# Table of Contents

**Figure S1. Precision, recall, and F1 scores for SNV and indel calls for ONT and Illumina small variant calls using PacBio assemblies as the truth set.**

ONT data (PMDV) or Illumina data (GATK) compared to GIAB or HPRC calls for 5 high-confidence samples (HG002/NA24385, HG003/NA24149, HG004/NA24143, HG00733, and HG02723) genome-wide in GIAB high-confidence regions for HG002 (GIAB.HG002.mask.incl.HP) and when excluding homopolymer regions in the GIAB high-confidence HG002 regions (GIAB.HG002.mask.excl.HP). Homopolymer regions were defined as any sequence of four identical nucleotides or more, plus one bp flanking each side of the sequence.

**Figure S2. Precision, recall, and F1 scores for SNV and indel calls for ONT small variant calls from both the Napu (PMDV) and alignment (Clair3) pipelines using Illumina data as the truth set.**
Variants for each sample called by PMDV and Clair3 are benchmarked to Illumina GATK variant calls without masking (A,B), in GIAB high-confidence regions only (C,D), and in GIAB high-confidence regions excluding homopolymers 4 nucleotides or greater +/− one nucleotide on either side (E,F).



|  | Precision | Recall | F1 | Precision | Recall | F1 |
|---|---|---|---|---|---|---|
| Mean | 0.969 | 0.993 | 0.981 | 0.970 | 0.995 | 0.982 |
| SD | 0.002 | 0.001 | 0.001 | 0.002 | 0.001 | 0.001 |
| Mean | 0.817 | 0.690 | 0.748 | 0.857 | 0.899 | 0.878 |
| SD | 0.007 | 0.024 | 0.017 | 0.004 | 0.010 | 0.006 |

**Figure S3. Percent of genome assembled by Flye and Shasta-Hapdup.**
Violin plot showing the percentage of genome assembled for each sample by Flye and Shasta-Hapdup. The *x*-axis represents the assembler, and 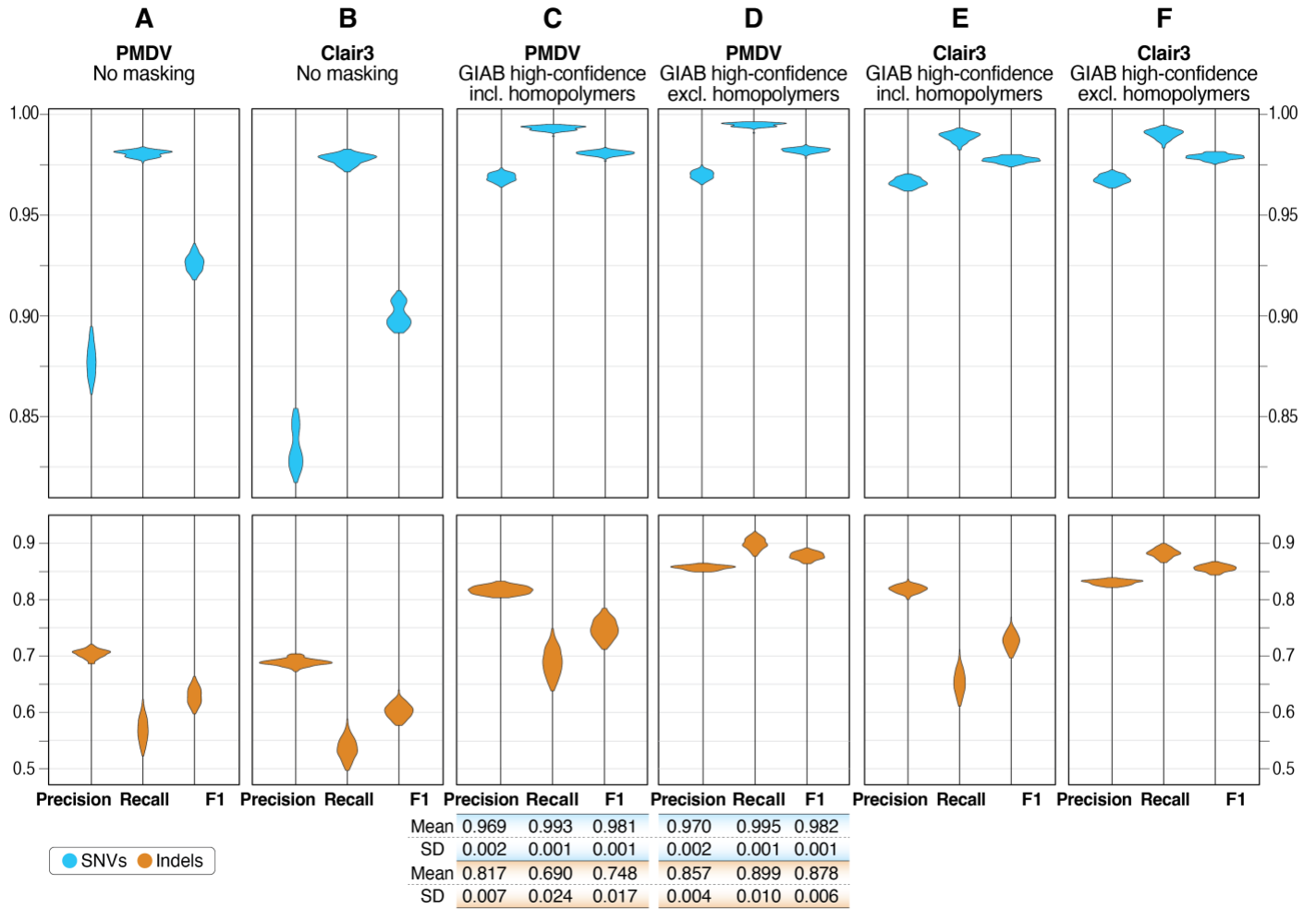the *y*-axis represents the percentage of GRCh38 reference genome covered by contigs for each assembly. The median values for Flye and Shasta-Hapdup assemblies for all 100 samples are 93.6% and 93.5%, respectively.
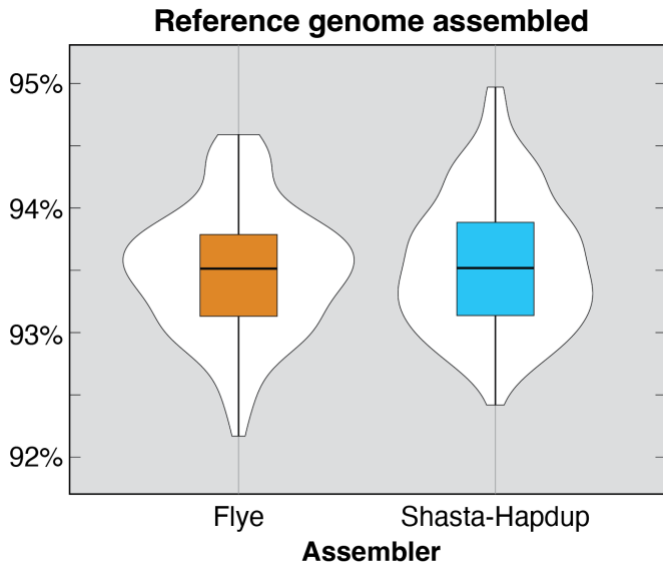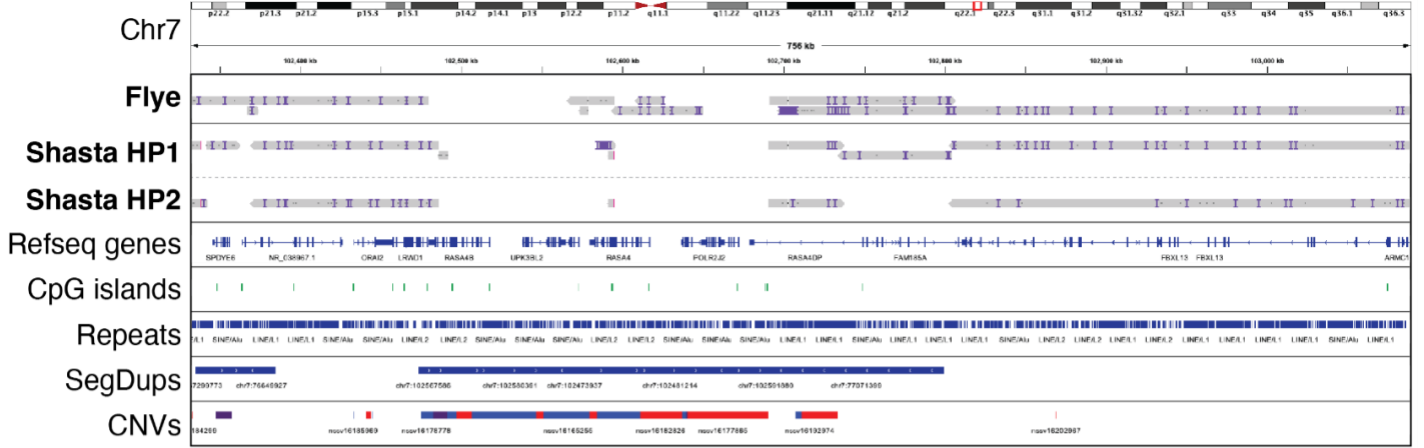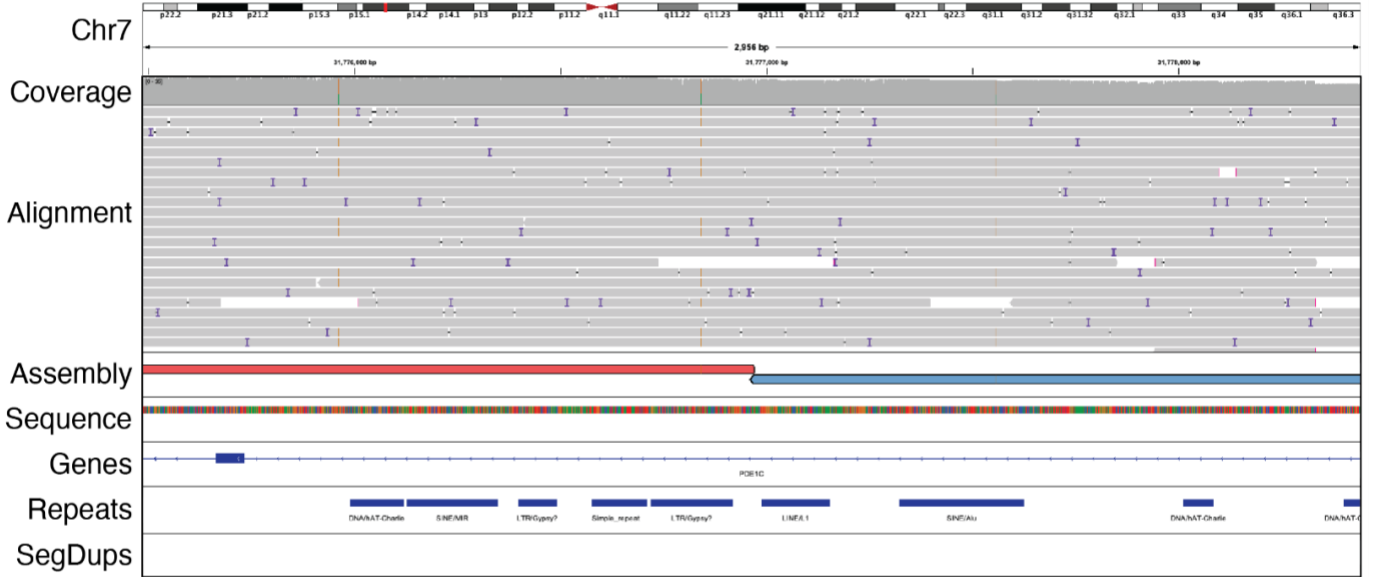
**Figure S4. IGV screenshots of Chromosome 7 assemblies** (*see next page*).
Among all contig breaks on Chromosome 7, 1.9% of those in Flye assemblies were breaks in nonrepetitive sequence that are not within 10 kbp of a segdup, similar to the genome-wide average of 2.9%. **A**. Example of a location where the presence of a segmental duplication is associated with breaks in assemblies in both Flye and Shasta-Hapdup for HG02373. **B**. Example of a region with a discontinuity in the Flye assembly contigs in GM19035. This break occurs in a region outside of a segmental duplication or a satellite repeat. The alignment and coverage tracks at the top suggest no aberrations or increased coverage that could indicate increased likelihood of an assembly break. IGV view is approximately 2,900 bp. **C**. Example of a position where Shasta-Hapdup assembled contigs break in both haplotype-resolved assemblies in HG01801. This break occurs in a region outside of a segmental duplication or satellite repeat. The alignment and coverage tracks at the top suggest no aberrations or increased coverage that could indicate increased likelihood of assembly breaks. IGV view is approximately 23 kbp.

**A** **HG02373** | Flye and Shasta-Hapdup

**B** **GM19035** | Flye
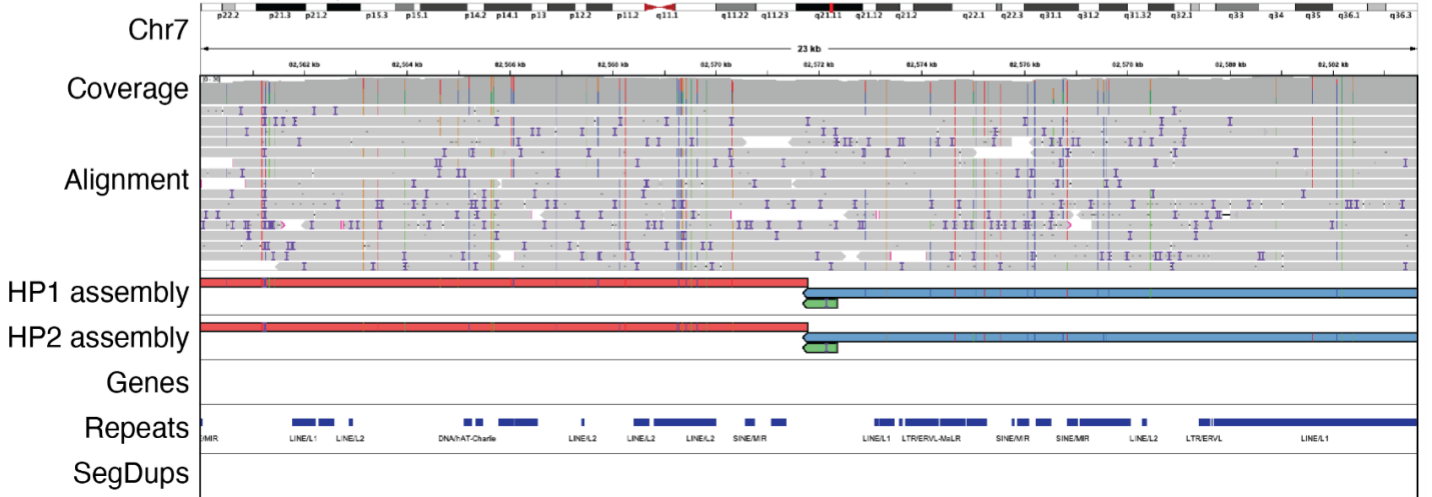
**C** **HG1801** | Shasta-Hapdup

**Figure S5. Pangenome principal component analysis** (*see next page*).
**A**. Graph-based principal component analysis of Chromosome 20 assemblies – including 100 Shasta-Hapdup samples, 44 HPRC samples, and the GRCh38, CHM13 (v2.0), and HG002 (v1.0.1) human reference genomes – revealed clear separation along the first principal component (PC1). The two clusters correspond to the HPRCy1 assemblies (right) and 100 Shasta-Hapdup assemblies (left). Because our assemblies cluster close to GRCh38, whose centromeric sequences were masked, we interpreted PC1 as representing centromeric diversity of the assemblies. Experience with ONT based assemblies suggests that the ONT-only assemblies are currently less able to assemble highly repetitive sequences compared to those which include HiFi data. **B**. Graph-based principal component analysis of euchromatic, non-centromeric fraction of the Chromosome 20 assemblies—including 100 Shasta-Hapdup samples, 44 HPRC samples, and the GRCh38, CHM13 (v2.0), and HG002 (v1.0.1) human reference genomes—revealed no distinct separation between the Shasta-Hapdup and HPRC samples. The 4 reference genomes cluster because their entire sequence is considered, including the heterochromatic regions. Data points were visually distinguished by color and shape according to their group. In general, our assemblies are slightly more spread in both PC1 and PC2, which possibly relates to higher base-level error rates in the ONT-only assemblies.
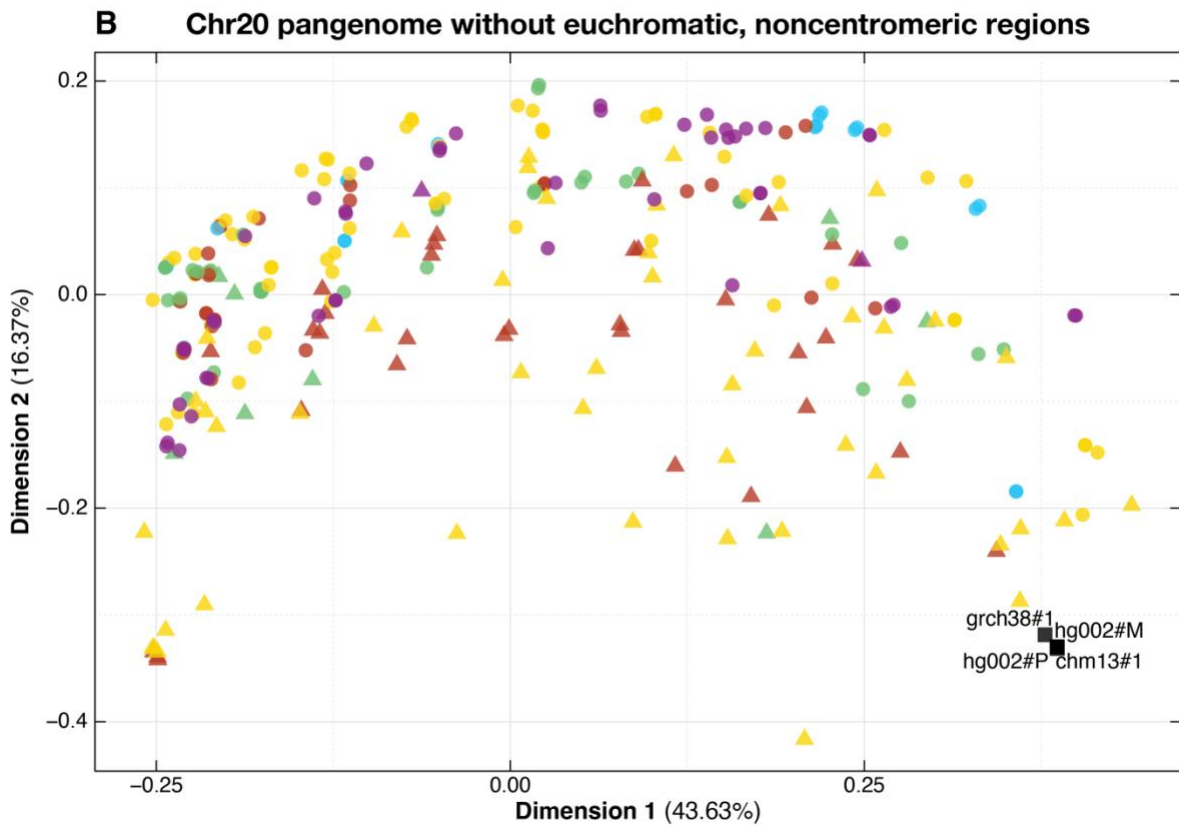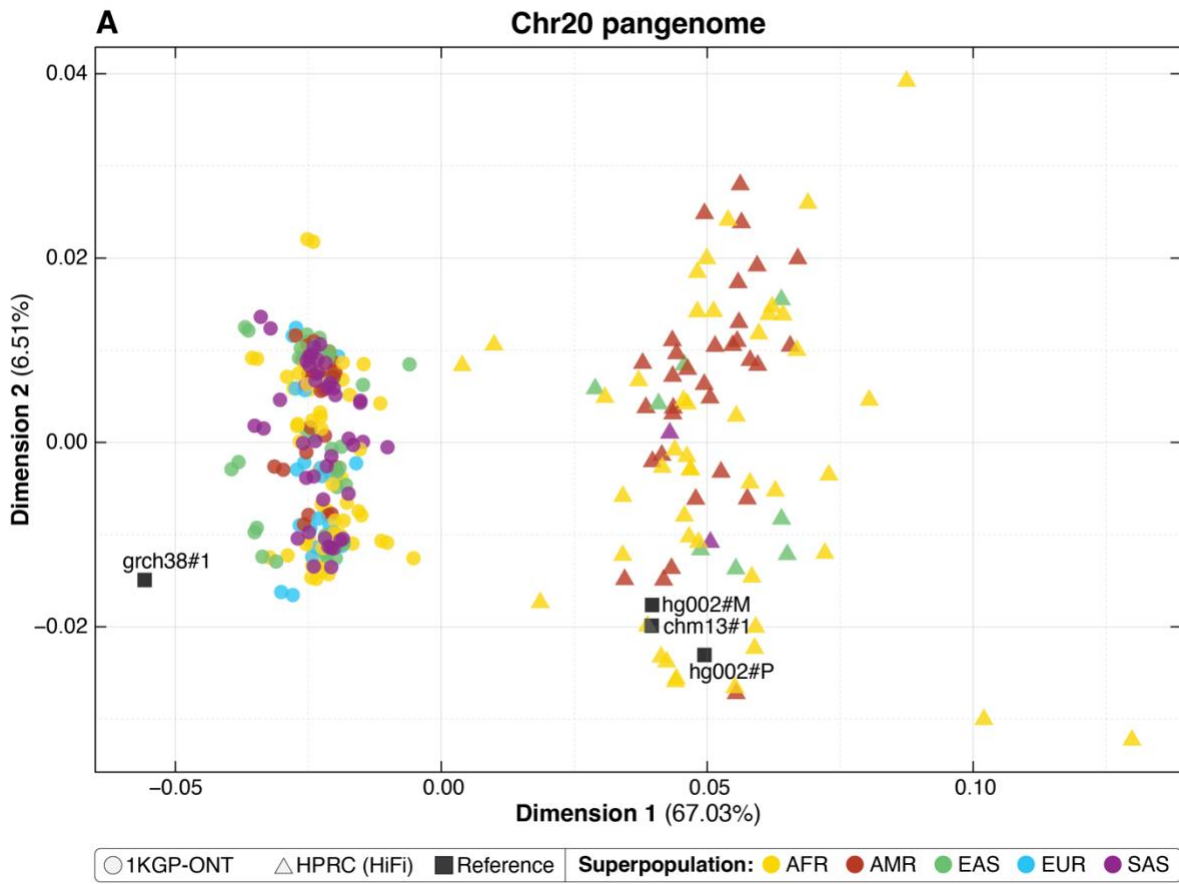
**A**

## Chr20 pangenome



○ 1KGP-ONT  △ HPRC (HiFi)  ■ Reference  | **Superpopulation:** ● AFR  ● AMR  ● EAS  ● EUR  ● SAS

**B**  ## Chr20 pangenome without euchromatic, noncentromeric regions

**Figure S6. Evaluation of L1HS elements in the first 100 assemblies.**
Comparison of the number of L1HS in the major populations in the 100 samples (standard boxplots) against the number identified in phased haplotypes from HG002 and HG005 from GIAB and the haploid CHM13 T2T assembly.
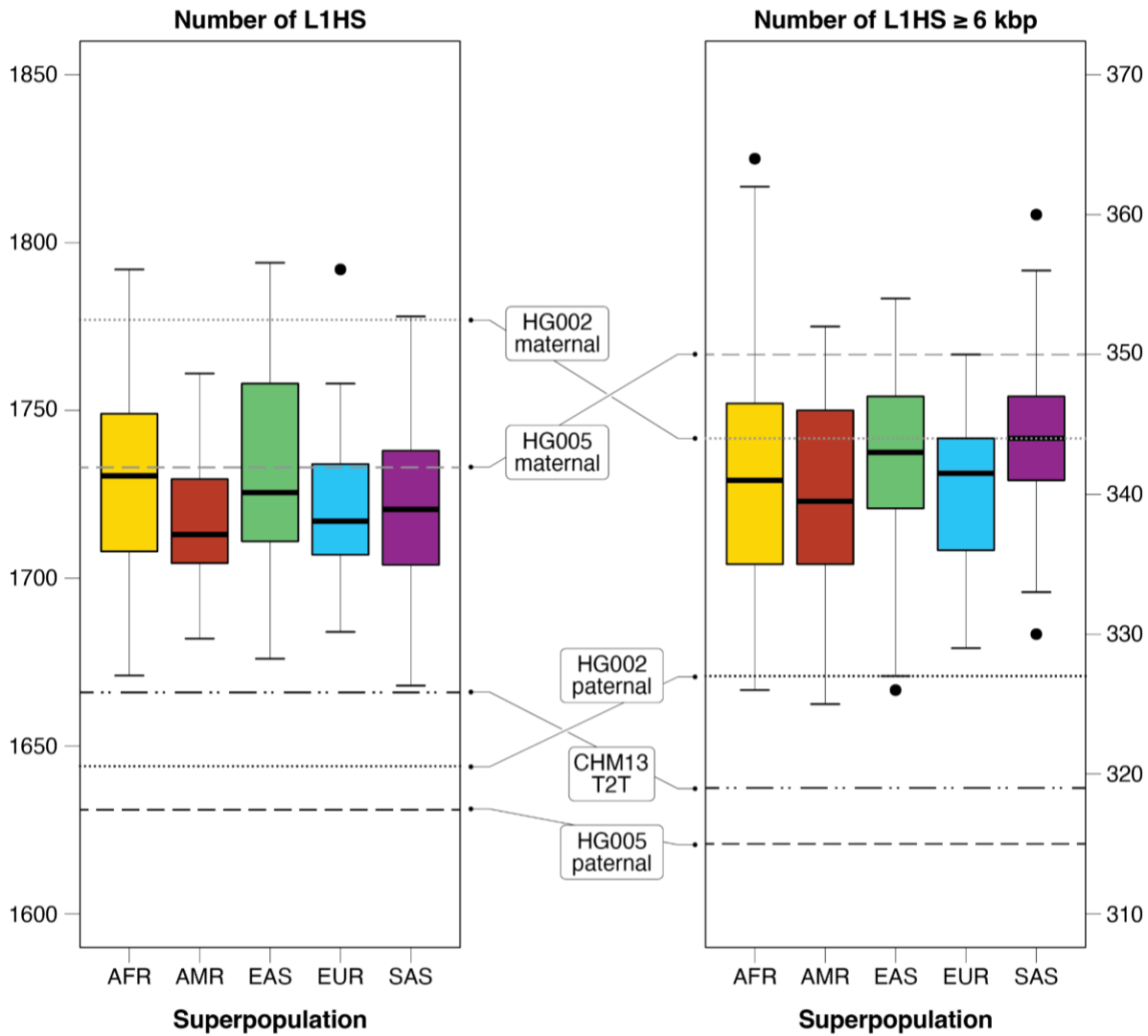
**Figure S7. Evaluation of per-sample insertions and deletions by type, chromosome, and SV caller.**
The number of insertion and deletion SVs was evaluated by chromosome for each of the five SV callers across all 100 samples. Analyses were conducted using the GRCh38 reference genome.
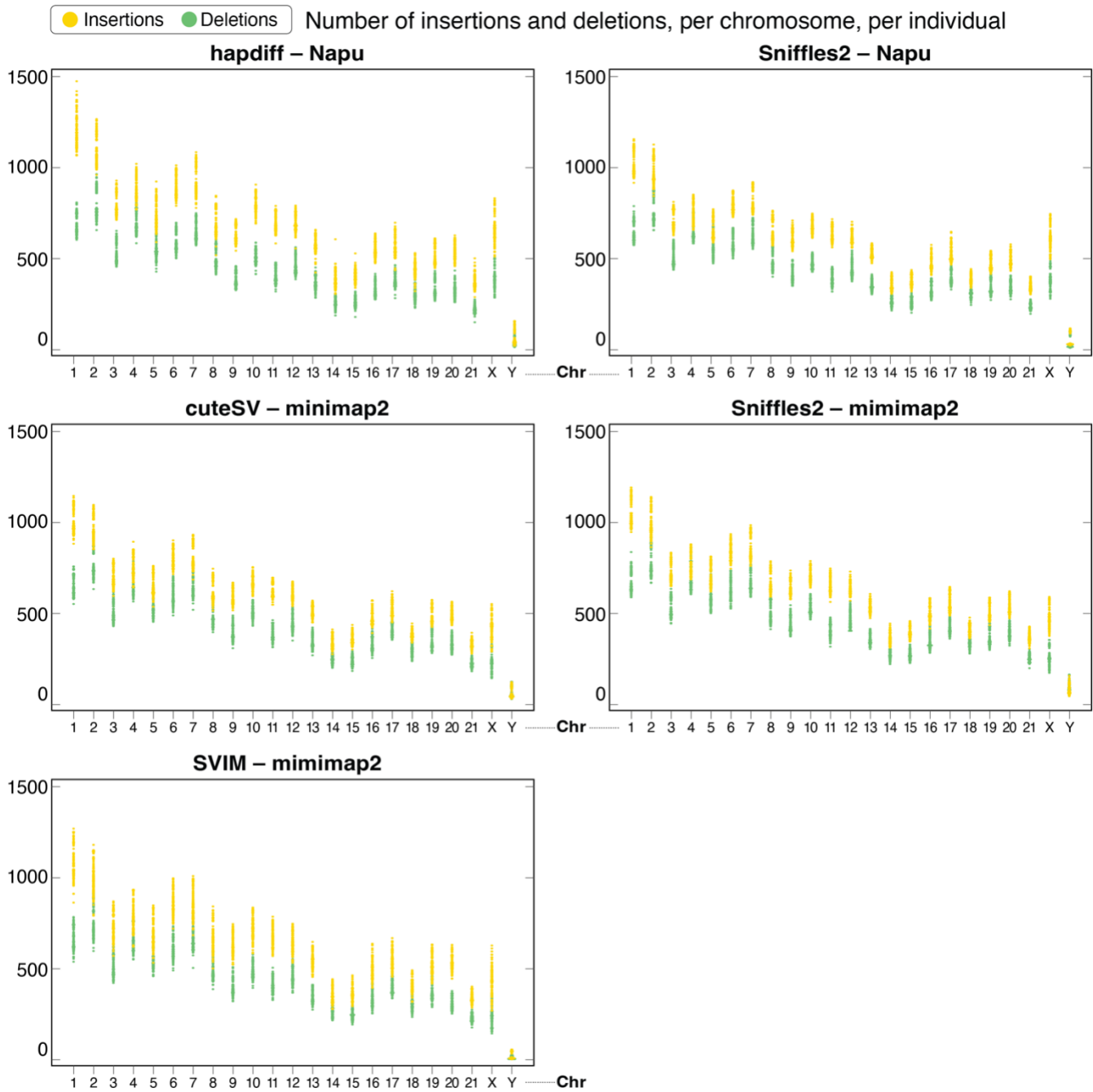
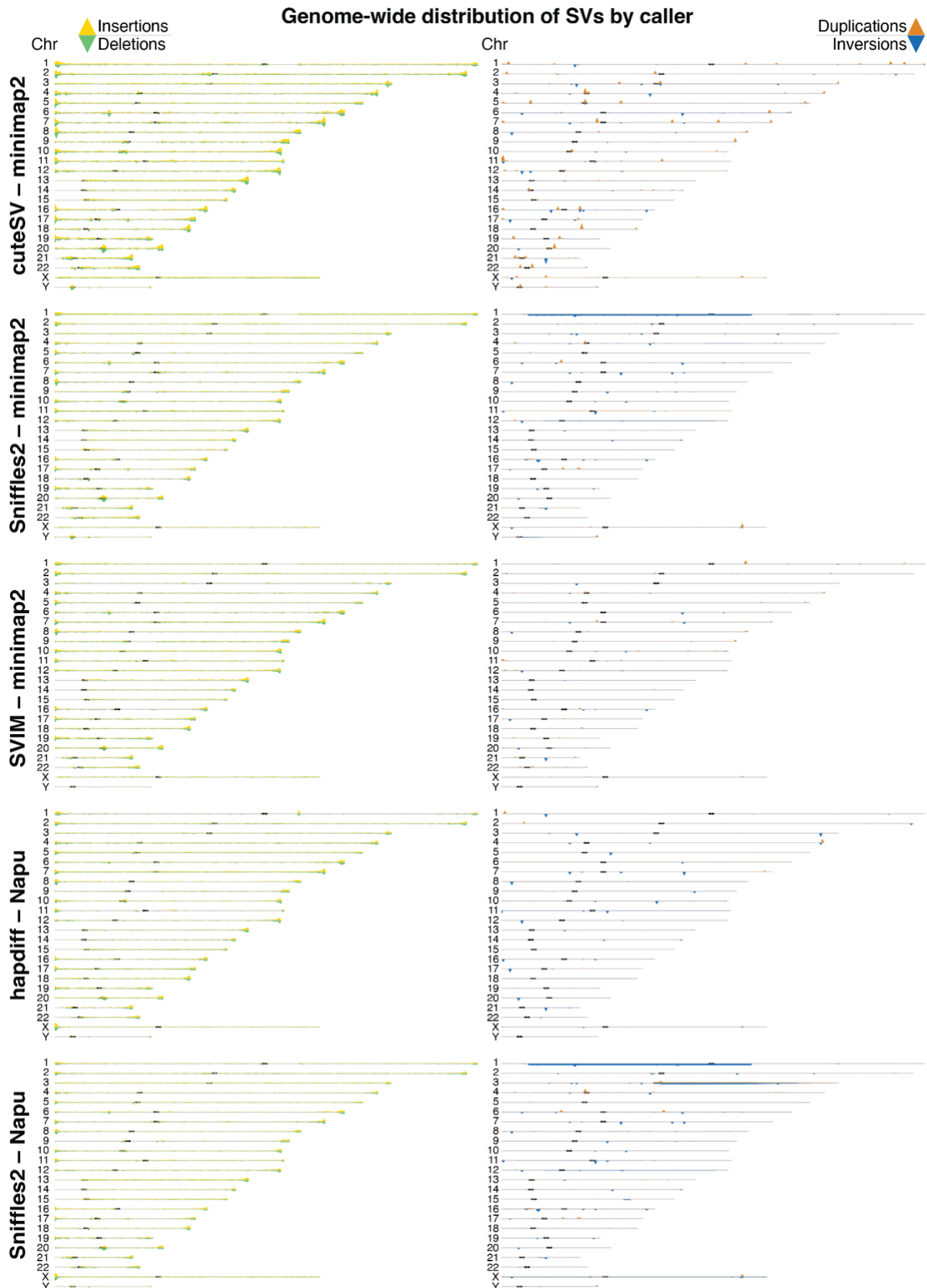**Figure S8. Distribution of SV events genome-wide, all populations.**

**Figure S9. Concordance of SVs by caller and evaluation of joint SV calls.**
**A,C**. Concordance of SV calls for all callers (A) and for only hapdiff (C) for insertion and deletion events. **B**.
Number of SVs per individual called by at least 1, 2, 3 ,4 and 5 (all) callers. Based on Jasmine merging of all 5
SV callers per sample (cuteSV–minimap2, Sniffles2–minimap2, SVIM–minimap2, hapdiff–Napu, Sniffles2–
minimap2). If an SV was called by at least one caller (in other words, called at all by any caller), it has a
minimum support of 1. If an SV was called in an individual in at least 2 callers, it has a minimum support of 2,
etc. 5 means that all 5 callers agree on the SV. We see the same expected distribution of AFR samples having
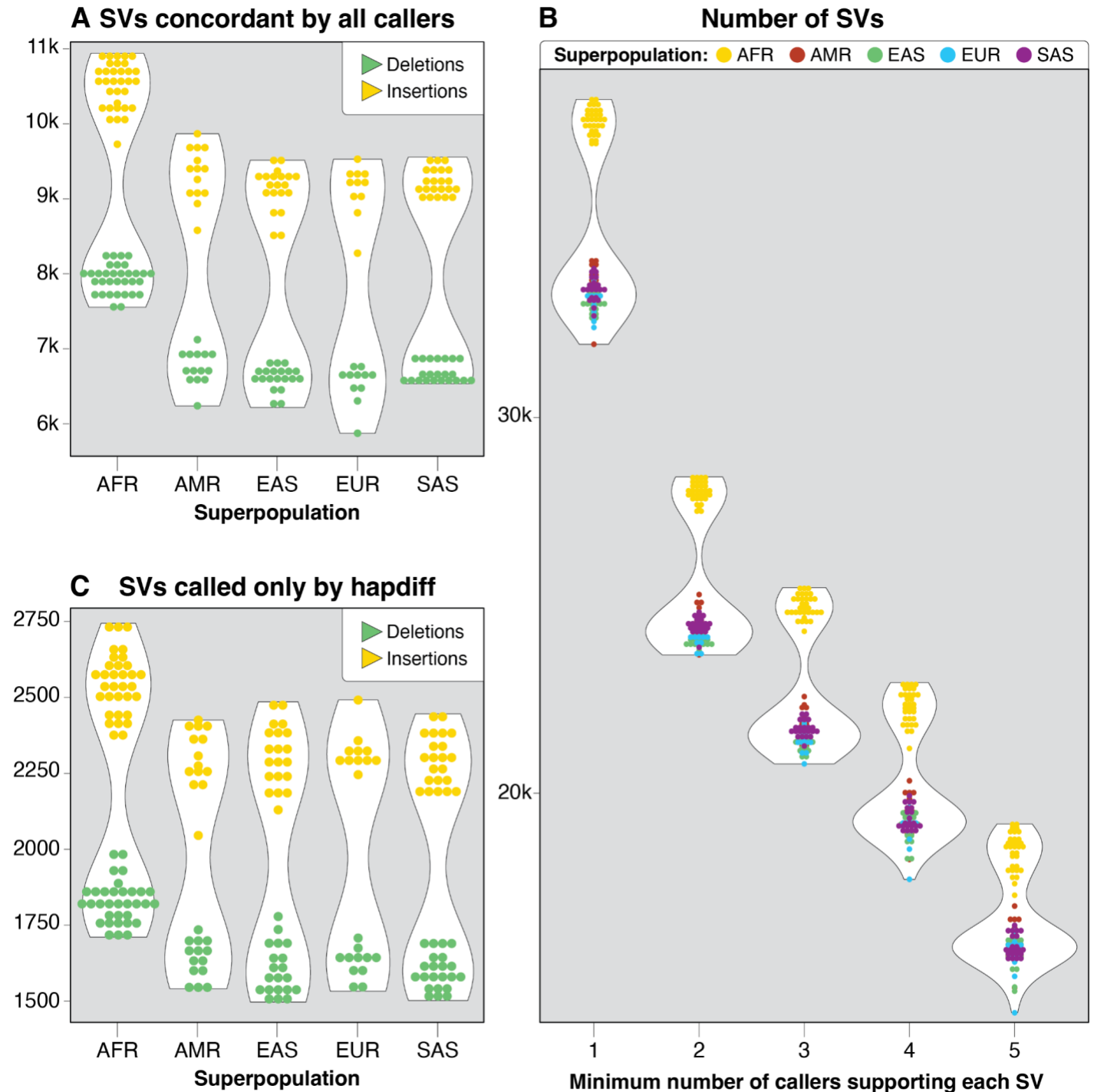more SVs in each support threshold.

**Figure S10. Example of false positive SVs called only by hapdiff.**

We evaluated SVs called only by hapdiff to determine whether they represented false positive calls. Among the 30 SVs visually evaluated, 28 were likely false positive, and these false-positive SV calls were either located within centromeric (6/30), telomeric (8/30), or other low-complexity regions (3/30), or they had allelic representations different from those observed with alignment-based methods (10/30) **A**. IGV screenshot shows a single deletion in an assembly represented as two deletions in the alignment. **B**. IGV screenshot shows a call within a centromeric region found in only one assembly contig.
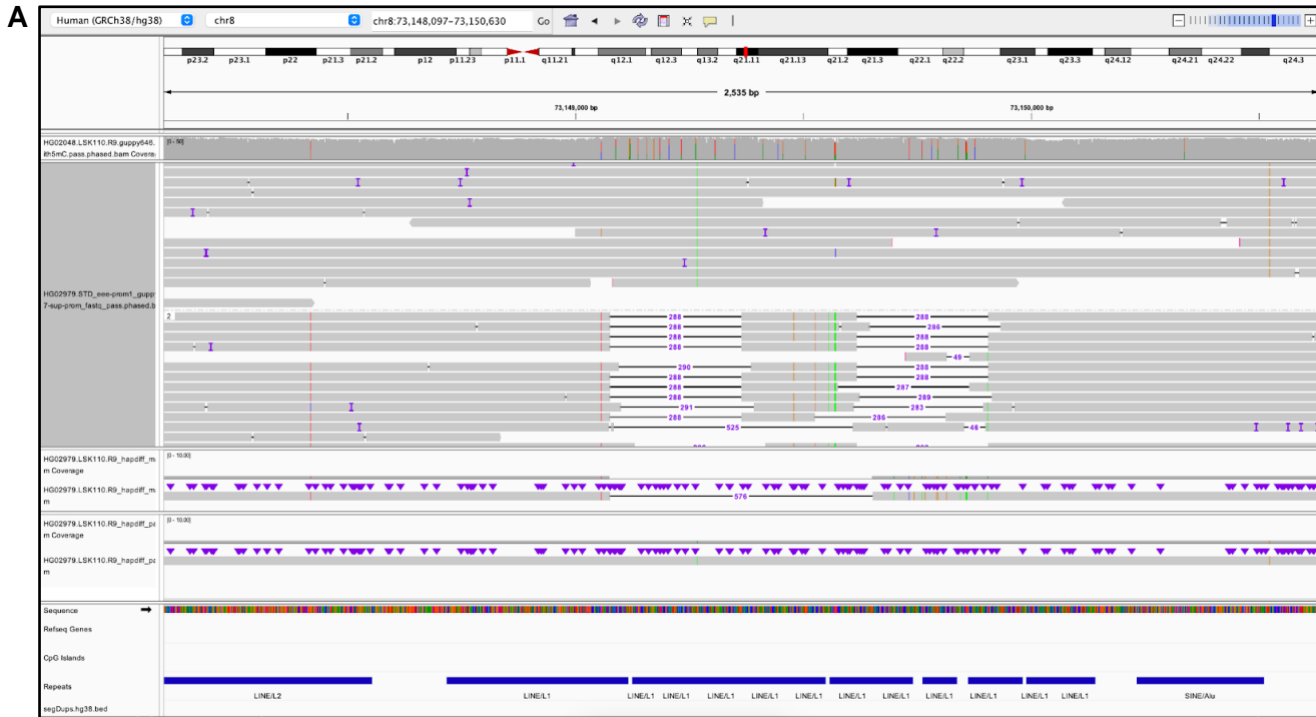
**Figure S11. Evaluating the functional impact of structural variants.**
**A**. Variant counts in the 65 samples common to both the 1000G-ONT and MAGE datasets, with heterozygous variants in orange (genotype 1) and homozygous variants in blue (genotype 2). **B**. Relationship between the best-fit slope (β) derived from OLS regression and gene-level $q$-values. "All" (blue) refers to the use of identified SVs in selected samples for the eQTL analysis. "Novel" (red) indicates the eQTL analysis conducted with SVs that were not reported in the previous long-read SV dataset of 31 samples by Kirsche *et al.* **C**. Integrative Genomics Viewer (IGV) screenshot of a 484-bp insertion upstream of the 5' TSS of the *ZNF79* gene. **D**. Distribution of genotypes and gene expression values in the 65 samples for the *ZNF79*-associated deletion, along with the q-value (0.014). **E**. IGV screenshot of an 81-bp deletion in the *NAPRT* gene. **F**. Distribution of genotypes and gene expression values in the 65 samples for the *NAPRT*-associated deletion, along with the $q$-value (0.00023). **G**. Proportion of genotypes that agree between short- and long-read SV calls among shared variants in the MAGE samples ($n$=29,147 variants, $n$=65 samples). Call consistency is separated by the genotype of the long-read call. **H**. Distribution of genotypes and gene expression values in the 731 short-read MAGE samples for the *GBP3*-associated deletion, along with the $q$-value (1.3e-27).
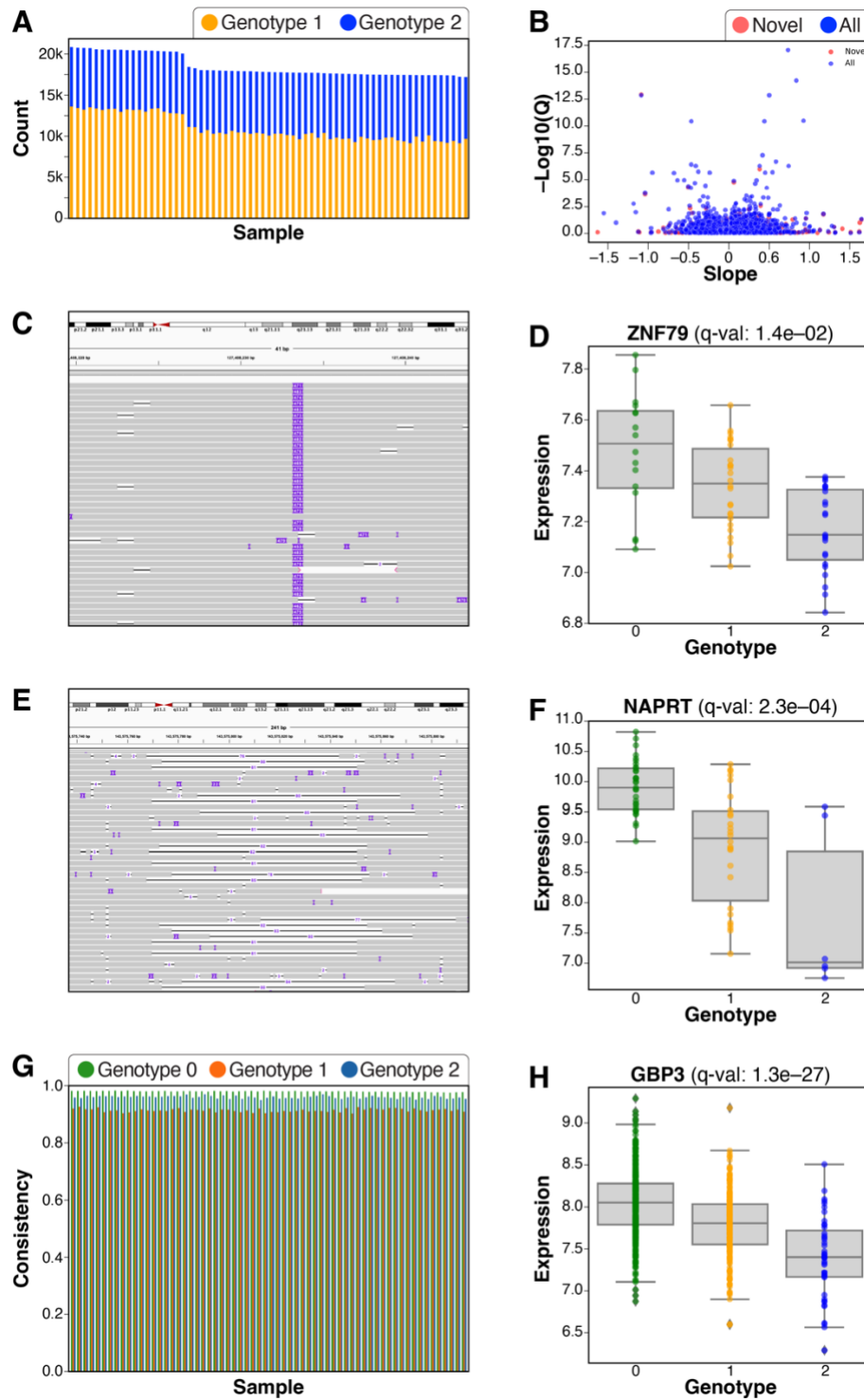
**Figure S12. Evaluation of SVs in exons of medically relevant genes.**
**A.** Number of samples harboring each INS/DEL intersecting an OMIM exon. **B.** Ideogram of where each SV that intersects with an OMIM exon is located in the GRCh38 genome.
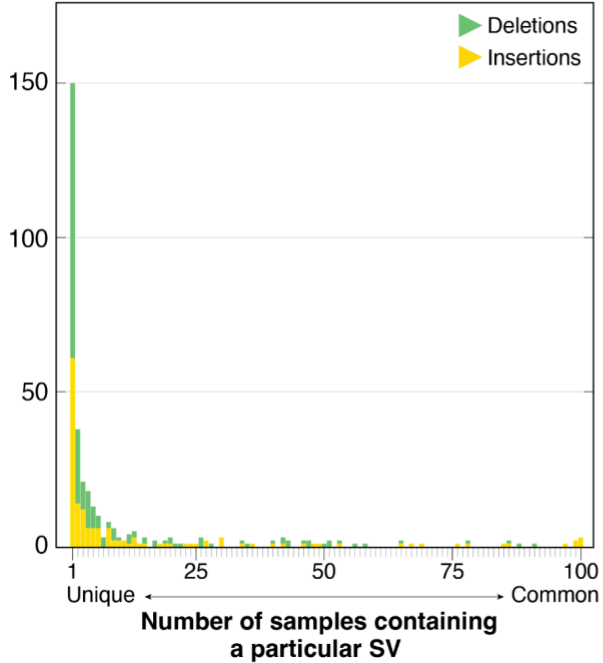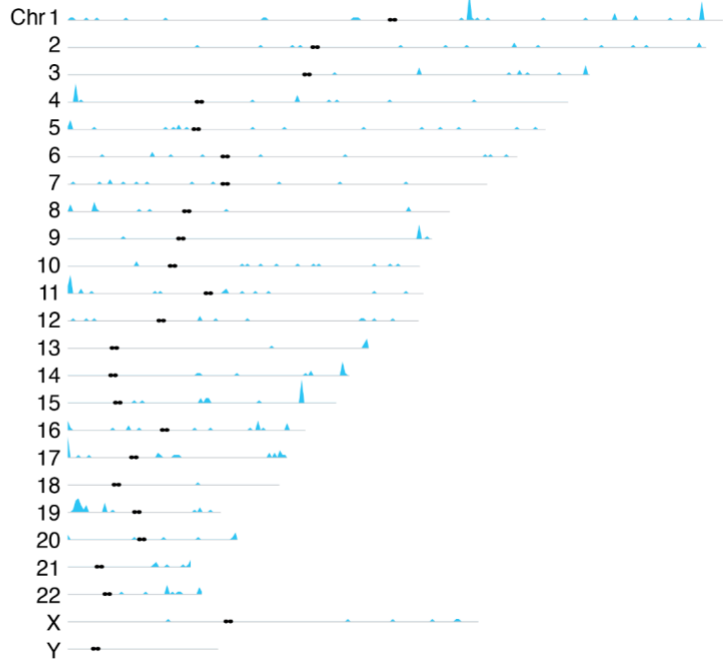
**Figure S13. IGV view of previously described deletions that segregate in the population.**
**A.** Phased IGV view of a 22,791-bp deletion in GM19035 that includes all or part of HBB, HBD, HBBP1, BGLT3, and HBG1. Variants in this region are associated with beta thalassemia (MIM: 613985), with this specific deletion known as Hemoglobin Kenya, with this individual likely being an asymptomatic carrier (Huisman *et al*. 1972). GM19035 is an individual from the Luhya population within the African superpopulation.
**B**. A 19,304-bp deletion (NC_000016.10:g.165401_184701del) that includes *HBA1* and *HBA2* in HG01812 (Dia Chinese/EAS) and HG00728 (Southern Han Chinese/EAS). This deletion is commonly referred to as the Southeast Asian deletion (--SEA) and is associated with alpha-thalassemia (MIM: 604131).
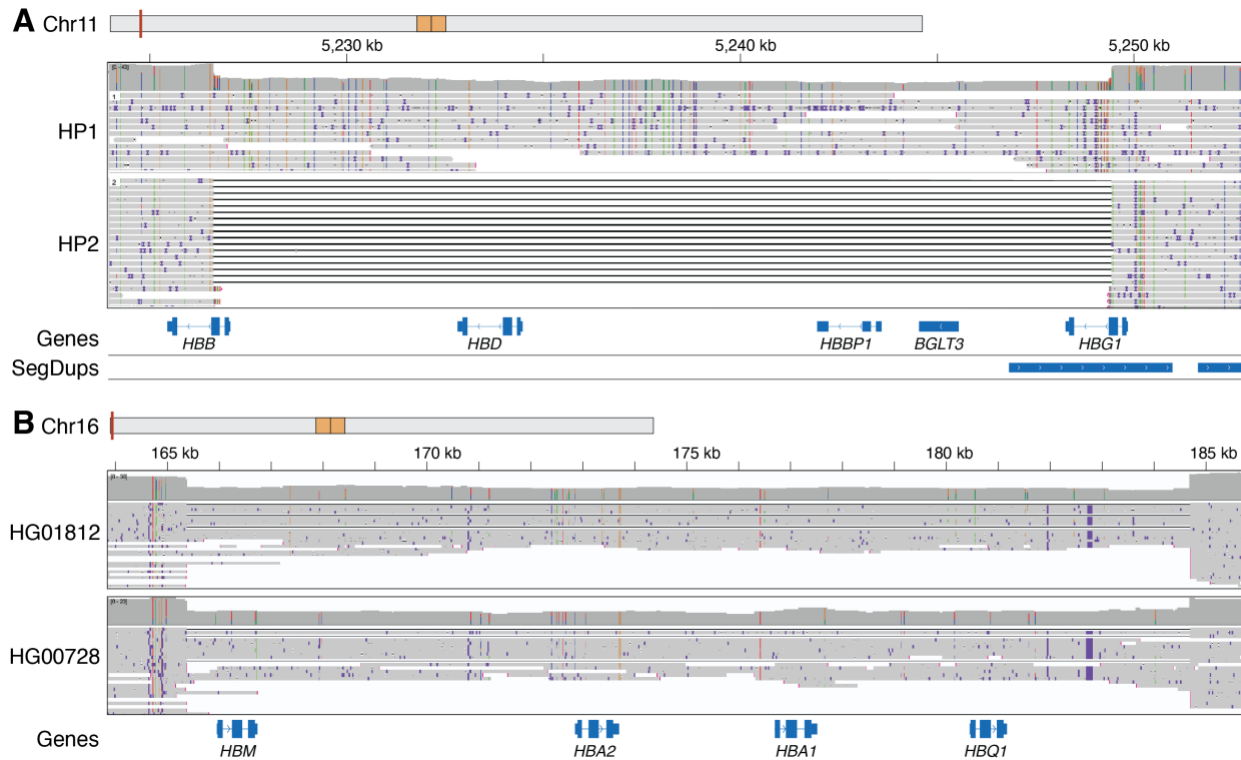
**Figure S14. Evaluation of X-linked SVs in 46,XY individuals.**
**A.** Screenshot of alignments and assemblies for an approximately 141-bp insertion in *RPGR* in GM18865. The top two tracks depict Coverage and Alignments. Track 3 and 4 show assembly-based calls derived from Shasta-Hapdup assemblies with 146-bp insertion on one haplotype and 141-bp insertion on the other. The bottom track depicts Flye assembly with a 141-bp insertion. **B.** Screenshot of short-read data from GM18865 and LRS data showing that the 141-bp insertion is not apparent in the short-read data.

**A**



**B**

**Figure S15. Phased IGV screenshots depicting long-read alignments of *CYP2D6* and *CYP2B6* star alleles.**

**A.** Phased IGV view from LRS data showing a *CYP2D6* full gene deletion on one haplotype (HP1) and a hybrid tandem arrangement (*\*36+\*10*) represented by an insertion on the second haplotype (HP2) in HG02396, compared to short-read whole genome sequence data from the same sample in which the complex nature of this event cannot be resolved. **B**. Long-read alignment across *CYP2B6* and *CYP2B7* for GM18871 from an individual from the Yoruba in Ibadan (Nigeria) who has the *CYP2B6\*18/\*29* diplotype. Short-read alignments are included for comparison. *CYP2B6\*29* is a hybrid deletion allele with a *CYP2B7*-derived portion (across exons 1–4) and a *CYP2B6*-derived portion (covering exons 5-9). Read alignment results in reads aligned to either the gene (*CYP2B6*) or the pseudogene (*CYP2B7*), resulting in portions with numerous spurious SNVs.
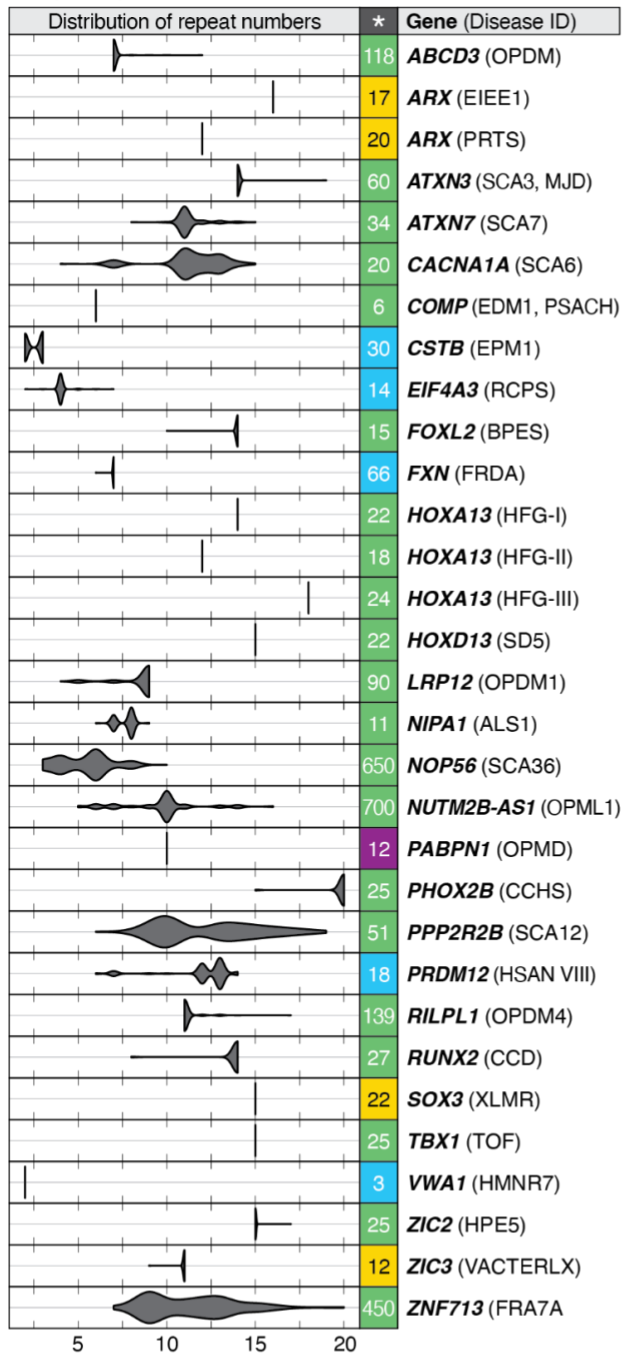
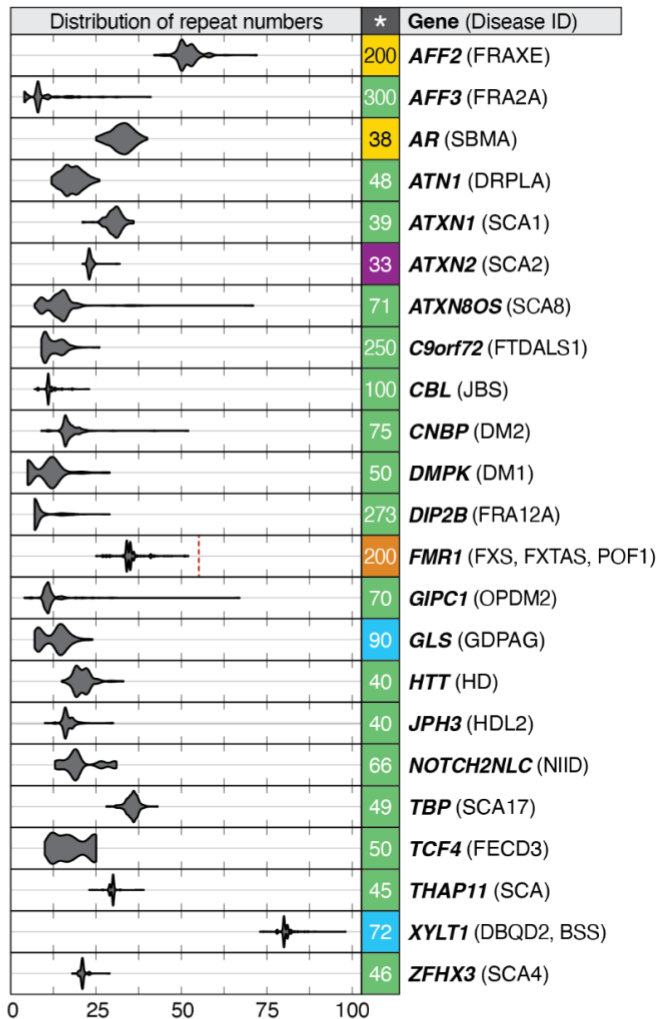**Figure S16. Repeat expansion plots for 66 disease-associated loci.**
Violin plots for 66 disease-associated repeat expansion loci listed at STRchive:
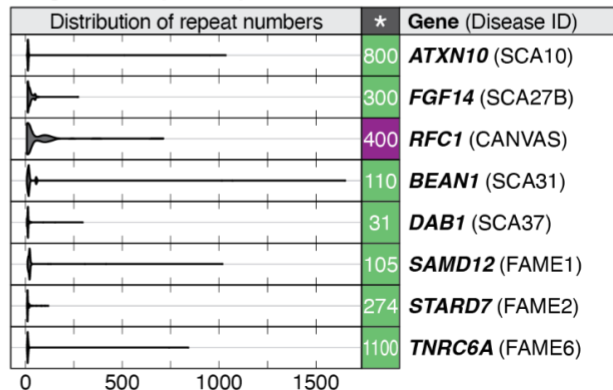https://github.com/hdashnow/STRchive. Motif counts were genotyped with vamos and plotted with R.
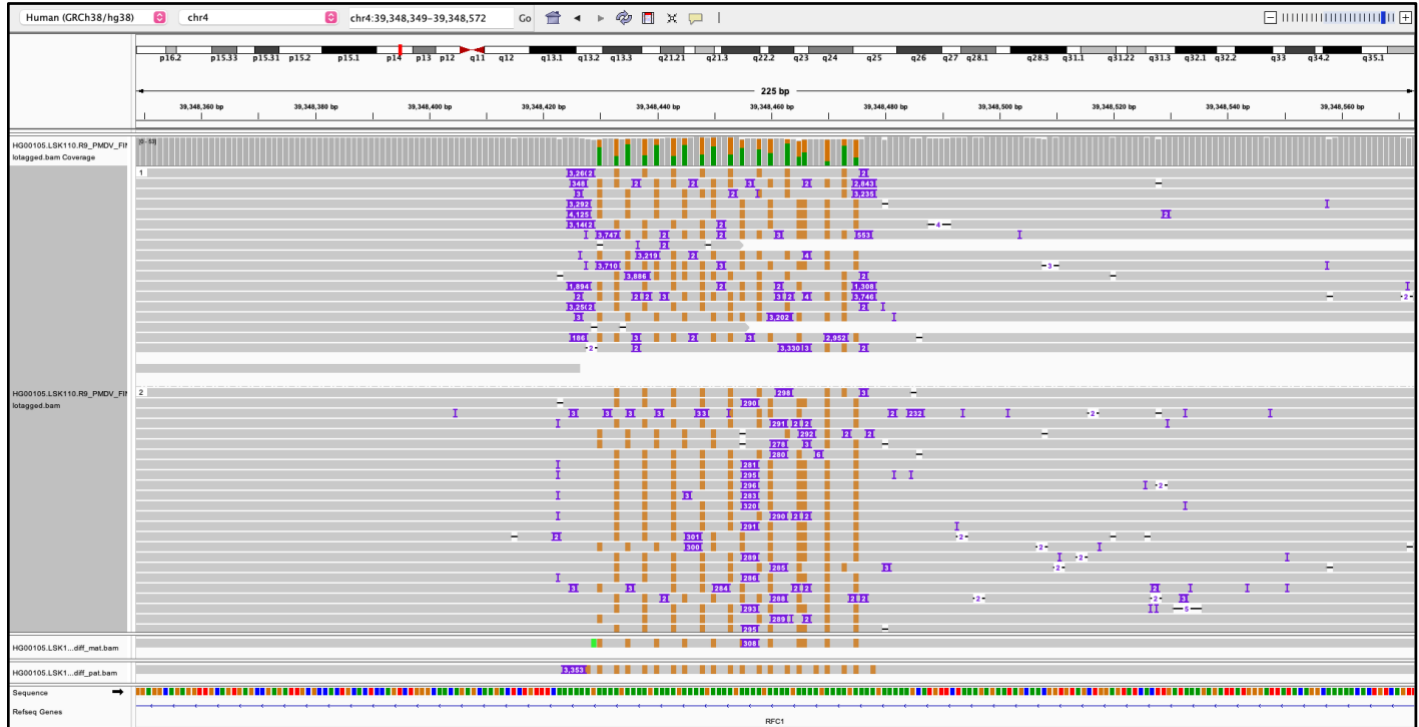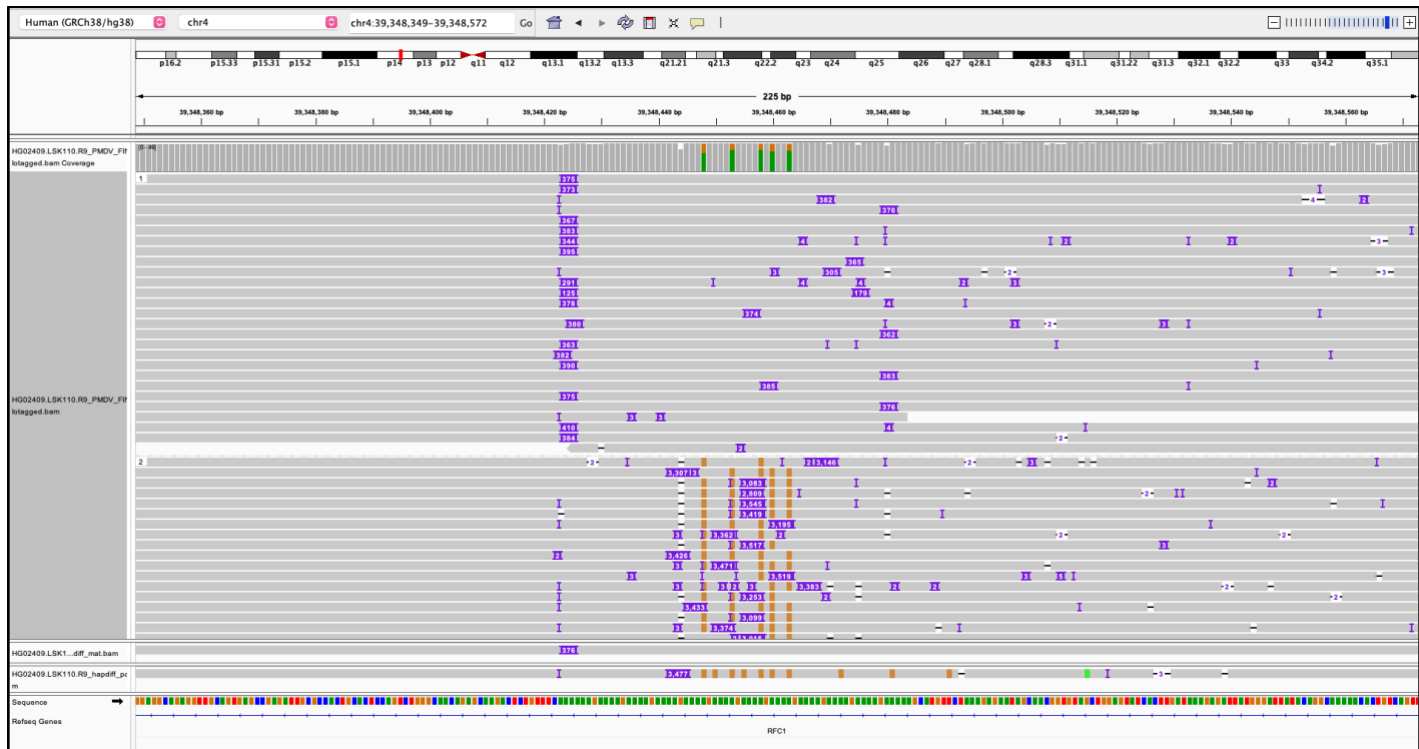
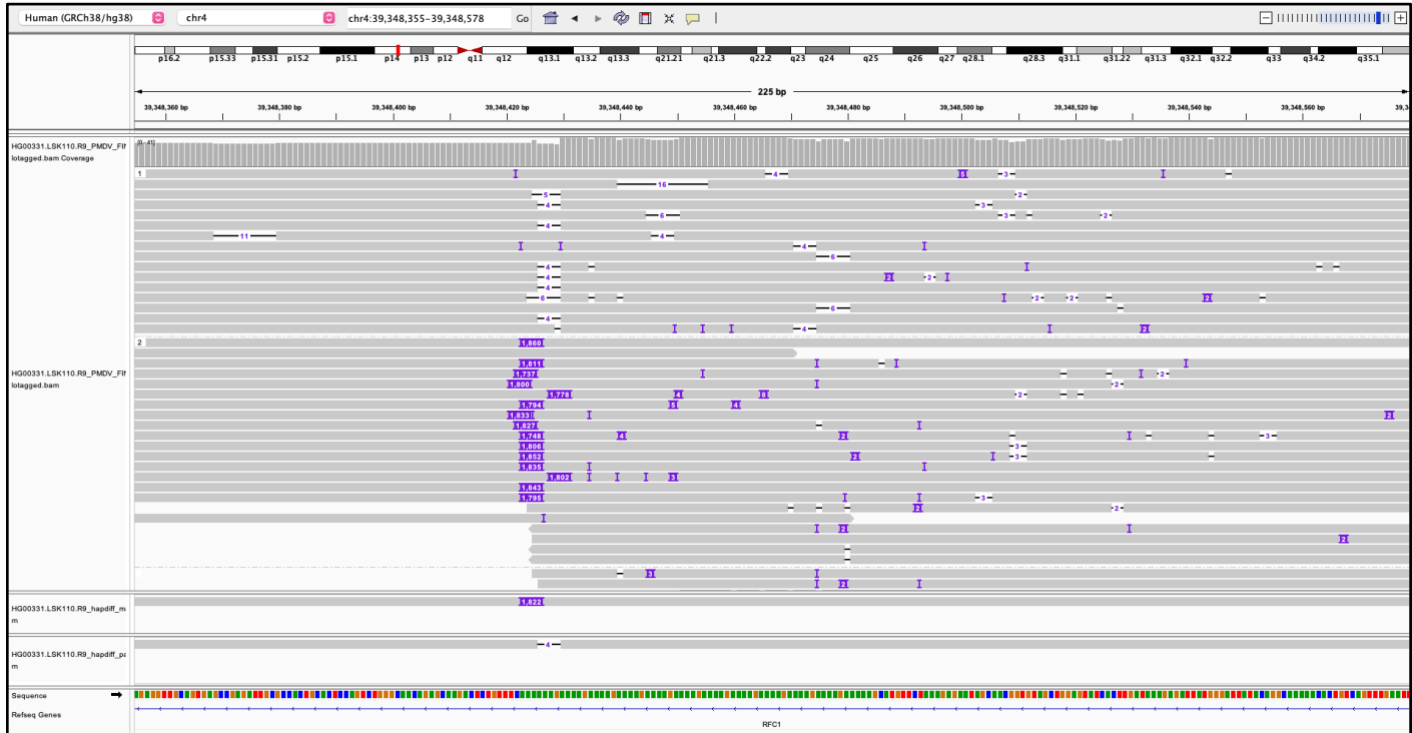**Figure S17. IGV screenshots of repeat expansions in *RFC1* observed in 5 samples.**

IGV screenshots of reads for the five samples with the largest alleles at the CANVAS-associated locus at *RFC1*. The PMDV haplotagged alignment and the hapdiff assemblies from the Napu pipeline for each sample are shown. The hapdiff assembly files from the Napu pipeline were used for genotyping with vamos.

**HG00105**



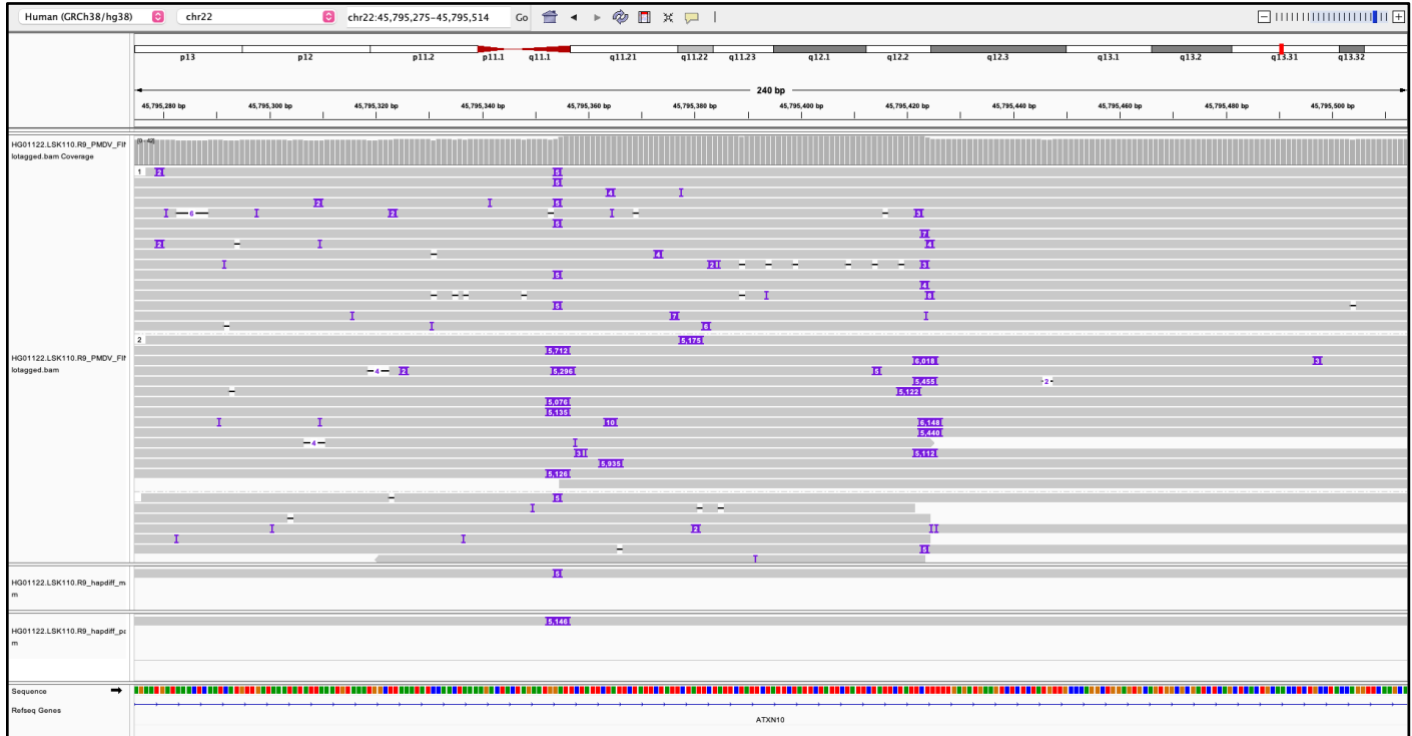**HG00209**

**HG01122**



**HG00331**

## HG01862

**Figure S18. IGV screenshots of *ATXN10* repeat expansions.**

IGV screenshots for the three samples at the SCA31-associated locus at *ATXN10* with over 250 motifs. The PMDV haplotagged alignment and the hapdiff assemblies from the Napu pipeline for each sample are shown. The hapdiff assembly files were used for genotyping with vamos.

**HG01122**
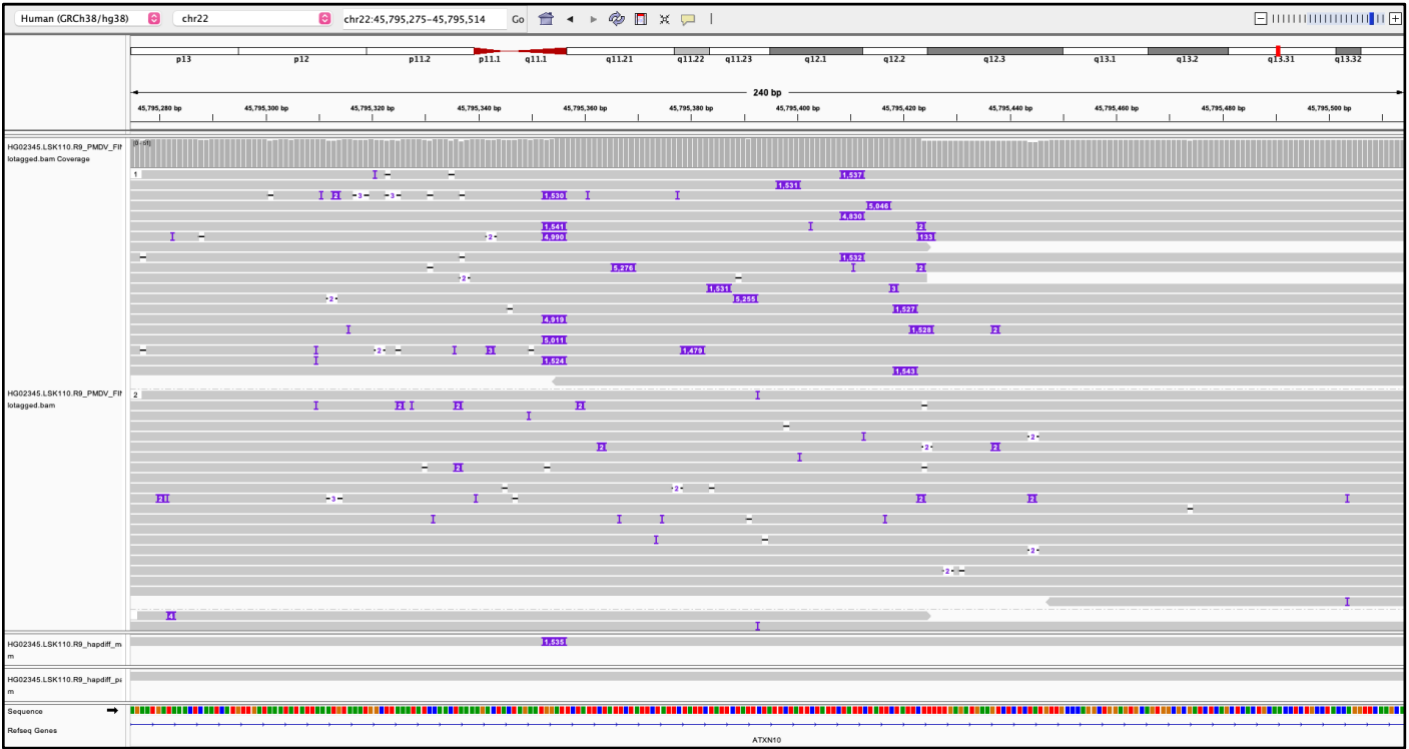


**HG02252**

## HG02345

**Figure S19. IGV screenshots of repeat expansion within *FGF14*.**

IGV screenshot for the sample at the SCA27B-associated locus at *FGF14* which is associated with autosomal dominant spinocerebellar ataxia 27B (SCA27B, MIM: 620174) with over 250 motifs which is within the range (250–400 repeat units) previously reported to be associated with reduced penetrance. The PMDV haplotagged alignment and the hapdiff assemblies from the Napu pipeline for each sample are shown. The hapdiff assembly files were used for genotyping with vamos.
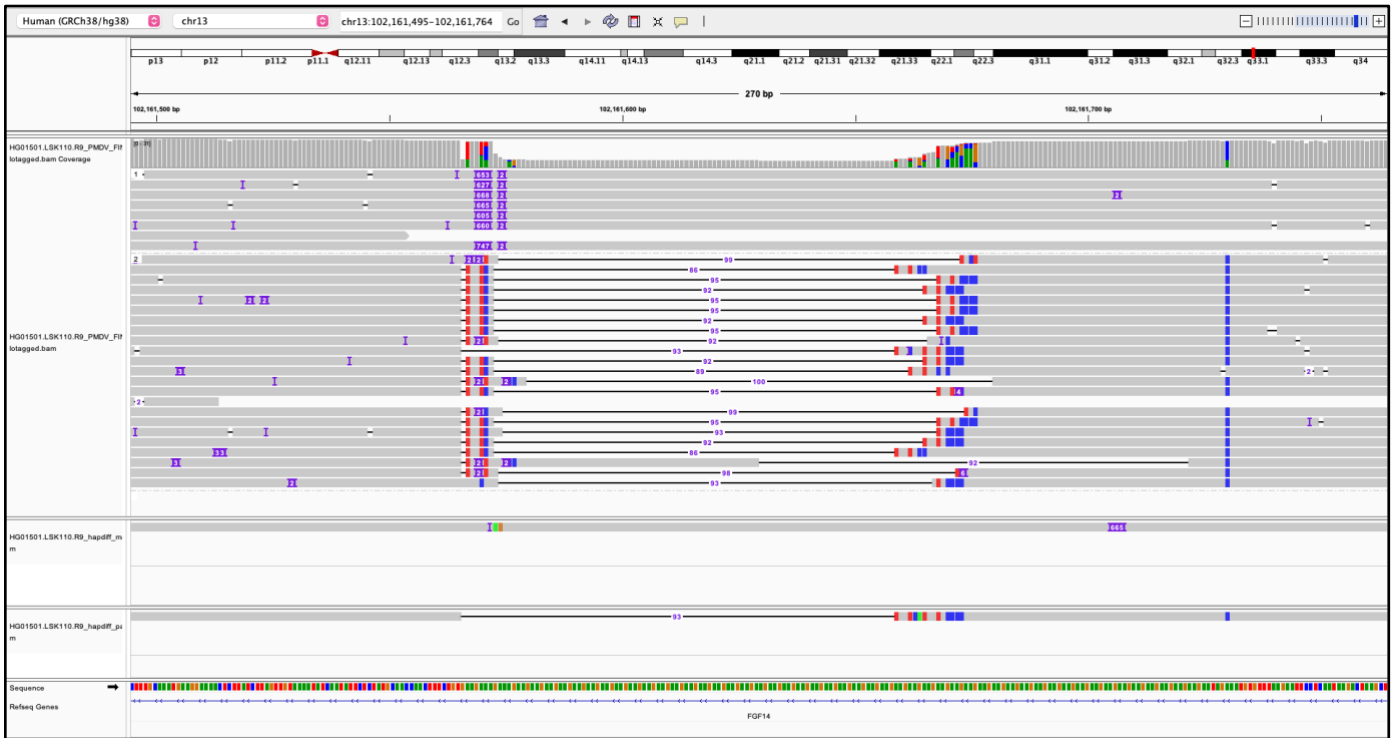
**HG01501**

**Figure S20. Principal component analysis using methylation.**
The average methylation frequency for 27,050 CpG islands for 76 samples was used for performing PCA. One 46,XX sample (GM18864) clusters with the 46,XY samples.
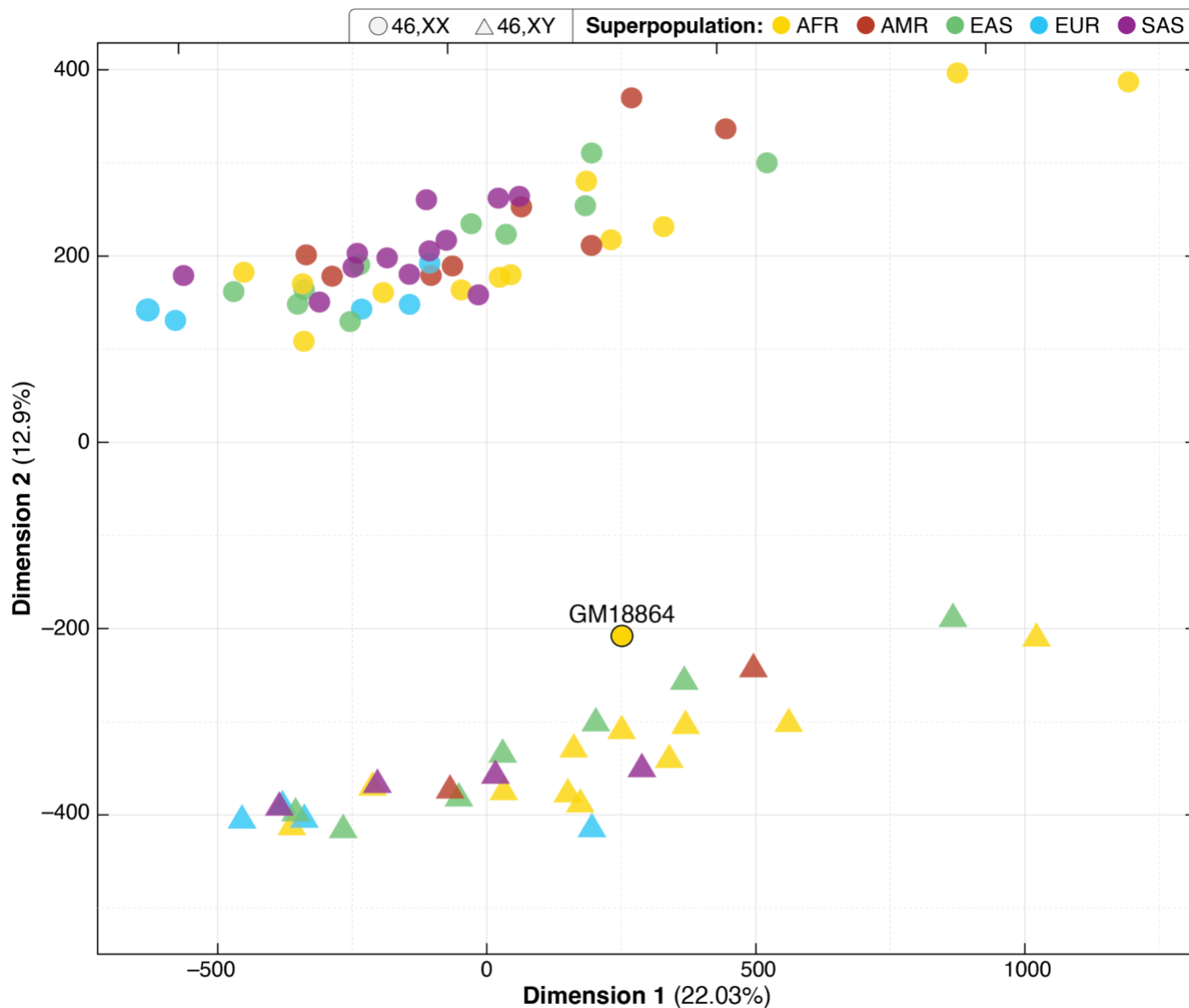
**Figure S21. Evaluating methylation differences at *SNURF-SNRPN*.**
**A**. Gain of methylation was observed in GM19473 (bottom), where both haplotypes are mostly methylated (red). HG03069 (top) is used as a control. **B**. Loss of methylation is seen in HP1 of HG00525 (top). HG04216 (bottom) is used as a control.
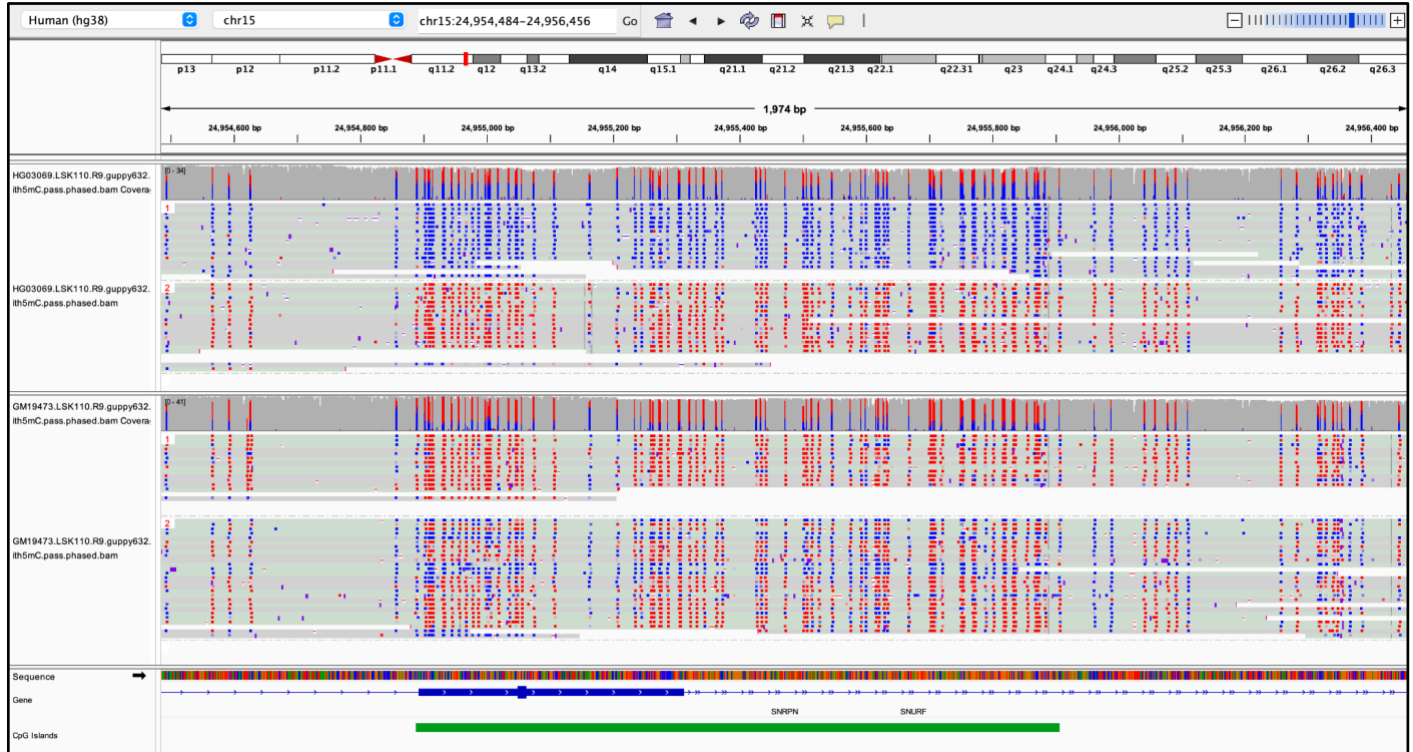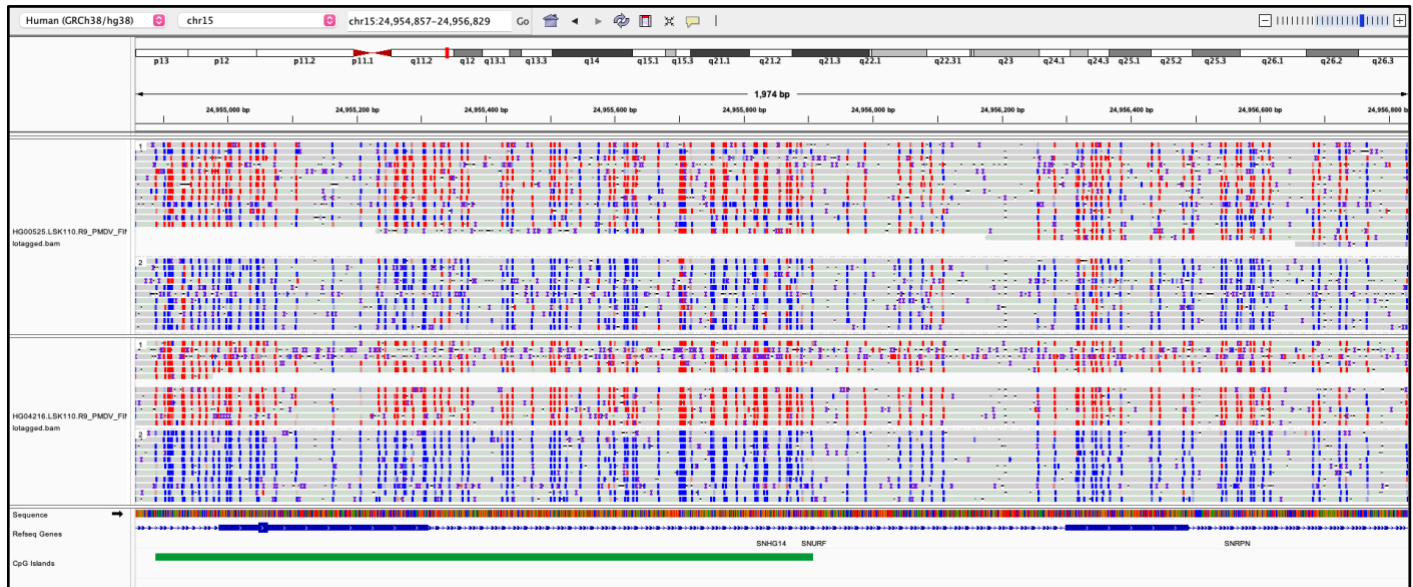
**A**



**B**

**Figure S22. Unique differentially methylated region identified by MeOW.**
One haplotype in HG02389 had increased methylation at an internal CpG site (3007) within an intron of
*SLC29A3*. Methylation frequency by haplotype is shown for HG02389 and one control (HG03022). Methylation
status is shown for individual reads for each haplotype from HG02389 only (red indicates a methylated CpG;
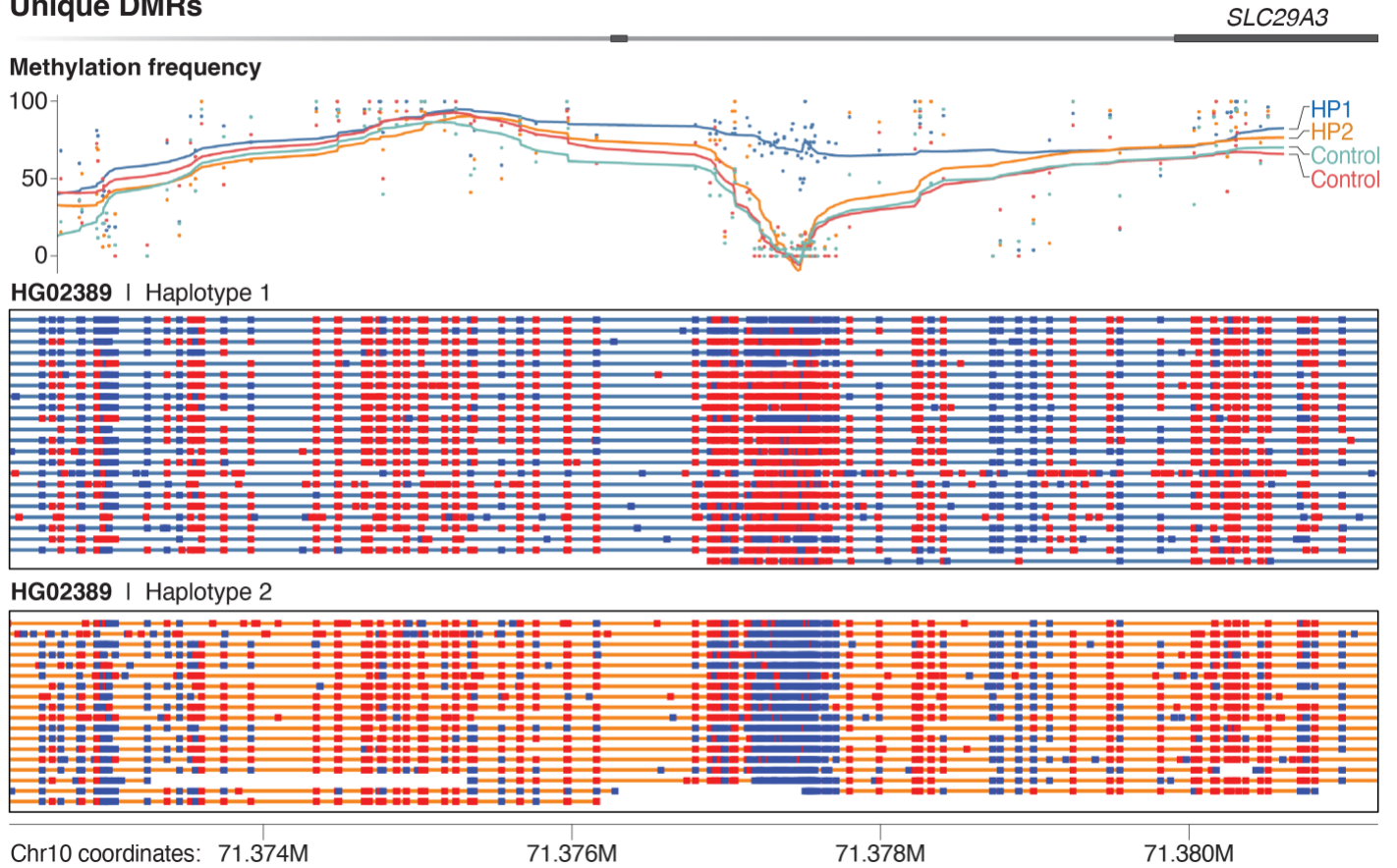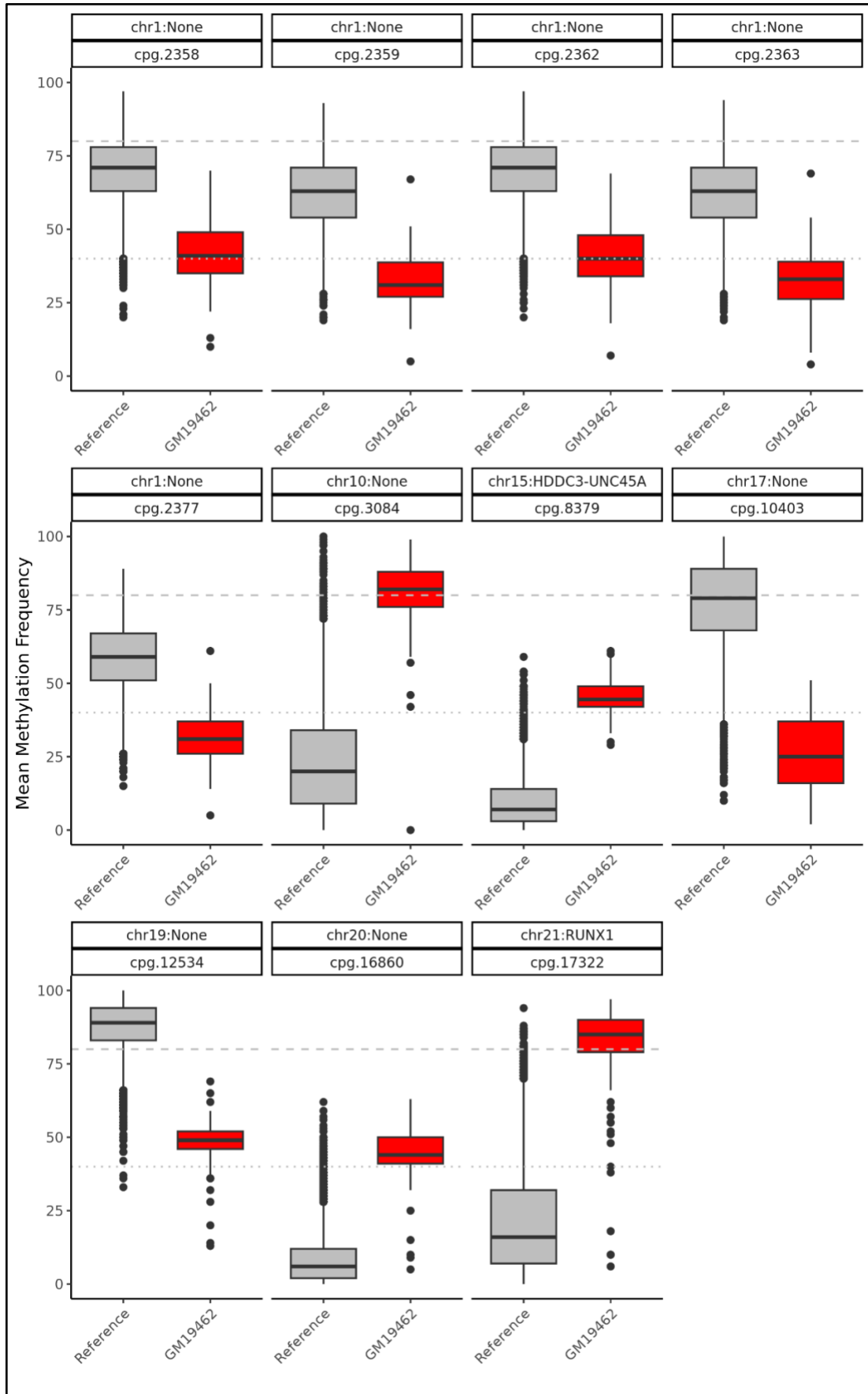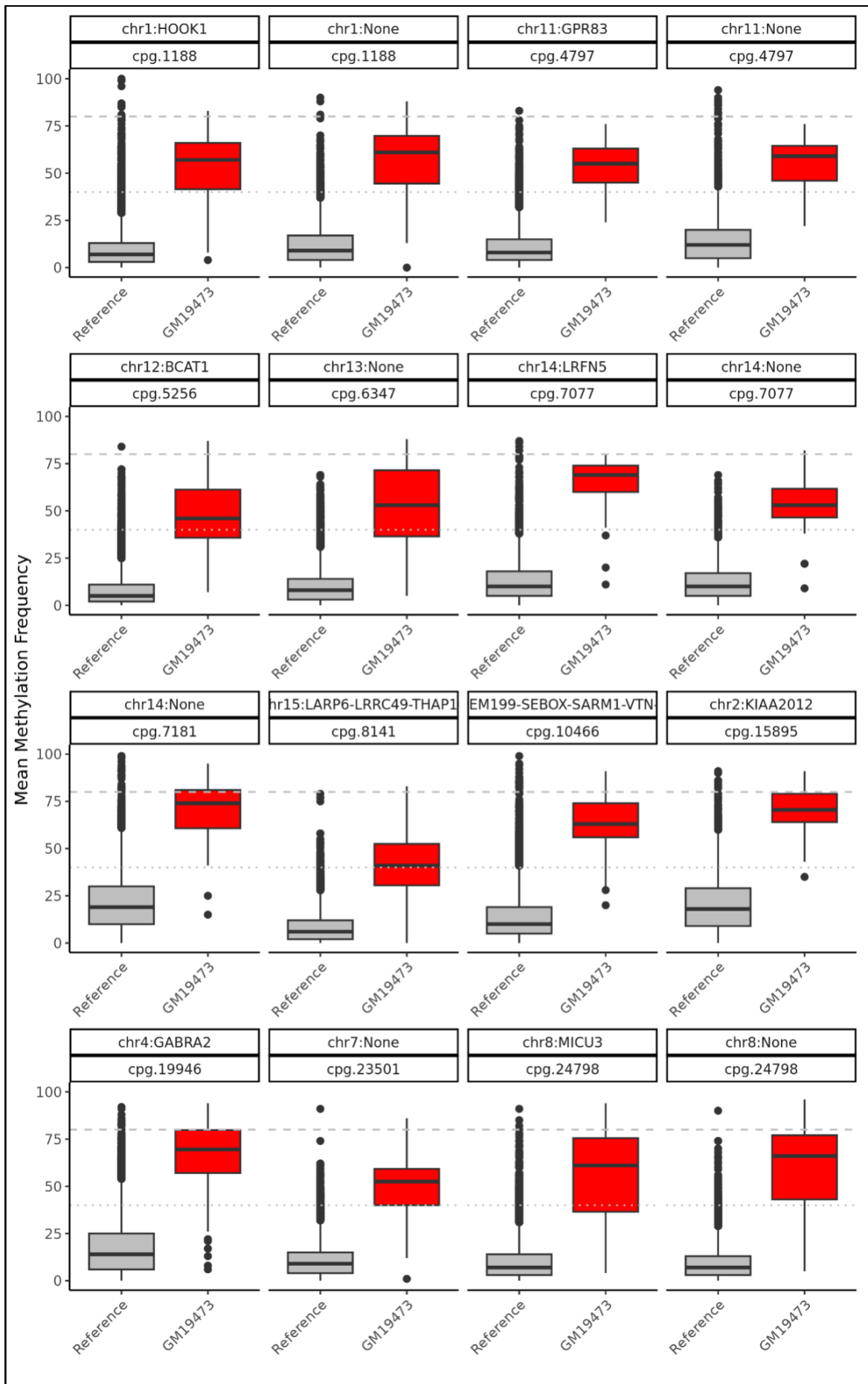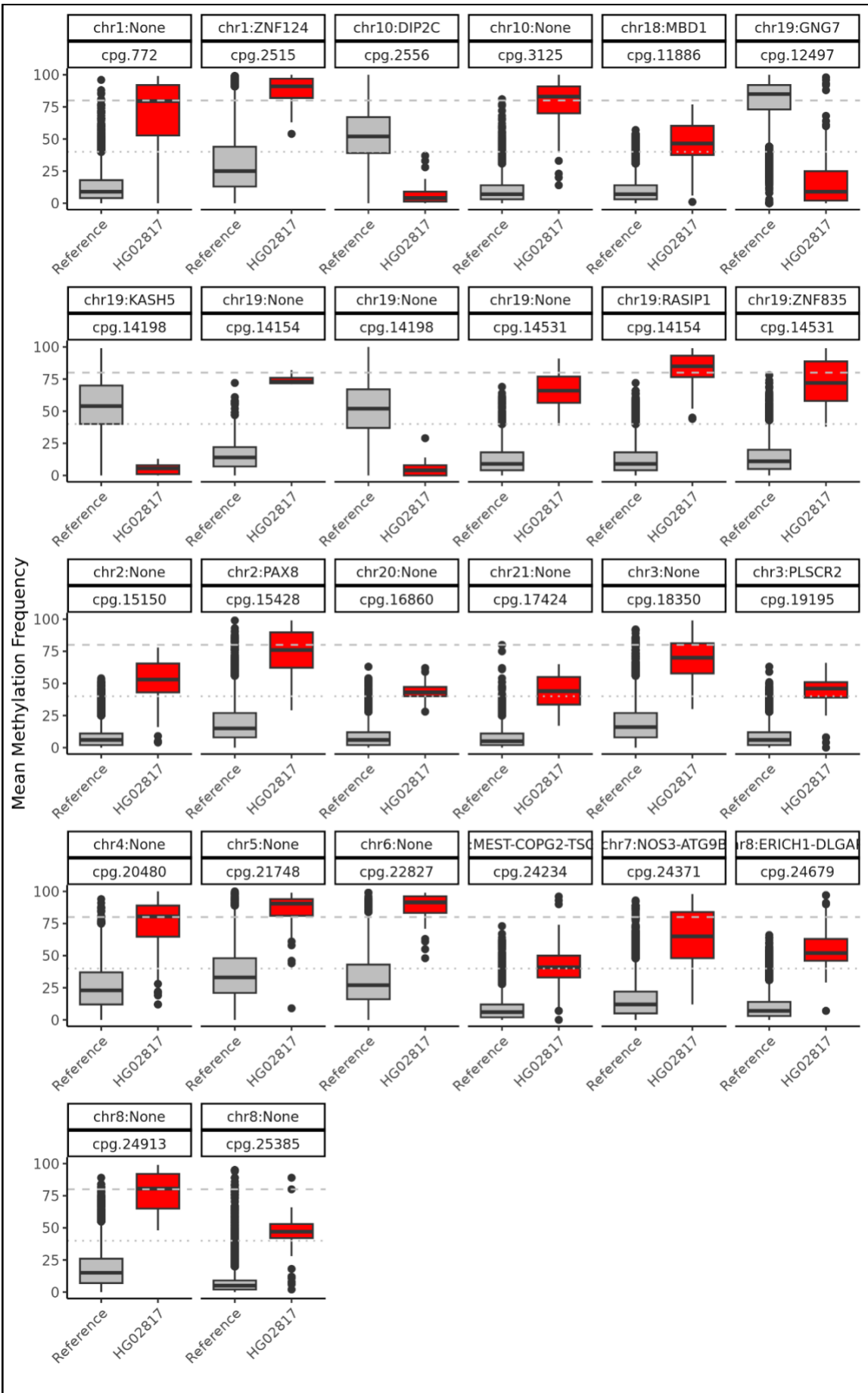blue indicates an unmethylated CpG).

**Figure S23. Output from MeOW showing DMRs within 3 samples that have 10 or more DMRs.**
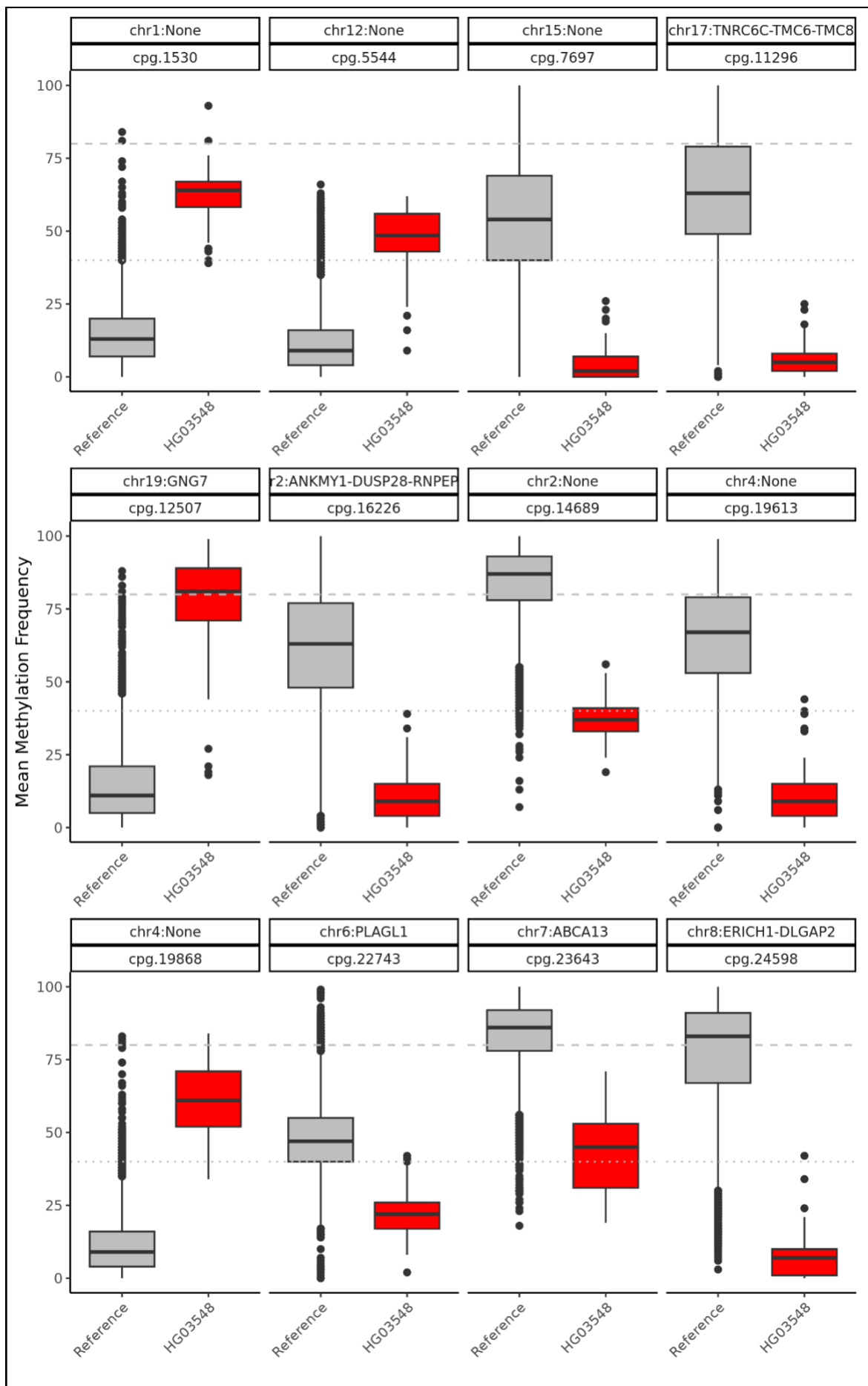
**Figure S24. Unique DMRs are associated with nearby changes in gene expression.**
Analysis of RNA-sequencing data from 15 samples of African ancestry focused on 85 genes within 10 kbp of DMRs identified by MeOW. Expression *Z*-score thresholds versus log odds ratios across all sample–gene pairs identify expression outliers.