M., M.H.W. and K.D.M.); HHMI Hanna Gray Fellowship (N.A.); NIH no. CA266339 (J.G. and T.P.); NIH no. GM147352 (G.A.L.); NIH nos. HG002939 and HG010136 (R.M.H. and J.M.S.); NIH no. HG009190 (P.W.H., A. Gershman and W.T.); NIH nos. HG010263, HG006620 and CA253481 and NSF no. DBI-1627442 (M.C.S.); NIH no. GM136684 (K.D.M.); NIH nos. HG011274 and HG010548 (K.H.M.); NIH nos. HG010961 and HG010040 (H.L.); NIH no. HG007234 (M.D.); NIH no. HG011758 (F.J.S.); NIH no. DAO47638 (E.G.); NIH no. GM124827 (M.A.W.); NIH no. GM133747 (R.C.M.); NIH no. CA240199 (R.J.O.); NIH nos. HG002385, HG010169 and HG010971 (E.E.E.); Stowers Institute for Medical Research (J.L.G. and T.P.); National Center for Biotechnology Information of the National Library of Medicine, NIH (F.T.-N. and T.D.M.); intramural funding at NIST (J.M.Z.); NIST no. 70NANB20H206 (M.J.); and NIH nos. HG010972 and WT22155/Z/20/Z and the European Molecular Biology Laboratory (J.A., P.F., C.G.G., L.H., T.H., S.E.H., F.J.M. and L.S.). RNA generation was supported by NIST no. 70NANB21H101 and NIH no. 1S10OD02587; the Ministry of Science and Higher Education of the Russian Federation, St. Petersburg State University, no. PURE 73023672 (I.A.A.); the Computation, Bioinformatics, and Statistics Predoctoral Training Program awarded to Penn State by the NIH (A.C.W.); and Achievement Rewards for College Scientists Foundation, The Graduate College at Arizona State University (A.M.T.O.). E.E.E. is an investigator for HHMI.

Article

Author contributions V.A.S. is retired from the Institute of Molecular Genetics. Assembly was carried out by S.N., S.K. and M.R. Validation was performed by A.R., S.K., M.A., A.V.B., G.F., A.F., A.M.M., J.M., A.M., L.F.P., D.P., F.J.S., K.S., P.M., J.M.Z. and K.D.M. ChrY haplogroups were determined by A.R. and A.C.W. Alignment was done by C.-S.C., M.D., R. Harris, M.R.V. and K.D.M. Satellite annotation was performed by N.A., I.A.A., G.A.L., F.R., V.A.S. and K.H.M. N.A., J.G. and T.P. carried out FISH. Repeat annotation was done by S.J.H., P.G.S.G., G.A.H., R.M.H., J.M.S. and R.J.O. Retro-elements were dealt with by R. Halabian and W.M. Non-B DNA was dealt with by M.H.W. and K.D.M. Gene annotation was undertaken by A.R., M.D., P.F., C.G.G., L.H., M.H., J.H., T.H., F.J.M., T.D.M., S.L.S., A.S. and F.T.-N. A.R., R. Harris, W.T.H., P.M., M.T. and K.D.M. dealt with ampliconic genes. Structural annotation was performed by A.R., M.C., H.L., P.M. and K.D.M. Epigenetic analysis was performed by A.R., P.W.H., A. Gershman, W.T. and A.M.W. Mappability was performed by A.M.T.O., M.A.W. and J.M.Z. Non-B DNA was dealt with by M.H.W. and K.D.M. Variants and liftover were carried out by A.R., D.J.T., S.K., J.A., N.-C.C., M.D., E.G., A. Guarracino, N.F.H., W.T.H., S.E.H., S.H., R.C.M., N.D.O., M.E.G.S., L.S., M.R.V., S.Z., J.M.Z., E.E.E. and A.M.P. A.R., S.L.S., B.P.W. and A.M.P. dealt with contamination. Data generation was carried out by M.J., R.K.K., A.P.L., J.K.L., C.M., B.M.M., K.M.M., H.E.O., F.J.S. and Y.Z. Data management was undertaken by A.R., M.D., M.J. and J.K.L. Computational resources were sourced by R.J.O., M.C.S. and A.M.P. A.R., S.N., M.C., S.J.H., D.J.T., N.A., I.A.A., N.-C.C., E.G., J.G., P.G.S.G., A. Guarracino, R. Halabian, W.M., J.M., T.P., F.R., S.L.S., J.M.S., A.M.T.O., A.C.W., M.A.W., S.Z., J.M.Z., E.E.E., R.J.O., M.C.S., K.H.M., K.D.M. and A.M.P. wrote the manuscript draft. A.R. and A.M.P. edited the manuscript, with the assistance of all authors. J.M.Z., E.E.E., R.J.O., M.C.S.,

K.H.M., K.D.M. and A.M.P. supervised the research. Conceptualization was the responsibility of A.R., S.N., M.C., E.E.E., K.H.M., K.D.M. and A.M.P.

Competing interests S.N. is now an employee of ONT. S.K. has received travel funding for speaking at events hosted by ONT. A.F. is an employee of DNAnexus. C.-S.C. is an employee of GeneDX Holdings Corp. N.-C.C. is an employee of Exai Bio. L.F.P. receives research support from Genetech. F.J.S. receives research support from Pacific Biosciences, ONT, Illumina and Genetech. K.S. is an employee of Google LLC and owns Alphabet stock as part of the standard compensation package. W.T. has two patents (nos. 8,748,091 and 8,394,584) licensed to ONT. E.E.E. is a scientific advisory board member of Variant Bio, Inc. The remaining authors declare no competing interests.

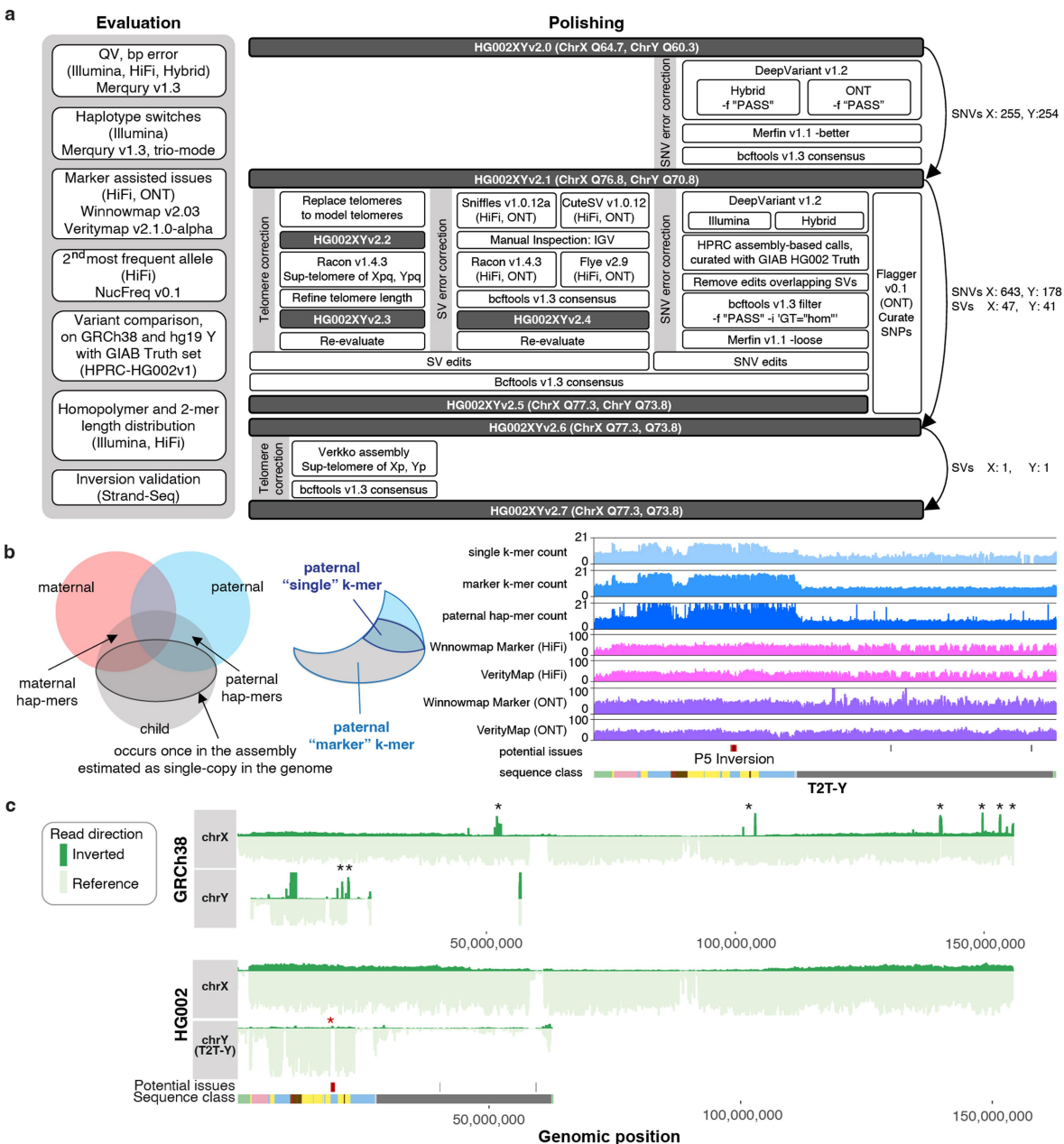
Additional information

Supplementary information The online version contains supplementary material available at <https://doi.org/10.1038/s41586-023-06457-y>.

Correspondence and requests for materials should be addressed to Adam M. Phillippy.

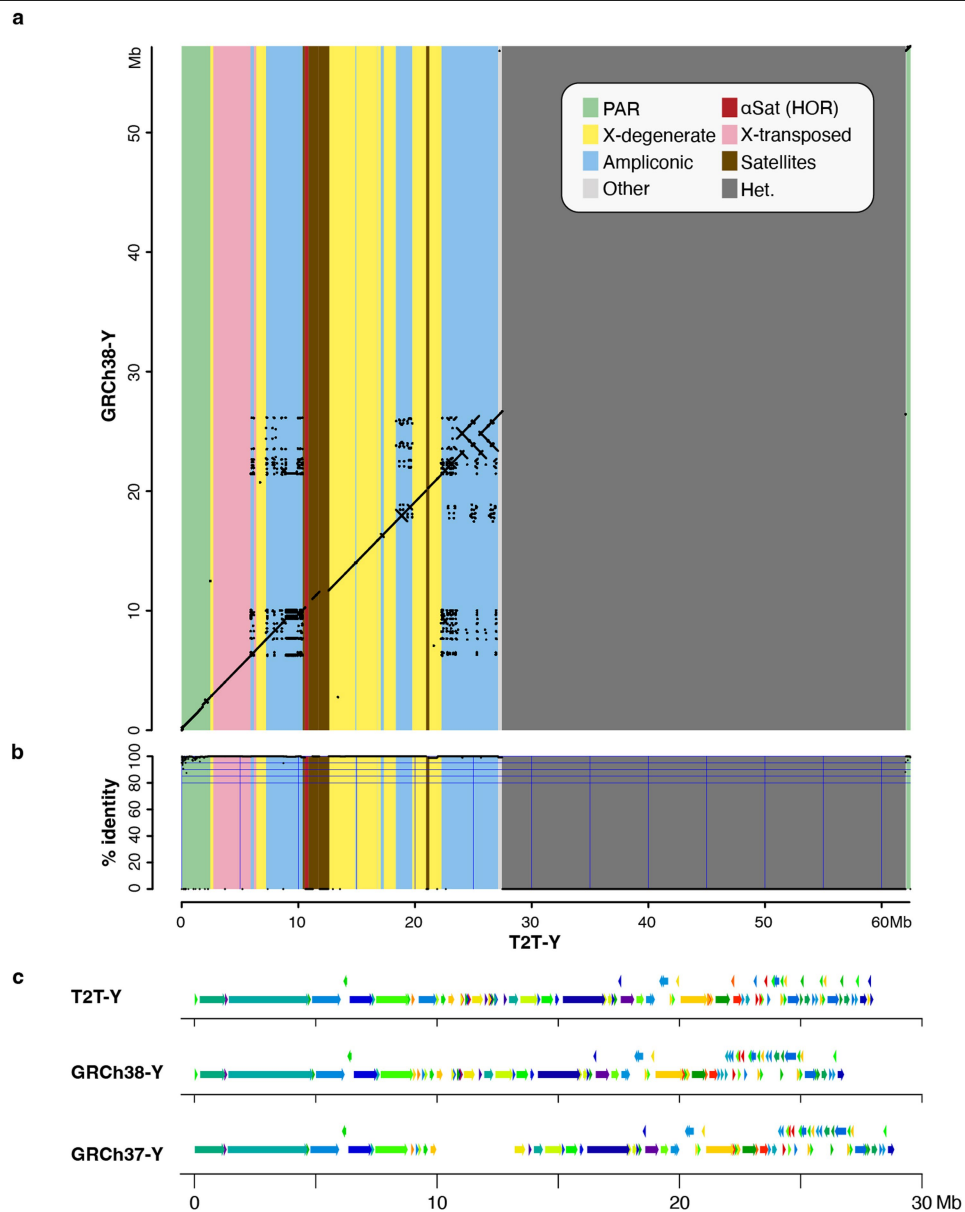
Peer review information *Nature* thanks John Lovell, Mikkel Heide Schierup and the other, anonymous, reviewer(s) for their contribution to the peer review of this work. Peer reviewer reports are available.

Reprints and permissions information is available at <http://www.nature.com/reprints>.



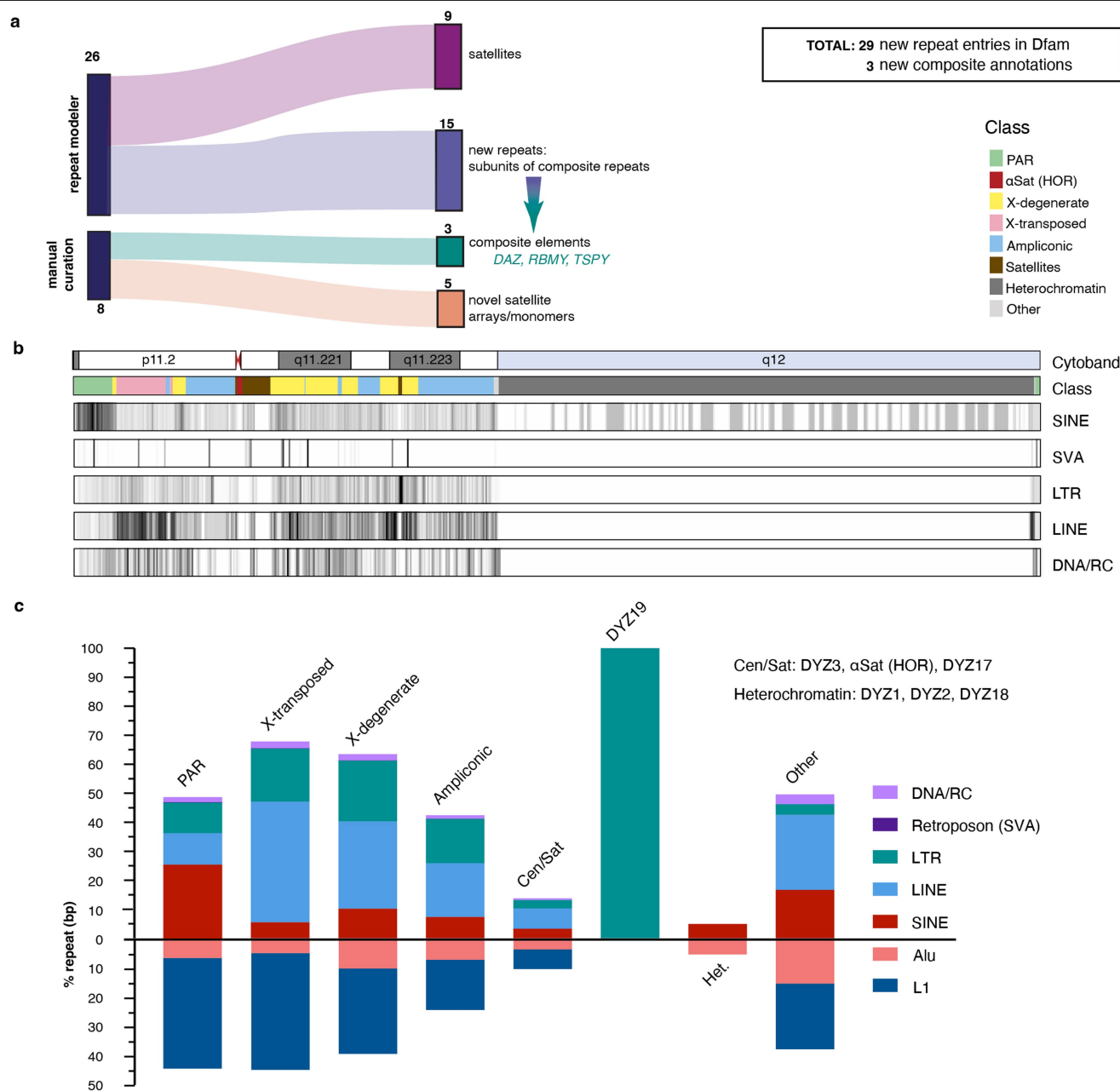
Extended Data Fig. 2 | Validation and polishing of the T2T-Y. a. Evaluation and polishing workflow performed on T2T-CHM13v1.1 autosomes + HG002 XY assemblies. **b.** Venn diagram of the k-mers from the parents and child. On the left, hap-mers¹⁸ represent haplotype specific k-mers inherited by the child. The darker outlined circle inside the child k-mers represent single-copy k-mers (k-mers occurring once in the assembly and single-copy in the child's genome). Right figure shows an example of the paternal specific, "single-copy" and "marker" k-mers. The marker set includes both multi-copy and single-copy k-mers specific to the paternal haplotype that were inherited by the child. Unlike polishing the nearly haploid CHM13 assembly¹⁷, both single-copy k-mers and marker k-mers were used for the marker-assisted alignments to HG002 XY. This helped align more reads within repetitive regions to the correct chromosome for evaluation during polishing. Right panel shows counts of the k-mers and coverage of HiFi and ONT reads using the marker-assisted Winnowmap2

alignment, in addition to alignments from VerityMap, which uses locally unique k-mers for anchoring the reads. **c.** Aggregated Strand-seq coverage profile across all 65 libraries on GRCh38-Y (top) and T2T-Y (bottom). Each bar represents read counts in every 20 kb bin supporting the reference in forward direction (light green) or reverse direction (dark green). Multiple spikes in reverse direction (black asterisks) in GRCh38-Y indicate inversion polymorphisms relative to HG002, likely due to differences between the haplogroups. Such spikes in coverage are not observed on T2T-X and T2T-Y, which confirm the structural and directional accuracy of the HG002 assemblies. A 3 kb inversion of the unique sequence between the P5 palindromic arms was identified as erroneous in T2T-Y (red asterisk), but was confirmed to be polymorphic in the population and left uncorrected in this version of the assembly.



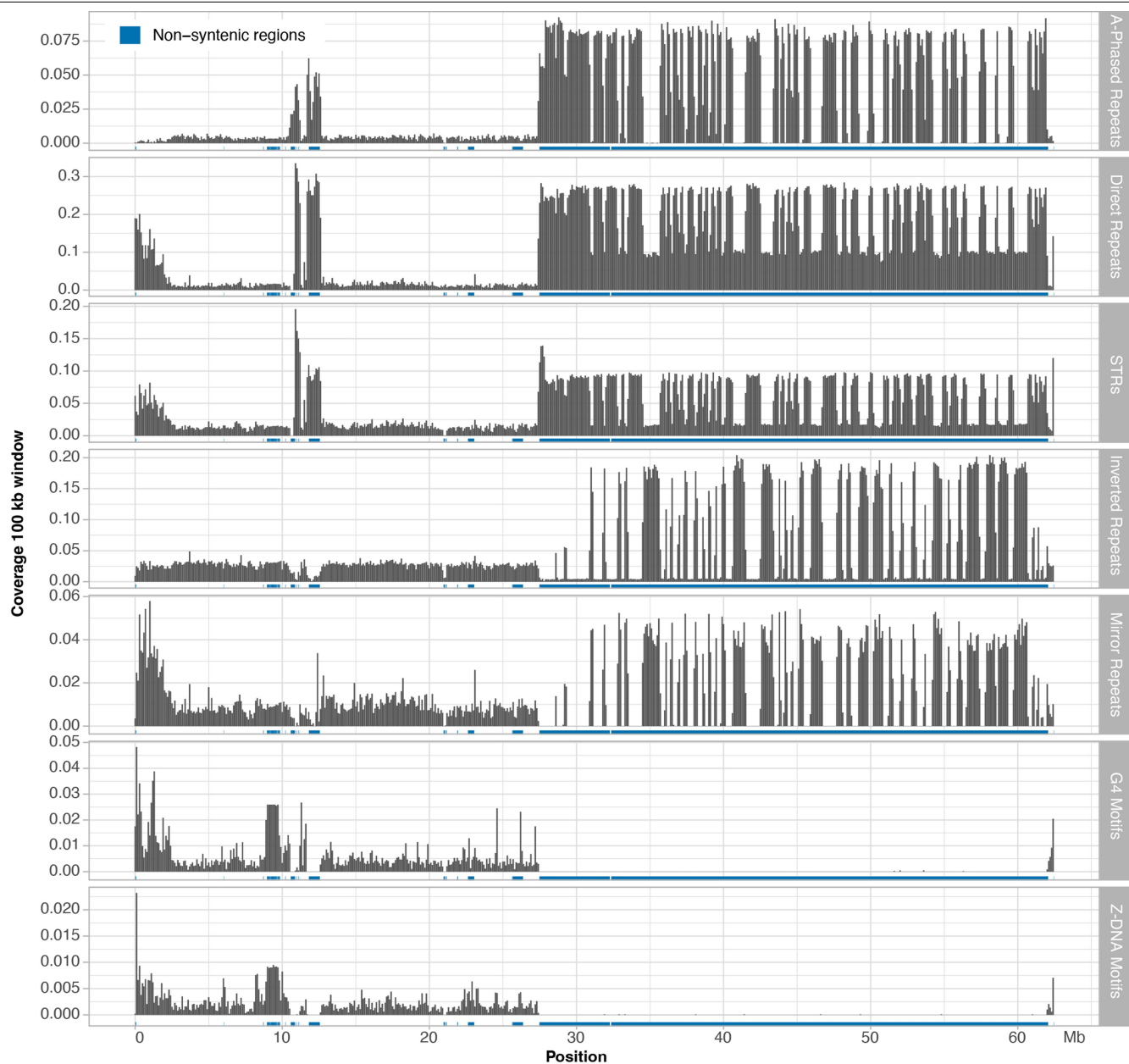
Extended Data Fig. 3 | Large structural differences between T2T-Y and previous GRCh Y assemblies. a-b. Ampliconic genes and X-degenerate sequences revealed from alignments between GRCh38-Y (Y-axis) and T2T-Y (X-axis). **a.** Dotplot generated using LastZ⁹³ after softmasking with

WindowMasker⁹⁴. **b.** Identity was computed from matches and mismatches over positions with alignments, excluding gaps. **c.** Structural differences revealed using PRG-TK⁹⁵ against GRCh38-Y and GRCh37-Y in the euchromatic region of the Y chromosome.



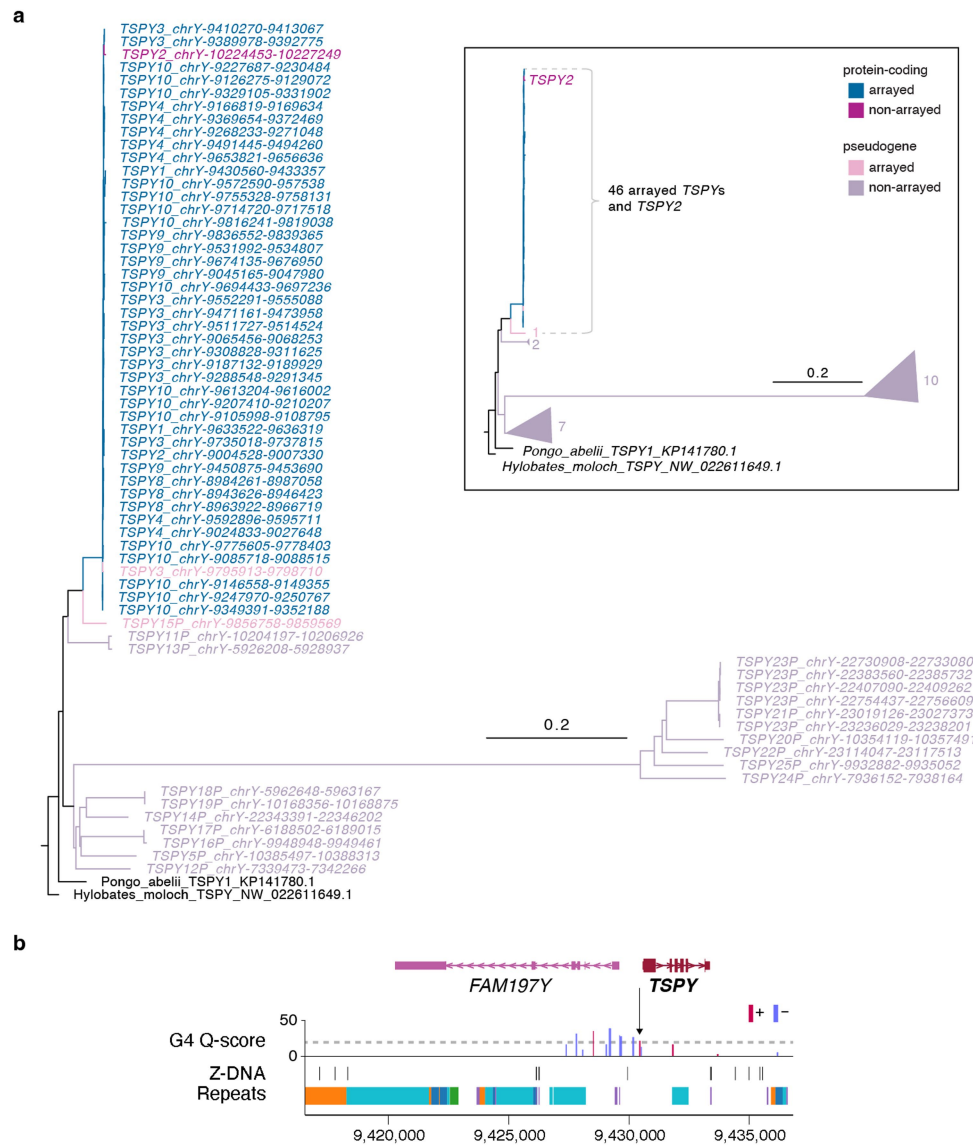
Extended Data Fig. 4 | Repeat discovery and annotation of T2T-Y.
a. Assembly completion allowed for a full assessment of repeats and resulted in the identification of previously unknown satellite arrays (predominantly in the PAR1) and subunit repeats that fall within one of three composite repeat units (*TSPY*, *RBMY*, *DAZ*). **b.** Ideogram of TE density (per 100 kb bin). This is an extension of Fig. 1 with non-SINEs expanded into separate TE classes (SVA, LTR, LINE, DNA/RC). Density scale ranges from low (white, zero) to high

(black, relative to total density) and sequence classes are denoted by color. **c.** Summary (in terms of base coverage per region) across all five TE classes and two specific families: *Alu*/SINE and L1/LINE. The satellites in **(b)** were kept separate as two categories; Cen/Sat as the left satellite block including alpha satellites and DYZ19, while all other categories were combined per sequence classes.



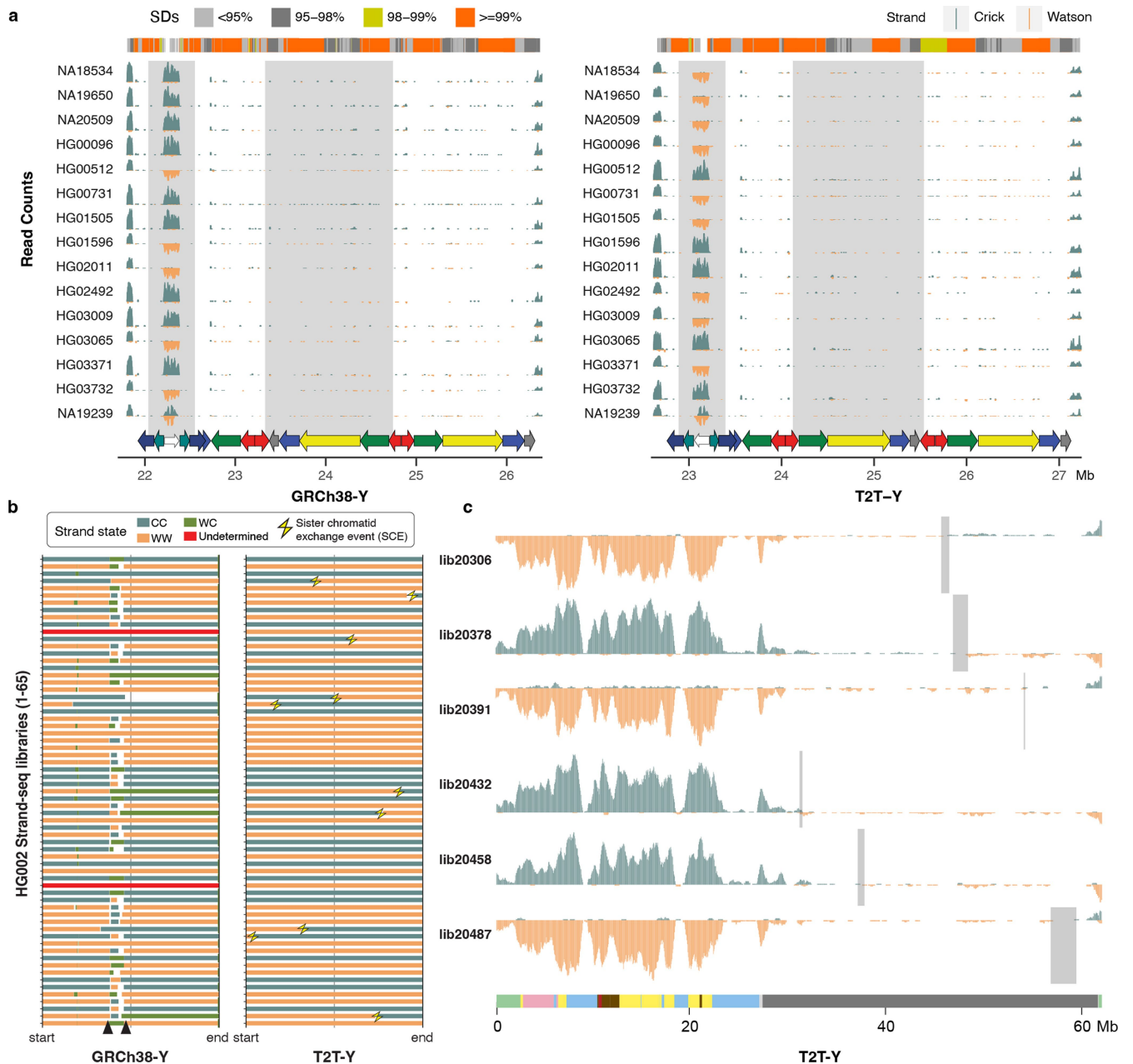
Extended Data Fig. 5 | Non-B DNA motifs along the T2T-Y. HSat3 on the Yq and satellite sequences around the centromere are more enriched with A-phased repeats, direct repeats and STRs, while HSat1B is more enriched with inverted

repeats and mirror repeats. Enrichment of non-B DNA sequences were also observed in the PAR region. Notably, the *TSPY* gene array is enriched for G4 and Z-DNA motifs, as shown in Extended Data Fig. 6b.



Extended Data Fig. 6 | Phylogenetic tree analysis of the ampliconic *TSPY* gene family and pattern of non-B DNA structure. a. Phylogenetic tree analysis using protein-coding *TSPY*s from a Sumatran Orangutan (*Pongo abelii*) and a Silvery gibbon (*Hylobates moloch*) as outgroups confirmed *TSPY2* (distal to the array) and *TSPY* copies within the array originated from the same branch, distinguished from the rest of the *TSPY* pseudogenes. Rectangular inset shows a cartoon representation of the simplified tree. Numbers next to the triangles

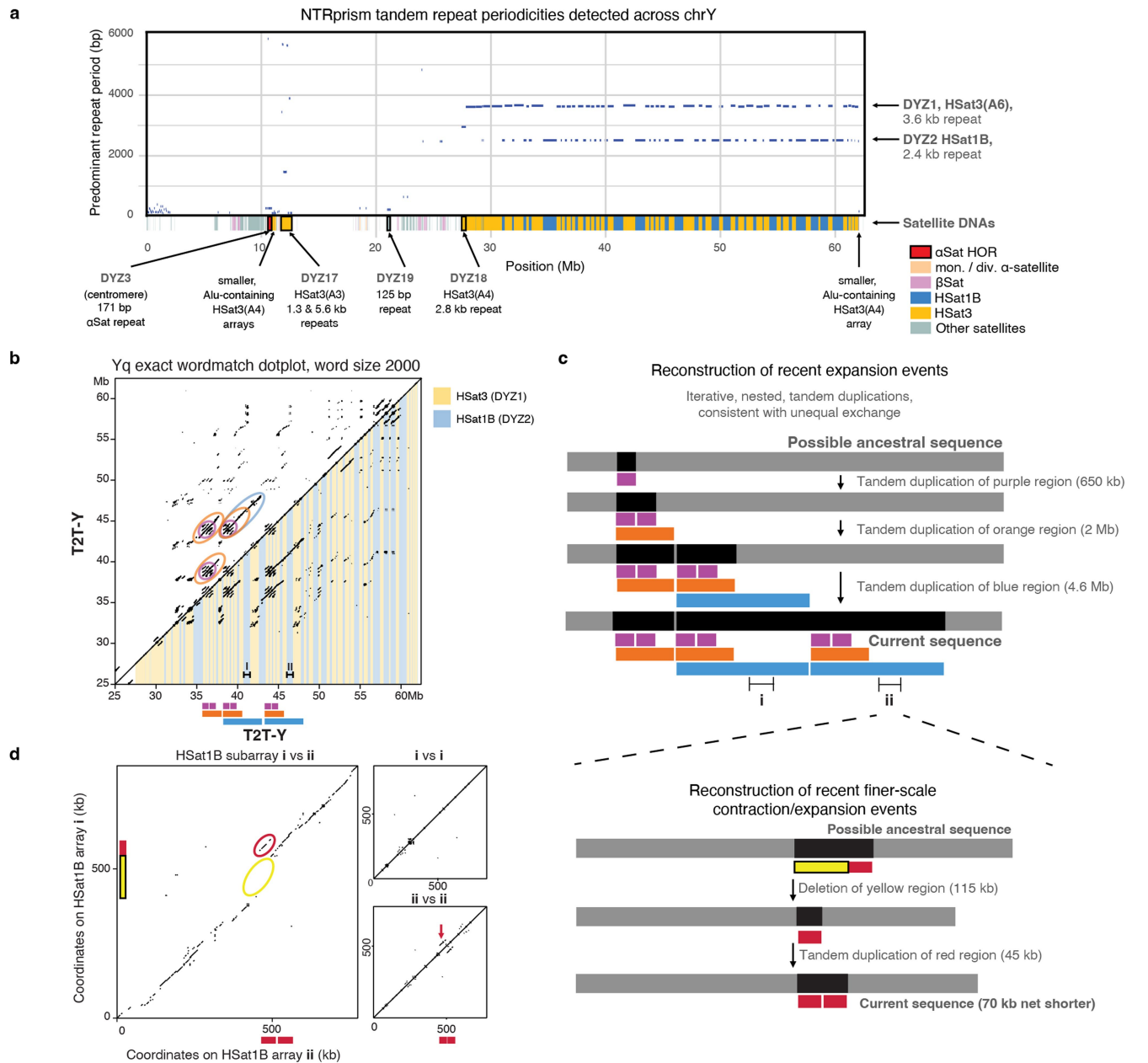
indicate the number of *TSPY* genes in the same branch. **b.** G4 and Z-DNA structures predicted for a typical *TSPY* copy inside the *TSPY* array. All *TSPY* copies in the array have the same signature, with one G4 peak present -500 bases upstream of the *TSPY* (arrow). Higher Quadron score¹²² (Q-score) indicates a more stable G4 structure, with scores over 19 considered stable (dotted line).



Extended Data Fig. 7 | Recurrent inversions identified with Strand-seq.

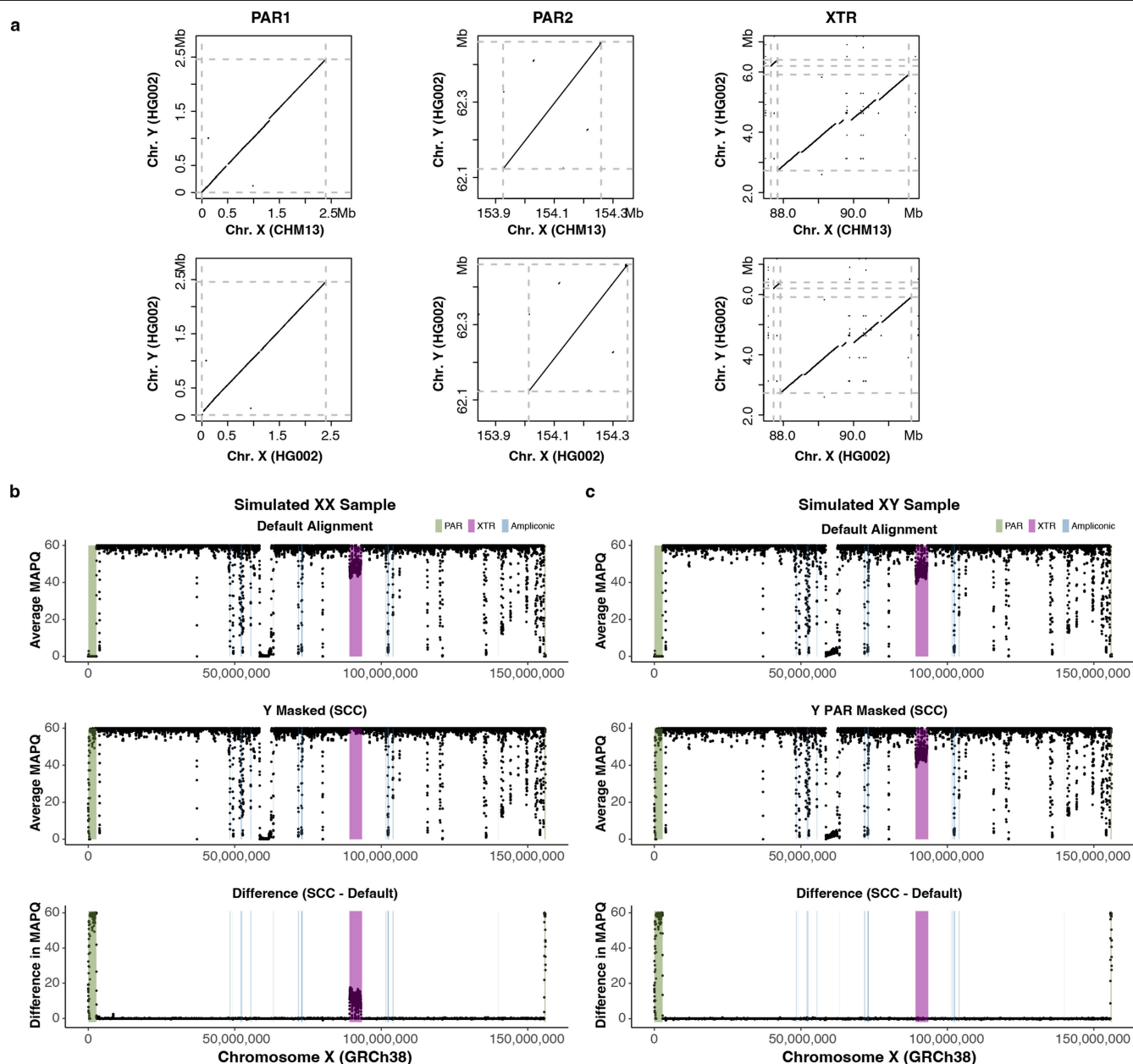
a. Five out of 15 individuals have the inverted variant as present in HG002 at the P3 palindrome (white arrow). Although inversions across P1–P2 (yellow and red arrows) are difficult to confirm with Strand-seq because of the high sequence similarity between the palindromic arms, different orientations are observable in these samples. **b.** Strand states for 65 Strand-seq libraries of HG002. Depending on the mappings of directional Strand-seq reads (+ reads: 'Crick', C, – reads: 'Watson', W), reference sequence was assigned in three states: WC, WW, and CC. WC, roughly equal mixture of plus and minus reads; WW, all reads mapped in minus orientation; CC, all reads mapped in plus orientation. Changes in strand state along a single chromosome are normally caused by a double-strand-break (DSBs) that occurred during DNA replication¹⁶⁰ in a random fashion and we refer to them as sister-chromatid-exchanges (SCEs, yellow

thunderbolts). Recurrent change in strand state over the same region in multiple Strand-seq cells indicates misassembly. Similarly, collapsed or incomplete assembly of a certain genomic region will result in a recurrent strand state change as observed for GRCh38-Y (black arrowheads). In contrast, T2T-Y shows strand state changes randomly distributed along each Strand-seq library with no evidence of misassembly or collapse. **c.** Strand-seq profile of selected libraries over T2T-Y summarized in bins (bin size: 500 kb, step size: 50 kb). Teal, Crick read counts; orange, Watson read counts. As ChrY is haploid, reads are expected to map only in Watson or Crick orientation. Light gray rectangles highlight regions where SCEs were detected in the heterochromatic Yq12 despite a lower coverage of Strand-seq reads. A modified breakpointR parameter was used (windowsize = 500000 minReads = 20) in order to refine detected SCEs presented in panel **b** and **c**.



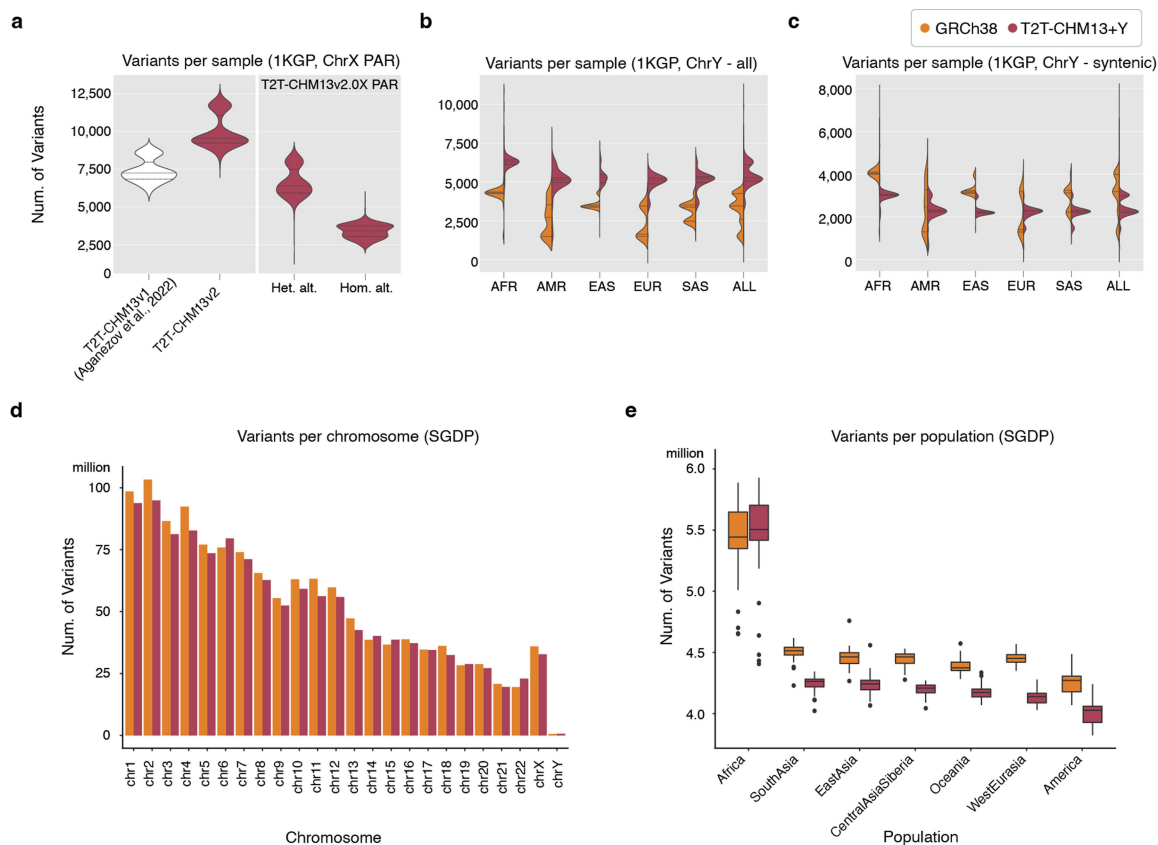
Extended Data Fig. 8 | Satellite annotation and recent expansion events in the Yq heterochromatin. **a.** A plot showing the top repeat periodicities detected by NTRprism⁴⁴ in 50 kb blocks tiled across T2T-Y, with centromeric satellite annotations overlaid on the X axis. Large arrays are labeled with their historic nomenclature¹, HSat subfamilies⁶¹, and predominant repeat periodicities. **b.** An exact 2000-mer match dotplot of the Yq region (a dot is plotted when an identical 2000 base sequence is found at positions X and Y). The lower triangle has DYZ1/DYZ2 annotations overlaid as yellow and blue bars, respectively. Circled patterns in the upper triangle correspond to

recent iterative duplication events, which are illustrated below the X axis. **c.** A reconstruction of a possible sequence of recent iterative duplications that could explain the observed dotplot patterns. **d.** A 2000-mer dotplot comparison of two ~800 kb HSat1B sub-arrays that were part of a recent large duplication event, along with self-self comparisons of the same arrays, revealing sites of more recent and smaller-scale deletions and expansions (annotated in yellow and red, with a possible sequence of events illustrated by the schematic on the right).



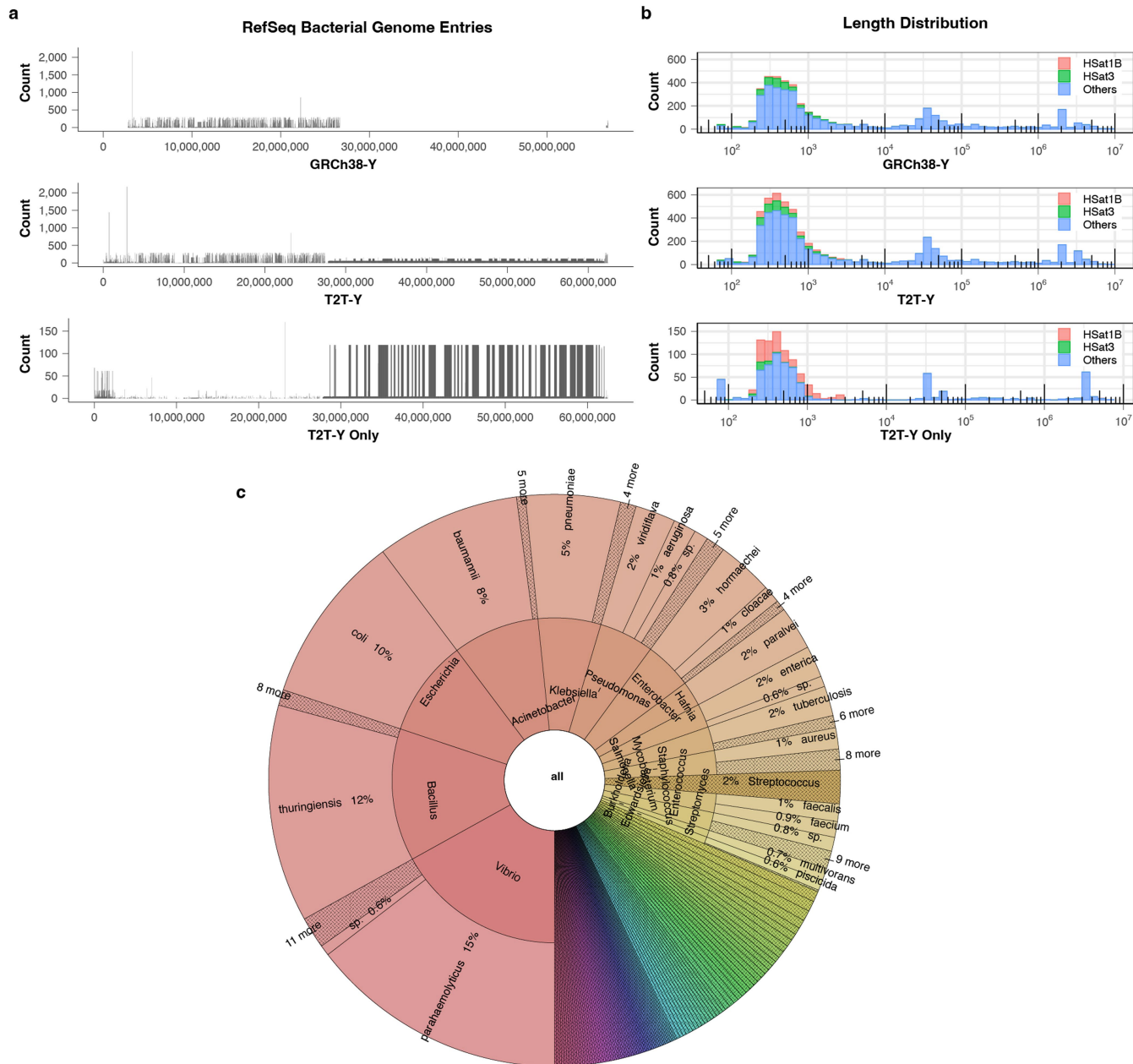
Extended Data Fig. 9 | Genomic similarity in PARs and XTR and improved MAPQ of the PARs through informed sex chromosome complement reference. **a.** Dotplots from LASTZ alignments of the CHM13-X, HG002-X, and HG002-Y (T2T-Y) over 96% sequence identity. Dashed gray lines represent the start and end of the approximate PARs or XTR boundaries. Disconnected diagonal lines indicate the presence of genomic diversity between each paired region. More genomic differences are observed in the PAR1 between the HG002-Y and CHM13-X. **b-c.** Average mapping quality (MAPQ) across

GRCh38-X from simulated reads of an XX (**b**) and XY (**c**) sample. Top, a default version of GRCh38 (with two copies of identical PARs on XY). Middle, a version of GRCh38 informed on the sex chromosome complement (SCC) of the sample (entire Y hard-masked for the XX sample vs. only PARs on the Y hard-masked for the XY sample). Bottom, the difference in average MAPQ between the SCC and default approaches. MAPQ was averaged in 50 kb windows, sliding 10 kb across the chromosome. A positive value means MAPQ score is higher with SCC reference alignment compared to default alignment.



Extended Data Fig. 10 | Number of variants called from 1KGP and SGDP individuals. **a.** More variants are called on the X-PARs when using the sex chromosome complement reference approach (calling variants in diploid mode on PARs) than the non-masked approach (calling variants in haploid mode on PARs). The 1KGP results for GRCh38-Y are from Aganezov et al.⁶⁶, which was performed on CHM13v1.0+GRCh38-Y. **b.** Num. of variants called from each 1KGP XY sample on chromosome GRCh38-Y and T2T-Y. **c.** Num. of variants called in the syntenic region between the two Ys. A large num. of additional variants are called on each sample attributed to the newly added, non-syntenic sequences on T2T-Y. Within the syntenic regions, a reduction in

the number of variants is observed for each population except for samples from R1 haplogroups as shown in Fig. 6c. **d.** Aggregated total number of variants for the 279 SGDP samples per chromosome. **e.** SGDP genome-wide counts of variants per-sample ($n = 279$) demonstrate increased variation in African samples regardless of reference. Each bar in the box plot represents the 1st, 2nd (median), and 3rd quartile of the number of variants in each population. Whiskers are bound to the $1.5 \times$ interquartile range. Data outside of the whisker ranges are shown as dots. For the SGDP samples, variants were called using T2T-CHM13+Y or GRCh38 as the reference. All variants shown in this figure were filtered for “high quality (PASS)”.



Extended Data Fig. 11 | Human contaminants in bacterial reference genomes. **a.** Number of distinct RefSeq accessions in every 10 kb window containing 64-mers of GRCh38-Y (top), T2T-Y (middle), and in T2T-Y only (bottom). Here, RefSeq sequences with more than 20 64-mers or matching over 10% of the Y chromosome are included. **b.** Length distribution of the

sequences from (a) in log scale. Majority of the shorter (<1 kb) sequences contain 64-mers found in HSat1B or HSat3. **c.** Number of bacterial RefSeq entries by strain identified to contain sequences of T2T-Y and not GRCh38-Y, visualized with Krona¹⁵⁸.

Reporting Summary

Nature Portfolio wishes to improve the reproducibility of the work that we publish. This form provides structure for consistency and transparency in reporting. For further information on Nature Portfolio policies, see our [Editorial Policies](#) and the [Editorial Policy Checklist](#).

Statistics

For all statistical analyses, confirm that the following items are present in the figure legend, table legend, main text, or Methods section.

n/a Confirmed

- ☒ ☐ The exact sample size (n) for each experimental group/condition, given as a discrete number and unit of measurement
- ☐ ☒ A statement on whether measurements were taken from distinct samples or whether the same sample was measured repeatedly
- ☒ ☐ The statistical test(s) used AND whether they are one- or two-sided
Only common tests should be described solely by name; describe more complex techniques in the Methods section.
- ☒ ☐ A description of all covariates tested
- ☒ ☐ A description of any assumptions or corrections, such as tests of normality and adjustment for multiple comparisons
- ☐ ☒ A full description of the statistical parameters including central tendency (e.g. means) or other basic estimates (e.g. regression coefficient) AND variation (e.g. standard deviation) or associated estimates of uncertainty (e.g. confidence intervals)
- ☒ ☐ For null hypothesis testing, the test statistic (e.g. F , t , r) with confidence intervals, effect sizes, degrees of freedom and P value noted
Give P values as exact values whenever suitable.
- ☒ ☐ For Bayesian analysis, information on the choice of priors and Markov chain Monte Carlo settings
- ☒ ☐ For hierarchical and complex designs, identification of the appropriate level for tests and full reporting of outcomes
- ☒ ☐ Estimates of effect sizes (e.g. Cohen's d , Pearson's r), indicating how they were calculated

Our web collection on [statistics for biologists](#) contains articles on many of the points above.

Software and code

Policy information about [availability of computer code](#)

Data collection

AWS command line interface (aws-cli) v1.14.36 was used to retrieve data from the s3://human-pangenomics bucket.

Data analysis

Custom codes developed for data analysis and visualization are available at <https://github.com/arangrhie/T2T-HG002Y>, https://github.com/snurk/sg_sandbox, and <https://github.com/schatzlab/t2t-chm13-chry>. Software and parameters used are stated in the Supplementary Methods with more details.

All sequencing data generated were downloaded from s3://human-pangenomics/. ONT R9 UL reads were re-basecalled with Guppy v6.1.2. Throughout all analysis, bwa v0.7.17, Winnowmap v2.03, samtools v1.9, minimap2 v2.17-26 and bedtools v2.20.0-v2.30.0 were used for aligning short or long reads.

Source code and scripts used for assembly graph construction, pruning, semi-automated repeat resolution and consensus can be found in https://github.com/snurk/sg_sandbox (commit ver. 19ee5e306f83f8eb5f5a6ac6a3477e2f925b375e), which also utilized GraphAligner v1.0.13, and later released in Verkko. HiFi gaps were patched with Flye v2.7-b1585. Note that after finishing the T2T-Y assembly, the entire assembly procedure was re-engineered and updated in the Verkko assembler.

Polishing was performed using codes and scripts in <https://github.com/arangrhie/T2T-Polish> v1.0 release.

We used Genomescope2, Meryl v1.3, Merqury v1.3, DeepVariant v1.2, Merfin v1.1, Racon v1.4.3 and v1.6.0, Flagger v0.1, bcftools v1.10.2, hap.py v0.3.14, Sniffles v1.0.12a, cuteSV v1.0.12, SURVIVOR v1.0.7, dipcall v0.3, Flye v2.9, VerityMap v2.1.0-alpha, NucFreq v0.1, and IGV v2.14.1 for polishing and validation. For Y haplogroup identification, we used yhaplo v1.1.2 and Y-SNP haplogroup hierarchy finder (<http://forensic.mc.ntu.edu.tw:9000/DNAToolWeb/YHGSearch>, based on ISOGG Tree 11 Jul 2020 Version 15.73 release).

For alignments between GRCh38 and HG002Y, we used rustybam v0.1.29, Saffire v0.2, Rhodonite v0.12, LASTZ v1.04.15 and PRG-TK v0.3.4. Gene annotation was performed with Cactus v2.0.5, Stringtie2 commit 647ab51, Liftoff v1.6.1~v1.6.3, BUSCO v4.1.4, Splign v2.1.0, ProSplign and Gnomon from NCBI C++ ToolKit r645952, MAFFT v7.475 (2020/Nov/23), and Exonerate v2.2.0.

Iso-Seq alignments were performed using uLTRA v0.0.4.1, deSALT v1.5.5, cDNA_cupcake v28.0.0.

Repeats were annotated using RepeatMasker v4.1.2-p1-v4.1.3, RepeatModeler v2.0.1, TRF v409.linux64, ULTRA v1.0, NTRprism v0.22-v1.0.0, kalign v3.3.2, HMMER with HMM v3.3.2, GraphPad Prism v9.1.0, and BLAT v36.5. Transduction analysis was performed with TSDfinder v1.0, RepeatMasker v4.1.2-p1, and BLAST v2.11.0. Non-B DNAs were annotated using nBMST commit 1c8f963 and Quadron commit 19047e3). Methylation analysis was performed with Nanopolish v0.13.2, modbam2bed v0.6.2, primrose v1.3.0, pbmm2 v1.9.0, pb-CpG-tools v1.1.0, Bismark v0.23.1dev, Sequence class annotation utilized Gepard v2.1 and LASTZ v1.04.00. TSPY copy number analysis used fastCN v0.2, mrsFAST v3.4.2, RepeatMasker v4.1.2-p1. Phylogenetic tree analysis of the TSPY and AluYs were performed with MAFFT v7.471, RAXML-NG v0.9.0, Geneious v2019.2.3, and FigTree v1.4.4. Short-read variant calling benchmark and pipeline to call variants from the 1KGP and SGDP samples are using GATK v4.2.1.0-v4.2.4, XYalign v1.15, bcftools v1.16, IGVtools v2.14.1. Source code is available at <https://github.com/schatzlab/t2t-chm13-chry>, released as v1.0.0. For creating CHM13-GRCh38 chain file, nf-LO v1.5.1, rustybam v0.1.29, paf2chain commit f68eeca, chaintools v0.1, SafFire commit aa16e43, wfmash commit a36ab5f, rustybam commit f68eeca were used. For lifting over dbSNP, ClinVar, and GWAS resources, GATK v4.1.1.9 and picard v2.23.3 was used. ENCODE data were generated with Bowtie2 v2.4.1, samtools v1.10, Picard v2.22.1, deepTools2 v3.4.3, and MACS2 v2.2.7.1. Genome Aggregation Database (gnomAD) release v3.1.2 was lifted over from GRCh38 to each assembly using CrossMap v0.6.1, and a nextflow pipeline utilizing Ensembl Variant Effect Predictor (VEP) <https://github.com/Ensembl/ensembl-vep/tree/release/108/nextflow> commit cb84684. Human Y chromosome contamination analysis was performed with MUMmer v4.0, Meryl v1.3, seqrequester (commit r95 fa5bdac1), and Kronatools v2.8.

For manuscripts utilizing custom algorithms or software that are central to the research but not yet described in published literature, software must be made available to editors and reviewers. We strongly encourage code deposition in a community repository (e.g. GitHub). See the Nature Portfolio [guidelines for submitting code & software](#) for further information.

Data

Policy information about [availability of data](#)

All manuscripts must include a [data availability statement](#). This statement should provide the following information, where applicable:

- Accession codes, unique identifiers, or web links for publicly available datasets
- A description of any restrictions on data availability
- For clinical datasets or third party data, please ensure that the statement adheres to our [policy](#)

The T2T-CHM13v2.0 (T2T-CHM13+Y) assembly, reference analysis set, complete list of resources including gene annotation, repeat annotation, epigenetic profiles, variant calling results from 1KGP and SGDP, gnomAD, ClinVar, GWAS, and dbSNP datasets are available for download at <https://github.com/marbl/CHM13>. The assembly is also available from NCBI and EBI with GenBank accession GCA_009914755.4. Annotation and associated resources are also browsable as “hs1” from the UCSC Genome Browser http://genome.ucsc.edu/cgi-bin/hgTracks?db=hub_3671779_hs1, the Ensembl Genome Browser <https://projects.ensembl.org/hprc/> (assembly name T2T-CHM13v2.0) and NCBI data-hub https://www.ncbi.nlm.nih.gov/data-hub/genome/GCF_009914755.1/. Potential assembly issues are listed and tracked at <https://github.com/marbl/CHM13-issues>. 1KGP and SGDP short read alignments and variant calls are available within AnVIL at https://anvil.terra.bio/#workspaces/anvil-datastorage/AnVIL_T2T_CHRY. Sequencing data used in this study is listed in Supplementary Table 1.

Custom codes developed for data analysis and visualization are available at <https://github.com/arangrhie/T2T-HG002Y>, https://github.com/snurk/sg_sandbox, and <https://github.com/schatzlab/t2t-chm13-chry> and deposited on Zenodo along with <https://github.com/marbl/CHM13> and <https://github.com/marbl/CHM13-issues> (<https://doi.org/10.5281/zenodo.8136598>). Software and parameters used are stated in the Supplementary Methods with more details.

Human research participants

Policy information about [studies involving human research participants and Sex and Gender in Research](#).

Reporting on sex and gender

Population characteristics

Recruitment

Ethics oversight

Note that full information on the approval of the study protocol must also be provided in the manuscript.

Field-specific reporting

Please select the one below that is the best fit for your research. If you are not sure, read the appropriate sections before making your selection.

☒ Life sciences ☐ Behavioural & social sciences ☐ Ecological, evolutionary & environmental sciences

For a reference copy of the document with all sections, see nature.com/documents/nr-reporting-summary-flat.pdf

Life sciences study design

All studies must disclose on these points even when the disclosure is negative.

Sample size	One sample was chosen for generating complete sequences of the X and Y chromosomes. To overcome challenges in genome assembly, a well characterized reference sample, HG002 (GM24385), was chosen based on availability of its previously characterized genomic variation data and its use as a National Institute of Standards and Technology (NIST) reference material.
Data exclusions	No data was excluded.
Replication	Replication has been performed for the Yq satellite (DYZ1 and DYZ2) painting. The results of this experiment were successfully replicated using two different sets of PCR probes. Fifteen large-field images containing at least 20 spreads were analyzed per condition. Assembly integrity was also successfully confirmed with PacBio HiFi and ONT reads obtained from HG002 and its paternal cell line HG003 (GM24149), sequenced at different time points.
Randomization	Randomization is not applicable to the assembly of a reference genome because the genomic material is derived from one cell line (GM24385).
Blinding	Blinding is not applicable to this study. The reference material has been fully consented for its genomic data release. Known genotypes from HG002 had to be used for validating the accuracy of the assembly.

Reporting for specific materials, systems and methods

We require information from authors about some types of materials, experimental systems and methods used in many studies. Here, indicate whether each material, system or method listed is relevant to your study. If you are not sure if a list item applies to your research, read the appropriate section before selecting a response.

Materials & experimental systems

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> Antibodies
<input type="checkbox"/>	<input checked="" type="checkbox"/> Eukaryotic cell lines
<input checked="" type="checkbox"/>	<input type="checkbox"/> Palaeontology and archaeology
<input checked="" type="checkbox"/>	<input type="checkbox"/> Animals and other organisms
<input checked="" type="checkbox"/>	<input type="checkbox"/> Clinical data
<input checked="" type="checkbox"/>	<input type="checkbox"/> Dual use research of concern

Methods

n/a	Involved in the study
<input checked="" type="checkbox"/>	<input type="checkbox"/> ChIP-seq
<input checked="" type="checkbox"/>	<input type="checkbox"/> Flow cytometry
<input checked="" type="checkbox"/>	<input type="checkbox"/> MRI-based neuroimaging

Eukaryotic cell lines

Policy information about [cell lines and Sex and Gender in Research](#)

Cell line source(s)	Genome in a bottle reference material HG002 cell lines (GM24385, GM26105 and GM27730) were purchased from Coriell Institute and used for sequencing and generating data used in this study. HG002 DNA is available as a reference material from NIST, and the associated cell lines have been previously consented for both research use and commercial redistribution. More details can be found at https://www.nist.gov/programs-projects/genome-bottle and https://www.coriell.org/1/NIGMS/Collections/NIST-Reference-Materials .
Authentication	Cell lines and DNA were obtained directly from authoritative sources (NIST, Coriell, and PGP) and the authenticity subsequently confirmed by comparing assembly-based variant calls to the HG002 GIAB truth set and karyotyping.
Mycoplasma contamination	The cell lines were not tested for mycoplasma contamination. The final product of this study (T2T-Y assembly) was tested for contamination, none was identified except the EBV used for immortalization, which was found as an external chromosomal component, as expected.
Commonly misidentified lines (See ICLAC register)	No commonly misidentified cell lines were used in the study.