

Punctuated bursts in human male demography inferred from 1,244 worldwide Y-chromosome sequences

G. David Poznik^{1,2,*}, Yali Xue^{3,*}, Fernando L. Mendez², Thomas F. Willems⁴, Andrea Massaia³, Melissa A. Wilson Sayres⁵, Qasim Ayub³, Shane A. McCarthy³, Apurva Narechania⁶, Seva Kashin⁷, Yuan Chen³, Ruby Banerjee³, Juan L. Rodriguez-Flores⁸, Maria Cerezo³, Haojing Shao⁹, Melissa Gymrek¹⁰, Ankit Malhotra¹¹, Sandra Louzada³, Rob Desalle¹², Graham R. S. Ritchie³, Eliza Cerveira¹¹, Tomas W. Fitzgerald³, Erik Garrison³, Anthony Marcketta¹³, David Mittelman¹⁴, Mallory Romanovitch¹¹, Chengsheng Zhang¹¹, Xiangqun Zheng-Bradley¹⁵, Goncalo R. Abecasis¹⁶, Steven A. McCarroll¹⁷, Paul Flicek¹⁵, Peter A. Underhill², Lachlan Coin⁹, Daniel R. Zerbino¹⁵, Fengtang Yang³, Charles Lee^{11,18}, Laura Clarke¹⁵, Adam Auton¹³, Yaniv Erlich¹⁹, Robert E. Handsaker⁷, The 1000 Genomes Project Consortium, Carlos D. Bustamante^{2,20} and Chris Tyler-Smith³

¹Program in Biomedical Informatics, Stanford University, Stanford, CA 94305, USA.

²Department of Genetics, Stanford University, Stanford, CA 94305, USA.

³The Wellcome Trust Sanger Institute, Wellcome Genome Campus, Hinxton, Cambridgeshire CB10 1SA, UK.

⁴Computational and Systems Biology Program, MIT, Cambridge, MA 02139, USA.

⁵School of Life Sciences and The Biodesign Institute, Arizona State University, Tempe, AZ, USA.

⁶Sackler Institute for Comparative Genomics, American Museum of Natural History, New York, NY, USA.

⁷Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA.

⁸Department of Genetic Medicine, Weill Cornell Medical College, New York, NY, USA.

⁹Institute for Molecular Bioscience, University of Queensland, St Lucia, QLD 4072, Australia.

¹⁰Harvard-MIT Division of Health Sciences and Technology, MIT, Cambridge, MA 02139, USA.

¹¹The Jackson Laboratory for Genomic Medicine, 10 Discovery Drive, Farmington, CT 06032, USA.

¹²Sackler Institute for Comparative Genomics, American Museum of Natural History, New York, NY, USA.

¹³Department of Genetics, Albert Einstein College of Medicine, 1301 Morris Park Avenue, Bronx, NY 10461, USA.

¹⁴Virginia Bioinformatics Institute and Department of Biological Sciences, Virginia Tech, VA 24061, USA.

¹⁵European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Trust Genome Campus, Hinxton, Cambridge CB10 1SD, UK.

¹⁶Department of Biostatistics, School of Public Health, University of Michigan, Ann Arbor, MI 48109, USA.

¹⁷Harvard Medical School, Boston, MA 02115, USA.

¹⁸Department of Graduate Studies – Life Sciences, Ewha Womans University, Ewhayeodae-gil, Seodaemun-gu, Seoul, South Korea 120-750

¹⁹New York Genome Center, New York, NY 10013, USA. Columbia University, New York, NY 10027, USA.

²⁰Department of Biomedical Data Science, Stanford University, Stanford, CA 94305, USA.

*These authors contributed equally to this work.

Correspondence should be addressed to C.D.B. (cdbustam@stanford.edu) or C.T.-S. (cts@sanger.ac.uk).

Abstract

We report the sequences of 1,244 human Y chromosomes randomly ascertained from 26 worldwide populations by the 1000 Genomes Project. We discovered more than 65,000 variants, including SNVs, MNVs, indels, STRs, and CNVs. Of these, CNVs contribute the greatest predicted functional impact. We constructed a calibrated phylogenetic tree based on binary SNVs and projected the more complex variants onto it, estimating the numbers of mutations for each class. Our phylogeny reveals bursts of extreme expansions in male numbers that have occurred independently among each of the five continental super-populations examined, at times of known migrations and technological innovations.

Main Text

Due to its male-specific inheritance and the absence of crossover for most of its length, which together link it completely to male phenotype and behavior, the Y chromosome bears a unique record of human history¹. Previous studies have demonstrated the value of full sequences for characterizing and calibrating the human Y-chromosome phylogeny^{2,3}. This work has led to insights into male demography, but further work is needed: to more comprehensively describe the range of Y-chromosome variation, including non-SNV classes of variation; to investigate the mutational processes operating in the different classes; and to understand the relative roles of selection⁴ and demography⁵ in shaping Y-chromosome variation. The role of demography has risen to prominence with reports of male-specific bottlenecks in several geographical areas after 10 thousand years ago (10 kya)⁵⁻⁷, at times putatively associated with the spread of farming⁵ or Bronze Age culture⁶. With improved calibration of the Y-SNV mutation rate⁸⁻¹⁰ and, consequently, more secure dating of relevant features of the Y-chromosome phylogeny, it is now possible to hone such interpretations.

We have conducted a comprehensive analysis of Y-chromosome variation using the largest extant sequence-based survey of global genetic variation, Phase 3 of the 1000 Genomes Project¹¹. With this, we have documented the extent of, and biological processes acting on, six types of genetic variation, and we have generated new insights into human male history.

Results

Dataset

Our dataset comprises 1,244 Y chromosomes sampled from 26 populations and sequenced to a median haploid coverage of 4.3×. Reads were mapped to the GRCh37 human reference assembly used by Phase 3 of the 1000 Genomes Project¹¹ and to the GRCh38 reference for our analysis of short tandem repeats (STRs). We used multiple haploid-tailored methods to call variants and generate callsets containing more than 65,000 variants of six types, including single nucleotide variants (SNVs) (**Supplementary Information [SI] 1.1**), multiple nucleotide variants (MNVs) (**SI 1.2**), short insertions/deletions (indels) (**SI 1.2**), copy-number variants (CNVs) (**SI 2**), and STRs (**SI 3**). We also identified karyotype variation that included one instance of 47,XXY and several mosaics of the karyotypes 46,XY and 45,X (**SI 2.2.3**). We applied stringent quality control to meet the Project's requirement of FDR < 5% for SNVs, indels and MNVs, and CNVs, and in our validation analysis with independent datasets, genotype concordance was greater than 99% for SNVs and was 86%–97% for the more complex variants (**Table 1**).

We generated six distinct callsets to construct a set of putative SNVs, which we input to a consensus genotype caller. In an iterative process, we leveraged the phylogeny to tune the final genotype calling strategy. We used similar methods for MNVs and indels, and we ran HipSTR¹² to call STRs.

We discovered CNVs from the sequence data using two approaches, GenomeSTRiP (SI 2.1.1) and CnvHitSeq (SI 2.1.2), and we validated calls using array comparative genomic hybridization (aCGH) (SI 2.2.1), supplemented by fluorescence *in situ* hybridization onto DNA fibres (fibre-FISH) in a few cases (SI 2.2.2). **Figure 1** illustrates a representative large deletion we discovered in a single individual using GenomeSTRiP (**Fig. 1b**). We validated its presence by aCGH (**Fig. 1c**) and ascertained its structure with fibre-FISH (**Fig. 1d**). Notably, the event that gave rise to this variant was not a simple recombination between the segmental duplication elements it partially encompasses (**Fig. 1a** and **Fig. 1d**).

Phylogeny

We identified each individual's Y-chromosome haplogroup and constructed a maximum-likelihood phylogenetic tree using 60,555 biallelic SNVs derived from 10.3 megabases of accessible DNA (**Fig. 2** and **Supplementary Fig. 14**) (SI 4). Our tree recapitulates and refines the expected structure^{2,3,5}, with all but two major haplogroups from A0 through T represented. The only haplogroups absent are M and S, both subgroups of K2b1 that are largely specific to New Guinea, which was not included in the 1000 Genomes Project. Notably, the branching patterns of several lineages suggest extreme expansions around 50–55 kya and also within the last few millennia. We investigated these later expansions in some detail and describe our findings in the “Demographic Analysis” section below.

When calibrated with a mutation rate of 0.76×10^{-9} mutations/site/year⁹, the time to the most recent common ancestor (TMRCA) of the tree is ~190 ky, but we consider the implications of alternative mutation rate estimates in the “Discussion” section. Of the clades resulting from the four deepest branching events, all but one are exclusive to Africa, and the TMRCA of all non-African lineages (i.e., the TMRCA of haplogroups DE and CF) is ~76 ky. We see a major increase in the number of lineages outside Africa ~50–55 kya, perhaps reflecting the geographic expansion and differentiation of Eurasian populations as they settled the vast expanse of these continents. Consistent with previous proposals¹³, a parsimonious interpretation of the phylogeny is that the predominant African haplogroup, E, arose outside the continent. This model of geographic segregation within the CT clade requires just one continental haplogroup exchange (E to Africa), rather than three (D, C, and F out of Africa). Furthermore, the timing of this putative return to Africa—between the emergence of E and its differentiation within Africa by 58 kya—is consistent with proposals, based on non-Y data, of abundant gene flow between Africa and nearby regions of Asia 50–80 kya¹⁴.

Three novel features of the phylogeny underscore the importance of South and Southeast Asia as likely locations where lineages currently distributed throughout Eurasia first diversified. First, we observed in a Vietnamese individual a rare F lineage that is an outgroup for the rest of the megahaplogroup. This sequence includes the derived allele for 147 SNVs shared by, and specific to, the 857 F chromosomes in our sample, but the lineage split off from rest of the group ~55 kya. This finding enabled us to define a new megagroup, GHIJK-M3658, whose subclades include the vast majority of the world's non-African males¹ (SI 4.4.6). Second, we identified in 12 South Asian individuals a new clade, here designated

“H0,” that split with the rest of haplogroup H ~51 kya. This new structure highlights the ancient diversity within the haplogroup and requires a more inclusive redefinition using, for example, the deeper SNV M2713, a G→A mutation at GRCh37 coordinate 6,855,809 (**SI 4.4.8**). Third, a lineage carried by a South Asian Telugu individual, HG03742, enabled us to refine early differentiation within the K2a clade ~50 kya. Using the high resolving power of the SNVs in our phylogeny, we determined that this lineage split off from the branch leading to haplogroups N and O (NO) not long after the ancestors of two individuals with well-known ancient DNA (aDNA) sequences did. Ust²-Ishim⁹ and Oase1¹⁵ lived in Western Siberia 43–47 kya and Romania 37–42 kya, respectively. Their Y chromosomes join HG03742 in sharing with haplogroup NO the derived T allele at M2308 (GRCh37 Y:7,690,182), and the modern sample shares just four additional mutations (**SI 4.4.11**) with the NO clade.

Mutations

To map each SNV to a branch (or branches) of the phylogeny, we first partitioned the tree into eight overlapping subtrees. Within each subtree, we provisionally assigned each SNV to the internal branch constituting the minimum superset of carriers of one allele or the other, designating the derived state to the allele specific to this clade. When no member of the clade bore the ancestral allele, we deemed the site compatible with the subtree and assigned the SNV to the branch (**SI 4.3, Supplementary Data File 5**). Most SNVs (94%) mapped to a single branch of the phylogeny, corresponding to a single mutation event during the Y-chromosome history captured by this tree. We projected the other variants onto the tree to infer the number of mutations associated with each (**Fig. 3a**).

Supplementary Figure 10 summarizes our workflow to count the number of independent mutation events associated with each CNV. We found that 39% of CNVs have mutated multiple times, a much higher proportion than SNVs (**Fig. 3a**). CNVs can arise by several different mutation mechanisms, one of which is homologous recombination between misaligned repeated sequences. This mechanism is particularly susceptible to recurrent mutations¹⁶ but, in comparing CNVs associated with repeated sequences to those that are not repeat-associated, we did not observe a significant difference in the proportion that have mutated multiple times (Mann-Whitney two-sided test). We did, however, observe that repeat-associated CNVs tend to be longer ($p = 0.01$).

We inferred more than six independent mutation events for each of three CNVs. One in particular stood out with 154 events. An apparent CNV hotspot spans a gene-free stretch of the chromosome’s long arm at GRCh37 Y:22,216,565–22,512,935. The region includes two arrays of long terminal repeat 12B (*LTR12B*) elements that together harbor 48 of the genome’s 211 copies (23%). In principle, our inference of numerous independent mutations could have been due to a “shadowing” effect from *LTR12B* elements elsewhere in the genome. That is, mismapping sequencing reads, and cross-hybridizing CGH probes, can lead to false inference of variation. But, in a phylogenetic analysis of all 211 *LTR12B* elements (**Supplementary Figure 11**), those within the putative CNV hotspot formed a pure monophyletic clade, demonstrating that the copy-number signal was genuine. The CNV has no predicted functional consequence.

Short tandem repeats (STRs) constituted the most mutable variant class, with a median of 16 mutations per locus and an average mutation rate of 3.9×10^{-4} mutations/generation. Assuming a generation time of 30 years, this equates to 1.3×10^{-5} mutations/year. Allele

length explains more than half the variance of the log mutation rate for uninterrupted STRs. Longer STRs mutate more rapidly, and, conditional on allele length, mutability decreases when the repeat structure is interrupted, with a general trend toward slower mutations rates for STRs with more interruptions (**Fig. 3b**).

Functional Impact

A small proportion of SNVs have a predicted functional impact (**SI 5**). Among 60,555 SNVs, we observed two singleton premature stop-codons, one each in *AMELY* and *USP9Y*, and one splice-site SNV that affects all known transcripts of *TBL1Y*. Among 94 missense SNVs with SIFT scores, all 30 deleterious variants are singletons or doubletons, while 17/64 tolerated variants are present at higher frequency ($p = 0.001$), underscoring the impact of purifying selection on variation at protein-coding genes. No STRs overlapped protein-coding regions, but, in contrast to the SNVs, a high proportion of CNVs have a predicted functional impact.

Twenty of 100 CNVs in our final callset overlap with 27 protein-coding genes from 17 of the 33 Y-chromosome gene families. In our analysis of 1000 Genomes autosomal data, we observed that the ratio of the proportion of deletions overlapping protein-coding genes to the proportion of duplications overlapping protein-coding genes is 0.84. Whereas on the autosomes deletions are less likely to overlap protein-coding genes than duplications are, as others have also reported¹⁷, we found the reverse to be true for the Y chromosome. Despite its haploidy, we calculated its ratio of proportions to be 1.5, indicating a surprising increased tolerance of gene loss, as compared with the diploid genes on autosomes.

Demographic Analysis

Given observed diversity levels of the autosomes, the X chromosome, and the mitochondrial genome (mtDNA) (**SI 6**), Y-chromosome diversity was reported to be lower than expected from simple population-genetic models that assume a Poisson-distributed number of offspring⁴, and the role of selection in this disparity is debated. We confirmed that Y-chromosome diversity in our sample is low and found that positing extreme male-specific bottlenecks in the last few millennia can lead to a good fit between modeled and observed relative diversity levels of the autosomes, the X chromosome, the Y chromosome, and the mtDNA (**SI 7.1**). Therefore, we conclude that Y diversity may be shaped primarily by neutral demographic processes.

To investigate punctuated bursts within the phylogeny and estimate growth rates, we modeled lineage growth as a rapid phase followed by a moderate phase and applied this model to lineages showing rapid expansions (**SI 7.2**), noting that such extreme expansions are seldom seen in the mtDNA phylogeny here (**SI 6**) or in other studies⁵. We examined 20 nodes of the tree whose branching patterns were well-fit by this model. These nodes were drawn from eight haplogroups and included at least one lineage from each of the five continental regions surveyed (**Fig. 4**). As the haplogroup expansions we report are among the most extreme yet observed in humans, we think it more likely than not that such events should correspond to historical processes that have also left archaeological footprints. Therefore, in what follows, we propose links between genetic and historical or archaeological data. We caution that, especially in light of as yet imperfect calibration, these connections remain unproven. But they are testable, for example using aDNA.

First, in the Americas, we observed expansion of Q1a-M3 at ~15 kya, the time of the initial colonization of the hemisphere¹⁸. This correspondence, based on one of the most thoroughly examined dates in human prehistory, attests to the suitability of the calibration we have chosen. Second, in sub-Saharan Africa, two independent E1b-M180 lineages expanded ~5 kya, a period before the numerical and geographical expansions of Bantu-speakers in whom E1b-M180 now predominates¹⁹. The presence of these lineages in non-Bantu-speakers (e.g., Yoruba, Esan) indicates an expansion pre-dating the Bantu migrations, perhaps triggered by the development of ironworking²⁰. Third, in Western Europe, related lineages within R1b-L11 expanded ~4.8–5.9 kya, most markedly around 4.8 and 5.5 kya. The earlier of these times, 5.5 kya, is associated with the origin of the Bronze Age Yamnaya culture. The Yamnaya have been linked by aDNA evidence to a massive migration from the Steppe, which may have replaced much of the previous European population^{21,22}, but the six Yamnaya with informative genotypes did not bear lineages descending from or ancestral to R1b-L11, so a Y-chromosome connection has not been established. The later time, 4.8 kya, coincides with the origins of the Corded Ware (Battle Axe) culture in Eastern Europe and the Bell-Beaker culture in Western Europe²³.

Potential correspondences between genetics and archaeology in South and East Asia have received less investigation. In South Asia, we detect eight lineage expansions dating to ~4.0–7.3 kya and involving haplogroups H1-M52, L-M11, and R1a-Z93. The most striking are expansions within R1a-Z93, ~4.0–4.5 kya. This time predates by a few centuries the collapse of the Indus Valley Civilization, associated by some with the historical migration of Indo-European speakers from the western steppes into the Indian sub-continent²⁴. There is a notable parallel with events in Europe, and future aDNA evidence may prove to be as informative as it has been in Europe. Finally, East Asia stands out from the rest of the Old World for its paucity of sudden expansions, perhaps reflecting a larger starting population or the coexistence of multiple prehistoric cultures wherein one lineage could rarely dominate. We observed just one notable expansion within each of the O2b-M176 and O3-M122 clades.

Discussion

The 1000 Genomes Project dataset provides a rich and unparalleled resource of Y-chromosome variation coupled with open access to DNA and cell lines that will facilitate diverse further investigations. By cataloging the phylogenetic position of ~60,000 SNVs, we have constructed a database of diagnostic variants with which one can assign Y-chromosome haplogroups to DNA samples. This resource is particularly valuable for SNP-chip design and for aDNA studies, in which sequencing coverage is often quite low, as exemplified by our reanalysis of the Ust'-Ishim and Oase1 Y chromosomes.

The variants we report have well-calibrated FDRs. Nevertheless, due to the modest sequencing coverage, data missingness was a principal concern. Small CNVs and long STRs are largely undetected, and low frequency variants in general, including SNVs, are under-represented. We therefore took great care to minimize the impact of missing variants. In particular, we designed the relevant downstream analyses to only use information from higher frequency, shared, variation, corresponding to mutations on internal branches of the tree.

Since many DNA samples were extracted from lymphoblastoid cells, another potential concern was variation that has arisen during cell culture²⁵. However, these false discoveries are inherently not shared. Therefore, the precautions we took to minimize the impact of

missingness also precluded in vitro mutations from influencing our findings. We discuss additional caveats on the mapping of SNVs to branches in **SI 4**.

Our findings illustrate unique properties of the Y chromosome. Foremost, the abundance of extreme male-lineage expansions underscores differences between male and female demographic histories. A caveat to our expansion analysis is that our inference method assumes that population structure did not affect the branching patterns immediately downstream of the particular phylogenetic node under investigation. This is reasonable, because population structure is unlikely when a very rapid expansion is in progress, but to accommodate this strong assumption, we limited all analyses to pruned internal subtrees short enough for it to hold. A second caveat regards the choice of calibration metric, which is relevant to the links we have suggested between expansions and historical or archaeological events. Present-day geographical distributions provide strong support for the correspondences we proposed for the initial peopling of most of Eurasia by fully modern humans 50–55 kya and for the first colonization of the Americas 15 kya. For later male-specific expansions, we should consider the consequences of alternative mutation rate estimates, where pedigree-based estimates integrated over a few centuries^{8,10,25} may be more relevant. The estimate from the largest set of mutations⁸ would lead to a decrease in expansion times by ~15%, increasing the precision of the correspondences proposed for E1b and R1a. For R1b, a 15% decrease would suggest an expansion postdating the Yamnaya migration, perhaps explaining better the distinction between the Yamnaya R1b chromosomes and the expanding R1b-L11 lineage. Either way, the lineage expansions seem to have followed innovations that may have elicited increased variance in male reproductive success²⁶, innovations such as metallurgy, wheeled transport, or social stratification and organized warfare. In each case, privileged male lineages could undergo preferential amplification for generations. We find that rapid expansions are not confined to unusual circumstances^{27,28}. Rather, they can dominate on a continental scale and do so in some of the populations most studied by medical geneticists. Inferences incorporating demography may benefit from taking these male-female differences into account.

Supplementary Information and Supplementary Data are available with the online version of the paper. The SI includes detailed descriptions of our methods for calling and validating each variant type and for conducting our phylogenetic, functional, and demographic analyses. Please see **SI 8** for details on how to access sequence read alignments (BAM files) and genotype calls (VCF files).

Acknowledgements We thank the 1000 Genomes Project sample donors for making this work possible and all Project members for their contributions. Figures were generated with FigTree²⁹ and ggplot2³⁰. Thanks to A. Martin for ADMIXTURE results. G.D.P. was supported by the National Science Foundation (NSF) Graduate Research Fellowship under grant number DGE-1147470 and by the National Library of Medicine training grant LM-007033. Work at The Wellcome Trust Sanger Institute (Q.A., R.B., M.C., Y.C., S.L., A.M., S.A.M., C.T.-S., Y.X., and F.Y.) was supported by Wellcome Trust grant number 098051. F.L.M. was supported by the National Institutes of Health (NIH) grant number 1R01GM090087, by NSF grant number DMS-1201234, and by a postdoctoral fellowship from the Stanford Center for Computational, Evolutionary and Human Genomics (CEHG). T.W. was supported by an AWS Education Grant, and the work of T.W., M.G. and Y.E was supported in part by an NIJ Award 2014-DN-BX-K089. M.C. is supported by a Fundacion Barrie Fellowship. M.G. was supported by a National Defense Science & Engineering Graduate Fellowship. G.R.S.R. was supported by the European Molecular Biology Laboratory and the Sanger Institute through an EBI-Sanger Postdoctoral Fellowship. X.Z.-B., P.F., D.R.Z. and L.C. were supported by Wellcome Trust grant number 085532 and by the European Molecular Biology Laboratory. C.L. was supported in part by NIH grant U41HG007497. Y.E. holds a Career Award at the Scientific Interface from the Burroughs Wellcome Fund. C.D.B. was supported by NIH grant number 5R01HG003229-09.

References

1. Jobling, M.A. & Tyler-Smith, C. The human Y chromosome: an evolutionary marker comes of age. *Nat Rev Genet* **4**, 598-612 (2003).
2. Wei, W. *et al.* A calibrated human Y-chromosomal phylogeny based on resequencing. *Genome Res* **23**, 388-95 (2013).
3. Poznik, G.D. *et al.* Sequencing Y chromosomes resolves discrepancy in time to common ancestor of males versus females. *Science* **341**, 562-5 (2013).
4. Wilson Sayres, M.A., Lohmueller, K.E. & Nielsen, R. Natural selection reduced diversity on human Y chromosomes. *PLoS Genet* **10**, e1004064 (2014).
5. Karmin, M. *et al.* A recent bottleneck of Y chromosome diversity coincides with a global change in culture. *Genome Res* **25**, 459-66 (2015).
6. Batini, C. *et al.* Large-scale recent expansion of European patrilineages shown by population resequencing. *Nat Commun* **6**, 7152 (2015).
7. Sikora, M.J., Colonna, V., Xue, Y. & Tyler-Smith, C. Modeling the contrasting Neolithic male lineage expansions in Europe and Africa. *Investig Genet* **4**, 25 (2013).
8. Helgason, A. *et al.* The Y-chromosome point mutation rate in humans. *Nat Genet* **47**, 453-7 (2015).
9. Fu, Q. *et al.* Genome sequence of a 45,000-year-old modern human from western Siberia. *Nature* **514**, 445-9 (2014).
10. Balanovsky, O. *et al.* Deep phylogenetic analysis of haplogroup G1 provides estimates of SNP and STR mutation rates on the human Y-chromosome and reveals migrations of Iranic speakers. *PLoS ONE* **10**, e0122968 (2015).
11. The 1000 Genomes Project Consortium. A global reference for human genetic variation. *Nature* **526**, 68-74 (2015).
12. Willems, T. HipSTR. (<https://github.com/tfwillems/HipSTR>, 2015).
13. Hammer, M.F. *et al.* Out of Africa and back again: nested cladistic analysis of human Y chromosome variation. *Mol Biol Evol* **15**, 427-41 (1998).
14. Groucutt, H.S. *et al.* Rethinking the dispersal of *Homo sapiens* out of Africa. *Evol Anthropol* **24**, 149-64 (2015).
15. Fu, Q. *et al.* An early modern human from Romania with a recent Neanderthal ancestor. *Nature* **524**, 216-219 (2015).
16. Zhang, F., Gu, W., Hurles, M.E. & Lupski, J.R. Copy number variation in human health, disease, and evolution. *Annu Rev Genomics Hum Genet* **10**, 451-81 (2009).
17. Sudmant, P.H. *et al.* Global diversity, population stratification, and selection of human copy-number variation. *Science* **349**, aab3761 (2015).
18. Raghavan, M. *et al.* Genomic evidence for the Pleistocene and recent population history of Native Americans. *Science* **349**, aab3884 (2015).
19. de Filippo, C., Bostoen, K., Stoneking, M. & Pakendorf, B. Bringing together linguistic and genetic evidence to test the Bantu expansion. *Proc Biol Sci* **279**, 3256-63 (2012).
20. Jobling, M., Hollox, E., Hurles, M., Kivisild, T. & Tyler-Smith, C. *Human Evolutionary Genetics, second edition*, (Garland Science, New York and London, 2014).
21. Allentoft, M.E. *et al.* Population genomics of Bronze Age Eurasia. *Nature* **522**, 167-72 (2015).
22. Haak, W. *et al.* Massive migration from the steppe was a source for Indo-European languages in Europe. *Nature* **522**, 207-11 (2015).

23. Harding, A.F. *European Societies in the Bronze Age*, (Cambridge University Press, Cambridge, UK, 2000).
24. Bryant, E.F. & Patton, L.L. (eds.). *The Indo-Aryan Controversy: Evidence and Inference in Indian History*, (Routledge, Abingdon, UK, 2005).
25. Xue, Y. *et al.* Human Y chromosome base-substitution mutation rate measured by direct sequencing in a deep-rooting pedigree. *Curr Biol* **19**, 1453–1457 (2009).
26. Betzig, L. Means, variances, and ranges in reproductive success: comparative evidence. *Evol Hum Behav* **33**, 309-317 (2012).
27. Zerjal, T. *et al.* The genetic legacy of the Mongols. *Am J Hum Genet* **72**, 717-21 (2003).
28. Balaresque, P. *et al.* Y-chromosome descent clusters and male differential reproductive success: young lineage expansions dominate Asian pastoral nomadic populations. *Eur J Hum Genet* **23**, 1413-22 (2015).
29. Rambaut, A. Figtree. (<http://tree.bio.ed.ac.uk/software/figtree/>, 2006).
30. Wickham, H. *ggplot2: Elegant Graphics for Data Analysis*, (Springer, New York, 2009).

Variant Type	Number	FDR (%)	Concordance (%)
SNVs	60,555	3.9	99.6
Indels & MNVs	1,427	3.6	96.4
CNVs	110	2.7	86
STRs	3,253	N/A	89–97

Table 1 | Y-chromosome variants discovered in 1,244 males. FDR, false discovery rate; Concordance, with independent genotype calls. CNVs considered are those computationally inferred using Genome STRiP. N/A, not available.

a

Segmental duplication in the human reference sequence

Y:17,986,738-17,995,460 Y:18,008,099-18,016,824

FISH probes

Custom PCR probes



BAC clone

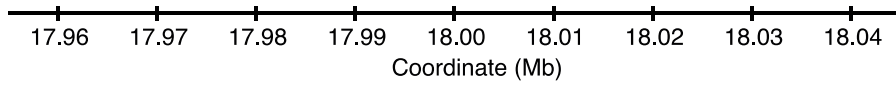


HG00183 deletion calls

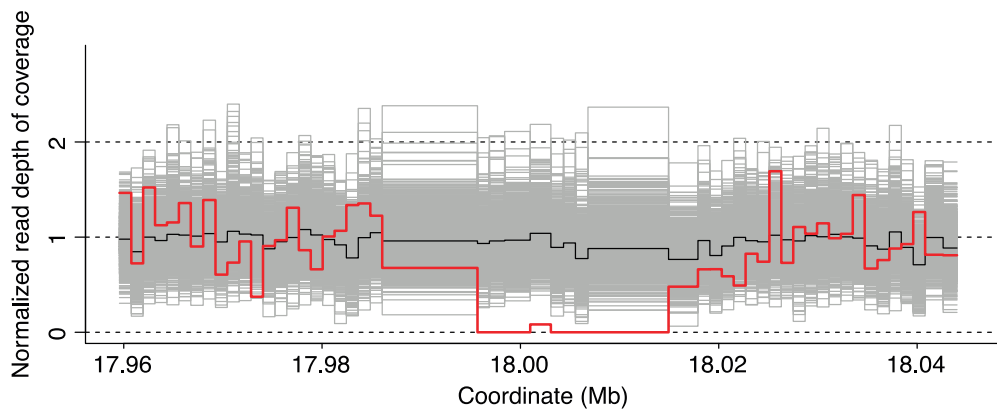
GenomeSTRiP



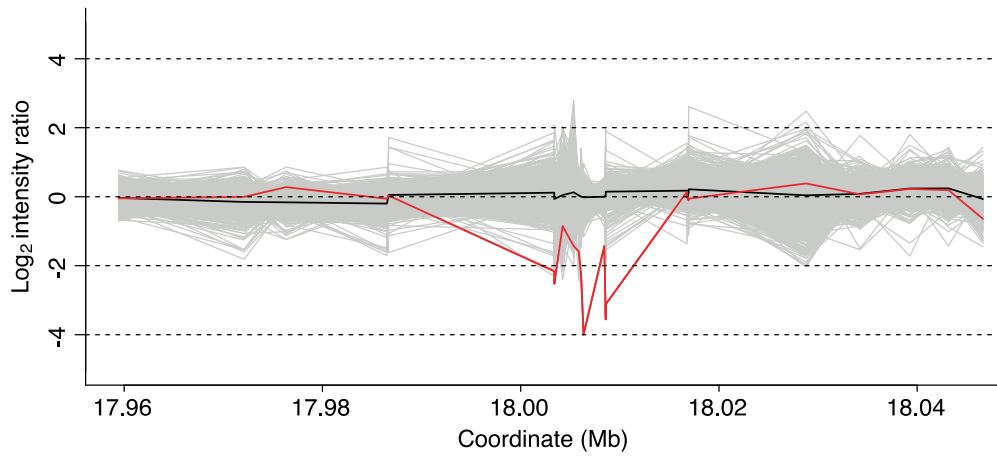
aCGH



b

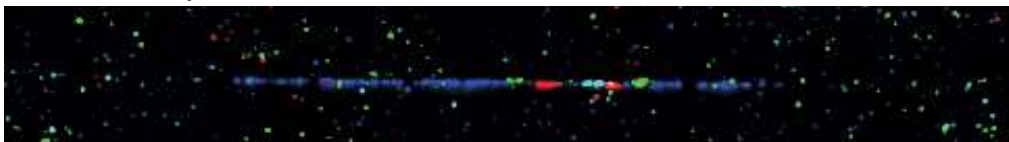


c



d

Reference sample: HG00096



Sample with deletion: HG00183



Figure 1 | Discovery and validation of a representative Y-chromosome CNV. a, The GRCh37 reference sequence contains an inverted segmental duplication within Y:17,986,738–18,016,824 (orange bars). We designed FISH probes to target the 3' termini of the two segments (red and green bars labeled “P1” and “P3,” respectively) and the unique region between them (light blue, “P2”). A fourth probe used reference sequence BAC clone RP11-12J24 (dark blue, “P4”). Unlabeled green and red bars indicate expected cross-hybridization, and black bars indicate CNV events called by Genome STRiP and aCGH, respectively. GenomeSTRiP called a 30-kb deletion that includes the duplicated segments and the unique spacer region, whereas aCGH lacks probes in the duplicated regions. **b,** Genome STRiP discovery plot. The red curve indicates the normalized read depth of HG00183, as compared to those for 1,232 other samples (grey) and the median (black curve). **c,** Validation by aCGH, indicating the \log_2 intensity ratio for HG00183 (red), versus 1,233 other samples (grey) and the median signal (black). **d,** Fibre-FISH validation using the probes illustrated in (a). The reference sample, HG00096, matches the human reference sequence, with green, red, light blue, red, and green hybridizations occurring in sequence. In contrast, we observed just one green and one red hybridization in HG00183, indicating the deletion of one copy of the segmental duplication and the central unique region. The consistent length scale across panels a–c does not apply to this panel, and the lengths of the BAC-clone hybridizations (dark blue) differ between the two samples solely due to the molecular combing process.

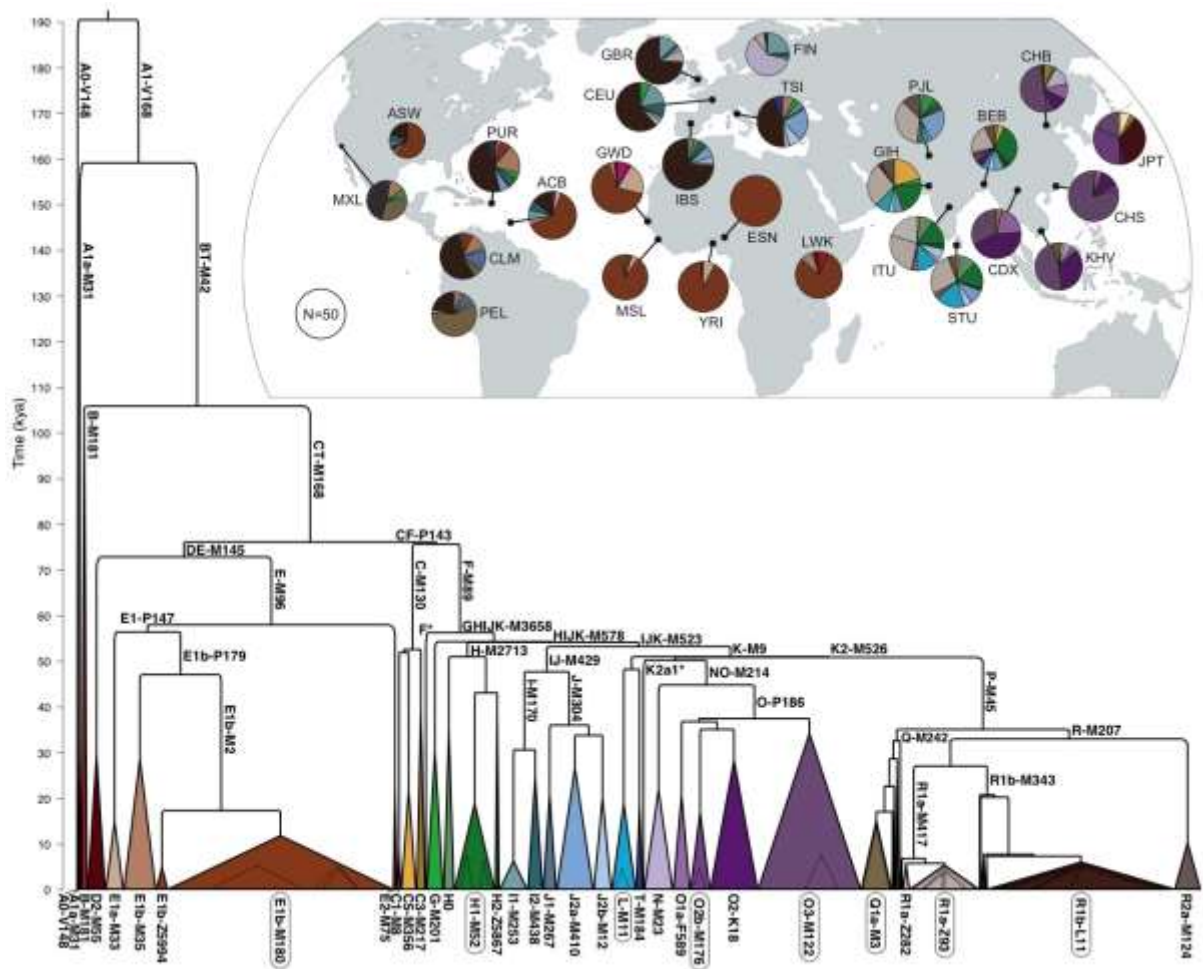


Figure 2 | Y-chromosome phylogeny and haplogroup distribution. Branch lengths are drawn proportional to the estimated times between successive splits, with the most ancient division occurring ~190 kya. Colored triangles represent the major clades, and the width of each base is proportional to one less than the corresponding sample size. We modeled expansions within eight of the major haplogroups (circled) (**Figure 4**), and dotted triangles represent the ages and sample sizes of the expanding lineages. (**Inset**) World map indicating, for each of the 26 populations, the geographic source, sample size, and haplogroup distribution.

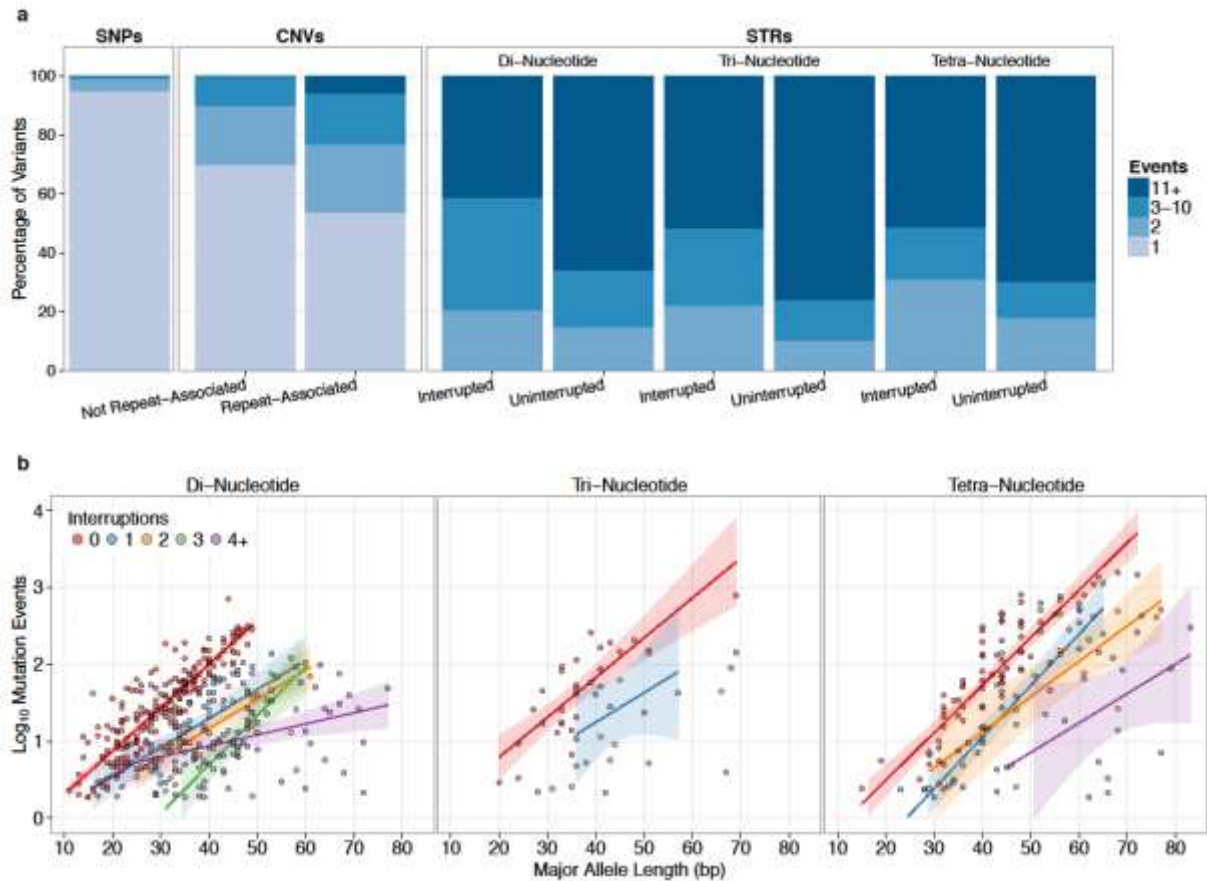


Figure 3 | Mutation events. **a**, Bar plots show the percentage of each variant type stratum associated with 1, 2, 3–10, or more mutations across the phylogeny. **b**, For STRs, scatter plots show the logarithm of the number of mutation events versus major allele length, stratified by motif length and the number of interruptions to the repeat structure. We have plotted regressions lines for categories with at least 10 data points, and we have omitted from the plots 44 STRs with motif lengths greater than four and 91 STRs whose mutation rate estimates were equal to the minimum threshold of 10^{-5} mutations/generation.

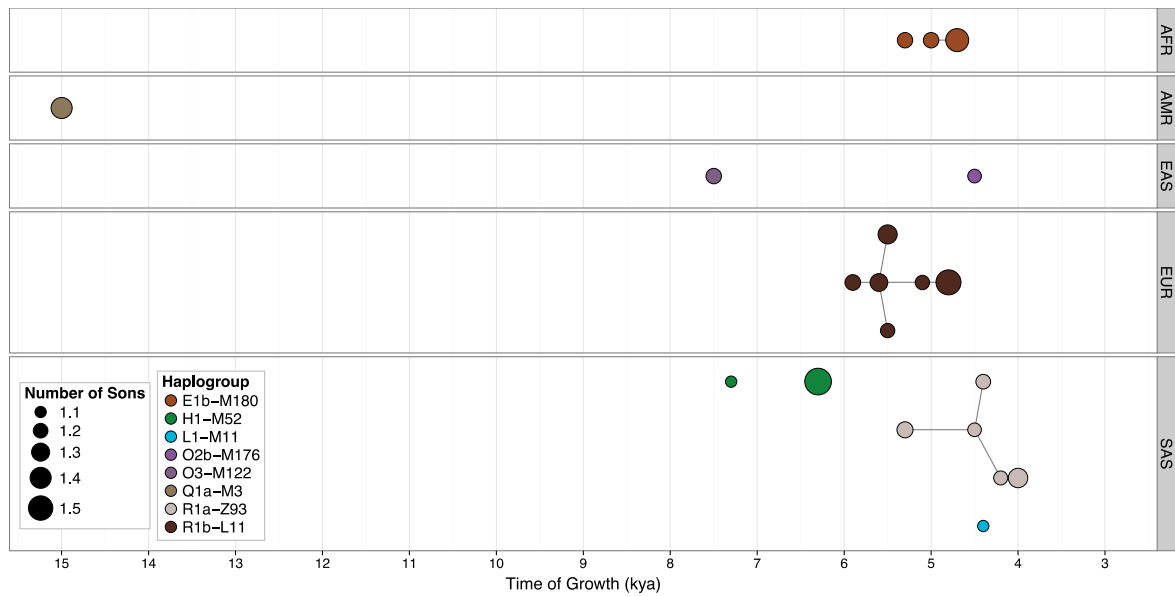


Figure 4 | Explosive male-lineage expansions of the last 15 thousand years. Each circle represents a phylogenetic node whose branching pattern suggests rapid expansion. The x -axis indicates the timings of the expansions, and circle radii reflect growth rates—the minimum number of sons per generation, as estimated by our two-phase growth model. Nodes are grouped by continental super-population (AFR, African; AMR, Admixed American; EAS, East Asian; EUR, European; SAS, South Asian) and colored by haplogroup. Line segments connect phylogenetically nested lineages.