

RESEARCH ARTICLE SUMMARY

HUMAN GENOMICS

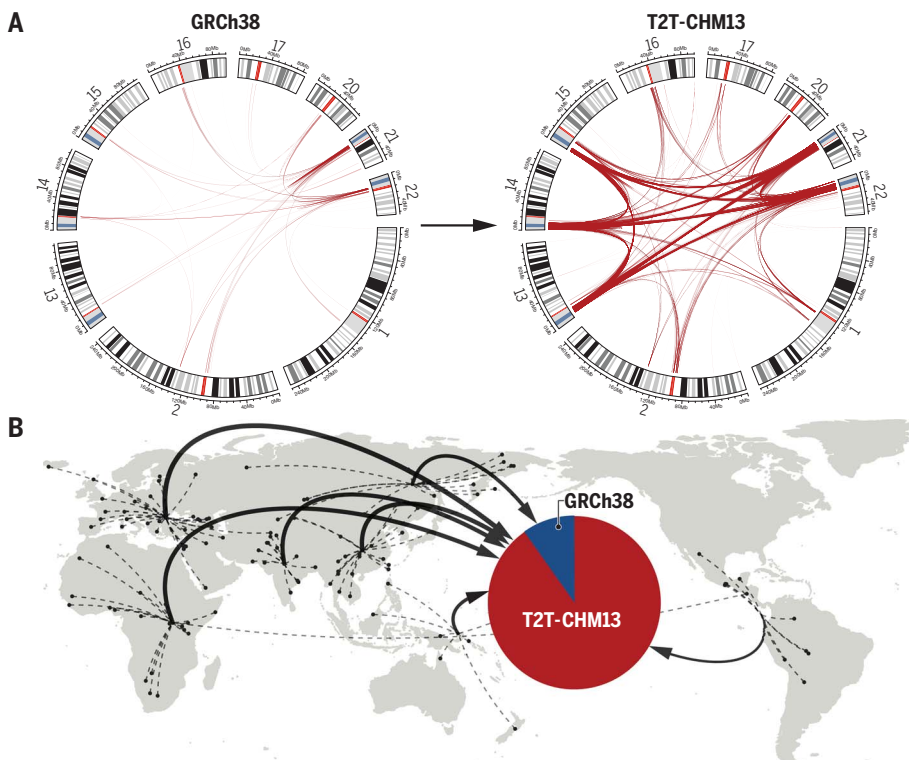
Segmental duplications and their variation in a complete human genome

Mitchell R. Vollger, Xavi Guitart, Philip C. Dishuck, Ludovica Mercuri, William T. Harvey, Ariel Gershman, Mark Diekhans, Arvis Sulovari, Katherine M. Munson, Alexandra P. Lewis, Kendra Hoekzema, David Porubsky, Ruiyang Li, Sergey Nurk, Sergey Koren, Karen H. Miga, Adam M. Phillippy, Winston Timp, Mario Ventura, Evan E. Eichler*

INTRODUCTION: Large, high-identity duplicated sequences—termed segmental duplications (SDs)—are frequently the last regions of genomes to be sequenced and assembled. While the human reference genome provided a road-map of the SD landscape, >50% of the remaining gaps correspond to regions of complex SDs.

RATIONALE: SDs are major sources of evolutionary gene innovations and contribute disproportionately to genetic variation within and between ape species. With the complete human genome (T2T-CHM13), researchers have the potential to identify genes and uncover patterns of human genetic variation.

RESULTS: We identified 51 million base pairs (Mbp) of additional human SD in T2T-CHM13 and now estimate that 7% of the human genome consists of SDs [(218 Mbp of 3.1 billion base pairs (Gbp)]. SDs make up two-thirds (45.1 of 68.1 Mbp) of acrocentric short arms, and these SDs are the largest in the human genome (see the figure, panel A). Additionally, 54% of acrocentric SDs are copy number variable or map to different chromosomes among the six individuals examined. A detailed comparison between the current reference genome (GRCh38) and T2T-CHM13 for SD content identifies 81 Mbp of previously unresolved or structurally variable SDs. Short-read whole-



More-complete segmental duplication content improves genotyping. (A) Increase (by a factor of 10) in the number of large (>10 kilo-base pairs) acrocentric segmental duplications (red) in T2T-CHM13 (right) compared with GRCh38 (left). (B) Copy number genotyping based on read-depth from 268 diverse human genomes across the globe shows that 90% of new SDs in T2T-CHM13 (red) are more likely to reflect human copy number when compared to GRCh38 (blue) irrespective of human population group considered.

genome sequence data from a diversity panel of 268 humans show that human copy number is nine times (59.26 versus 6.55 Mbp) more likely to match T2T-CHM13 rather than GRCh38, including 119 protein-coding genes (see the figure, panel B). Using long-read-sequencing data from 25 human haplotypes, we investigated patterns of human genetic variation identifying significant increases in structural and single-nucleotide diversity. We identified gene-rich regions (e.g., *TBC1D3*) that vary by hundreds of kilo-base pairs and gene copy number between individuals showing some of the highest genome-wide structural heterozygosity (85 to 90%). Our analysis identified 182 candidate protein-coding genes as well as the complete sequence for structurally variable gene models that were previously unresolved. Among these is the complete gene structure of lipoprotein A (*LPA*), including the expanded kringle IV repeat domain. Reduced copies of this domain are among the strongest genetic associations with cardiovascular disease, especially among African Americans, and sequencing of multiple human haplotypes identified not only copy number variation but also other forms of rare coding variation potentially relevant to disease risk. Finally, we compared global methylation and expression patterns between duplicated and unique genes. Transcriptionally inactive duplicate genes are more likely to map to hypomethylated genomic regions; however, specifically over the transcription start site we observe an increase in methylation, suggesting that as many as two-thirds of duplicated genes are epigenetically silenced. Additionally, SD genes show a high degree of concordance between methylation profiles and transcription levels, allowing us to define the actively transcribed members of high-identity gene families that are otherwise indistinguishable by coding sequence.

CONCLUSION: A complete human genome provides a more comprehensive understanding of the organization, expression, and regulation of duplicated genes. Our analysis reveals underappreciated patterns of human genetic diversity and suggests characteristic features of methylation and gene regulation. This resource will serve as a critical baseline for improved gene annotation, genotyping, and previously unknown associations for some of the most dynamic regions of our genome. ■

The list of author affiliations is available in the full article online.
*Corresponding author. Email: eee@gs.washington.edu
Cite this article as M. R. Vollger et al., *Science* 376, eabj6965 (2022). DOI: 10.1126/science.abj6965

S READ THE FULL ARTICLE AT
<https://doi.org/10.1126/science.abj6965>

RESEARCH ARTICLE

HUMAN GENOMICS

Segmental duplications and their variation in a complete human genome

Mitchell R. Vollger¹, Xavi Guitart¹, Philip C. Dishuck¹, Ludovica Mercuri², William T. Harvey¹, Ariel Gershman³, Mark Diekhans⁴, Arvis Sulovari¹, Katherine M. Munson¹, Alexandra P. Lewis¹, Kendra Hoekzema¹, David Porubsky¹, Ruiyang Li¹, Sergey Nurk⁵, Sergey Koren⁵, Karen H. Miga⁴, Adam M. Phillippy⁵, Winston Timp³, Mario Ventura², Evan E. Eichler^{1,6*}

Despite their importance in disease and evolution, highly identical segmental duplications (SDs) are among the last regions of the human reference genome (GRCh38) to be fully sequenced. Using a complete telomere-to-telomere human genome (T2T-CHM13), we present a comprehensive view of human SD organization. SDs account for nearly one-third of the additional sequence, increasing the genome-wide estimate from 5.4 to 7.0% [218 million base pairs (Mbp)]. An analysis of 268 human genomes shows that 91% of the previously unresolved T2T-CHM13 SD sequence (68.3 Mbp) better represents human copy number variation. Comparing long-read assemblies from human ($n = 12$) and nonhuman primate ($n = 5$) genomes, we systematically reconstruct the evolution and structural haplotype diversity of biomedically relevant and duplicated genes. This analysis reveals patterns of structural heterozygosity and evolutionary differences in SD organization between humans and other primates.

Genomic duplications have long been recognized as important sources of structural change and gene innovation (1, 2). In humans, the most recent and highly identical sequences [$>90\%$ and >1 kilo-base pairs (kbp)]—referred to as segmental duplications (SDs) (3)—promote meiotic unequal crossover events that contribute to recurrent rearrangements associated with $\sim 5\%$ of developmental delay and autism (4). These same SDs are reservoirs for human-specific genes that have been important in increasing synaptic density and the expansion of the frontal cortex since humans diverged from other ape lineages (5–8). SDs are also enriched ~ 10 -fold for normal copy number variation, although most of this genetic diversity has yet to be fully characterized or associated with human phenotypes (9, 10). SD length (frequently >100 kbp), sequence identity, and extensive structural diversity among human haplotypes have hampered our ability to characterize these regions at a genomic level. This is because sequence reads have been insuffi-

ciently long and human haplotypes too structurally diverse to resolve duplicate copies or distinguish allelic variants.

One of the first human whole-genome sequence (WGS) assembly drafts created with Sanger sequencing technology was almost devoid of SDs and their underlying genes (11, 12). Similarly, bacterial artificial chromosome (BAC)-based approaches to assembling the human genome from different haplotypes led to many misjoins, creating de facto gaps that took years to resolve (13). Although combining WGS- and BAC-based data for early sequencing of human genomes provided a roadmap of the SD landscape (14), more than 50% of the gaps within the human reference genome have corresponded to regions of complex SDs.

The development of genomic resources (15–17), including BAC libraries and long-read sequence data from complete hydatidiform moles (CHM, which represent a single human haplotype), was motivated in large part by efforts to resolve the organization of these regions and concomitantly complete the human reference genome. The CHM13 cell line has the advantage of originating from a single haplotype and predominantly a single ancestral group (European) (18) in contrast to the GRCh38 reference, which is a composite representation of multiple human haplotypes and ancestries (19). These resources, combined with advances in long-read technologies, have produced the gapless human genome assembly T2T-CHM13 (20). We use this genome assembly to present a complete view of SDs in a human genome and highlight their importance in

advancing our understanding of genetic diversity, evolution, and disease in humans.

SD content and organization

We characterized the SD content of the T2T-CHM13 v1.0 assembly by sequence read-depth and pairwise sequence alignments ($>90\%$ and >1 kbp) (21). Our analysis of the assembly identifies 208 Mbp of nonredundant segmentally duplicated sequences within chromosome-level scaffolds (including 15.6 Mbp of SD located on chrY, which is included from GRCh38), compared with just 167 Mbp in the current reference (GRCh38) (Table 1 and Fig. 1). This raises the percent estimate of the human genome that is segmentally duplicated from 5.4 to 6.7%. However, five SD-related gaps remained in the initial assembly of the female CHM13 genome (T2T-CHM13 v1.0). Each corresponded to a cluster of tandemly repeated ribosomal DNA (rDNA) genes on each acrocentric chromosome where we confirmed long-read sequence pileups consistent with unresolved SDs. The estimated amount of missing rDNA sequence was calculated by Nurk *et al.* by using both digital droplet PCR (22) and a whole-genome Illumina coverage analysis (20). Assuming a canonical repeat length of 45 kbp for the rDNA molecule (23, 24), the total amount of missing sequences was approximated at ~ 10 Mbp and ~ 200 copies of unresolved rDNA sequence (20). These findings are consistent with the subsequent specialized assembly of the rDNA released as part of the T2T-CHM13 v1.1 assembly. Including this estimate, the overall SD content of the human genome is now 7.0% (6.7% not including rDNA; see Table 1 for statistics breakdown by SD type) and is likely to increase as more complete genomes of diverse origins are sequenced and assembled.

One-third (81.3 Mbp) (25) of SD sequence in T2T-CHM13 is wholly uncharacterized in GRCh38 (16.5 Mbp) or differs in copy number and structure (64.8 Mbp) (25). Most of these involve large, high-identity SDs. For example, there is a 70% increase (41,285 versus 24,280) in the number of SD pairs and a doubling of the number of bases in pairwise alignments with $>95\%$ identity (Fig. 1C). Among these previously unresolved or variable SDs, 13,258 (35.0 Mbp) map to the acrocentric short arms of chromosomes 13, 14, 15, 21, and 22 (Fig. 1B and Table 1), which are unassembled in the GRCh38. These SDs do not correspond to rDNA duplications but represent other segments predominantly shared among acrocentric ($n = 5332$ alignments) and nonacrocentric chromosomes ($n = 5500$ alignments; table S1). In particular, the pericentromeric regions of chromosomes 1, 3, 4, 7, 9, 16, and 20 show the most extensive SD homology with acrocentric DNA (Fig. 1B). Non-rDNA acrocentric SDs are 1.74 times as long as all other SDs (N50: 74,704 versus 42,842) and significantly longer

¹Department of Genome Sciences, University of Washington School of Medicine, Seattle, WA, USA.

²Department of Biology, University of Bari, Aldo Moro, Bari 70125, Italy. ³Department of Molecular Biology and Genetics, Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD, USA. ⁴UC Santa Cruz Genomics Institute, University of California Santa Cruz, Santa Cruz, CA, USA. ⁵Genome Informatics Section, Computational and Statistical Genomics Branch, National Human Genome Research Institute, National Institutes of Health, Bethesda, MD, USA.

⁶Howard Hughes Medical Institute, University of Washington, Seattle, WA, USA.

*Corresponding author. Email: eee@gs.washington.edu

Table 1. Summary statistics of segmental duplications in T2T-CHM13 and GRCh38. Mbp, number of nonredundant Mbp of SD; peri, within 5 Mbp of the heterochromatin surrounding the centromere; telo, within 500 kbp of the telomere; acro, within the short arms of the acrocentric chromosomes. Difference: SD content difference between T2T-CHM13 v1.0 and GRCh38. Previously unresolved or structurally variable: Sequence in T2T-CHM13 that does not have 1 Mbp of synteny with GRCh38. GRCh38 contains 149,690,719 bp of gap sequence included in the reported number of Gbp.

Assembly	Gbp	% SD	SD (Mbp)	# SDs	inter (Mbp)	# inter	intra (Mbp)	# intra	acro (Mbp)	# acro	peri (Mbp)	# peri	telo (Mbp)	# telo
T2T-CHM13 v1.0*	3.114	6.665	207.563	41289	121.113	30484	142.958	10805	35.106	13264	88.606	24985	10.975	4998
GRCh38	3.114	5.372	167.297	24280	83.556	16348	120.710	7932	6.624	1407	53.944	10606	8.926	1529
Difference	0.000	1.293	40.266	17009	37.556	14136	22.248	2873	28.482	11857	34.662	14379	2.049	3469
Previously unresolved or structurally variable	0.240	33.885	81.338	25161	61.873	20579	54.932	4582	35.039	13258	54.037	19607	5.616	4005
T2T-CHM13 v1.0* + rDNA estimate	3.114	6.987	217.598	66042	131.148	49213	152.993	16829	45.141	38017	98.641	49738	10.975	4998

*The version of T2T-CHM13 that was used (v1.0) included chrY from GRCh38.

(P value < 0.01, one-sided Wilcoxon rank-sum test) than any other defined SD category in the human genome (intrachromosomal, interchromosomal, pericentromeric, and telomeric; fig. S1).

We annotated all T2T-CHM13 SDs with DupMasker (26), which defines ancestral evolutionary units of duplication on the basis of mammalian outgroups and a repeat graph (27). Focusing on duplcons that carry genes or duplicated portions of genes, we identified 30 duplcons that show the greatest copy number change between T2T-CHM13 and GRCh38. These 30 genic SDs represent regions where gene annotation is most likely to change; all predicted differences favor an increase in copy number for the T2T-CHM13 assembly (Fig. 1D and table S2).

We also compared the number of SDs more directly by defining syntenic regions (5 Mbp) between GRCh38 and T2T-CHM13 (25). Of the 15 windows with the largest increase, nine mapped to the acrocentric short arms and six were in pericentromeric regions (fig. S1 and table S3). In particular, the intervals between the centromeric satellite and secondary constrictions (qh regions) on chromosomes 1, 9, and 16 show a 4.6-fold increase in the number of SDs (5254 versus 1141) and the most differences in organization compared with GRCh38. SDs in these regions are almost exclusively interchromosomal and depleted for intrachromosomal duplications (figs. S2 and S3).

Validation and heteromorphic variation

Because the acrocentric short arms as well as the qh regions on chromosomes 1, 9, and 16 were either previously unresolved or showed the most considerable differences in terms of SD content, we focused first on validating their organization. We mapped available end-sequence

data from a human fosmid genome library (28) to the T2T-CHM13 assembly and selected nine distinct clones as probes (Fig. 2A) to confirm the patterns of high-identity (>95%) SDs (25). All 30 of the distinct duplication predictions from the T2T-CHM13 SDs were corroborated by fluorescence in situ hybridization (FISH) against chromosomal metaphases of the CHM13 cell line (Fig. 2, B and C, and table S4).

FISH also revealed nine additional signals not originally predicted by our SD analysis (fig. S4). However, we were able to identify lower identity duplications that confirmed seven of these sites, leading to an overall concordance of 95% (37 of 39) between FISH and the T2T-CHM13 SD assembly content. We extended this analysis to five additional human cell lines of diploid origin because both pericentromeric and acrocentric portions of chromosomes have been shown to be cytogenetically heteromorphic (29–31). In total, we identified 61 distinct cytogenetic locations of which 28 (46%) were fixed whereas 33 (54%) were variable in their presence or absence on specific homologs (both acrocentric and pericentromeric regions of the human genome) (fig. S4). Of the 61 FISH signals, all but six were observed in more than one of the six human cell lines indicating that such heteromorphic variation is common and prevalent.

We found a correlation (Pearson's correlation coefficient, $r = 0.96$) between genome-wide copy number variation from the assembly and Illumina read-depth data generated from the same CHM13 source (25). Because SDs frequently map to the breakpoints of inversion polymorphisms (28, 32, 33), we validated 65 inversions relative to GRCh38 with single-cell DNA template strand sequencing (Strand-seq analysis) of the T2T-CHM13 assembly (figs. S5 and S6)

(25). Although 32 of these represent known human polymorphisms, 33 have not been observed in six previously analyzed human genomes (32). However, by analysis of Strand-seq data from one additional human haplotype (CHM1), we further confirmed 30 of these inversions (i.e., present in CHM1 and CHM13), suggesting that at least 95.4% (62 of 65) represent true large-scale human inversion polymorphisms (fig. S5). Consistent with previous literature (34), inversions associated with SDs ($n = 30$) are significantly longer than those not associated with SDs (P value < 0.01, one-sided Wilcoxon rank-sum test) and are polymorphic among humans (fig. S6). One notable example is an inversion polymorphism mapping to chromosome 1q21. It is a complex event consisting of two inversions (262.3 kbp and 2.26 Mbp) originally predicted by Sanders and colleagues (33) but our sequence analysis shows a relocation of 767.6 kbp of genic sequences (Fig. 2D). The large inversion (chr1:146,350,000 to 148,610,000) is flanked by the core duplcon—the *NBPF* gene family—and in combination with the other rearrangements changes the order of human-specific genes *NOTCH2NLA*, *-B*, and *-C*, which have been implicated in expansion of the frontal cortex (8, 35). As a final test, we resolved this region in eight additional human haplotypes (25), all of which support the T2T-CHM13 configuration with one exception (CHM1), which is consistent with the GRCh38 configuration (fig. S7).

Single-nucleotide and copy number variation within SDs

The high quality and single haplotype nature of both the T2T-CHM13 and GRCh38 reference genomes provides an opportunity to compare the genome-wide pattern of single-nucleotide

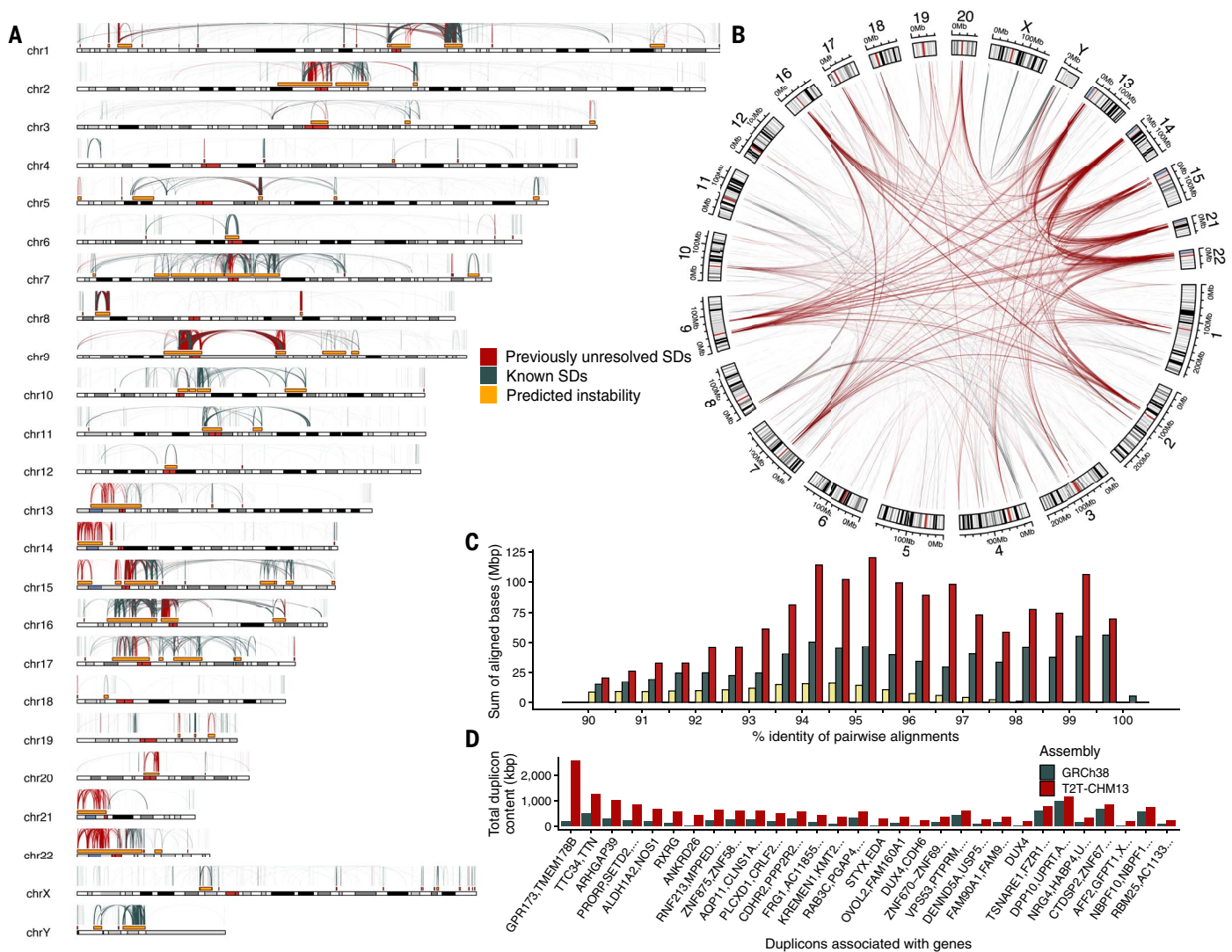


Fig. 1. Segmental duplication (SD) content of the T2T-CHM13 genome.

(A) The pattern of previously unresolved or structurally variant intrachromosomal duplications in T2T-CHM13 (red) compared with known duplications in GRCh38 (blue-gray). These predict hotspots of genomic instability (gold) flanked by large (>10 kbp), high-identity (>95%), interspersed (>50 kbp) SDs. (B) Circos plot highlighting previously unresolved interchromosomal SDs (red) and showing the preponderance of previously unresolved SDs mapping to pericentromeric and

acrocentric regions. (C) A histogram comparing SD content in different human reference genomes, including the sum of bases in pairwise SD alignments stratified by their percent identity for the celera (yellow, Sanger-based); GRCh38 (blue-gray, BAC-based); and T2T-CHM13 (red, long-read) assemblies. (D) The 30 genic duplicons (ancestral repeat units) with the greatest copy number difference between GRCh38 and T2T-CHM13 as determined by DupMasker (table S2). The 30 largest differences all exhibit increase in T2T-CHM13.

variation in regions that have typically been excluded from most analyses because of their repetitive nature. We aligned GRCh38 to T2T-CHM13 and retained only regions deemed to be “syntenic” on the basis of an unambiguous one-to-one correspondence between both reference genomes and at least 1 Mbp of aligned sequence (25).

Most unique regions of the genome (2693 Mbp) could be compared, whereas only 60% (124 Mbp) of the SDs within T2T-CHM13 exhibited a clear orthologous relationship between the two human reference genomes. As expected, the X chromosome and the region corresponding to the major histocompatibility complex (MHC) are the least and most diver-

gent, respectively (Fig. 3A), as a result of the slower rate of evolution for the female X and the deep coalescence of MHC.

Notably, SD sequences are significantly more diverged than unique sequences (P value < 0.001, one-sided Mann-Whitney U test) (fig. S8). Comparing only syntenic regions (25) between GRCh38 and T2T-CHM13, we estimate the single-nucleotide variant (SNV) density to be 0.95 SNVs/kbp for unique regions of the genome when compared with SD regions where density rises to 1.47 SNVs per kbp (table S5). This 50% increase could be a result of an increased mutation rate of SDs (e.g., a result of the action of interlocus gene conversion), or a deeper average coalescence of duplicated sequences. Another

possible explanation for this observation is erroneous alignment of paralogous instead of allelic sequences; however, we believe this is unlikely given the requirement of at least 1 Mbp of continuous, one-to-one, best alignment between GRCh38 and T2T-CHM13 (25).

As part of this analysis, we also identified regions that structurally differ or are absent from GRCh38 when compared with the T2T-CHM13 assembly. Using 1-Mbp LASTZ alignments (25), we identified 126 nonsyntenic regions for a total of 240 Mbp (N50 length of 12.7 Mbp; fig. S9). Of these, 33.9% (81.34 of 240 Mbp) overlapped SD regions. Using sequence read depth (25) from 268 human genomes [(Simons Genome Diversity Project (SGDP)],

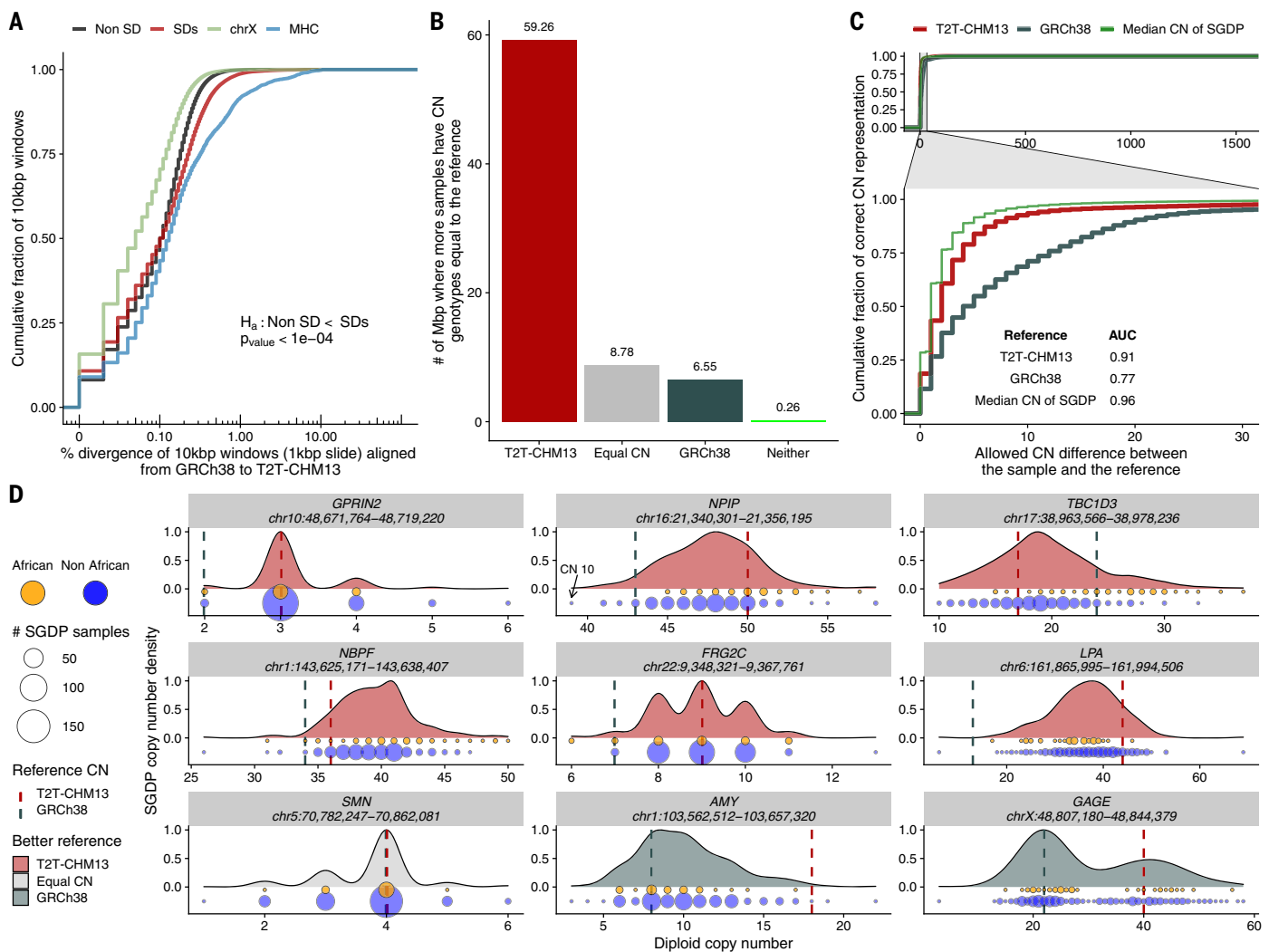


Fig. 3. SD single-nucleotide and copy number variation. (A) Sequence divergence (percent in 10-kbp bins) on the basis of syntenic alignments between GRCh38 and T2T-CHM13 for SDs (red), and unique genomic regions (black). SD regions show considerably more divergence when compared with unique sequences (black) and chromosome X (blue) but less than that of the MHC regions (green). (B) Copy number of SD regions that were previously unresolved or structurally different in T2T-CHM13 compared with GRCh38 on the basis of 268 human genomes from the Simons Genome Diversity Project (SGDP). The histogram shows the number of Mbp, in which more samples support the copy number of the given assembly [T2T-CHM13 (red), GRCh38 (blue), neither (green), or both equally (equal copy number, gray)]. (C) Empirical cumulative

distribution showing how many samples genotype correctly with either GRCh38 or T2T-CHM13 as a function of the allowed difference between sample and reference copy number. The inset shows the area under the curve (AUC) calculation for both references allowing a maximum copy number difference of 30. The green curve shows an in silico reference made using the median copy number of the SGDP samples at each site. (D) Genic copy number variation. Copy number variation in nine gene families are shown (generated with SGDP) and distribution is colored according to which reference better reflects the median copy number; GRCh38 generally underestimates copy number (vertical lines) and Africans (orange) tend to show higher copy number than non-Africans (blue); circle size indicates number of samples.

reference for copy number variation irrespective of population group (fig. S11).

Structural variation and massive evolutionary changes in the human lineage

Advances in long-read genome assembly (47, 48) enable sequence resolution of complex structural variation associated with SDs at the haplotype level (49). We generated or used existing high-fidelity (HiFi) sequence data from 12 human and five nonhuman primate genomes to understand both the structural diversity and evolution of specific SD regions.

To guide the selection of candidate regions for analysis, we constructed a hifiasm assembly of a chimpanzee genome (note that the cell line used to make the sequencing data for this genome assembly was originally made from a now deceased chimpanzee individual, Clint), compared it with the T2T-CHM13 assembly, and searched for regions of substantial structural difference between the lineages. We focused first on the largest regions of insertion on the human lineage before sub-selecting those regions that contain genes of biomedical or evolutionary importance (tables

S8 and S9). We restricted the analysis to insertions >50 kbp in length and selected 10 loci for a more detailed analysis, including genes associated with the expansion of the human frontal cortex (tables S8 and S9 and fig. S12). Assemblies of additional haplotypes recapitulated the structural organization of T2T-CHM13 for eight of the 10 loci, whereas evidence for the structural organization of GRCh38 was only found in five of the 10 loci (25). Overall, 73% of the human haplotype assemblies were successfully reconstructed (table S8); however, the fraction of human haplotypes resolved at each

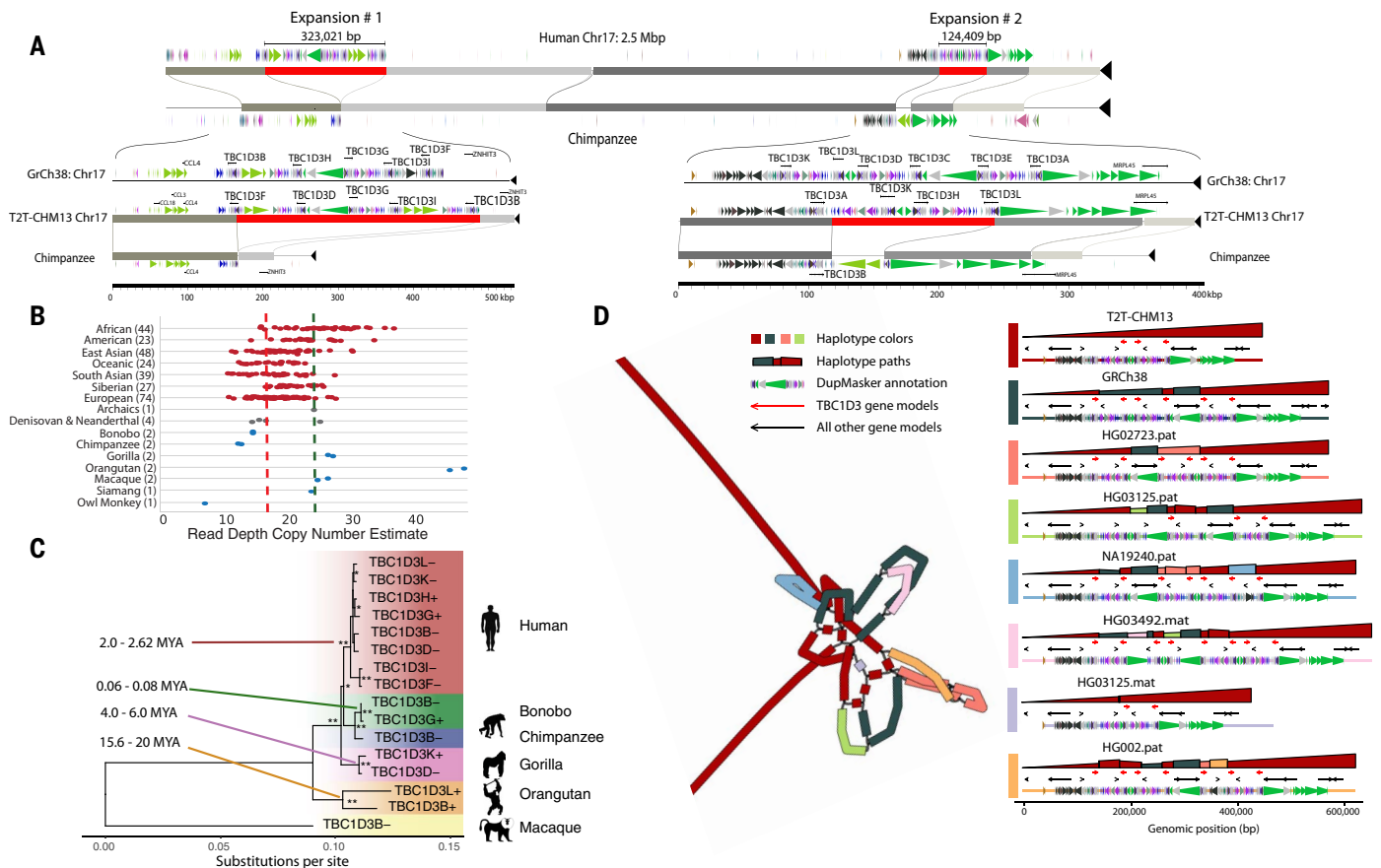


Fig. 4. Human-specific expansion of *TBC1D3* compared with nonhuman primates. (A) Regions of homology between human T2T-CHM13's chromosome 17 (top) and a HiFi assembly of the chimpanzee genome (bottom). Red blocks represent regions of human-specific expansion, including *TBC1D3* duplications. Colored arrows above and below the homologous sequence represent unique ancestral units (duplicons) identified by DupMasker. Inset plots for both expansion sites are included below with the gene models identified by Liftoff (94). (B) Copy number (diploid) estimates from an Illumina read-depth analysis of SGDP, ancient hominids, and nonhuman primates for a *TBC1D3* paralog (table S14). Copy number estimates include five pseudogenes not included in the phylogeny, explaining the higher counts observed. The T2T-CHM13 copy number and GRCh38 copy number are represented by red and blue lines, respectively. (C) Phylogeny of *TBC1D3* copies at these two expansion sites as well as nonhuman primate copies. Single asterisks at nodes indicate bootstrap values $\geq 70\%$ whereas double asterisks indicate 100%. The data illustrate a human-specific expansion as well as several

locus varied considerably depending on the size and complexity of the region (fig. S13). For example, in the case of the 8.9-Mbp region corresponding to *NOTCH2NL* and *SRGAP2B/2D*, we recovered only 37.5% of human haplotypes (table S8 and fig. S7). Similarly, we resolved only six haplotypes (from a potential of 24 haplotypes) for the 3.4-Mbp region harboring the *SMN1* and *SMN2* loci (fig. S14).

Among the haplotypes that could be resolved, we found a high degree of structural heterozygosity among human genomes (25) with 249 kbp differing on average when com-

pared with T2T-CHM13 (table S10). In some cases the structural changes are simple, such as ~ 12 kbp insertion or deletion of *CYP2D6*, which contributes to differential drug metabolism activity in addition to other human disease susceptibilities (50–56) (fig. S15). In other cases the patterns of structural variation are complex, involving hundreds of kilobase pairs of inserted or deleted gene-rich sequences along with large-scale inversion events that alter gene order for specific human haplotypes (see *ARHGAP11A/B*; fig. S16, and *NOTCH2NLA/B*; fig. S7). Furthermore, the spinal muscular atrophy

independent expansions in macaques, gorillas, and orangutans. Using a macaque sequence as an outgroup, we estimate the human-specific expansion to be ~ 2.3 million years ago (MYA). (D) Variation in human haplotypes across the first *TBC1D3* expansion site: a graph representation (rGFA, left) of the locus where colors indicate the source genome for the sequence, and on the right the path for each haplotype-resolved assembly through the graph. The top row for each haplotype composed of large polygons represents an alignment comparing the haplotype-resolved sequence (horizontal) against the graph (vertical), and color represents the source haplotype for the vertical sequence. For example, a single large red triangle indicates there is a one-to-one alignment between T2T-CHM13 and the haplotype. Structural variants can be identified from discontinuities in height (deletion), changes between colors (insertion), or changes in the direction of the polygon (inversion). Below is shown the gene of interest (red arrow) and other genic content in the region (black arrow). Colored bars show ancestral duplication segments (duplicons) that compose the larger duplication blocks.

(SMA) locus containing *SMN1* and *SMN2*—one of the most difficult regions to finish as part of the Human Genome Project on chromosome 5 (57)—shows a distinct structure for all seven assembled haplotypes including GRCh38. Some haplotypes not only show increases in *SMN2* copy number (fig. S14)—a genetic modifier of SMA (58)—but also potential functional differences in the organization and composition of *SMN2*. Because *SMN2* serves as a target for small-molecule drug therapy that improves splice-site efficiency compensating for the loss of *SMN1* in SMA patients (59), this level of

sequence resolution suggests practical utility for disease risk assessment and treatment of patients.

Of particular interest is the *TBCID3* gene family (44) (Fig. 4 and figs. S17 and S18), the protein products of which modulate epidermal growth factor receptor signaling and trafficking (60); further, their duplication in humans has been associated with expansion of the human prefrontal cortex as evidenced by mouse transgenic experiments (7). A comparison with chimpanzees (Fig. 4A) shows two massive genomic expansions in the human lineage (323.0 and 124.4 kbp). Both the high sequence identity (99.6%) and sequence read-depth comparisons of *TBCID3* copy number are consistent with expansion occurring in the human lineage after divergence from chimpanzees (Fig. 4B).

We extended this analysis to other non-human primates by generating HiFi assemblies for bonobos, gorillas, orangutans, and macaques. We identified *TBCID3* homologs in each species and constructed a maximum likelihood phylogeny by using intronic or noncoding sequences flanking the gene (Fig. 4C). The analysis reveals recurrent and independent expansions of *TBCID3* in orangutan, gorilla, and macaque species at different time points during primate evolution, with the most recent expansions occurring 2 million and 2.6 million years ago. However, these estimates assume that there has not been substantial interlocus gene conversion, which may not be the case.

Complete sequencing of human *TBCID3* haplotypes reveals notable structural diversity (Fig. 4D) with *TBCID3* copy number ranging from three to 14 *TBCID3* copies at expansion site 1 and two to nine copies at expansion site 2. In total, approximately one-third of human expansion site 2 shows large-scale structural variation, and we identify >1.8 Mbp of duplicated sequence and >650 kbp of inverted sequence across the 18 haplotypes (including GRCh38). We estimate the structural heterozygosity of this locus to be 90.1% with 14 of 18 haplotypes showing structurally distinct duplication configurations (fig. S18). Similarly, *TBCID3* expansion site 1 is 87.6% heterozygous with 14 of 22 haplotypes displaying unique structures corresponding to copy number differences in the *TBCID3* gene family (fig. S17). Using orthogonal Oxford Nanopore Technologies (ONT) ultra-long-read sequencing, we validated these complex patterns of structural variation in a subset of the samples investigated here (25) (figs. S19 and S20). To better represent the structural genetic variation at this locus, we used a graph-based representation (61), which identified two *TBCID3* genes as common among all human haplotypes examined thus far (*TBCID3B* at site 1 and *TBCID3A* at site 2).

Additional gene models and variable duplicate genes

We identified 182 candidate genes that were previously unresolved or nonsyntenic (25) in the T2T-CHM13 genome assembly (compared with GRCh38) with open reading frames and multiple exons (table S11). Of these, 91% (166) corresponded to SD gene families (Fig. 5A). Many of these represent expanded tandem duplications (e.g., *GAGE* gene family members on the X chromosome) or large interspersed duplications (e.g., beta-defensin locus) adding additional copies of nearly identical genes to the human genome (Fig. 5A).

We searched for evidence that these copy number polymorphic or structurally variant regions were transcribed by aligning long-read transcript sequencing data and searching for perfect matches (25). We constructed a database of 44.2 million full-length cDNA transcripts derived from 31 human tissue samples and compared them with both the GRCh38 and T2T-CHM13 human genome references. For those 182 previously unresolved protein-coding genes where an unambiguous assignment could be made, 36% (65 of 182, >20 Iso-Seq reads) were confirmed as expressed, and 23 of them showed that more reads mapped better to T2T-CHM13 when compared with GRCh38 (Fig. 5B).

Overall, across the entire genome 12% of full-length transcripts exhibited at least 0.2% higher alignment identity when mapped against T2T-CHM13, whereas 8% aligned better to GRCh38. These results are consistent with the notion that T2T-CHM13 is more complete, but that in some cases both assemblies capture different structurally variant haplotypes associated with genes. In addition to entirely new genes, we identified several gene models that were previously incomplete—many of which encode proteins with large tandem repeat domains (ZNF, LPA, Mucin; Fig. 5C). Among these is the complete gene structure of the kringle IV domain of the lipoprotein A gene. Reduced copies of this domain have some of the strongest genetic associations with cardiovascular disease, especially among African Americans (37–40, 62). Sequencing of multiple human haplotypes not only identified length variation but also other forms of rare coding variants potentially relevant for disease risk (Fig. 5D).

SD methylation and transcription

Because methylation is an important consideration in regulating gene transcription, we took advantage of the signal inherent in ultra-long ONT data (63–65) to investigate the CpG methylation status of SD genes within the CHM13 genome (25). Using hierarchical clustering, we found that SD blocks are generally either methylated or unmethylated as an entire block; (fig. S21 and Fig. 6A). Specifically, we found that 452 SD blocks (127.7 Mbp) flanked by unique

sequences are hypermethylated in contrast to 222 hypomethylated SD blocks (52.1 Mbp). Methylation status does not appear to be driven by genomic location, e.g., proximity to centromeres, acrocentric short arms, or telomeres (Fig. 6A).

Using full-length transcript data from CHM13, we compared methylation and transcription status of duplicated genes (25). If we stratify genes by their number of full-length transcripts, we observe distinct methylation patterns for transcribed and nontranscribed SD genes (Fig. 6B). For highly transcribed SD genes (genes without at least one exon overlapping with SD sequence) and unique genes, the gene body and flanking sequence are generally hypermethylated with a pronounced dip near the transcription start site (TSS) and promoter (66). By contrast, nontranscribed genes show moderate to low methylation across the gene body and flanking sequence.

Restricting the analysis to genes mapping within SDs, we find that transcriptionally silenced duplicate genes are more likely (10,000 permutations, $P = 0.0018$) to map to hypomethylated regions of SD sequences (Fig. 6A) when compared with transcribed duplicate genes. Additionally, in untranscribed SD genes, we observe a statistically significant (P value < 0.001, one-sided Mann-Whitney U test) increase in TSS methylation (6.6% increase) when compared with unique genes where the TSS is more likely to be depleted for methylation (8.2% decrease).

One important consideration in this analysis is the presence of a CpG island within 1500 bp of the promoter (67). In our analysis of CHM13, for example, unexpressed unique genes have a low CpG count, consistent with a lack of CpG islands (fig. S22). If we repeat the same analysis on SD genes, we find that the unexpressed SD genes exist with and without CpG islands (fig. S22). In total, these observations suggest a process of epigenetic silencing for a subset of duplicate genes through general demethylation of the gene body but hypermethylation of promoter regions. On the basis of these observed signatures, we investigated whether these epigenetic features coincided with actively transcribed members of duplicate gene families.

We investigated a recently duplicated hominid gene family (*NPIPA*) (68) where sufficient paralogous sequence differences exist to unambiguously assign full-length transcripts to specific loci. Although promoter and TSS signatures are less evident at the individual gene level, the gene body methylation signal appears diagnostic (Fig. 6C). *NPIPA1* and *NPIPA9*, for example, are the most transcriptionally active and show demonstrably distinct methylation patterns providing an epigenetic signature to distinguish transcriptionally active loci associated with high-identity gene families that are otherwise largely indistinguishable. We

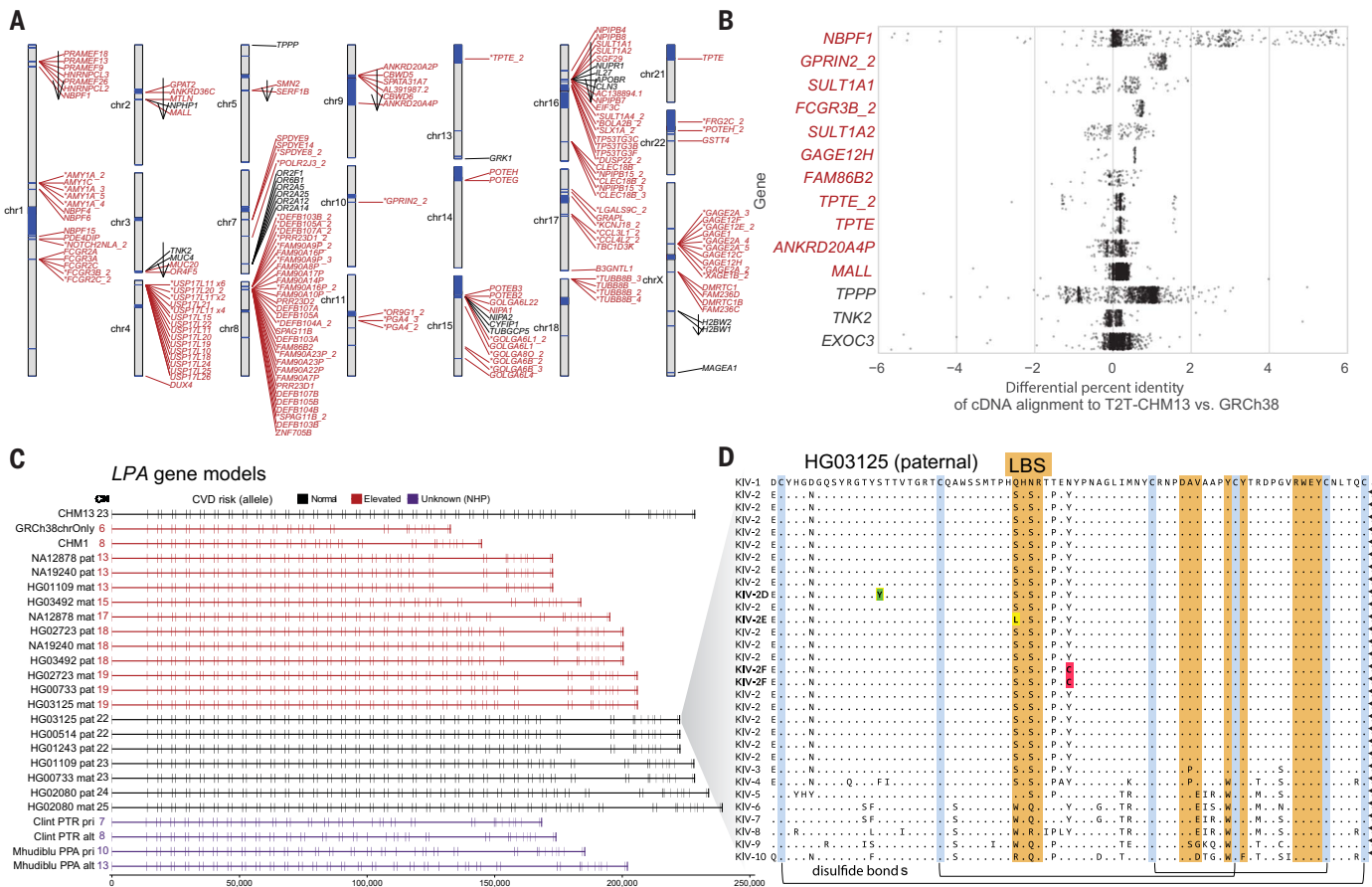


Fig. 5. Genic variation in previously unresolved SD regions of T2T-CHM13.

(A) Ideogram showing the previously unresolved or nonsyntenic gene models [open reading frames (ORFs) with >200 bp of coding sequence and multiple exons] in the T2T-CHM13 assembly as predicted by LiftOff. Previously unresolved genes mapping to SDs (red) are indicated with an asterisk if predicted to be an expansion in the gene family relative to GRCh38 (25). Arrows indicate inverted regions. Most unique genes mapping to nonsyntenic regions (black) are the result of an inversion (arrow). (B) Percent improvement in mapping of CHM13 Iso-Seq reads in candidate duplicated genes (red) mapping to nonsyntenic

regions of the T2T-CHM13 assembly. Positive values identify Iso-Seq reads aligning better to T2T-CHM13 than GRCh38. (C) Gene models of *LPA* with ORF generated from haplotype-resolved HiFi assemblies. The double-exon repeat in these gene models encode for the kringle IV subtype 2 domain of the *LPA* protein. Highlighted in red are haplotypes with reduced kringle IV subtype 2 repeats predicted to increase risk of cardiovascular disease. (D) Amino acid variation in the kringle IV subtype 2 repeat in the paternal haplotype of HG01325 identifies a previously unknown set of amino acid substitutions including rare variants: Ser42Leu in the active site, Ser24Tyr, and Tyr49Cys.

show that this trend also holds for other gene families with high copy numbers (fig. S23).

Discussion

This work provides a comprehensive view of the organization of SDs in the human genome. The T2T-CHM13 reference adds a chromosome's worth (81 Mbp) of SDs, increasing the human genome average from 5.4 to 7.0% and nearly doubling the number of SD pairwise relationships (24 thousand versus 41 thousand) and thus predicts regions of genomic instability as a result of their potential to drive unequal crossing-over events during meiosis.

By every metric, T2T-CHM13 improves our representation of the structure of the human genome. This includes sequence-based organization of the short arms of chromosomes 13, 14, 15, 20, and 21, where we find that SDs

account for more sequence (34.6 Mbp) than either heterochromatic satellite (26.7 Mbp) or rDNA (10 Mbp). Acrocentric SDs are almost twice as large when compared with non-acrocentric regions likely because of ectopic exchange events occurring among the short arms, which associate more frequently during the formation of the nucleolus (69).

Notably, nearly half of the acrocentric SDs involve duplications with nonacrocentric pericentromeric regions of chromosomes 1, 3, 4, 7, 9, 16, and 20. These duplicated islands of euchromatic-like sequences within acrocentric DNA are much more extensive than previously thought but have been shown to be transcriptionally active (70). Although the underlying mechanism for their formation is unknown, it is noteworthy that three of the nonacrocentric regions have large secondary constriction sites (chromosomes 1q, 9q, and 16q) composed

almost entirely of heterochromatic satellites (HSAT2 and HSAT3) (fig. S2). These particular SD blocks thus are bracketed by large tracts of heterochromatic satellites, and such configurations may make them particularly prone to double-strand breakage events (71) promoting interchromosomal duplications (fig. S3) between acrocentric and nonacrocentric chromosomes.

The T2T-CHM13 reference, along with resources from other human genomes, provides a baseline for investigating more complex forms of human genetic variation. For example, this complete reference sequence facilitates the design of sequence-anchored probes to systematically discover and characterize SD heteromorphic variation where chromosome organization differs among individuals (Fig. 2). Such chromosomal heteromorphisms have been traditionally investigated cytogenetically and are thought to be clinically benign (29–31). However,

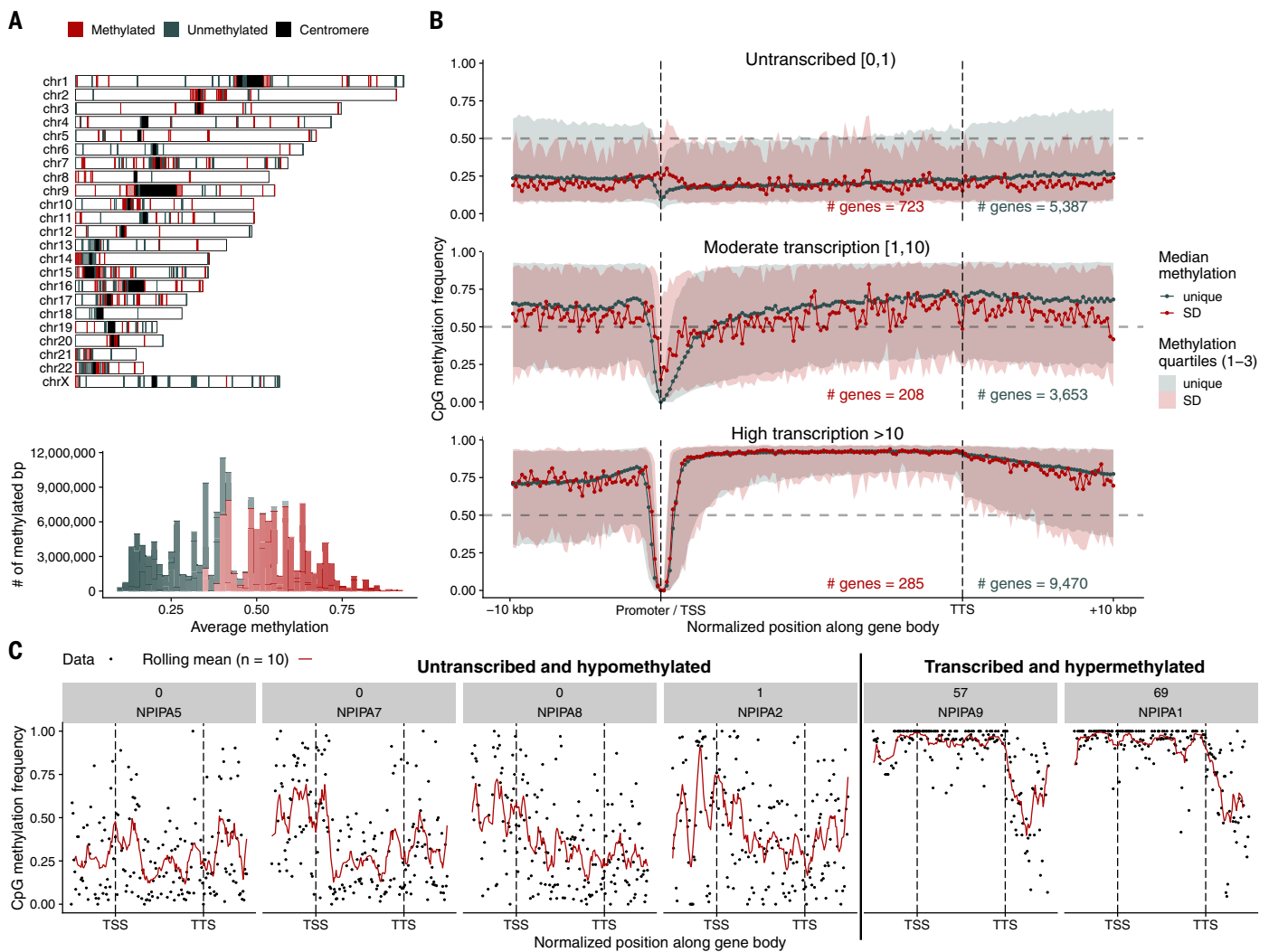


Fig. 6. SD methylation and gene transcription. (A) Methylated (red) or unmethylated (blue-gray) SD blocks in the CHM13 genome on the basis of processing ONT data. The histogram shows the distribution of average methylation across these regions. (B) Median methylation signal of SD (red) and unique (blue-gray) genes stratified by their Iso-Seq expression levels in CHM13. The filled intervals represent the 25 and 75 percentiles of the observed data.

Vertical lines indicate the position of the transcription start site (TSS) and the transcription termination site (TTS). (C) Methylation signal across the recently duplicated *NPIPA* gene family in CHM13, showing increased methylation in transcriptionally active copies. Black points are individual methylation calls, and the red line is a rolling mean across 10 methylation sites. The labels in gray show the number of CHM13 Iso-Seq transcripts and the gene name.

recent work indicates that these large-scale variants are associated with infertility through increasing sperm aneuploidy, decreasing rates of embryonic cleavage (in vitro fertilization, IVF), and increasing miscarriages (72–79). Distinguishing between fixed and heteromorphic acrocentric SDs will facilitate such research as well as the characterization of breakpoints associated with Robertsonian translocations—the most common form of human translocation (80).

At a finer-grained level, the T2T-CHM13 reference and the use of long reads from other human genomes provides access to other complex forms of variation involving duplicated gene families. Short-read copy number variation analyses and single-nucleotide polymor-

phism microarrays have long predicted that SDs are enriched 10-fold for copy number variation, but the structural differences underlying these regions—as well as their functional consequences—have remained elusive (10, 81). We reveal elevated levels of human genetic variation in genes important for human neurodevelopment (*TBCID3*) and disease (*LPA*, *SMN*). Even between just two genomes (GRCh38 and CHM13), we find that 37% (81 Mbp) of SD bases are uncharacterized or structurally variable, and that this predicts 182 copy number variable genes between two human haplotypes (table S6).

In cases such as *TBCID3*, we find that most human haplotypes vary structurally (64 to 78%). Different individuals carry different com-

plements and arrangements of the *TBCID3* gene family. The potential ramifications of this considerable expansion in humans versus chimpanzees and of such high structural heterozygosity among humans are notable given this gene's purported role in expansion of the frontal cortex (7). Similarly, we were able to reconstruct the complete structure of the *LPA* gene model in multiple human genotypes. Although *LPA* is only a single gene, variability in the tandemly repeated 5.2-kbp protein-encoding kringle IV domain underlies one of the most important genetic risk factors for cardiovascular disease. Sequence resolution of structural variation—as well as underlying amino acid differences—allow us to predict previously uncharacterized risk alleles for disease

(Fig. 5). Sequence-resolved structural variation improves genotyping and tests of selection (49, 82, 83), providing a path forward for understanding the disease and evolutionary implications of these complex forms of genetic variation.

Finally, the T2T-CHM13 reference—coupled with other long-read datasets—enables genome-wide functional characterization of recently duplicated genes. Both gene annotation and large-scale efforts to characterize the regulatory landscape of the human genome have typically excluded repetitive regions, including the 859 human genes that map to high-identity SDs (84, 85). This is because the underlying short-read sequencing limits conventional RNA-seq or Chip-seq data from being assigned unambiguously to specific duplicated genes.

In this study, we used long-read full-length transcript data (Iso-Seq) (86) with long-read methylation data from ONT sequencing of the same genome to simultaneously investigate epigenetic and transcriptional data against a fully assembled reference genome. The long-read data from the same haploid source facilitated the unambiguous assignment of these functional readouts, allowing us to correlate methylation and transcript abundance. Our initial analyses suggest that a large fraction of duplicate genes are in fact epigenetically silenced (characterized by hypermethylation of the promoter and hypomethylation of the gene body) and that this epigenetic mark may be used to predict actively transcribed loci even when genes are virtually identical (Fig. 6 and fig. S23). Although more human genomes and diverse tissues need to be investigated to assess the implications of this observation, it is clear that phased genome assemblies (49) with long-read functional readouts such as methylation (65), transcription, or Fiber-seq (87, 88) provide a powerful approach to understanding the regulatory landscape of duplicated and copy number polymorphic genes in the human genome.

However, there are several remaining challenges: First, not all human haplotypes corresponding to specific duplicated regions could be fully sequence-resolved with automated assembly of long-read HiFi sequencing technology. Most of the 250 unresolved regions of phased human genomes generated solely with HiFi long reads correspond to some of the largest and most variable duplicated regions of the human genome (49). For example, only 25% of *SMN1* and *SMN2* haplotypes were fully resolved by HiFi assembly, and these unresolved loci are predicted to carry some of the most complex structural variation patterns. In comparison, the T2T-CHM13 assembly used both accurate HiFi and ultra-long ONT data, and future assembly methods that combine these technologies will likely be critical for

diploid T2T assembly and the complete characterization of SD haplotypes (18, 86).

Another important challenge going forward will be accurate representation of these more complex forms of human genetic variation, including functional annotation where linear representations may be insufficient. Although a more complex pangenome reference graph could overcome these limitations, it is unclear how this will be achieved in practice or how it will be adopted by genomics and clinical communities. This highlights the importance of not only the construction of a pangenome reference but development of necessary tools that will distinguish paralogous and orthologous sequences within duplications to allow for comparison between haplotypes with different SD architectures. The work currently underway by the Human Pangenome Reference Consortium (HPRC), Human Genome Structural Variation Consortium (HGSVC), and Telomere-to-Telomere (T2T) Consortium will be key to developing these methods and completing our understanding of SDs and their role in human genetic variation.

Materials and methods summary

SDs in T2T-CHM13 were identified using SEDEF (21) after repeat masking with Tandem Repeats Finder (TRF) (89) and RepeatMasker (90). Syntenic one-to-one alignments were determined using halSynteny (91). Copy number prediction based on short-read data was performed with WSSD (3, 14) and mrsFAST (92), and regions of comparable copy number were determined with the changepoint package in R (93). To generate gene annotations we used the tools LiftOff (94) and GffRead (95). Fosmid probes were selected from the ABC10 library (28, 96) and two-color FISH was performed to experimentally validate acrocentric SDs (97–99). All assemblies (table S12) with the exception of T2T-CHM13 and GRCh38 were assembled with hifiasm v0.12 (47) using default parameters. Assembly validation of *TBCID3* was performed using sample-matched ONT data by checking the consistency of read alignments to the assemblies (100, 101). Phylogenetic analysis of *TBCID3* was performed with MAFFT and RAXML (102–105). Assembling pangenome graphs for select loci was performed with minigraph (61). Methylation analysis was performed using the methods and data described in Gershman *et al.* (106) using Winnowmap2 and Nanopolish for mapping and methylation calling (65, 107). Data visualization and figures (with the exception of Miroppeats) (108) were primarily made in R making use of GenomicRanges (109), Tidyverse (110), karyoploteR (111), and circize (112). Pipelines used for large-scale data analysis were constructed with Snakemake (113–115). Detailed descriptions of materials and methods are available in the supplementary materials (25).

REFERENCES AND NOTES

1. S. Ohno, U. Wolf, N. B. Atkin, Evolution from fish to mammals by gene duplication. *Hereditas* **59**, 169–187 (1968). doi: [10.1111/j.1601-5223.1968.tb02169.x](https://doi.org/10.1111/j.1601-5223.1968.tb02169.x); pmid: [5662632](https://pubmed.ncbi.nlm.nih.gov/5662632/)
2. S. Ohno, *Evolution by Gene Duplication* (Springer, 1970).
3. J. A. Bailey, A. M. Yavor, H. F. Massa, B. J. Trask, E. E. Eichler, Segmental duplications: Organization and impact within the current human genome project assembly. *Genome Res.* **11**, 1005–1017 (2001). doi: [10.1101/gr187101](https://doi.org/10.1101/gr187101); pmid: [11381028](https://pubmed.ncbi.nlm.nih.gov/11381028/)
4. G. M. Cooper *et al.*, A copy number variation morbidity map of developmental delay. *Nat. Genet.* **43**, 838–846 (2011). doi: [10.1038/ng.909](https://doi.org/10.1038/ng.909); pmid: [21841781](https://pubmed.ncbi.nlm.nih.gov/21841781/)
5. M. Y. Dennis *et al.*, Evolution of human-specific neural SRGAP2 genes by incomplete segmental duplication. *Cell* **149**, 912–922 (2012). doi: [10.1016/j.cell.2012.03.033](https://doi.org/10.1016/j.cell.2012.03.033); pmid: [22559943](https://pubmed.ncbi.nlm.nih.gov/22559943/)
6. M. Florio *et al.*, Evolution and cell-type specificity of human-specific genes preferentially expressed in progenitors of fetal neocortex. *eLife* **7**, e32332 (2018). doi: [10.7554/eLife.32332](https://doi.org/10.7554/eLife.32332); pmid: [29561261](https://pubmed.ncbi.nlm.nih.gov/29561261/)
7. X. C. Ju *et al.*, The hominoid-specific gene *TBCID3* promotes generation of basal neural progenitors and induces cortical folding in mice. *eLife* **5**, e18197 (2016). doi: [10.7554/eLife.18197](https://doi.org/10.7554/eLife.18197); pmid: [27504805](https://pubmed.ncbi.nlm.nih.gov/27504805/)
8. I. T. Fiddes *et al.*, Human-specific *NOTCH2NL* genes affect notch signaling and cortical neurogenesis. *Cell* **173**, 1356–1369.e22 (2018). doi: [10.1016/j.cell.2018.03.051](https://doi.org/10.1016/j.cell.2018.03.051); pmid: [29856954](https://pubmed.ncbi.nlm.nih.gov/29856954/)
9. P. H. Sudmant *et al.*, 1000 Genomes Project, Diversity of human copy number variation and multicopy genes. *Science* **330**, 641–646 (2010). doi: [10.1126/science.1197005](https://doi.org/10.1126/science.1197005); pmid: [21030649](https://pubmed.ncbi.nlm.nih.gov/21030649/)
10. P. H. Sudmant *et al.*, Global diversity, population stratification, and selection of human copy-number variation. *Science* **349**, aab3761 (2015). doi: [10.1126/science.aab3761](https://doi.org/10.1126/science.aab3761); pmid: [26249230](https://pubmed.ncbi.nlm.nih.gov/26249230/)
11. J. C. Venter *et al.*, The sequence of the human genome. *Science* **291**, 1304–1351 (2001). doi: [10.1126/science.1058040](https://doi.org/10.1126/science.1058040); pmid: [11181995](https://pubmed.ncbi.nlm.nih.gov/11181995/)
12. X. She *et al.*, Shotgun sequence assembly and recent segmental duplications within the human genome. *Nature* **431**, 927–930 (2004). doi: [10.1038/nature03062](https://doi.org/10.1038/nature03062); pmid: [15496912](https://pubmed.ncbi.nlm.nih.gov/15496912/)
13. E. S. Lander *et al.*, Initial sequencing and analysis of the human genome. *Nature* **409**, 860–921 (2001). doi: [10.1038/35057062](https://doi.org/10.1038/35057062); pmid: [11237011](https://pubmed.ncbi.nlm.nih.gov/11237011/)
14. J. A. Bailey *et al.*, Recent segmental duplications in the human genome. *Science* **297**, 1003–1007 (2002). doi: [10.1126/science.1072047](https://doi.org/10.1126/science.1072047); pmid: [12169732](https://pubmed.ncbi.nlm.nih.gov/12169732/)
15. E. Eichler, U. Surti, R. Ophoff, Proposal for Construction a Human Haploid BAC library from Hydatidiform Mole Source Material (2002). www.genome.gov/Pages/Research/Sequencing/BACLibrary/HydatidiformMoleBAC021203.pdf.
16. D. Fredman *et al.*, Complex SNP-related sequence variation in segmental genome duplications. *Nat. Genet.* **36**, 861–866 (2004). doi: [10.1038/ng1401](https://doi.org/10.1038/ng1401); pmid: [15247918](https://pubmed.ncbi.nlm.nih.gov/15247918/)
17. M. J. Chaisson *et al.*, Resolving the complexity of the human genome using single-molecule sequencing. *Nature* **517**, 608–611 (2015). doi: [10.1038/nature13907](https://doi.org/10.1038/nature13907); pmid: [25383537](https://pubmed.ncbi.nlm.nih.gov/25383537/)
18. K. H. Miga *et al.*, Telomere-to-telomere assembly of a complete human X chromosome. *Nature* **585**, 79–84 (2020). doi: [10.1038/s41586-020-2547-7](https://doi.org/10.1038/s41586-020-2547-7); pmid: [32663838](https://pubmed.ncbi.nlm.nih.gov/32663838/)
19. R. E. Green *et al.*, A draft sequence of the Neandertal genome. *Science* **328**, 710–722 (2010). doi: [10.1126/science.1188021](https://doi.org/10.1126/science.1188021); pmid: [20448178](https://pubmed.ncbi.nlm.nih.gov/20448178/)
20. S. Nurk *et al.*, The complete sequence of a human genome. *Science* **376**, 44–53 (2022). doi: [10.1126/science.abj6987](https://doi.org/10.1126/science.abj6987)
21. I. Numanagic *et al.*, Fast characterization of segmental duplications in genome assemblies. *Bioinformatics* **34**, i706–i714 (2018). doi: [10.1093/bioinformatics/bty586](https://doi.org/10.1093/bioinformatics/bty586); pmid: [30423092](https://pubmed.ncbi.nlm.nih.gov/30423092/)
22. A. D. Bell, C. L. Usher, S. A. McCarroll, Analyzing Copy Number Variation with Droplet Digital PCR. *Methods Mol. Biol.* **1768**, 143–160 (2018). doi: [10.1007/978-1-4939-7778-9_9](https://doi.org/10.1007/978-1-4939-7778-9_9); pmid: [29717442](https://pubmed.ncbi.nlm.nih.gov/29717442/)
23. I. L. Gonzalez, J. E. Sylvester, Complete sequence of the 43-kb human ribosomal DNA repeat: Analysis of the intergenic spacer. *Genomics* **27**, 320–328 (1995). doi: [10.1006/geno.1995.1049](https://doi.org/10.1006/geno.1995.1049); pmid: [7557999](https://pubmed.ncbi.nlm.nih.gov/7557999/)
24. J. H. Kim *et al.*, Variation in human chromosome 21 ribosomal RNA genes characterized by TAR cloning and long-read

- sequencing. *Nucleic Acids Res.* **46**, 6712–6725 (2018). doi: [10.1093/nar/gky442](https://doi.org/10.1093/nar/gky442); pmid: [29788454](https://pubmed.ncbi.nlm.nih.gov/29788454/)
25. Supplementary materials and methods for segmental duplications and their variation in a complete human genome.
 26. Z. Jiang, R. Hubley, A. Smit, E. E. Eichler, DupMasker: A tool for annotating primate segmental duplications. *Genome Res.* **18**, 1362–1368 (2008). doi: [10.1101/gr.078477.108](https://doi.org/10.1101/gr.078477.108); pmid: [18502942](https://pubmed.ncbi.nlm.nih.gov/18502942/)
 27. Z. Jiang *et al.*, Ancestral reconstruction of segmental duplications reveals punctuated cores of human genome evolution. *Nat. Genet.* **39**, 1361–1368 (2007). doi: [10.1038/ng.2007.9](https://doi.org/10.1038/ng.2007.9); pmid: [17922013](https://pubmed.ncbi.nlm.nih.gov/17922013/)
 28. J. M. Kidd *et al.*, Mapping and sequencing of structural variation from eight human genomes. *Nature* **453**, 56–64 (2008). doi: [10.1038/nature06862](https://doi.org/10.1038/nature06862); pmid: [18451855](https://pubmed.ncbi.nlm.nih.gov/18451855/)
 29. M. K. Bhasin, Bhasin, Human population cytogenetics: A review. *Int. J. Hum. Genet.* **5**, 83–152 (2005). doi: [10.1080/09723757.2005.11885918](https://doi.org/10.1080/09723757.2005.11885918)
 30. L. Y. F. Hsu *et al.*, Chromosomal polymorphisms of 1, 9, 16, and Y in 4 major ethnic groups: A large prenatal study. *Am. J. Med. Genet.* **26**, 95–101 (1987). doi: [10.1002/ajmg.1320260116](https://doi.org/10.1002/ajmg.1320260116); pmid: [3812584](https://pubmed.ncbi.nlm.nih.gov/3812584/)
 31. J. C. K. Barber, Euchromatic heteromorphism or duplication without phenotypic effect? *Prenat. Diagn.* **14**, 323–324 (1994). doi: [10.1002/pd.1970140418](https://doi.org/10.1002/pd.1970140418); pmid: [8066046](https://pubmed.ncbi.nlm.nih.gov/8066046/)
 32. M. J. P. Chaisson *et al.*, Multi-platform discovery of haplotype-resolved structural variation in human genomes. *Nat. Commun.* **10**, 1784 (2019). doi: [10.1038/s41467-018-08148-z](https://doi.org/10.1038/s41467-018-08148-z); pmid: [30992455](https://pubmed.ncbi.nlm.nih.gov/30992455/)
 33. A. D. Sanders *et al.*, Characterizing polymorphic inversions in human genomes by single-cell sequencing. *Genome Res.* **26**, 1575–1587 (2016). doi: [10.1101/gr.201160.115](https://doi.org/10.1101/gr.201160.115); pmid: [27472961](https://pubmed.ncbi.nlm.nih.gov/27472961/)
 34. D. Porubský *et al.*, Recurrent inversion toggling and great ape genome evolution. *Nat. Genet.* **52**, 849–858 (2020). doi: [10.1038/s41588-020-0646-x](https://doi.org/10.1038/s41588-020-0646-x); pmid: [32541924](https://pubmed.ncbi.nlm.nih.gov/32541924/)
 35. I. K. Suzuki *et al.*, Human-Specific NOTCH2NL Genes Expand Cortical Neurogenesis through Delta/Notch Regulation. *Cell* **173**, 1370–1384.e16 (2018). doi: [10.1016/j.cell.2018.03.067](https://doi.org/10.1016/j.cell.2018.03.067); pmid: [29856955](https://pubmed.ncbi.nlm.nih.gov/29856955/)
 36. S. Mallick *et al.*, The Simons Genome Diversity Project: 300 genomes from 142 diverse populations. *Nature* **538**, 201–206 (2016). doi: [10.1038/nature18964](https://doi.org/10.1038/nature18964); pmid: [27654912](https://pubmed.ncbi.nlm.nih.gov/27654912/)
 37. R. Clarke *et al.*, Genetic variants associated with Lp(a) lipoprotein level and coronary disease. *N. Engl. J. Med.* **361**, 2518–2528 (2009). doi: [10.1056/NEJMoa0902604](https://doi.org/10.1056/NEJMoa0902604); pmid: [20032323](https://pubmed.ncbi.nlm.nih.gov/20032323/)
 38. S. Coassin *et al.*, A comprehensive map of single-base polymorphisms in the hypervariable LPA kringle IV type 2 copy number variation region. *J. Lipid Res.* **60**, 186–199 (2019). doi: [10.1194/jlr.M090381](https://doi.org/10.1194/jlr.M090381); pmid: [30413653](https://pubmed.ncbi.nlm.nih.gov/30413653/)
 39. F. Kronenberg, G. Utermann, Lipoprotein(a): Resurrected by genetics. *J. Intern. Med.* **273**, 6–30 (2013). doi: [10.1111/j.1365-2796.2012.02592.x](https://doi.org/10.1111/j.1365-2796.2012.02592.x); pmid: [22998429](https://pubmed.ncbi.nlm.nih.gov/22998429/)
 40. K. Schmidt, A. Noureen, F. Kronenberg, G. Utermann, Structure, function, and genetics of lipoprotein (a). *J. Lipid Res.* **57**, 1339–1359 (2016). doi: [10.1194/jlr.R067314](https://doi.org/10.1194/jlr.R067314); pmid: [27074913](https://pubmed.ncbi.nlm.nih.gov/27074913/)
 41. J. R. Gum *et al.*, Molecular cloning of cDNAs derived from a novel human intestinal mucin gene. *Biochem. Biophys. Res. Commun.* **171**, 407–415 (1990). doi: [10.1016/0006-291X\(90\)91408-K](https://doi.org/10.1016/0006-291X(90)91408-K); pmid: [2393399](https://pubmed.ncbi.nlm.nih.gov/2393399/)
 42. W. S. Pratt *et al.*, Multiple transcripts of MUC3: Evidence for two genes, MUC3A and MUC3B. *Biochem. Biophys. Res. Commun.* **275**, 916–923 (2000). doi: [10.1006/bbrc.2000.3406](https://doi.org/10.1006/bbrc.2000.3406); pmid: [10973822](https://pubmed.ncbi.nlm.nih.gov/10973822/)
 43. K. Kyo, T. Muto, H. Nagawa, G. M. Lathrop, Y. Nakamura, Associations of distinct variants of the intestinal mucin gene MUC3A with ulcerative colitis and Crohn's disease. *J. Hum. Genet.* **46**, 5–20 (2001). doi: [10.1007/s100380170118](https://doi.org/10.1007/s100380170118); pmid: [11289722](https://pubmed.ncbi.nlm.nih.gov/11289722/)
 44. C. A. Paulding, M. Ruvolo, D. A. Haber, The Tre2 (USP6) oncogene is a hominoid-specific gene. *Proc. Natl. Acad. Sci. U.S.A.* **100**, 2507–2511 (2003). doi: [10.1073/pnas.0437015100](https://doi.org/10.1073/pnas.0437015100); pmid: [12604796](https://pubmed.ncbi.nlm.nih.gov/12604796/)
 45. S. Cantislilleris *et al.*, An evolutionary driver of interspersed segmental duplications in primates. *Genome Biol.* **21**, 202 (2020). doi: [10.1186/s13059-020-02074-4](https://doi.org/10.1186/s13059-020-02074-4); pmid: [32778141](https://pubmed.ncbi.nlm.nih.gov/32778141/)
 46. T. Marques-Bonet, E. E. Eichler, The evolution of human segmental duplications and the core duplcon hypothesis. *Cold Spring Harb. Symp. Quant. Biol.* **74**, 355–362 (2009). doi: [10.1101/sqb.2009.74.011](https://doi.org/10.1101/sqb.2009.74.011); pmid: [19717539](https://pubmed.ncbi.nlm.nih.gov/19717539/)
 47. H. Cheng, G. T. Concepcion, X. Feng, H. Zhang, H. Li, Haplotype-resolved de novo assembly using phased assembly graphs with hifiasm. *Nat. Methods* **18**, 170–175 (2021). doi: [10.1038/s41592-020-01056-5](https://doi.org/10.1038/s41592-020-01056-5); pmid: [33526886](https://pubmed.ncbi.nlm.nih.gov/33526886/)
 48. S. Nurk *et al.*, HiCanu: Accurate assembly of segmental duplications, satellites, and allelic variants from high-fidelity long reads. *Genome Res.* **30**, 1291–1305 (2020). doi: [10.1101/gr.263566.120](https://doi.org/10.1101/gr.263566.120); pmid: [32801147](https://pubmed.ncbi.nlm.nih.gov/32801147/)
 49. P. Ebert *et al.*, Haplotype-resolved diverse human genomes and integrated analysis of structural variation. *Science* **372**, eabf7117 (2021). doi: [10.1126/science.abf7117](https://doi.org/10.1126/science.abf7117); pmid: [33632895](https://pubmed.ncbi.nlm.nih.gov/33632895/)
 50. L. Bertilsson, M. L. Dahl, P. Dalén, A. Al-Shurbaji, Molecular genetics of CYP2D6: Clinical relevance with focus on psychotropic drugs. *Br. J. Clin. Pharmacol.* **53**, 111–122 (2002). doi: [10.1046/j.0306-5251.2001.01548.x](https://doi.org/10.1046/j.0306-5251.2001.01548.x); pmid: [11851634](https://pubmed.ncbi.nlm.nih.gov/11851634/)
 51. W. Hammer, F. Sjöqvist, Plasma levels of monomethylated tricyclic antidepressants during treatment with imipramine-like compounds. *Life Sci.* **6**, 1895–1903 (1967). doi: [10.1016/0024-3205\(67\)90218-4](https://doi.org/10.1016/0024-3205(67)90218-4); pmid: [6052684](https://pubmed.ncbi.nlm.nih.gov/6052684/)
 52. B. Alexanderson, D. A. Evans, F. Sjöqvist, Steady-state plasma levels of nortriptyline in twins: Influence of genetic factors and drug therapy. *BMJ* **4**, 764–768 (1969). doi: [10.1136/bmj.4.5686.764](https://doi.org/10.1136/bmj.4.5686.764); pmid: [5391106](https://pubmed.ncbi.nlm.nih.gov/5391106/)
 53. R. C. Skoda, F. J. Gonzalez, A. Demierre, U. A. Meyer, Two mutant alleles of the human cytochrome P-450db1 gene (P450C2D1) associated with genetically deficient metabolism of debrisoquine and other drugs. *Proc. Natl. Acad. Sci. U.S.A.* **85**, 5240–5243 (1988). doi: [10.1073/pnas.85.14.5240](https://doi.org/10.1073/pnas.85.14.5240); pmid: [2899325](https://pubmed.ncbi.nlm.nih.gov/2899325/)
 54. M. L. Dahl, I. Johansson, M. P. Palmertz, M. Ingelman-Sundberg, F. Sjöqvist, Analysis of the CYP2D6 gene in relation to debrisoquine and desipramine hydroxylation in a Swedish population. *Clin. Pharmacol. Ther.* **51**, 12–17 (1992). doi: [10.1038/clpt.1992.2](https://doi.org/10.1038/clpt.1992.2); pmid: [1346258](https://pubmed.ncbi.nlm.nih.gov/1346258/)
 55. A. Gaedigk, M. Blum, R. Gaedigk, M. Eichelbaum, U. A. Meyer, Deletion of the entire cytochrome P450 CYP2D6 gene as a cause of impaired drug metabolism in poor metabolizers of the debrisoquine/sparteine polymorphism. *Am. J. Hum. Genet.* **48**, 943–950 (1991). pmid: [1673290](https://pubmed.ncbi.nlm.nih.gov/1673290/)
 56. I. Johansson *et al.*, Inherited amplification of an active gene in the cytochrome P450 CYP2D6 locus as a cause of ultrarapid metabolism of debrisoquine. *Proc. Natl. Acad. Sci. U.S.A.* **90**, 11825–11829 (1993). doi: [10.1073/pnas.90.24.11825](https://doi.org/10.1073/pnas.90.24.11825); pmid: [7903454](https://pubmed.ncbi.nlm.nih.gov/7903454/)
 57. J. Schmutz *et al.*, The DNA sequence and comparative analysis of human chromosome 5. *Nature* **431**, 268–274 (2004). doi: [10.1038/nature02919](https://doi.org/10.1038/nature02919); pmid: [15372022](https://pubmed.ncbi.nlm.nih.gov/15372022/)
 58. M. E. Butchbach, Copy Number Variations in the Survival Motor Neuron Genes: Implications for Spinal Muscular Atrophy and Other Neurodegenerative Diseases. *Front. Mol. Biosci.* **3**, 7 (2016). doi: [10.3389/fmolb.2016.00007](https://doi.org/10.3389/fmolb.2016.00007); pmid: [27014701](https://pubmed.ncbi.nlm.nih.gov/27014701/)
 59. T. W. Bebee, J. T. Gladman, D. S. Chandler, Splicing regulation of the Survival Motor Neuron genes and implications for treatment of spinal muscular atrophy. *Front. Biosci.* **15**, 1191 (2010). doi: [10.2741/3670](https://doi.org/10.2741/3670); pmid: [20515750](https://pubmed.ncbi.nlm.nih.gov/20515750/)
 60. M. J. Wainzselbaum *et al.*, The hominoid-specific oncogene TBC1D3 activates Ras and modulates epidermal growth factor receptor signaling and trafficking. *J. Biol. Chem.* **283**, 13233–13242 (2008). doi: [10.1074/jbc.M800234200](https://doi.org/10.1074/jbc.M800234200); pmid: [18319245](https://pubmed.ncbi.nlm.nih.gov/18319245/)
 61. H. Li, X. Feng, C. Chu, The design and construction of reference pangenome graphs with minigraph. *Genome Biol.* **21**, 265 (2020). doi: [10.1186/s13059-020-02168-z](https://doi.org/10.1186/s13059-020-02168-z); pmid: [33066802](https://pubmed.ncbi.nlm.nih.gov/33066802/)
 62. F. Kronenberg, Human Genetics and the Causal Role of Lipoprotein(a) for Various Diseases. *Cardiovasc. Drugs Ther.* **30**, 87–100 (2016). doi: [10.1007/s10557-016-6648-3](https://doi.org/10.1007/s10557-016-6648-3); pmid: [26896185](https://pubmed.ncbi.nlm.nih.gov/26896185/)
 63. J. Quick *et al.*, Real-time, portable genome sequencing for Ebola surveillance. *Nature* **530**, 228–232 (2016). doi: [10.1038/nature16996](https://doi.org/10.1038/nature16996); pmid: [26840485](https://pubmed.ncbi.nlm.nih.gov/26840485/)
 64. N. J. Loman, J. Quick, J. T. Simpson, A complete bacterial genome assembled de novo using only nanopore sequencing data. *Nat. Methods* **12**, 733–735 (2015). doi: [10.1038/nmeth.3444](https://doi.org/10.1038/nmeth.3444); pmid: [26076426](https://pubmed.ncbi.nlm.nih.gov/26076426/)
 65. J. T. Simpson *et al.*, Detecting DNA cytosine methylation using nanopore sequencing. *Nat. Methods* **14**, 407–410 (2017). doi: [10.1038/nmeth.4184](https://doi.org/10.1038/nmeth.4184); pmid: [28218898](https://pubmed.ncbi.nlm.nih.gov/28218898/)
 66. M. P. Ball *et al.*, Targeted and genome-scale strategies reveal gene-body methylation signatures in human cells. *Nat. Biotechnol.* **27**, 361–368 (2009). doi: [10.1038/nbt.1533](https://doi.org/10.1038/nbt.1533); pmid: [19329998](https://pubmed.ncbi.nlm.nih.gov/19329998/)
 67. S. Saxonov, P. Berg, D. L. Brutlag, A genome-wide analysis of CpG dinucleotides in the human genome distinguishes two distinct classes of promoters. *Proc. Natl. Acad. Sci. U.S.A.* **103**, 1412–1417 (2006). doi: [10.1073/pnas.051010103](https://doi.org/10.1073/pnas.051010103); pmid: [16432200](https://pubmed.ncbi.nlm.nih.gov/16432200/)
 68. M. E. Johnson *et al.*, Positive selection of a gene family during the emergence of humans and African apes. *Nature* **413**, 514–519 (2001). doi: [10.1038/35097067](https://doi.org/10.1038/35097067); pmid: [11586358](https://pubmed.ncbi.nlm.nih.gov/11586358/)
 69. N. Arnheim, in *Evolution of Genes and Proteins*, M. Nei, R. K. Koehn, Eds. (Sinauer, 1983) pp. 38–61.
 70. R. Lyle *et al.*, Islands of euchromatin-like sequence and expressed polymorphic sequences within the short arm of human chromosome 21. *Genome Res.* **17**, 1690–1696 (2007). doi: [10.1101/jr.6675307](https://doi.org/10.1101/jr.6675307); pmid: [17895424](https://pubmed.ncbi.nlm.nih.gov/17895424/)
 71. S. Luke, R. S. Verma, R. A. Conte, T. Mathews, Molecular characterization of the secondary constriction region (qh) of human chromosome 9 with pericentric inversion. *J. Cell Sci.* **103**, 919–923 (1992). doi: [10.1242/jcs.103.4.919](https://doi.org/10.1242/jcs.103.4.919); pmid: [1487504](https://pubmed.ncbi.nlm.nih.gov/1487504/)
 72. J. C. Barber *et al.*, Duplications of proximal 16q flanked by heterochromatin are not euchromatic variants and show no evidence of heterochromatin position effect. *Cytogenet. Genome Res.* **114**, 351–358 (2006). doi: [10.1159/000094225](https://doi.org/10.1159/000094225); pmid: [16954678](https://pubmed.ncbi.nlm.nih.gov/16954678/)
 73. F. I. Sahin *et al.*, Chromosome heteromorphisms: An impact on infertility. *J. Assist. Reprod. Genet.* **25**, 191–195 (2008). doi: [10.1007/s10815-008-9216-3](https://doi.org/10.1007/s10815-008-9216-3); pmid: [18461436](https://pubmed.ncbi.nlm.nih.gov/18461436/)
 74. M. Codina-Pascual *et al.*, Behaviour of human heterochromatic regions during the synapsis of homologous chromosomes. *Hum. Reprod.* **21**, 1490–1497 (2006). doi: [10.1093/humrep/del028](https://doi.org/10.1093/humrep/del028); pmid: [16484310](https://pubmed.ncbi.nlm.nih.gov/16484310/)
 75. A. O. Caglayan, I. Ozyazgan, F. Demiryilmaz, M. T. Ozgun, Are heterochromatin polymorphisms associated with recurrent miscarriage? *J. Obstet. Gynaecol. Res.* **36**, 774–776 (2010). doi: [10.1111/j.1447-0756.2010.01207.x](https://doi.org/10.1111/j.1447-0756.2010.01207.x); pmid: [20666944](https://pubmed.ncbi.nlm.nih.gov/20666944/)
 76. P. F. Madon, A. S. Athalye, F. R. Parikh, Polymorphic variants on chromosomes probably play a significant role in infertility. *Reprod. Biomed. Online* **11**, 726–732 (2005). doi: [10.1016/S1472-6483\(10\)61691-4](https://doi.org/10.1016/S1472-6483(10)61691-4); pmid: [16417737](https://pubmed.ncbi.nlm.nih.gov/16417737/)
 77. S. Minocherhomji *et al.*, A case-control study identifying chromosomal polymorphic variations as forms of epigenetic alterations associated with the infertility phenotype. *Fertil. Steril.* **92**, 88–95 (2009). doi: [10.1016/j.fertnstert.2008.05.071](https://doi.org/10.1016/j.fertnstert.2008.05.071); pmid: [18692838](https://pubmed.ncbi.nlm.nih.gov/18692838/)
 78. Y. Hong, Y. W. Zhou, J. Tao, S. X. Wang, X. M. Zhao, Do polymorphic variants of chromosomes affect the outcome of in vitro fertilization and embryo transfer treatment? *Hum. Reprod.* **26**, 933–940 (2011). doi: [10.1093/humrep/deq333](https://doi.org/10.1093/humrep/deq333); pmid: [21266453](https://pubmed.ncbi.nlm.nih.gov/21266453/)
 79. P. Kalantari, H. Sepehri, F. Behjati, Z. O. Ashtiani, M. T. Akbari, Chromosomal studies in infertile men. *Tsitol. Genet.* **35**, 50–54 (2001). pmid: [11944328](https://pubmed.ncbi.nlm.nih.gov/11944328/)
 80. E. S. Wilch, C. C. Morton, Historical and Clinical Perspectives on Chromosomal Translocations. *Adv. Exp. Med. Biol.* **1044**, 1–14 (2018). doi: [10.1007/978-981-13-0593-1_1](https://doi.org/10.1007/978-981-13-0593-1_1); pmid: [29956287](https://pubmed.ncbi.nlm.nih.gov/29956287/)
 81. D. P. Locke *et al.*, BAC microarray analysis of 15q11-q13 rearrangements and the impact of segmental duplications. *J. Med. Genet.* **41**, 175–182 (2004). doi: [10.1136/jmg.2003.013813](https://doi.org/10.1136/jmg.2003.013813); pmid: [14985376](https://pubmed.ncbi.nlm.nih.gov/14985376/)
 82. J. Ebler *et al.*, Pangenome-based genome inference. bioRxiv 2020.11.11.378133v1 [Preprint] Cold Spring Harbor Laboratory (2020); doi: [10.1101/2020.11.11.378133](https://doi.org/10.1101/2020.11.11.378133)
 83. P. Hsieh *et al.*, Adaptive archaic introgression of copy number variants and the discovery of previously unknown human genes. *Science* **366**, eaax2083 (2019). doi: [10.1126/science.aax2083](https://doi.org/10.1126/science.aax2083); pmid: [31624180](https://pubmed.ncbi.nlm.nih.gov/31624180/)
 84. A. Battle, C. D. Brown, B. E. Engelhardt, S. B. Montgomery, GTEx Consortium; Laboratory, Data Analysis & Coordinating Center (LDACC)—Analysis Working Group; Statistical Methods groups—Analysis Working Group; Enhancing GTEx (eGTEx) groups; NIH Common Fund; NIH/NCI; NIH/NHGRI; NIH/NIMH; NIH/NIDA; Biospecimen Collection Source Site—NDRI; Biospecimen Collection Source Site—RPCI; Biospecimen Core Resource—VARI; Brain Bank Repository—University of Miami Brain Endowment Bank; Leidos Biomedical—Project Management; ELSI Study; Genome Browser Data Integration & Visualization—EBI; Genome Browser Data Integration & Visualization—UCSC

- Genomics Institute, University of California Santa Cruz; Lead analysts; Laboratory, Data Analysis & Coordinating Center (LDACC); NIH program management; Biospecimen collection; Pathology; eQTL manuscript working group. Genetic effects on gene expression across human tissues. *Nature* **550**, 204–213 (2017). doi: [10.1038/nature24277](https://doi.org/10.1038/nature24277); PMID: [29022597](https://pubmed.ncbi.nlm.nih.gov/29022597/)
85. M. L. Dougherty *et al.*, Transcriptional fates of human-specific segmental duplications in brain. *Genome Res.* **28**, 1566–1576 (2018). doi: [10.1101/gr.237610.118](https://doi.org/10.1101/gr.237610.118); PMID: [30228200](https://pubmed.ncbi.nlm.nih.gov/30228200/)
 86. G. A. Logsdon *et al.*, The structure, function and evolution of a complete human chromosome 8. *Nature* **593**, 101–107 (2021). doi: [10.1038/s41586-021-03420-7](https://doi.org/10.1038/s41586-021-03420-7); PMID: [33828295](https://pubmed.ncbi.nlm.nih.gov/33828295/)
 87. A. B. Stergachis, B. M. Debo, E. Haugen, L. S. Churchman, J. A. Stamatoyannopoulos, Single-molecule regulatory architectures captured by chromatin fiber sequencing. *Science* **368**, 1449–1454 (2020). doi: [10.1126/science.aaz1646](https://doi.org/10.1126/science.aaz1646); PMID: [32587015](https://pubmed.ncbi.nlm.nih.gov/32587015/)
 88. N. J. Abdulhay *et al.*, Massively multiplex single-molecule oligonucleosome footprinting. *elife* **9**, e59404 (2020). doi: [10.7554/eLife.59404](https://doi.org/10.7554/eLife.59404)
 89. G. Benson, Tandem repeats finder: A program to analyze DNA sequences. *Nucleic Acids Res.* **27**, 573–580 (1999). doi: [10.1093/nar/27.2.573](https://doi.org/10.1093/nar/27.2.573); PMID: [9862982](https://pubmed.ncbi.nlm.nih.gov/9862982/)
 90. A.F.A. Smit, R. Hubley, P. Green, RepeatMasker Open-4.0. (2015). www.repeatmasker.org.
 91. K. Krasheninnikova, M. Diekhans, J. Armstrong, A. Dievskii, B. Paten, S. O'Brien, halSynteny: A fast, easy-to-use conserved synteny block construction method for multiple whole-genome alignments. *Gigascience* **9**, g1aa047 (2020). doi: [10.1093/gigascience/g1aa047](https://doi.org/10.1093/gigascience/g1aa047); PMID: [32463100](https://pubmed.ncbi.nlm.nih.gov/32463100/)
 92. F. Hach *et al.*, mrsFAST: A cache-oblivious algorithm for short-read mapping. *Nat. Methods* **7**, 576–577 (2010). doi: [10.1038/nmeth0810-576](https://doi.org/10.1038/nmeth0810-576); PMID: [20676076](https://pubmed.ncbi.nlm.nih.gov/20676076/)
 93. H. Killick, F. Eckley, P. Lee, changepoint: R package version 0.4. (2016). <https://cran.r-project.org/web/packages/changepoint/changepoint.pdf>.
 94. A. Shumate, S. L. Salzberg, Liftoff: Accurate mapping of gene annotations. *Bioinformatics* **37**, 1639–1643 (2021). doi: [10.1093/bioinformatics/btaa1016](https://doi.org/10.1093/bioinformatics/btaa1016); PMID: [33320174](https://pubmed.ncbi.nlm.nih.gov/33320174/)
 95. G. Pertea, M. Pertea, G. F. F. Utilities, GffRead and GffCompare. *F1000 Res.* **9** (2020) [version 2; peer review: 3 approved]. doi: [10.12688/f1000research.23297.2](https://doi.org/10.12688/f1000research.23297.2); PMID: [32489650](https://pubmed.ncbi.nlm.nih.gov/32489650/)
 96. S. F. Altschul, W. Gish, W. Miller, E. W. Myers, D. J. Lipman, Basic local alignment search tool. *J. Mol. Biol.* **215**, 403–410 (1990). doi: [10.1016/S0022-2836\(05\)80360-2](https://doi.org/10.1016/S0022-2836(05)80360-2); PMID: [2231712](https://pubmed.ncbi.nlm.nih.gov/2231712/)
 97. M. F. Cardone *et al.*, Independent centromere formation in a capricious, gene-free domain of chromosome 13q21 in Old World monkeys and pigs. *Genome Biol.* **7**, R91 (2006). doi: [10.1186/gb-2006-7-10-r91](https://doi.org/10.1186/gb-2006-7-10-r91); PMID: [17040560](https://pubmed.ncbi.nlm.nih.gov/17040560/)
 98. A. D. Sanders, E. Falconer, M. Hills, D. C. J. Spierings, P. M. Lansdorp, Single-cell template strand sequencing by Strand-seq enables the characterization of individual homologs. *Nat. Protoc.* **12**, 1151–1176 (2017). doi: [10.1038/nprot.2017.029](https://doi.org/10.1038/nprot.2017.029); PMID: [28492527](https://pubmed.ncbi.nlm.nih.gov/28492527/)
 99. Standing Committee on Human Cytogenetic Nomenclature, *ISCN 1995: An International System for Human Cytogenetic Nomenclature: Recommendations of the International Standing Committee on Human Cytogenetic Nomenclature* (Karger Medical and Scientific Publishers, 1995).
 100. A. R. Quinlan, I. M. Hall, BEDTools: A flexible suite of utilities for comparing genomic features. *Bioinformatics* **26**, 841–842 (2010). doi: [10.1093/bioinformatics/btq033](https://doi.org/10.1093/bioinformatics/btq033); PMID: [20110278](https://pubmed.ncbi.nlm.nih.gov/20110278/)
 101. H. Li, Minimap2: Pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094–3100 (2018). doi: [10.1093/bioinformatics/bty191](https://doi.org/10.1093/bioinformatics/bty191); PMID: [29750242](https://pubmed.ncbi.nlm.nih.gov/29750242/)
 102. K. Katoh, K. Misawa, K. Kuma, T. Miyata, MAFFT: A novel method for rapid multiple sequence alignment based on fast Fourier transform. *Nucleic Acids Res.* **30**, 3059–3066 (2002). doi: [10.1093/nar/gk41436](https://doi.org/10.1093/nar/gk41436); PMID: [12136088](https://pubmed.ncbi.nlm.nih.gov/12136088/)
 103. A. Stamatakis, RAxML version 8: A tool for phylogenetic analysis and post-analysis of large phylogenies. *Bioinformatics* **30**, 1312–1313 (2014). doi: [10.1093/bioinformatics/btu033](https://doi.org/10.1093/bioinformatics/btu033); PMID: [24451623](https://pubmed.ncbi.nlm.nih.gov/24451623/)
 104. M. P. Berger, P. J. Munson, A novel randomized iterative strategy for aligning multiple protein sequences. *Comput. Appl. Biosci.* **7**, 479–484 (1991). doi: [10.1093/bioinformatics/7.4.479](https://doi.org/10.1093/bioinformatics/7.4.479); PMID: [1747779](https://pubmed.ncbi.nlm.nih.gov/1747779/)
 105. O. Gotoh, Optimal alignment between groups of sequences and its application to multiple sequence alignment. *Comput. Appl. Biosci.* **9**, 361–370 (1993). doi: [10.1093/bioinformatics/9.3.361](https://doi.org/10.1093/bioinformatics/9.3.361); PMID: [8324637](https://pubmed.ncbi.nlm.nih.gov/8324637/)
 106. A. Gershman *et al.*, Epigenetic Patterns in a Complete Human Genome. *Science* **376**, eabj5089 (2022). doi: [10.1126/science.abj5089](https://doi.org/10.1126/science.abj5089)
 107. C. Jain *et al.*, Weighted minimizer sampling improves long read mapping. *Bioinformatics* **36** (Suppl. 1), i111–i118 (2020). doi: [10.1093/bioinformatics/btaa435](https://doi.org/10.1093/bioinformatics/btaa435); PMID: [32657365](https://pubmed.ncbi.nlm.nih.gov/32657365/)
 108. J. D. Parsons, Miropeats: Graphical DNA sequence comparisons. *Comput. Appl. Biosci.* **11**, 615–619 (1995). PMID: [8808577](https://pubmed.ncbi.nlm.nih.gov/8808577/)
 109. M. Lawrence *et al.*, Software for computing and annotating genomic ranges. *PLOS Comput. Biol.* **9**, e1003118 (2013). doi: [10.1371/journal.pcbi.1003118](https://doi.org/10.1371/journal.pcbi.1003118); PMID: [23950696](https://pubmed.ncbi.nlm.nih.gov/23950696/)
 110. H. Wickham *et al.*, Welcome to the tidyverse. *J. Open Source Softw.* **4**, 1686 (2019). doi: [10.21105/joss.01686](https://doi.org/10.21105/joss.01686)
 111. B. Gel, E. Serra, karyoploteR: An R/Bioconductor package to plot customizable genomes displaying arbitrary data. *Bioinformatics* **33**, 3088–3090 (2017). doi: [10.1093/bioinformatics/btx346](https://doi.org/10.1093/bioinformatics/btx346); PMID: [28575171](https://pubmed.ncbi.nlm.nih.gov/28575171/)
 112. Z. Gu, L. Gu, R. Eils, M. Schlesner, B. Brors, circize: Implements and enhances circular visualization in R. *Bioinformatics* **30**, 2811–2812 (2014). doi: [10.1093/bioinformatics/btu393](https://doi.org/10.1093/bioinformatics/btu393); PMID: [24930139](https://pubmed.ncbi.nlm.nih.gov/24930139/)
 113. J. Köster, S. Rahmann, Snakemake—A scalable bioinformatics workflow engine. *Bioinformatics* **28**, 2520–2522 (2012). doi: [10.1093/bioinformatics/bts480](https://doi.org/10.1093/bioinformatics/bts480); PMID: [22908215](https://pubmed.ncbi.nlm.nih.gov/22908215/)
 114. J. Köster, S. Rahmann, Snakemake—a scalable bioinformatics workflow engine. *Bioinformatics* **34**, 3600 (2018). doi: [10.1093/bioinformatics/bty350](https://doi.org/10.1093/bioinformatics/bty350); PMID: [29788404](https://pubmed.ncbi.nlm.nih.gov/29788404/)
 115. F. Molder *et al.*, Sustainable data analysis with Snakemake. *F1000 Res.* **10**, 33 (2021). doi: [10.12688/f1000research.29032.2](https://doi.org/10.12688/f1000research.29032.2); PMID: [34035898](https://pubmed.ncbi.nlm.nih.gov/34035898/)
 116. M. R. Vollger, Assemblies and data generated for “Segmental duplications and their variation in a complete human genome” (2021). <https://zenodo.org/record/4726156>.
 117. M. R. Vollger, mrVollger/Data-Analysis-for-SDs-in-T2T-CHM13, Zenodo (2021); <https://zenodo.org/record/5498994>.
 118. M. R. Vollger, Version v0.3, Interconnected snakemake workflows for annotation and analysis of assemblies, Zenodo (2021); <https://zenodo.org/record/5499093>.

ACKNOWLEDGMENTS

The authors thank T. Brown and A. Lo for help in editing this manuscript. **Funding:** This work was supported in part by the Intramural Research Program of the National Human Genome Research Institute, National Institutes of Health (S.N., S.K., and A.M.P.), grants from the U.S. National Institutes of Health (NIH grants 5R01HG002385, 5U01HG010971, and 1U01HG010973 to E.E.E.; 1R01HG011274 to K.H.M.; 5R01HG009190 to W.T.; and U41HG007234 to M.D.), and a grant from Futuro in Ricerca (2010-RBFR103CE3 to M.V.). E.E.E. is an investigator of Howard Hughes Medical Institute. **Author contributions:** Identification of SDs in T2T-CHM13 and analysis: M.R.V. PacBio genome sequence generation: K.M.M., A.P.L., and K.H. FISH experiments and analysis: L.M., M.V., M.R.V., and E.E.E. Iso-Seq analysis: P.C.D., M.R.V., and R.L. *TBC1D3* analysis: X.G. and M.R.V. Copy number analysis: M.R.V. and W.T.H. inversion analysis: D.P. and M.R.V. T2T-CHM13 assembly generation: S.N., S.K., and A.M.P. Refinement of SD annotations near centromeres: K.H.M. and M.R.V. UCSC browser: M.D., W.T.H., and M.R.V. Methylation analysis: M.R.V., A.G., W.T., and E.E.E. Analysis of regions with genomic instability: M.R.V. and A.S. Organization of tables: M.R.V., P.C.D., and X.G. Organization of supplementary material: M.R.V. Display items: M.R.V., X.G., and P.C.D. Manuscript writing: M.R.V., E.E.E., and X.G. with input from all authors. **Competing interests:** E.E.E. is an SAB member of Variant Bio, Inc; S.K. and K.H.M. received travel funds to speak at events hosted by Oxford Nanopore Technologies; W.T. has licensed two patents to Oxford Nanopore Technologies (US 8748091 and 8394584). **Data and materials availability:** CHM13hTERT cells were obtained for research use via a material transfer agreement with the University of Pittsburgh. PacBio HiFi data has been deposited into NCBI Sequence Read Archive (SRA) under the following accessions: SRX7897688, SRX7897687, SRX7897686, and SRX7897685 for CHM13; SRR14407677 and SRR14407676 for CHM1; SRR10382244, SRR10382245, SRR10382248, and SRR10382249 for HG002; PRJNA540705 for NA12878; PRJEB36100 for HG00733 and HG00514; ERX4787609, ERX4787607, ERX4787606, ERX4782632, and ERX4781730 for NA19240; PRJNA701308 for HG01109, HG01243, HG02080, HG02723, HG03125, and HG03492; and PRJNA659034 and PRJNA691628 for all nonhuman primate samples. The human lymphoblastoid cell lines GM24385, GM19240, HG00514, and HG00733 used in the FISH experiments were obtained from Coriell. The T2T-CHM13 v1.0 assembly can be found on NCBI (GCA_009914755.2) and all associated read data were uploaded to SRA under the BioProject identifier PRJNA686988. Table S13 contains the accession information for all Iso-Seq data used in the paper. The canonical rDNA unit used to estimate copy number can be found at the NCBI nucleotide repository (KY962518.1). Human and nonhuman primate genome assemblies, SD annotations, methylation data, and Liftoff gene models can be found at Zenodo (116). Code for Snakemake pipelines, data analysis, and figure generation are also available at Zenodo (117, 118).

SUPPLEMENTARY MATERIALS

[science.org/doi/10.1126/science.abj6965](https://doi.org/10.1126/science.abj6965)
Extended Materials and Methods
Figs. S1 to S25
Tables S1 to S14

26 May 2021; accepted 13 December 2021
[10.1126/science.abj6965](https://doi.org/10.1126/science.abj6965)

Segmental duplications and their variation in a complete human genome

Mitchell R. VollgerXavi GuitartPhilip C. DishuckLudovica MercuriWilliam T. HarveyAriel GershmanMark DiekhansArvis SulovariKatherine M. MunsonAlexandra P. LewisKendra HoekzemaDavid PorubskyRuiyang LiSergey NurkSergey KorenKaren H. MigaAdam M. PhillippyWinston TimpMario VenturaEvan E. Eichler

Science, 376 (6588), eabj6965. • DOI: 10.1126/science.abj6965

View the article online

<https://www.science.org/doi/10.1126/science.abj6965>

Permissions

<https://www.science.org/help/reprints-and-permissions>

Use of this article is subject to the [Terms of service](#)

Science (ISSN) is published by the American Association for the Advancement of Science. 1200 New York Avenue NW, Washington, DC 20005. The title *Science* is a registered trademark of AAAS.

Copyright © 2022 The Authors, some rights reserved; exclusive licensee American Association for the Advancement of Science. No claim to original U.S. Government Works