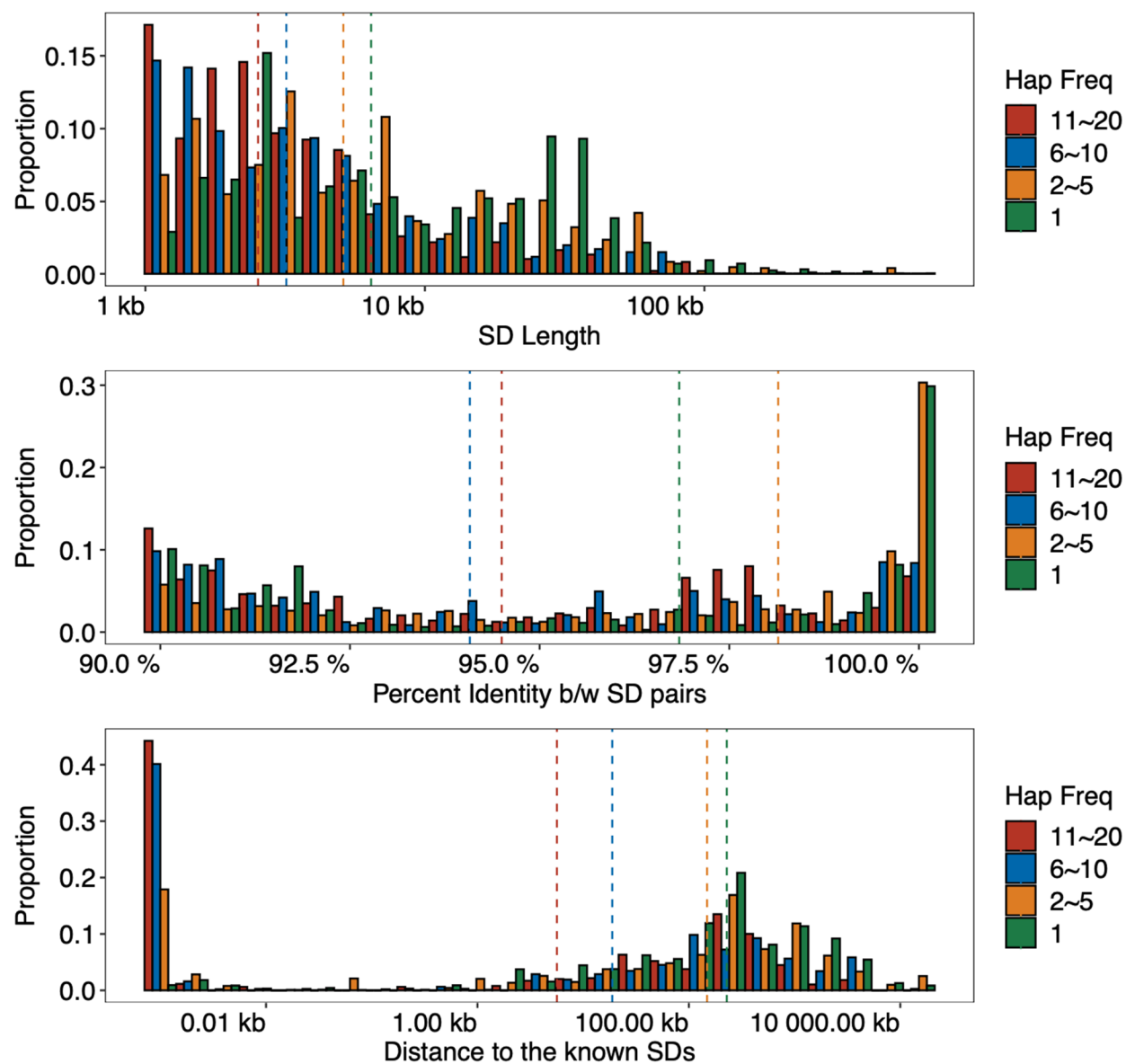


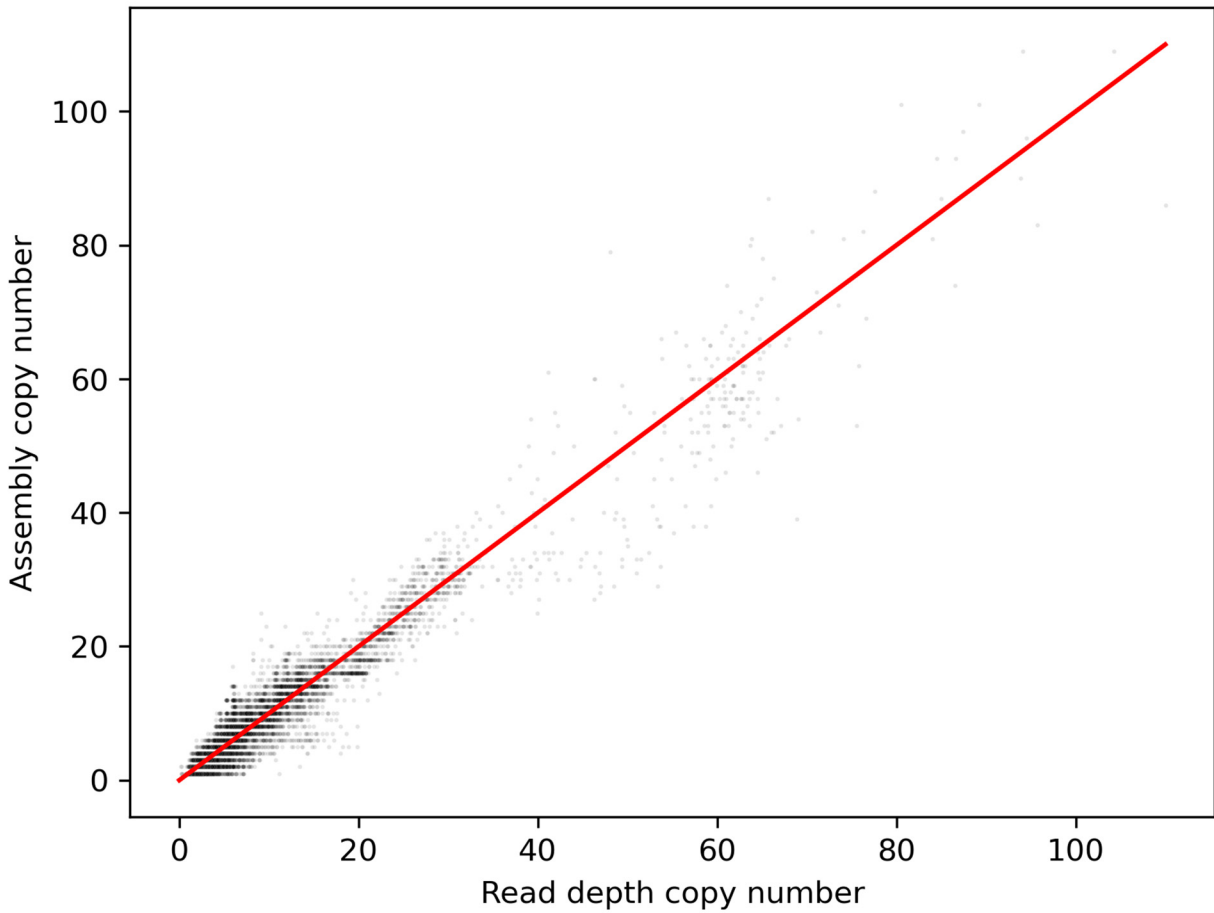
Structural polymorphism and diversity of human segmental duplications

In the format provided by the
authors and unedited

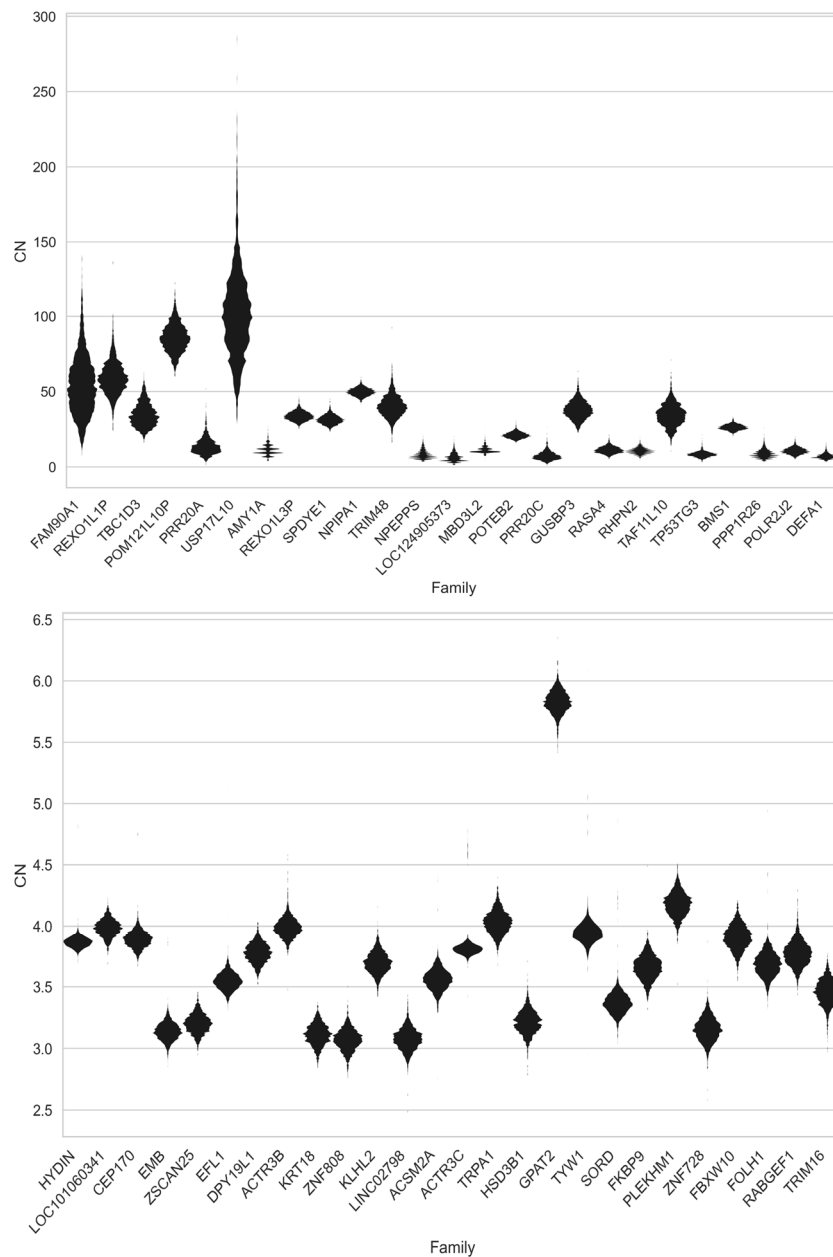
SUPPLEMENTARY INFORMATION & FIGURES for Jeong et al.



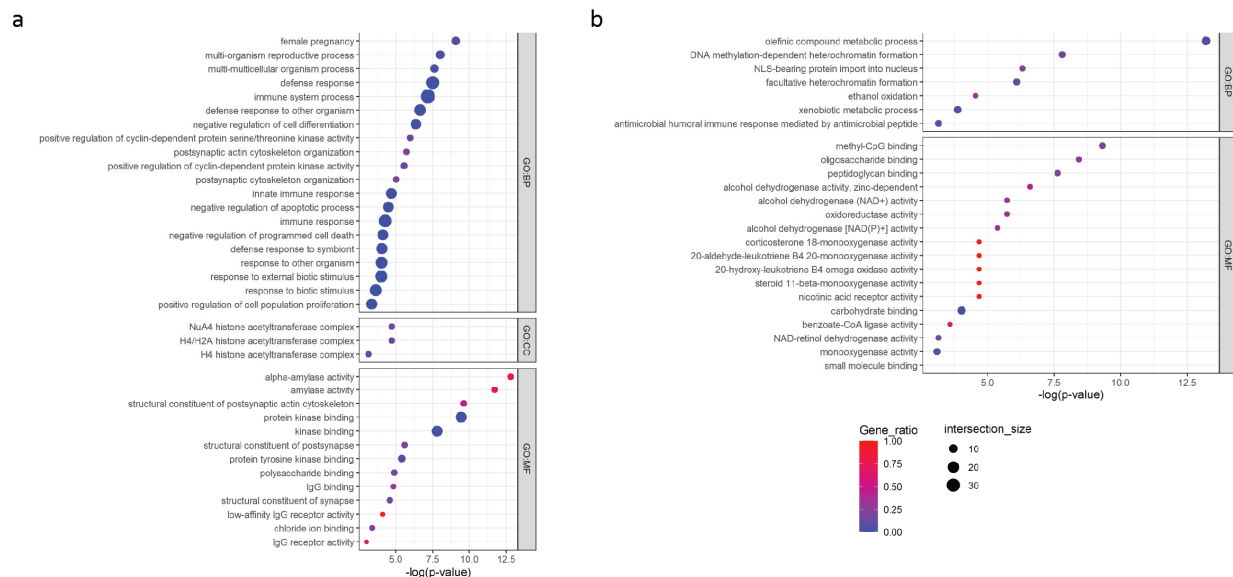
Supplementary Figure 1. Histogram comparing the sequence identity and length of polymorphic segmental duplications (SDs) at different haplotype frequencies.



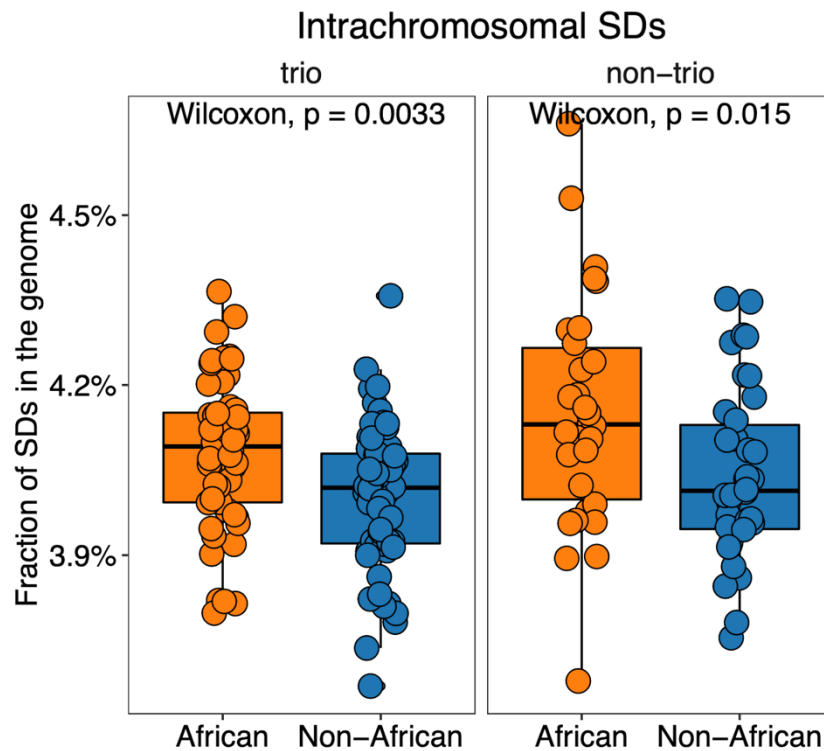
Supplementary Figure 2. Read-depth-based copy numbers estimated with fastCN compared to assembled copy number for each sample, summed between the two haplotypes. Each point represents the copy number estimates for a gene family in a sample ($R^2 = 0.94$).



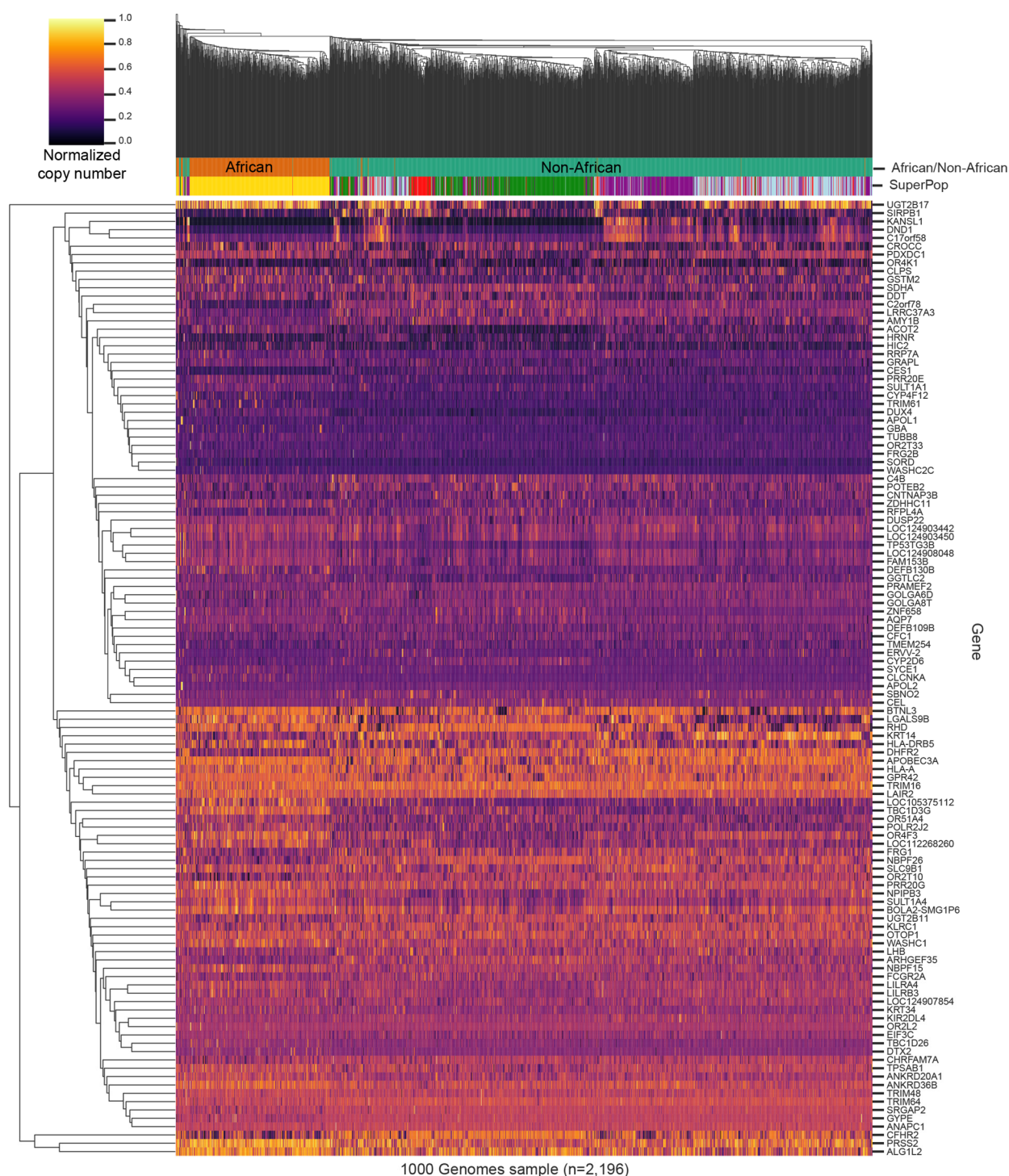
Supplementary Figure 3. Copy number distribution of high- and low-variance gene families. The read-depth copy number of gene families with highly variable (above) and nearly fixed copy number (below) are displayed. Gene families are selected and ordered by variance, requiring an average diploid copy number greater than three.



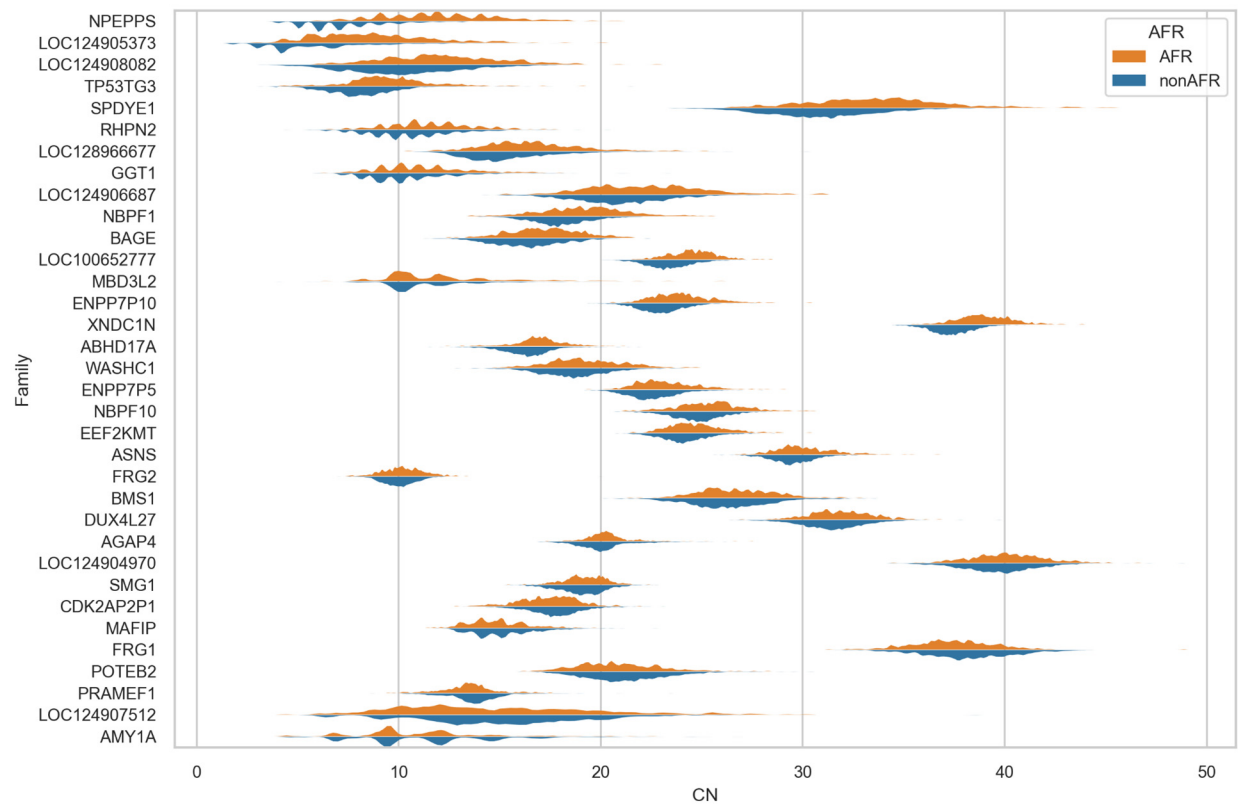
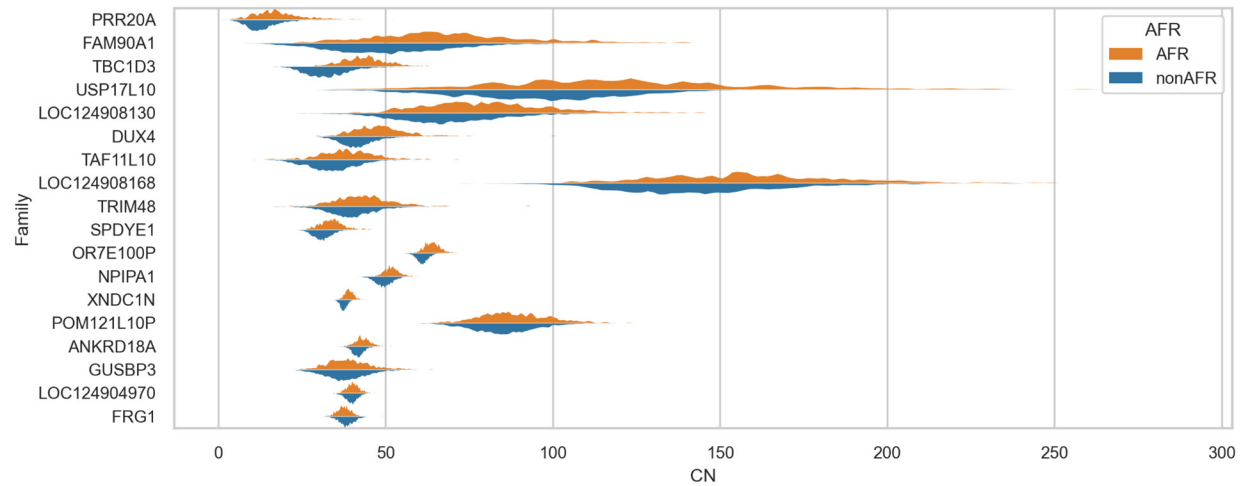
Supplementary Figure 4. Gene ontology enrichment of the (a) top 100 variable gene families (n = 358) and (b) invariable genes (n = 115). The x-axis indicates negative log transformed adjusted p-value. The number of intersecting genes is indicated by the size of the circle and the gene ratio represents the number of intersect/term size. To test statistical enrichment of gene ontology, Fisher's one-tailed test was performed. Multiple test corrections were done by the default g:SCS method of the g:profiler package, accounting for ontologies that are not independent from one another, which is less stringent test than Bonferroni correction or Benjamini-Hochberg FDR.

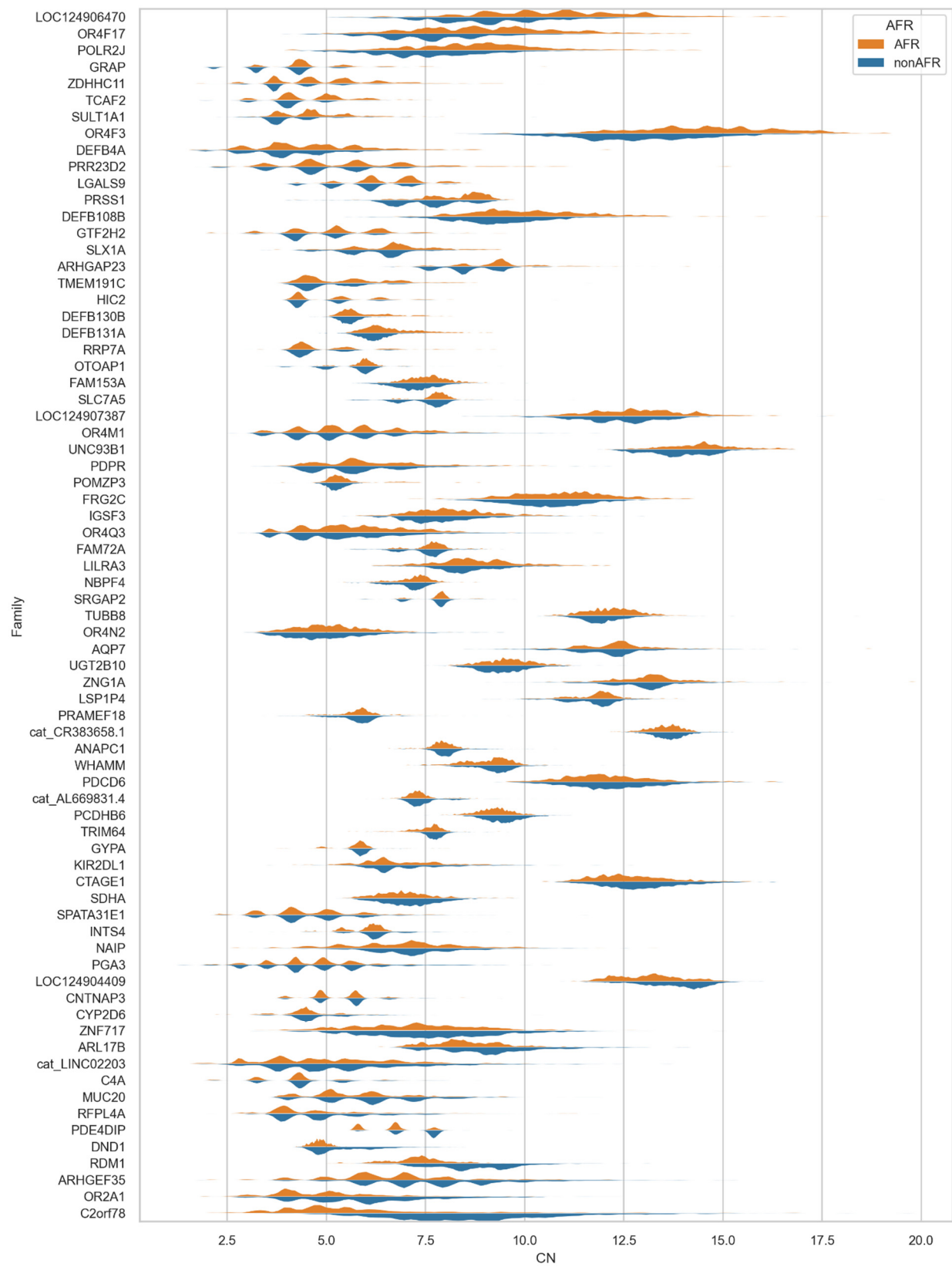


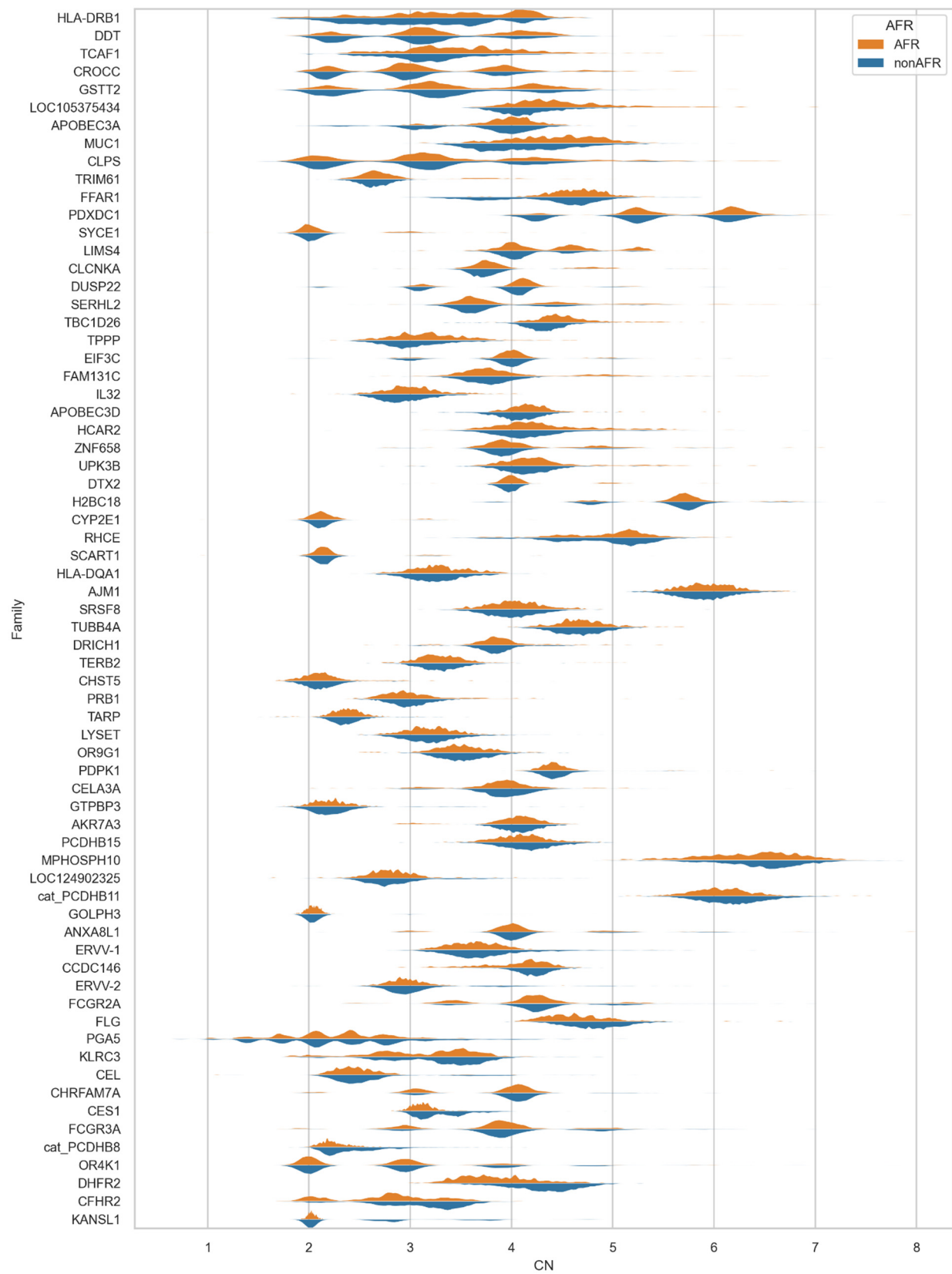
Supplementary Figure 5. Comparison of SD content in trio (n = 102) vs. non-trio (n = 68) genomes. Both datasets show a significant excess of SD content in African samples. The box plot ranges represent the interquartile range (first and third quartile), and the median is indicated by a horizontal line in each box. The whisker indicates the datapoints within 1.5*interquartile range. Two-tailed Wilcox ranked sum test was performed for statistical analyses.

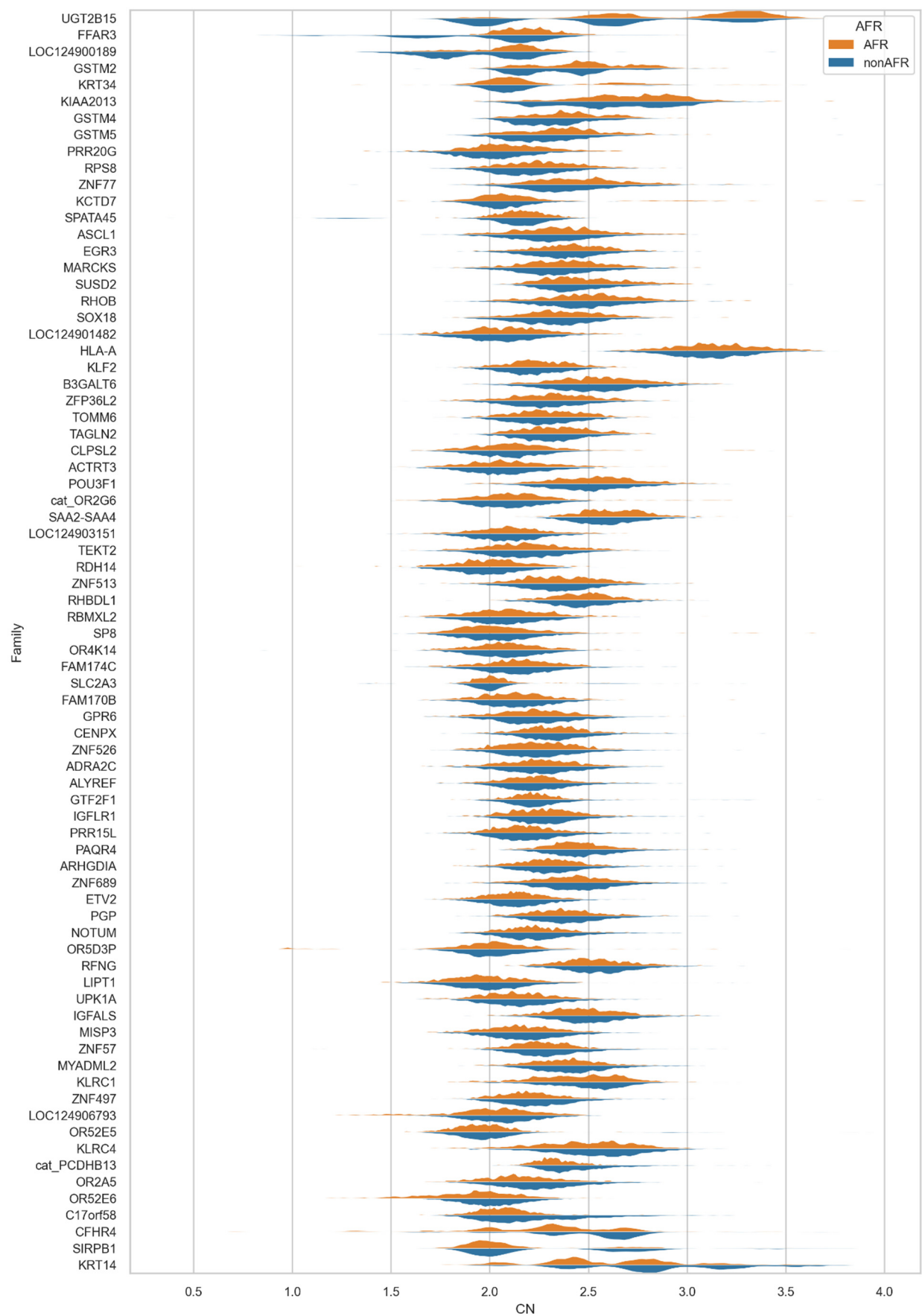


Supplementary Figure 6. Population-stratified genic copy number in 2,196 unrelated individuals from the 1000 Genomes Project. Gene copy number values are centered on the mean for each gene and scaled by unit variance to range from 0-1. One paralog per gene family and duplication block is shown. 73/115 deduplicated population-stratified genes have higher mean copy number in the African group as compared to the non-African group ($p=0.002$, two-tailed Wilcoxon ranked sum test). Superpopulations as described in the 1000 Genomes Project are shown above (Africa: gold, East Asia: green, South Asia: purple, Europe: blue, the Americas: red). Copy number estimates by population group (below).



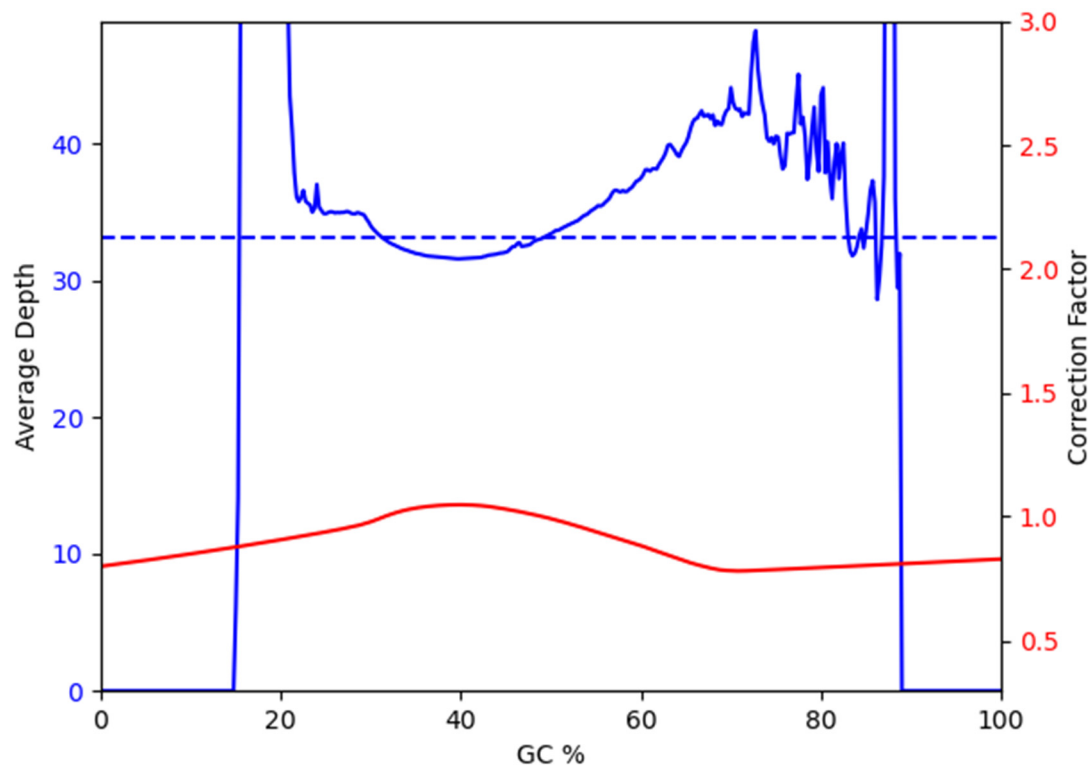






Supplementary Figure 7. Population-stratified gene copy number in 1000

Genomes Project. To validate population differences in gene copy number, we repeated the analysis on an unrelated subset of high-coverage Illumina data from the 1000 Genomes Project, excluding samples identified as related by published pedigree or with cryptic relationships (2nd degree or less) as identified by Somalier (Pedersen, Genome Medicine 2020). Individuals from the highly-admixed ACB, ASW, and PUR populations were also excluded, leaving $n=2,196$.



Supplementary Figure 9. Adjustment factor for gene copy number estimation.

The read depth of k-mers from decomposed T2T genome assembly is shown as a function of GC composition (blue). Even in a finished genome where there is no experimental or technical error, k-mer read depth is not uniform, since as GC and AT content increases so too does the number of low-complexity k-mers mapping elsewhere. Based on this, we estimated an adjustment factor required by fastCN (red line) to correct for this bias.

SUPPLEMENTARY NOTE

Full list of all consortium members:

Human Genome Structural Variation Consortium (HGSVC)

Adam M. Phillippy¹, Alexander T. Dilthey^{2,3}, Arda Söylev^{2,4}, Arvis Sulovari⁵, Benedict Paten⁶, Bida Gu⁷, Carolyn A. Paisie⁸, Charles Lee⁸, Chen-Shan Chin⁹, Chong Li^{10,11}, Christine R. Beck^{8,12}, David Porubsky⁵, DongAhn Yoo⁵, Evan E. Eichler^{5,13}, Feyza Yilmaz⁸, Gianni V. Martino^{14,15,16}, Glenn Hickey⁶, Glennis A. Logsdon^{5,17}, Haoyu Cheng¹⁸, Heng Li^{19,20}, Hufsah Ashraf^{2,4}, Jan O. Korbel²¹, Jana Ebler^{2,4}, Jiaqi Li^{22,23}, Jonathan Crabtree²⁴, Katherine M. Munson⁵, Keisuke K. Oshima¹⁷, Kendra Hoekzema⁵, Keon Rabbani⁷, Likhitha Surapaneni²⁵, Lisbeth A. Guethlein²⁶, Marc Jan Bonder^{27,28}, Mark Gerstein^{22,23}, Mark J.P. Chaisson⁷, Mark Loftus^{14,15}, Matthew Jensen^{22,23}, Michael C. Zody²⁹, Mike E. Talkowski^{30,31,32}, Mikko Rautiainen³³, Mir Henglin^{2,4}, Miriam K. Konkel^{14,15}, Nicholas R Pollock³⁴, Olanrewaju Austine-Orimoloye²⁵, Parithi Balachandran⁸, Patrick Hasenfeld²¹, Paul J. Norman^{34,35}, Peter A. Audano⁸, Peter Ebert^{2,36}, Pille Hallast⁸, Qihui Zhu^{8,37}, Ryan E. Mills³⁸, Sarah E. Hunt²⁴, Scott E. Devine²⁴, Sergey Koren¹, Stephan Scholz³, Timofey Prodanov^{2,4}, Tobias Marschall^{2,4}, Tobias Rausch²¹, Vasiliki Tsapalou²¹, Weichen Zhou³⁸, William T. Harvey⁵, Xinghua Shi^{10,11}, Xuefang Zhao^{30,31}, Ying Zhou^{19,20}, Youngjun Kwon⁵, Yunzhe Jiang^{22,23}, Yuwei Song³⁹, Zechen Chong³⁹

1. Genome Informatics Section, Center for Genomics and Data Science Research, National Human Genome Research Institute, National Institutes of Health, Bethesda, MD, USA
2. Center for Digital Medicine, Heinrich Heine University, Düsseldorf, Germany
3. Institute of Medical Microbiology and Hospital Hygiene, Medical Faculty, Heinrich Heine University, Düsseldorf, Germany
4. Institute for Medical Biometry and Bioinformatics, Medical Faculty and University Hospital Düsseldorf, Heinrich Heine University, Düsseldorf, Germany
5. Department of Genome Sciences, University of Washington School of Medicine, Seattle, WA, USA
6. UC Santa Cruz Genomics Institute, University of California, Santa Cruz, CA, USA
7. Department of Quantitative and Computational Biology, University of Southern California, Los Angeles, CA, USA
8. The Jackson Laboratory for Genomic Medicine, Farmington, CT, USA
9. Foundation of Biological Data Sciences, Belmont, CA, USA
10. Temple University, Department of Computer and Information Sciences, College of Science and Technology, Philadelphia, PA, USA
11. Temple University, Institute for Genomics and Evolutionary Medicine, Philadelphia, PA, USA
12. The University of Connecticut Health Center, Farmington, CT, USA
13. Howard Hughes Medical Institute, University of Washington, Seattle, WA, USA
14. Clemson University, Department of Genetics & Biochemistry, Clemson, SC, USA
15. Center for Human Genetics, Clemson University, Greenwood, SC, USA
16. Medical University of South Carolina, College of Graduate Studies, Charleston, SC, USA

17. Perelman School of Medicine, University of Pennsylvania, Department of Genetics, Epigenetics Institute, Philadelphia, PA, USA
18. Department of Biomedical Informatics and Data Science, Yale School of Medicine, New Haven, CT, USA
19. Department of Data Science, Dana-Farber Cancer Institute, Boston, MA, USA
20. Department of Biomedical Informatics, Harvard Medical School, Boston, MA, USA
21. European Molecular Biology Laboratory (EMBL), Genome Biology Unit, Heidelberg, Germany
22. Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, CT, USA
23. Program in Computational Biology and Bioinformatics, Yale University, New Haven, CT, USA
24. Institute for Genome Sciences, University of Maryland School of Medicine, Baltimore, MD, USA
25. European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge, United Kingdom
26. Department of Structural Biology, School of Medicine, Stanford University, Stanford, CA, USA
27. Department of Genetics, University of Groningen, University Medical Center Groningen, Groningen, the Netherlands; Oncode Institute, Utrecht, The Netherlands
28. Division of Computational Genomics and Systems Genetics, German Cancer Research Center, Heidelberg, Germany
29. New York Genome Center, New York, NY, USA
30. Program in Medical and Population Genetics, Broad Institute of MIT and Harvard, Cambridge, MA, USA
31. Center for Genomic Medicine, Massachusetts General Hospital, Boston, MA, USA
32. Department of Neurology, Harvard Medical School, Boston, MA, USA
33. Institute for Molecular Medicine Finland (FIMM), University of Helsinki, Helsinki, Finland
34. Department of Biomedical Informatics, University of Colorado School of Medicine, Aurora, CO, USA
35. Department of Immunology and Microbiology, University of Colorado School of Medicine, Aurora, CO, USA
36. Core Unit Bioinformatics, Medical Faculty and University Hospital Düsseldorf, Heinrich Heine University, Düsseldorf, Germany
37. Stanford Health Care, Palo Alto, CA, USA
38. Department of Computational Medicine & Bioinformatics, University of Michigan, MI, USA
39. Department of Biomedical Informatics and Data Science, Heersink School of Medicine, University of Alabama, Birmingham, AL, USA