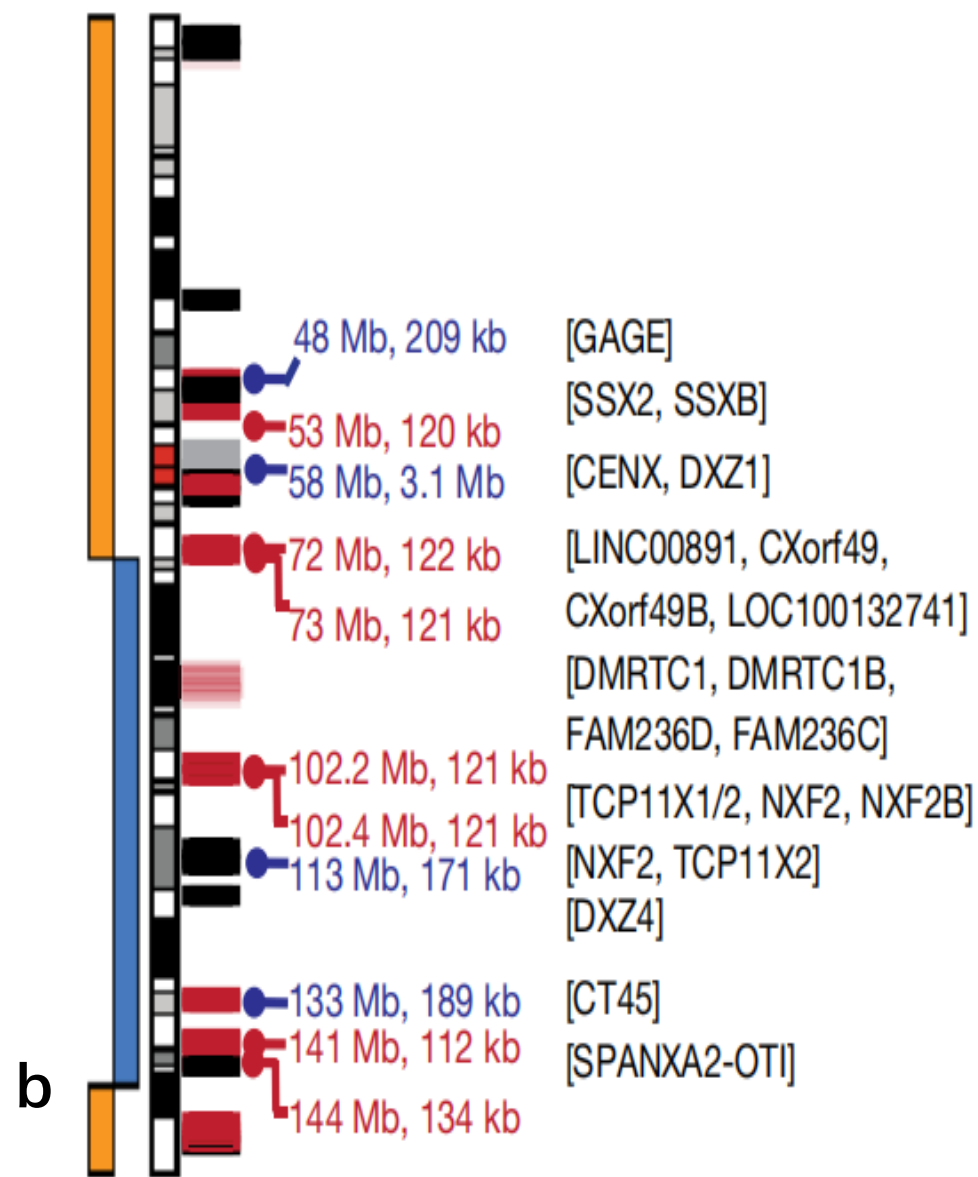
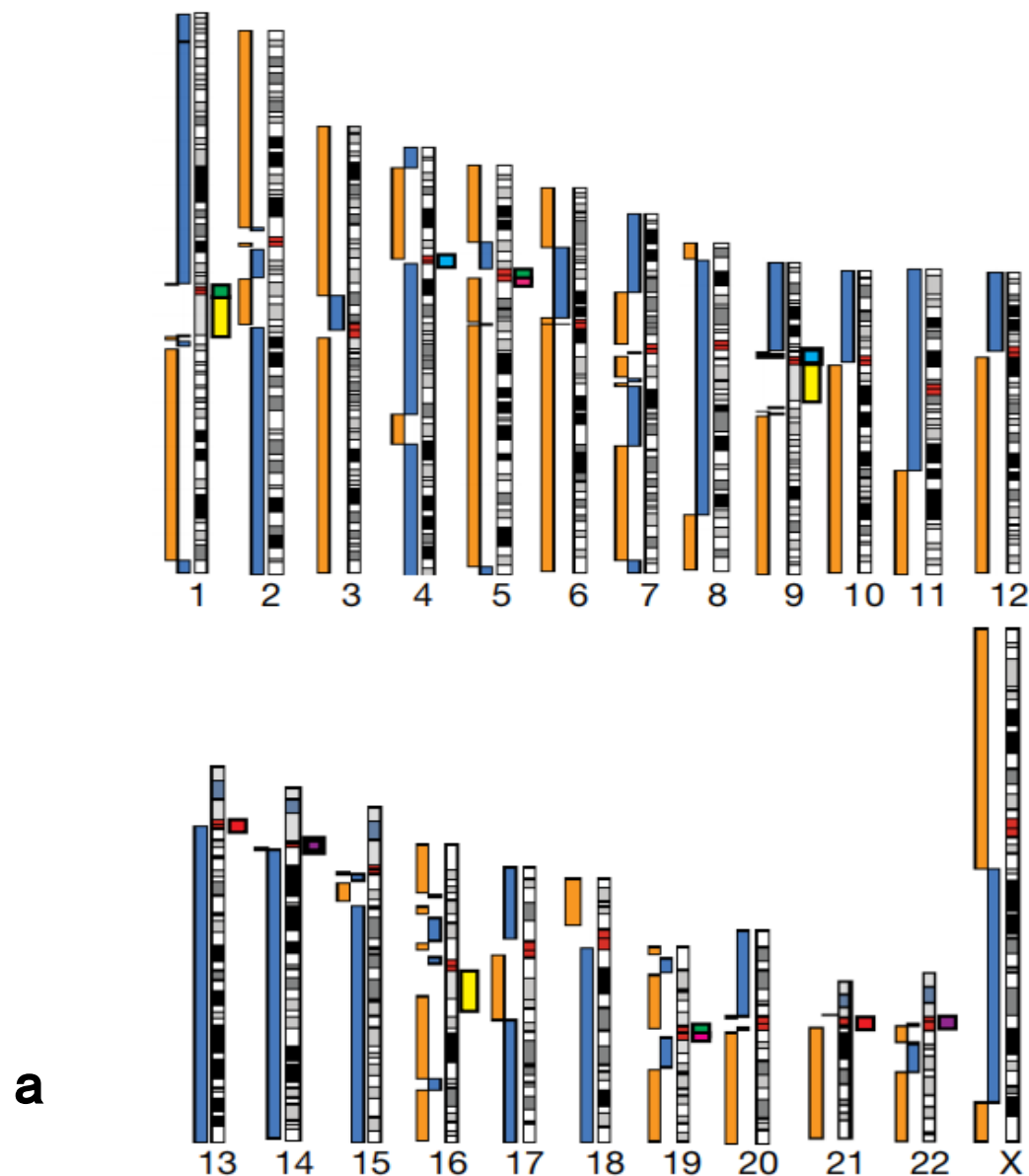
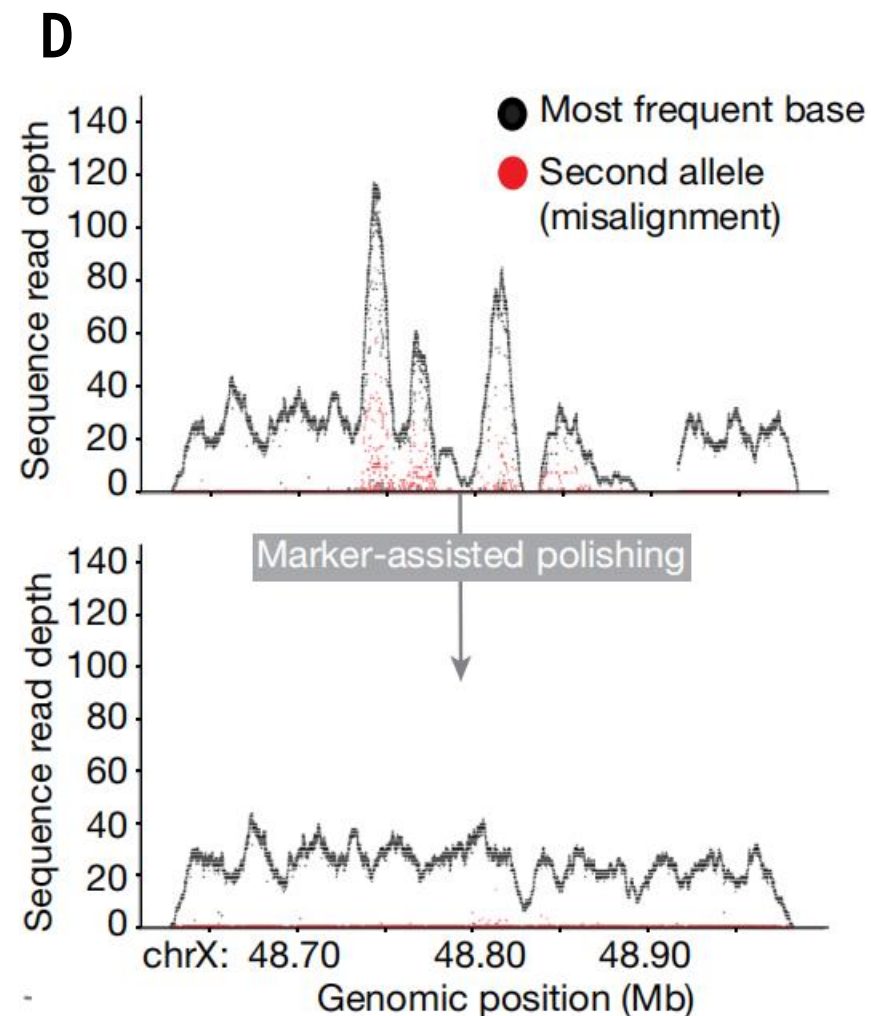
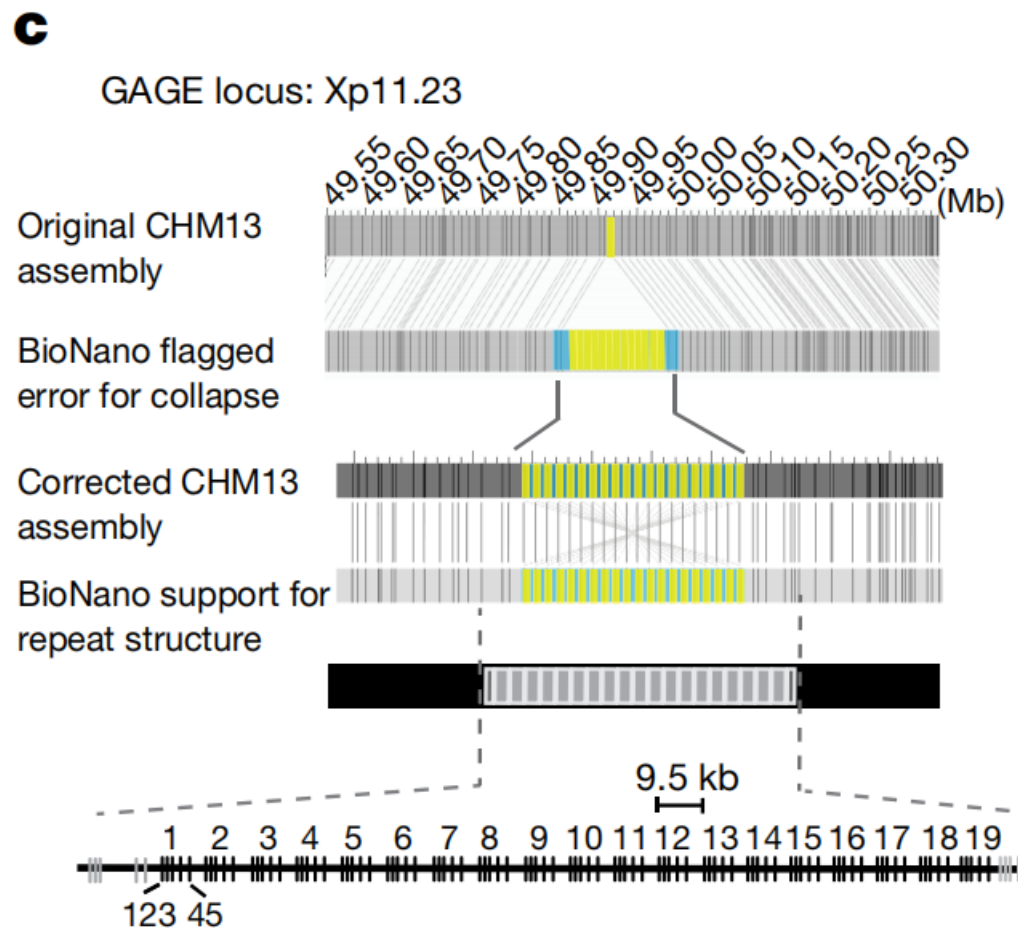


BGI 华大

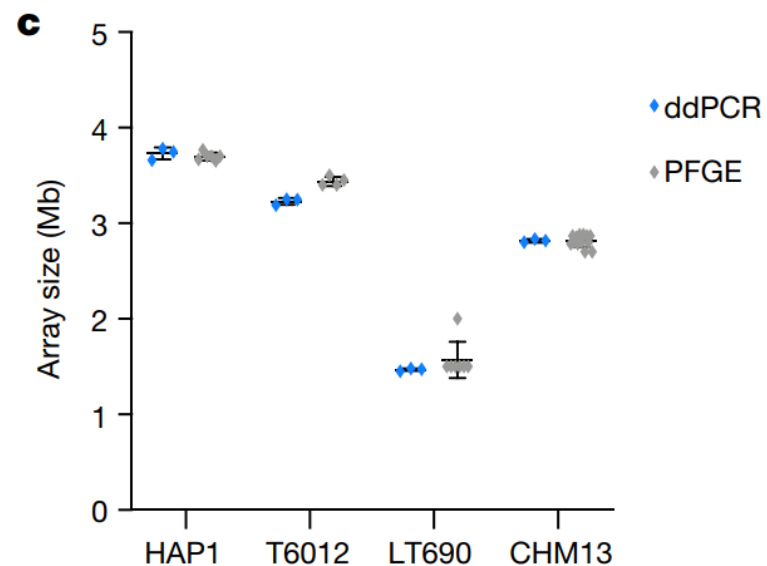
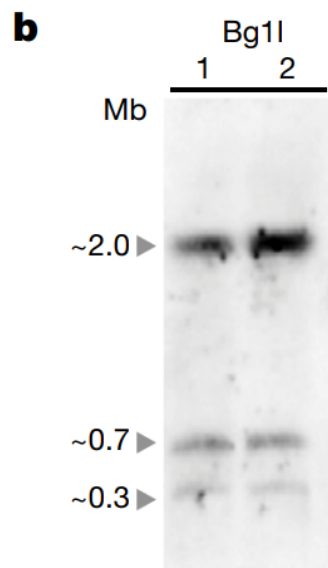
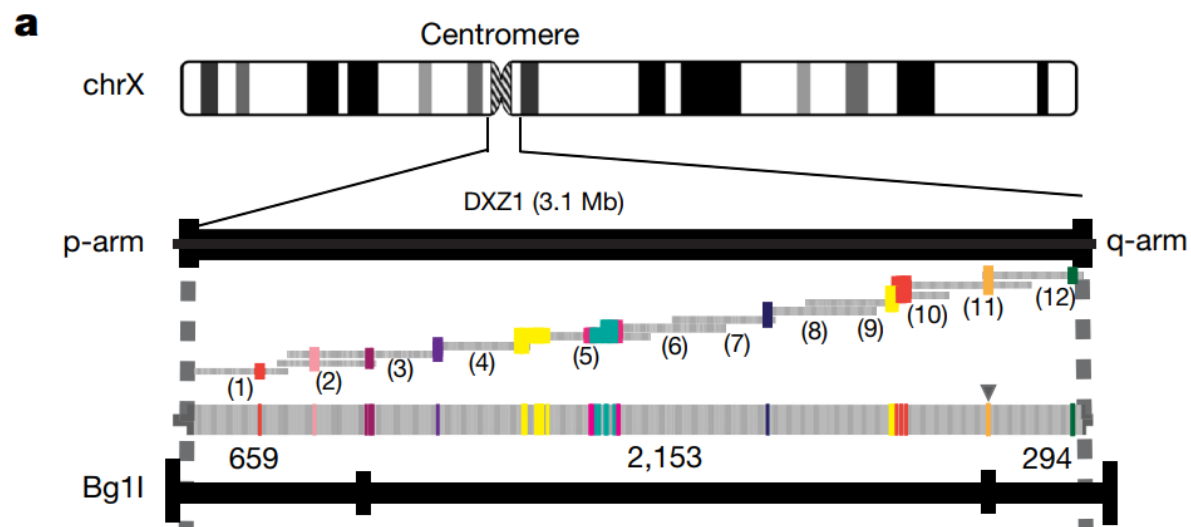
Telomere-to-telomere assembly of a complete human X chromosome

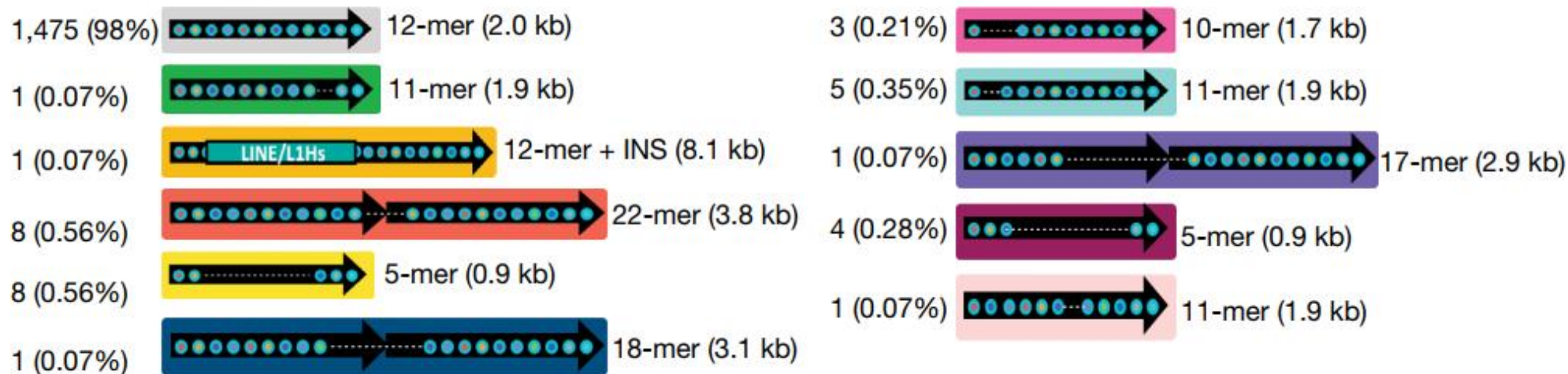
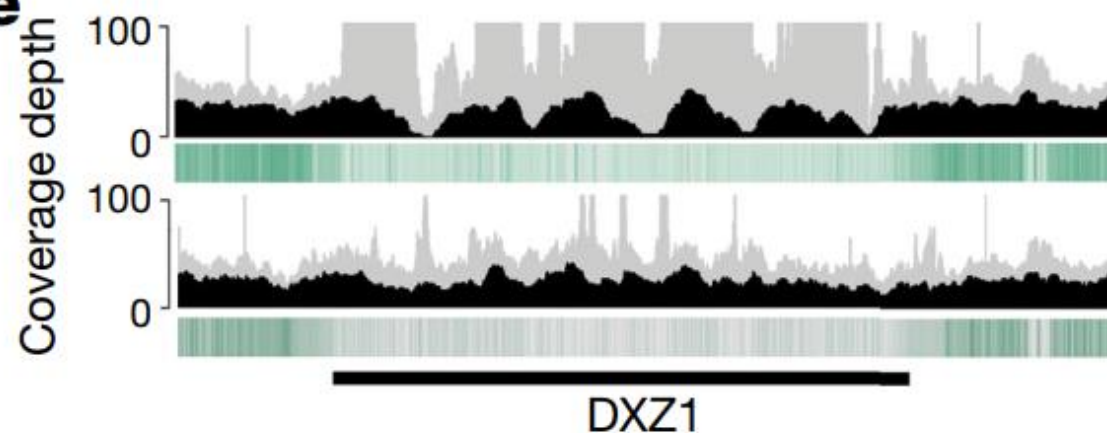
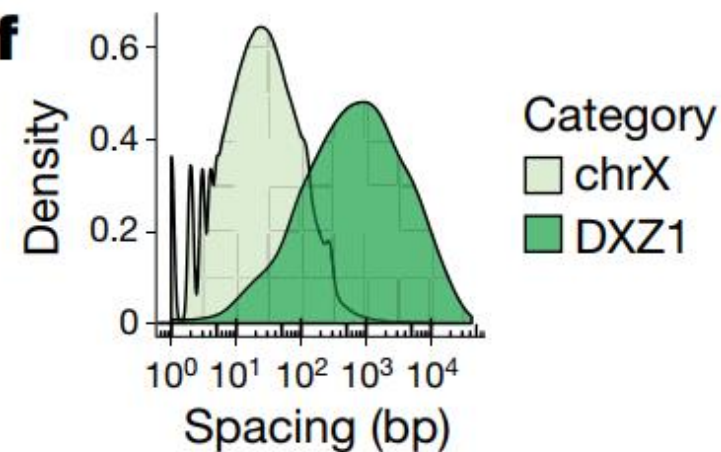


BGI华大 whole-genome assembly and validation

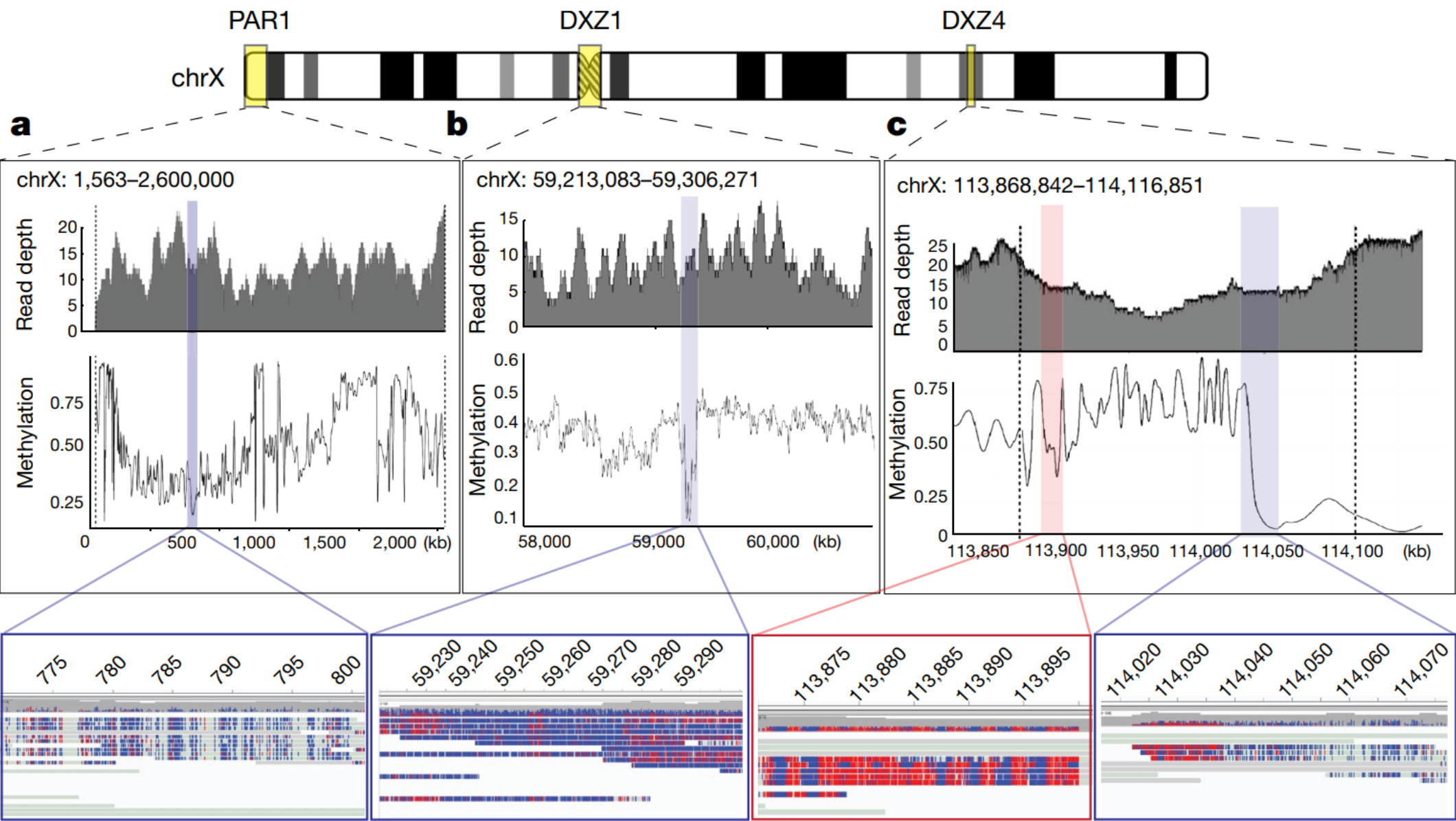


A finished human X chromosome



d**e****f**

BGI华大 Chromosome-wide DNA methylation maps



现有成果

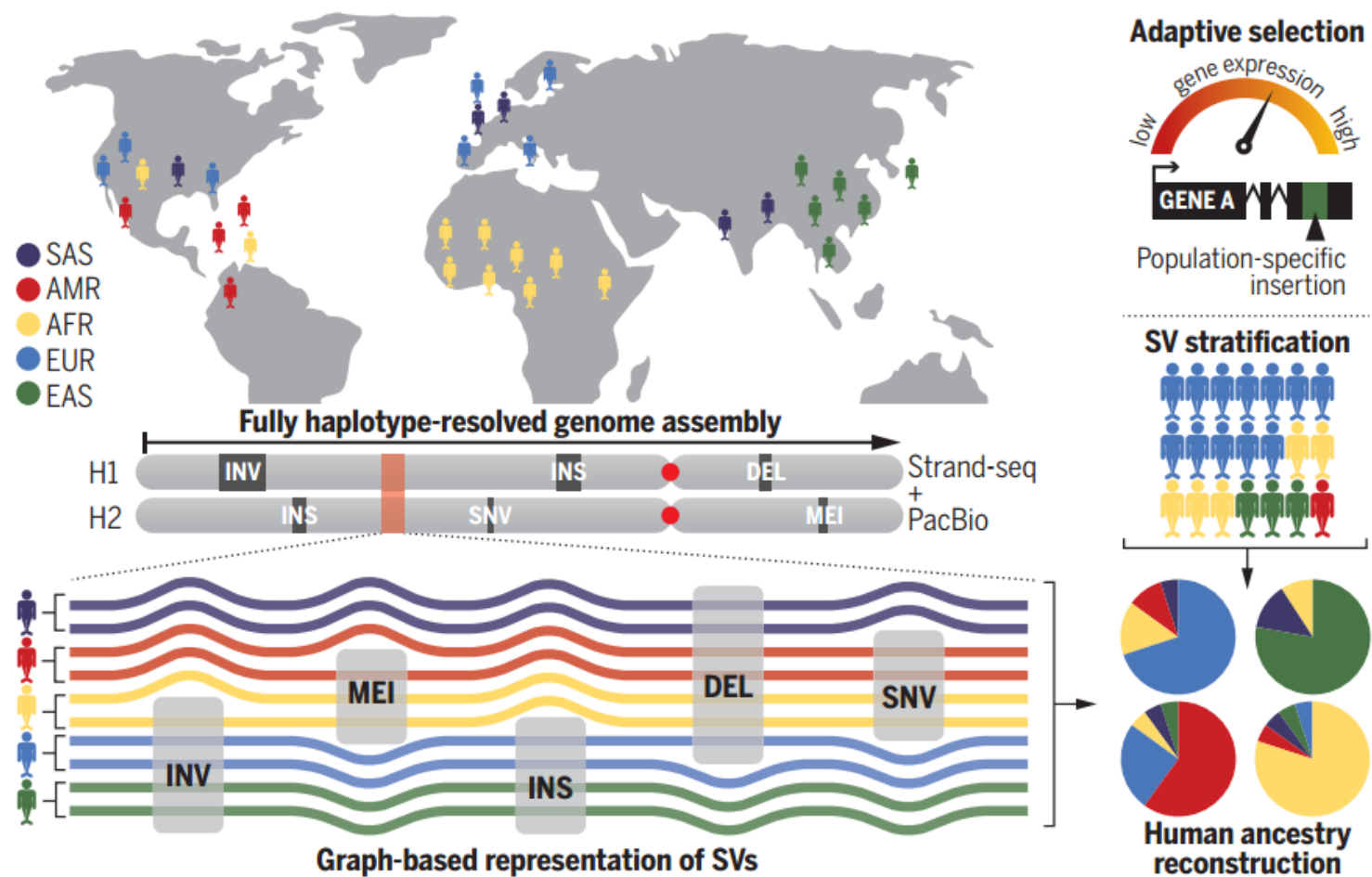
- 1.同时重建了其他几条染色体
- 2.提供了使用现有技术完成整个人类基因组的技术方案

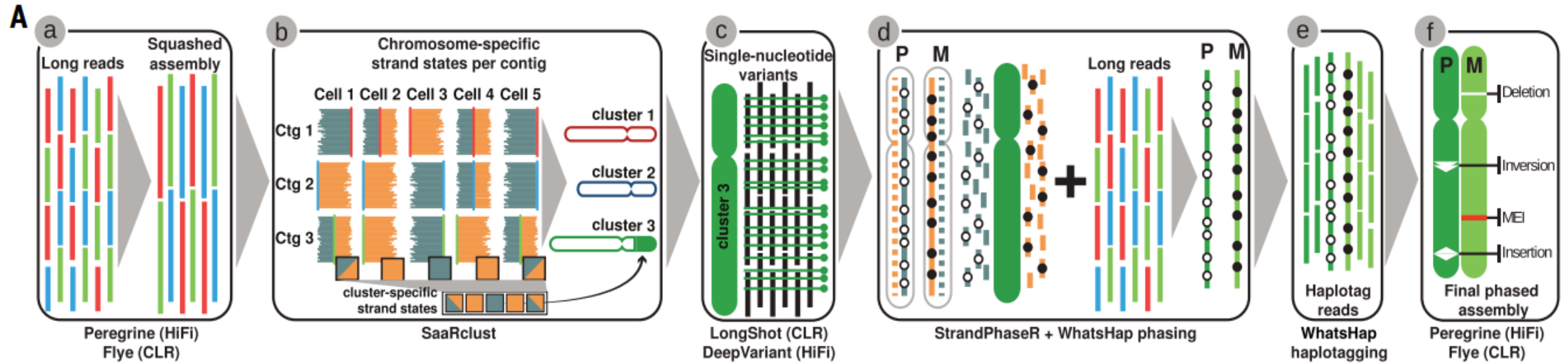
仍需克服的挑战

- 1.应用于二倍体样本将需要对潜在的单倍型进行定相，以避免混合复杂结构变异的区域
- 2.重复区域比X染色体更大的着丝粒卫星DNA将需要开发其他方法

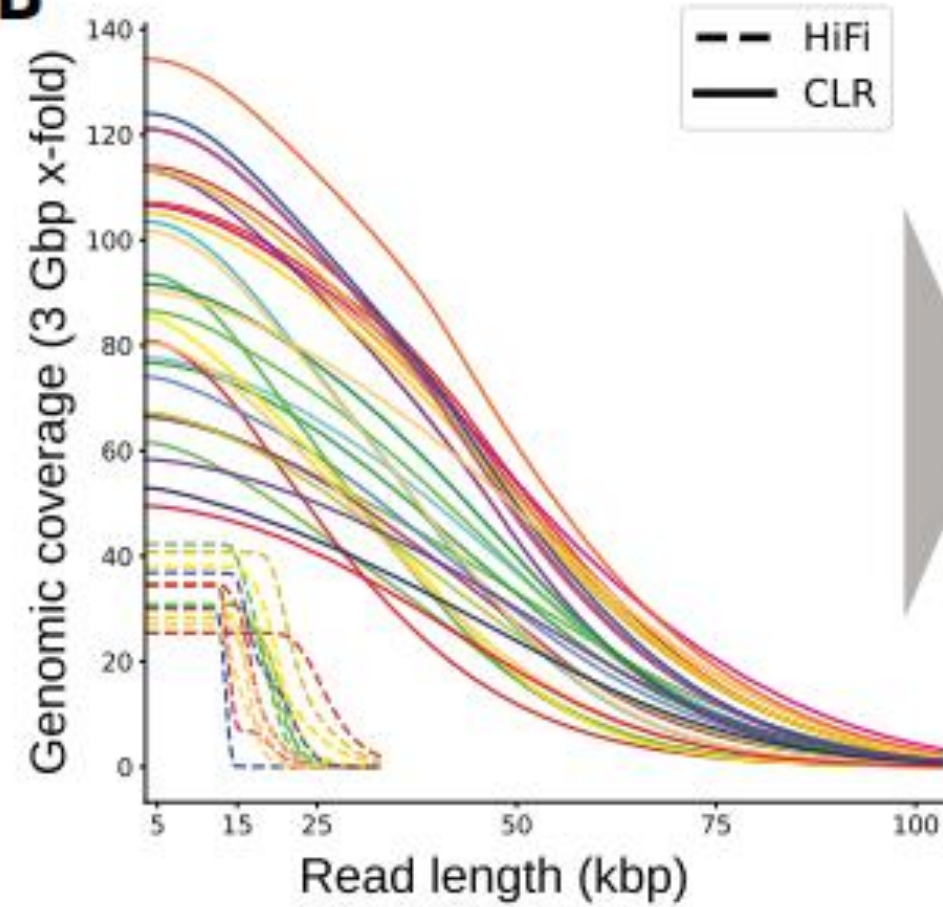
BGI 华大

Haplotype-resolved diverse human genomes
and integrated analysis of structural variation

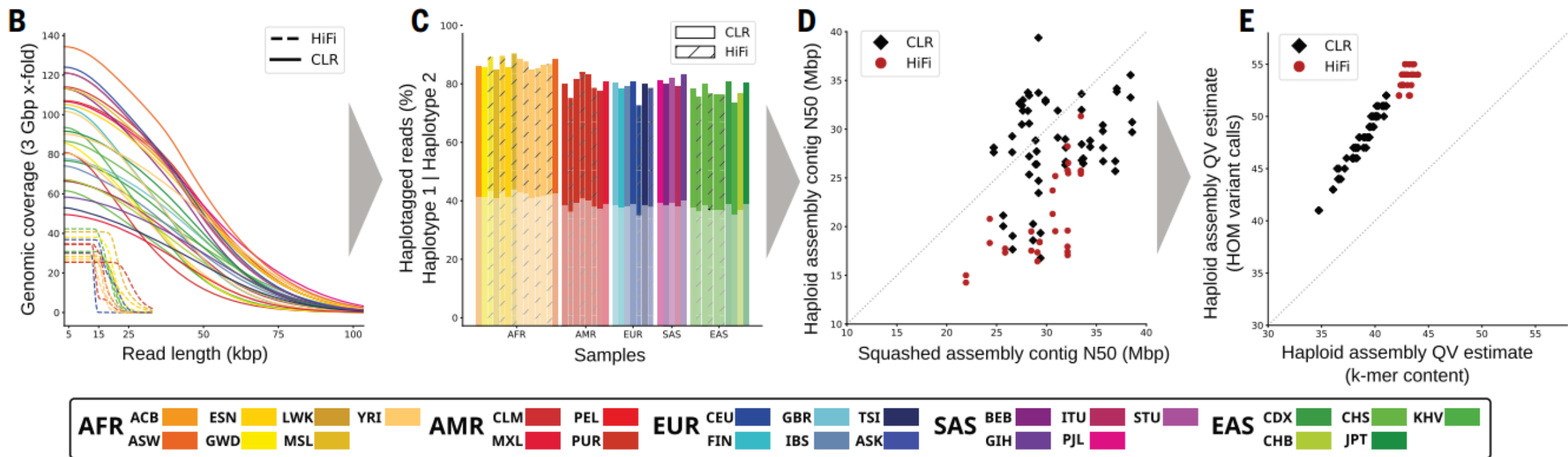




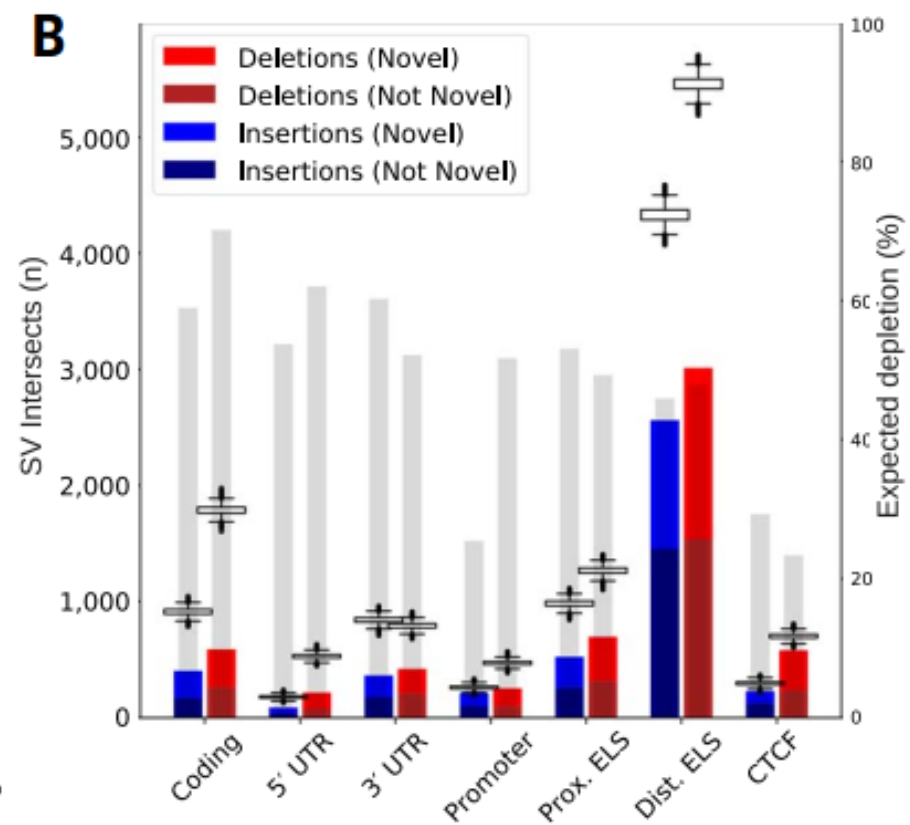
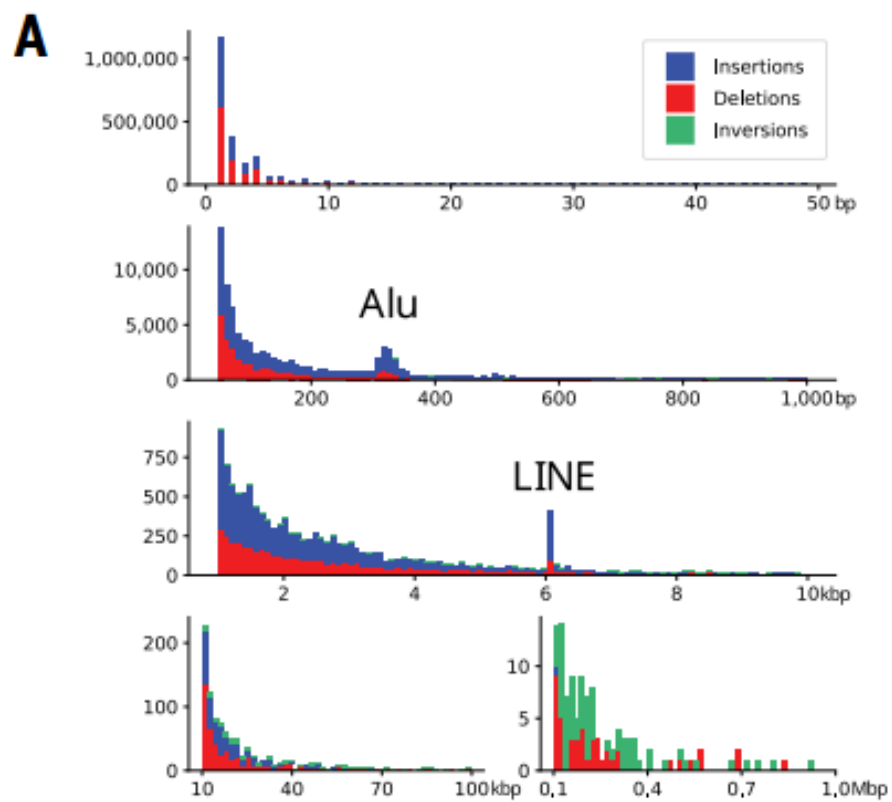
The Human Genome Structural Variation Consortium (HGSVC) recently developed a method for phased genome assembly that combines long-read PacBio whole-genome sequencing (WGS) and Strand-seq data to produce fully phased diploid genome assemblies without dependency on parent-child trio data

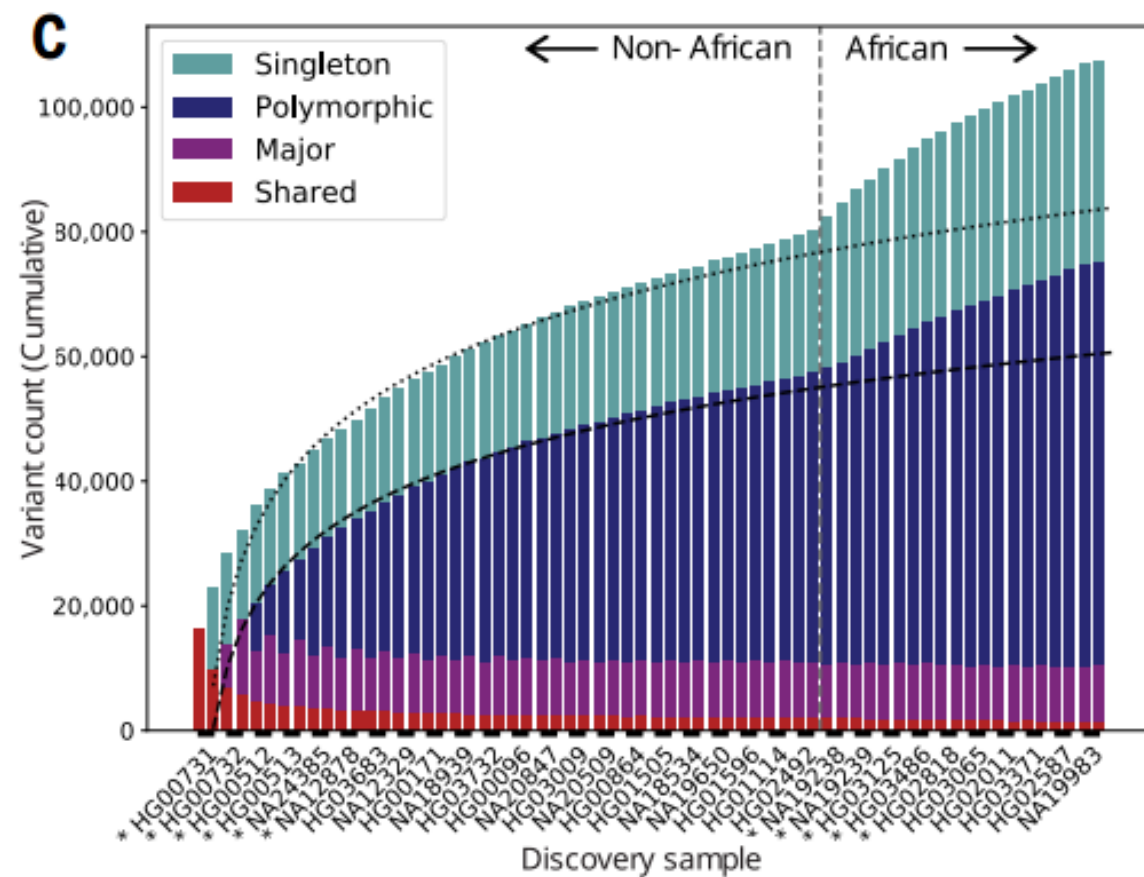
B

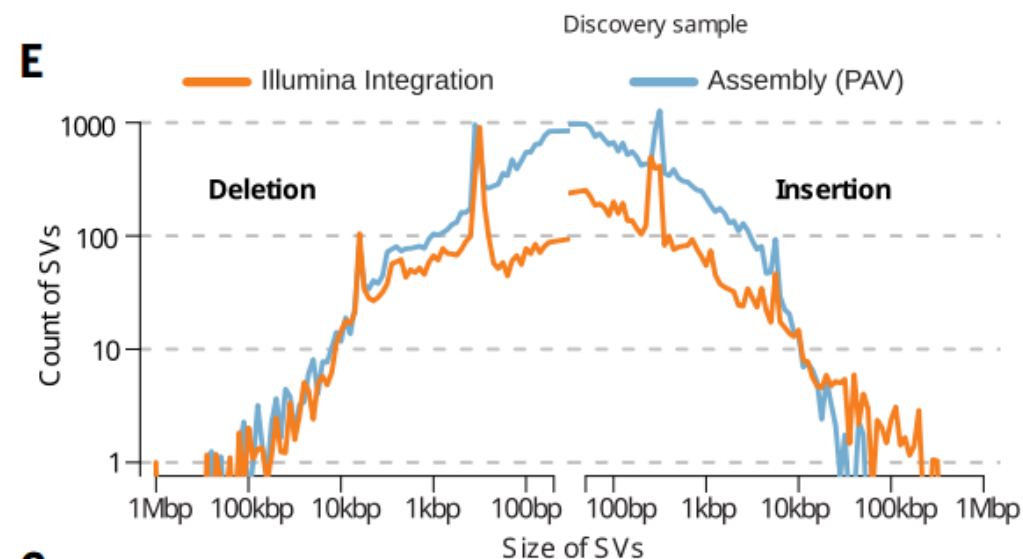
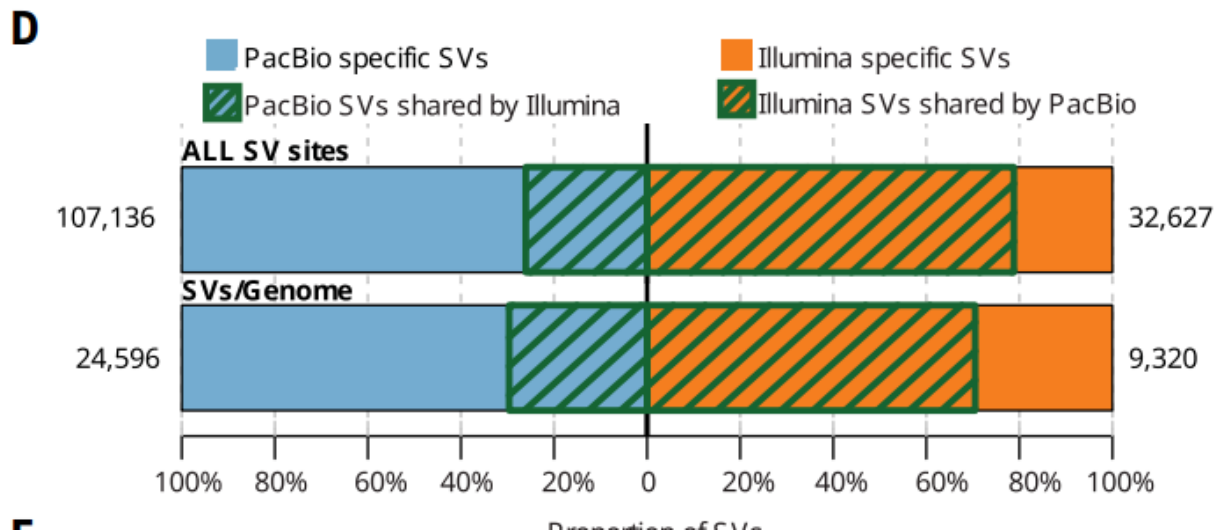
Unintegrated reads (%)



(B) Genomic coverage (y axis) as a function of the long-read length (x axis). (C) Fraction of reads that can be assigned ("haplotagged") to either haplotype 1 (semitransparent) or haplotype 2 for HiFi (hatched) and CLR (solid) datasets. (D) Contig-level N50 values for squashed (x axis) and haploid assemblies (y axis) for CLR (black diamonds) and HiFi (red circles) samples. (E) Haploid assembly QV estimates computed from distinct and shared k-mers (x axis) based on homozygous Illumina variant calls (y axis). Samples are colored according to the 1000GP population color scheme (15), with the exception of the added Ashkenazim individual NA24385/HG002 (Coriell family ID 3140) (ASK; dark blue).

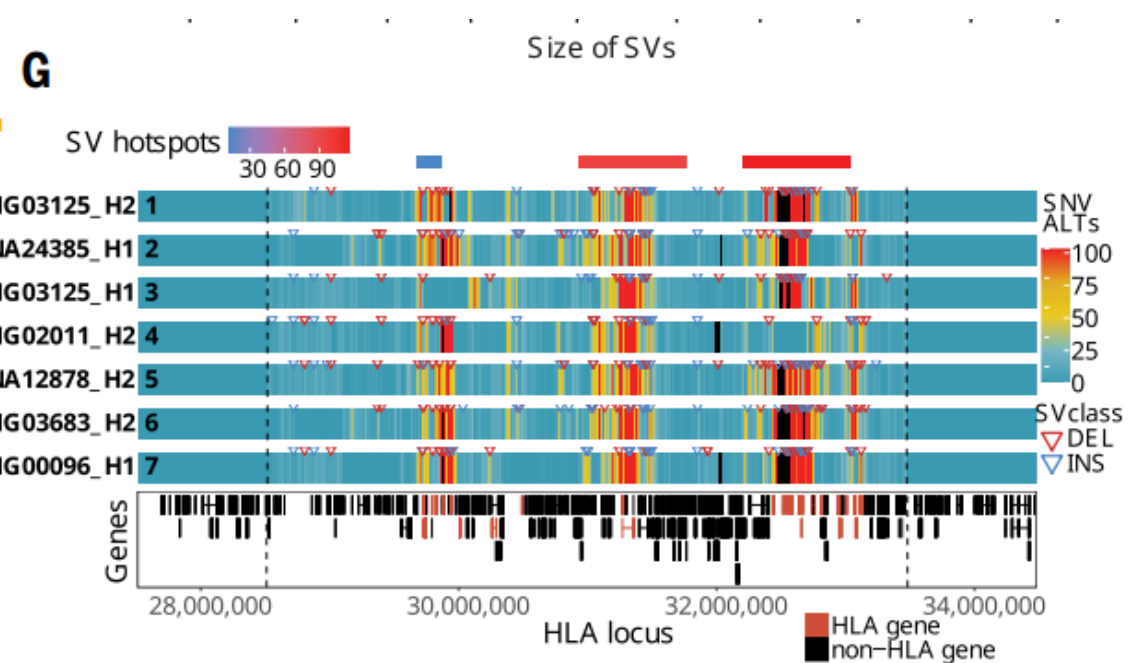
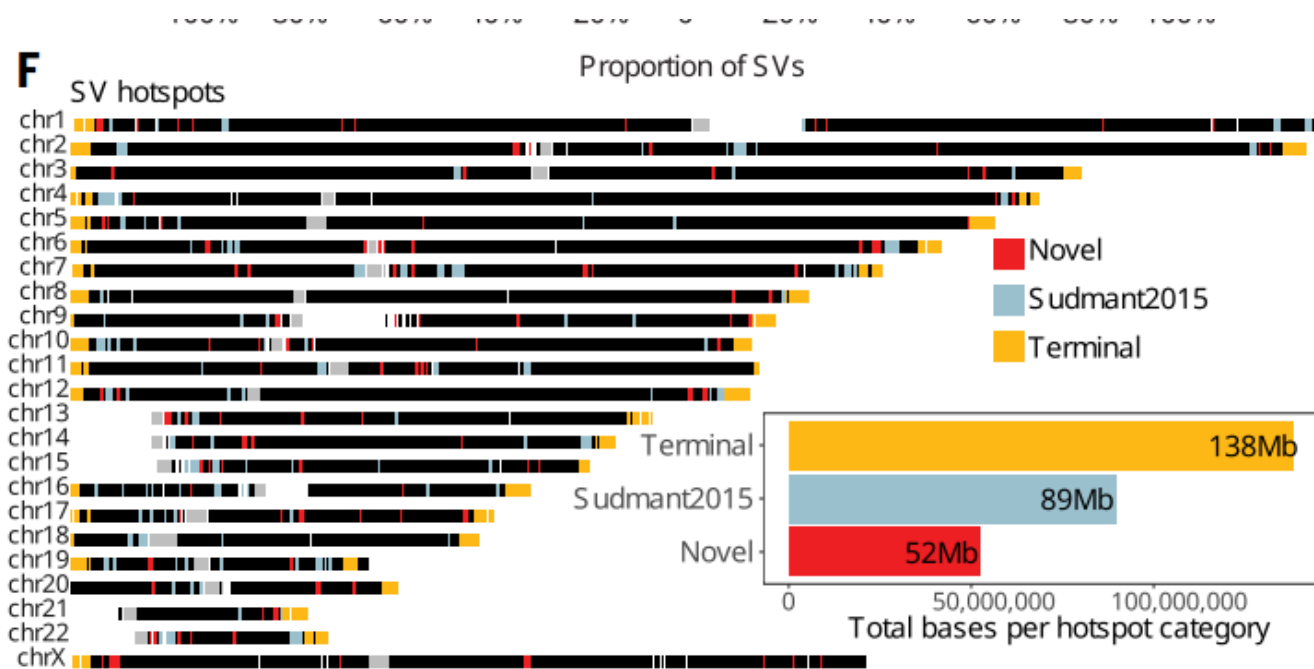






(D) Overlap between SVs detected by PacBio long-read assemblies and Illumina short-read alignments on 31 matched samples (NA24835, HG00514, HG00733, and NA19240 excluded). Top bar shows overall SV sites across 31 samples, and the bottom bar displays the average count of SVs per sample, with green stripes representing concordant SV calls between technologies.

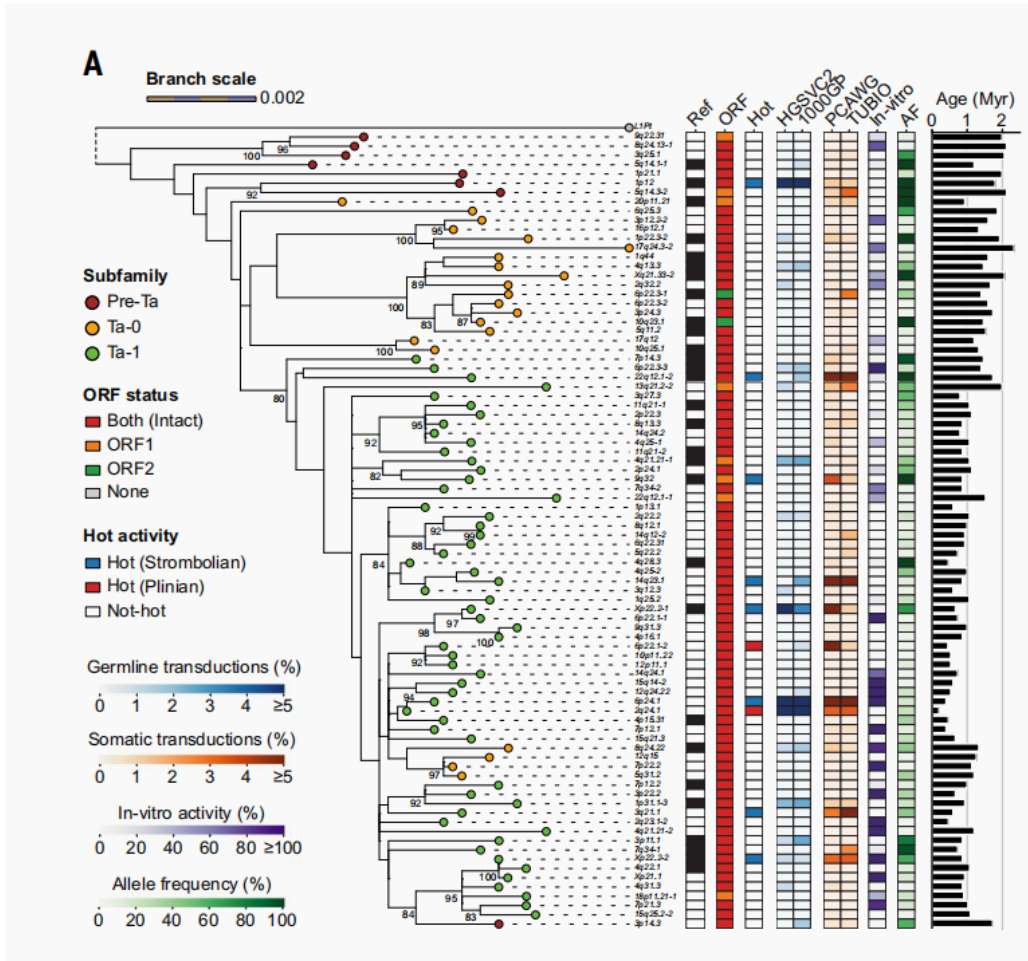
(E) Length distribution of SVs detected with PacBio long-read assemblies and Illumina short-read alignments across all 31 matched samples.



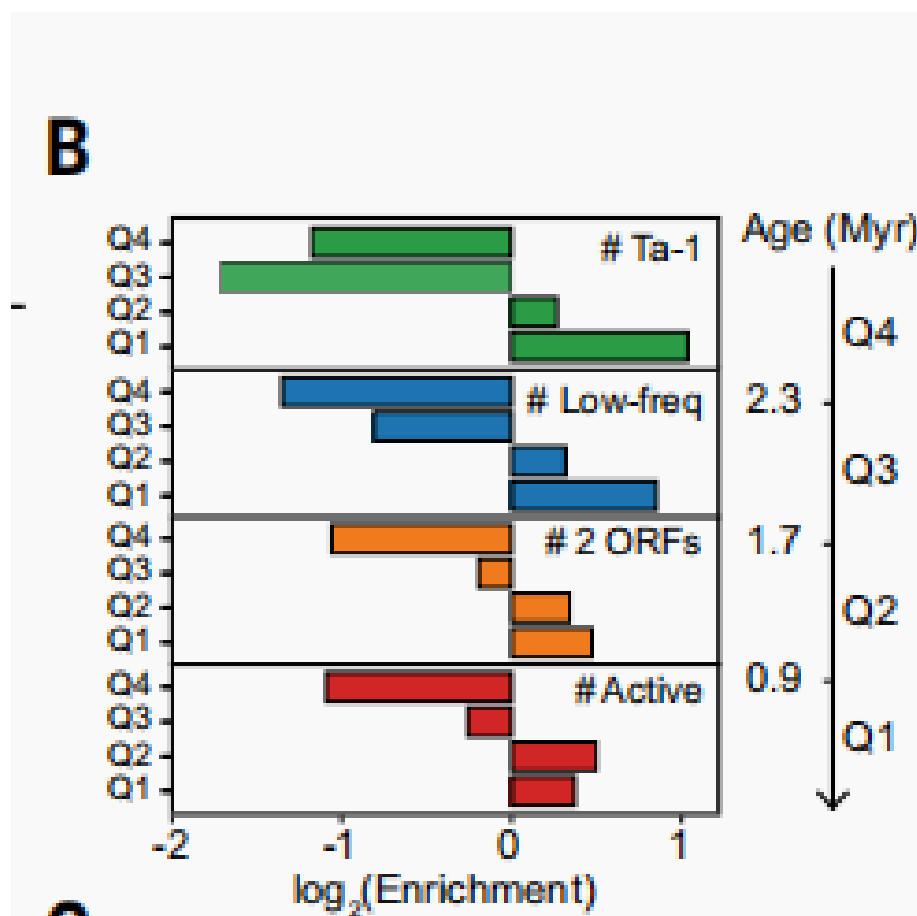
(F) Genome-wide distribution of SV hotspots divided in three categories: last 5 Mbp of chromosomes (yellow), overlapping (light blue), and previously unidentified (red) when compared with short-read SV analysis of 1000GP (23). (Inset) The total sequence length is represented by each hotspot category.

(G) Heatmap of seven selected SV haplotypes for 4-Mbp MHC region (chr6: 28,510,120 to 33,480,577; dashed lines) comparing regions of high-SNV (red) and low-diversity (blue) regions based on the number of alternate SNVs compared with the reference (GRCh38; alignment bin size 10 kbp, step 1 kbp). Phased SV insertions (blue open arrowheads) and deletions (red open arrowheads) are mapped above each haplotype. The most diverse regions correspond to SV hotspots (red and blue bars top row) and cluster with HLA genes (red bottom track)

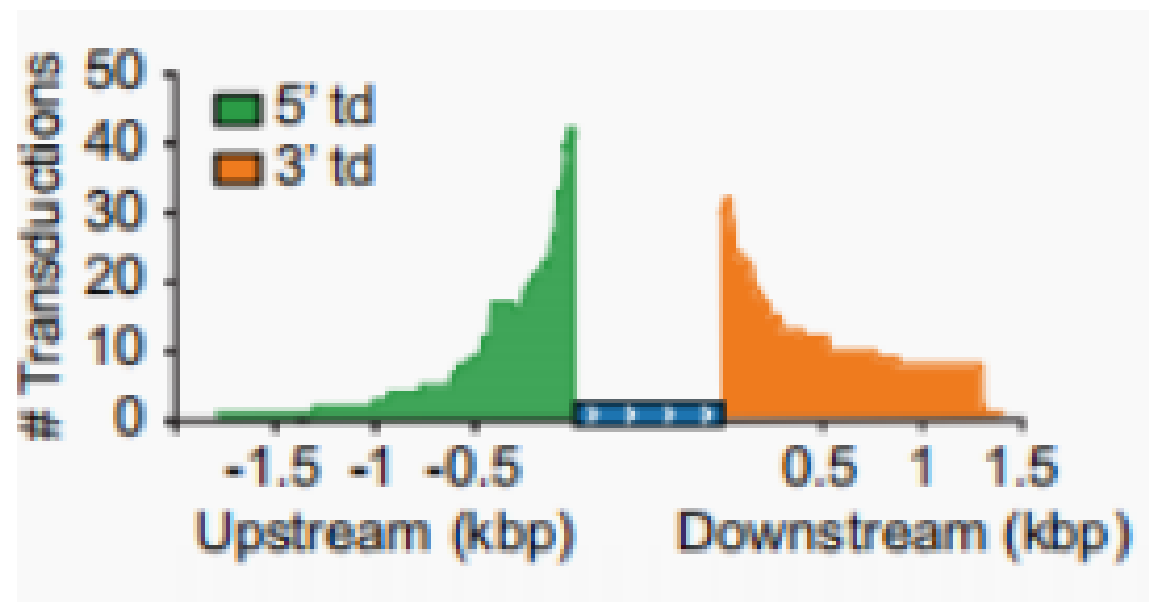
ME | S



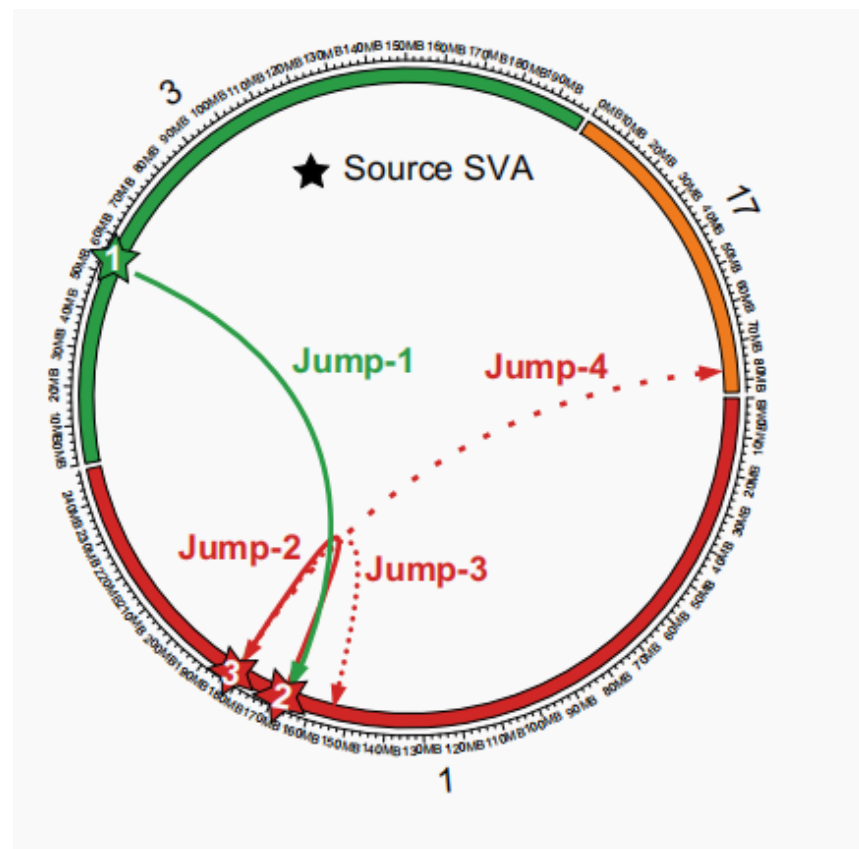
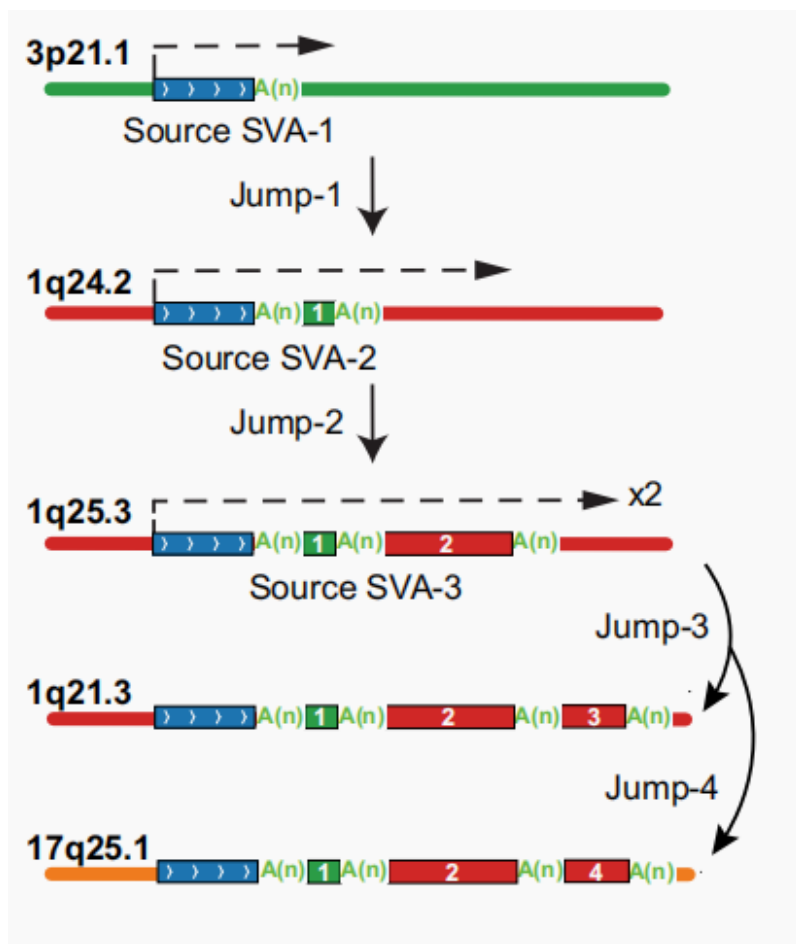
进化年龄与L1的亚家族、活性水平和等位基因频率等特征相关。



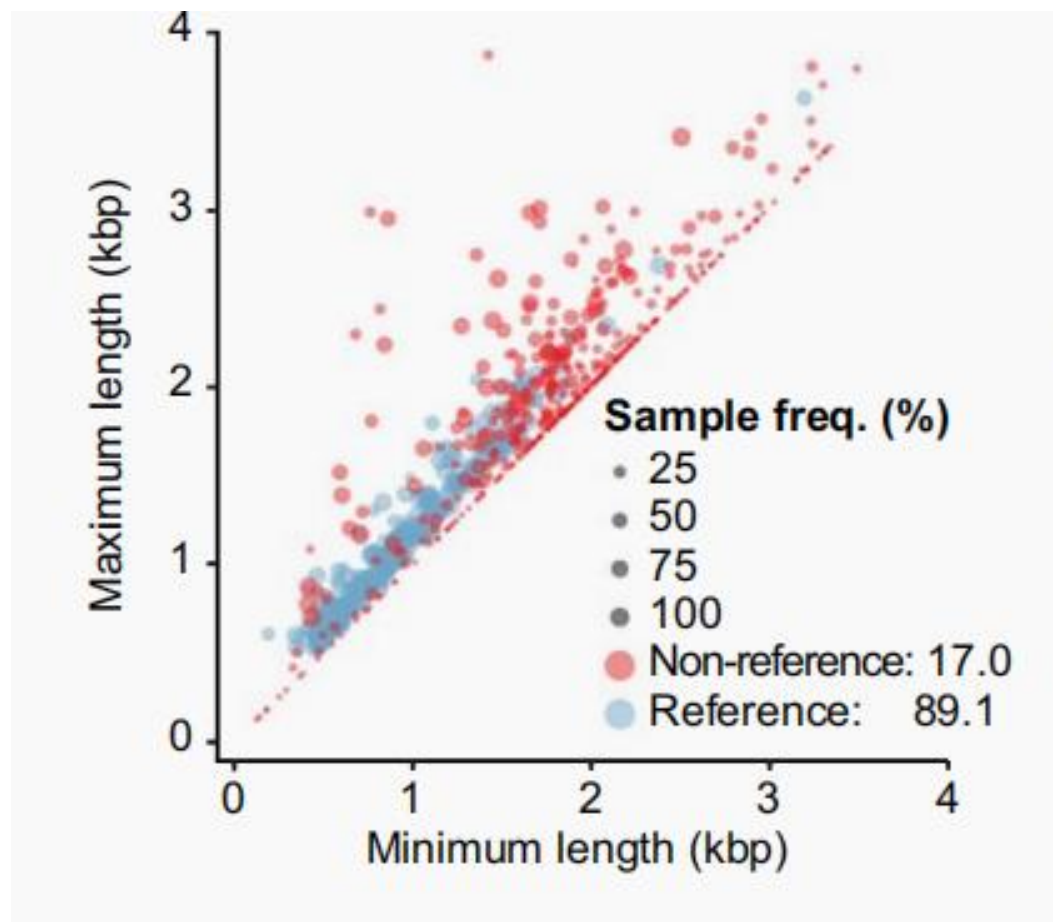
基于侧翼序列分析的5'和3' SVA介导的转导 (td) 的大小分布和数量



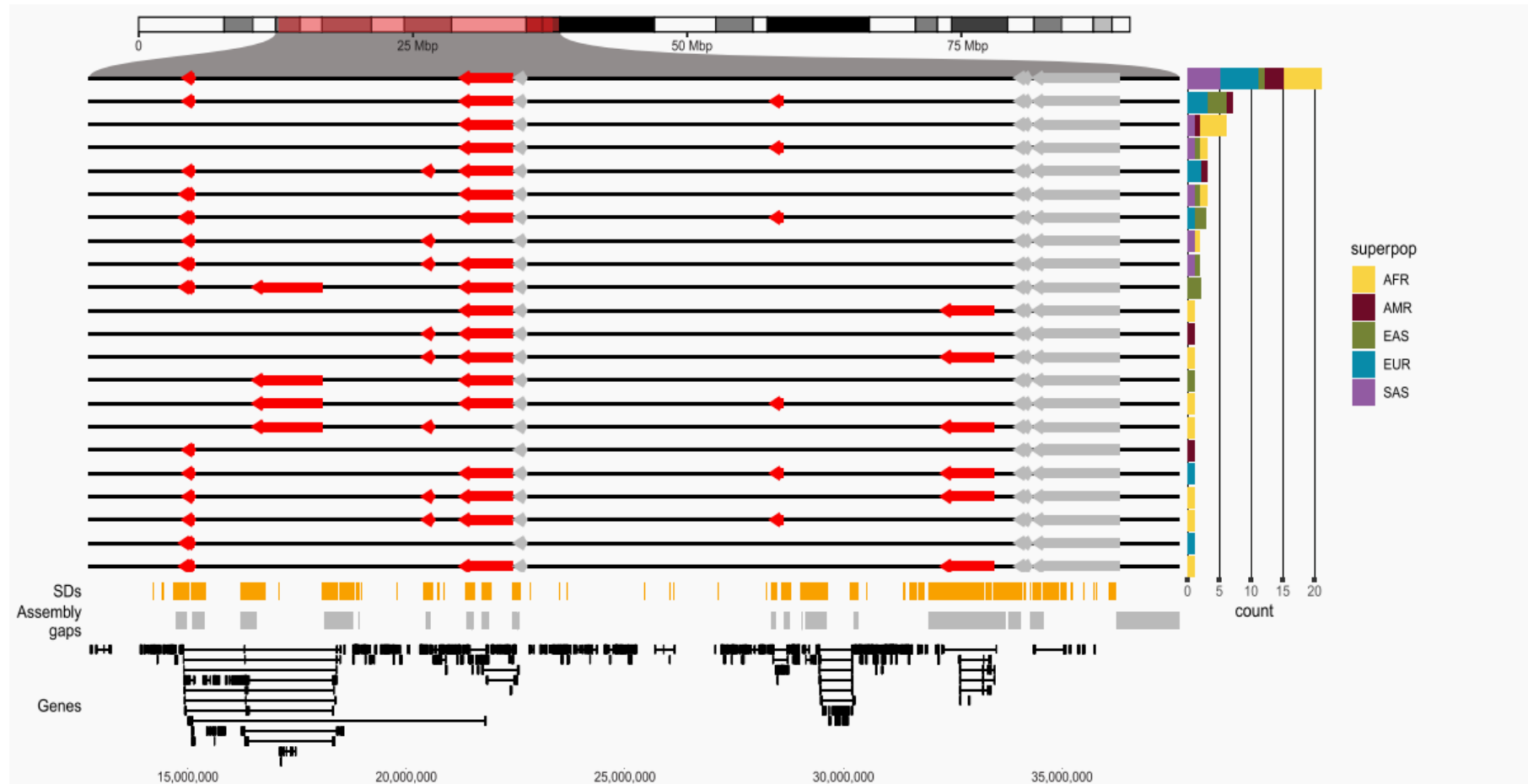
由SVA介导的一系列转导事件



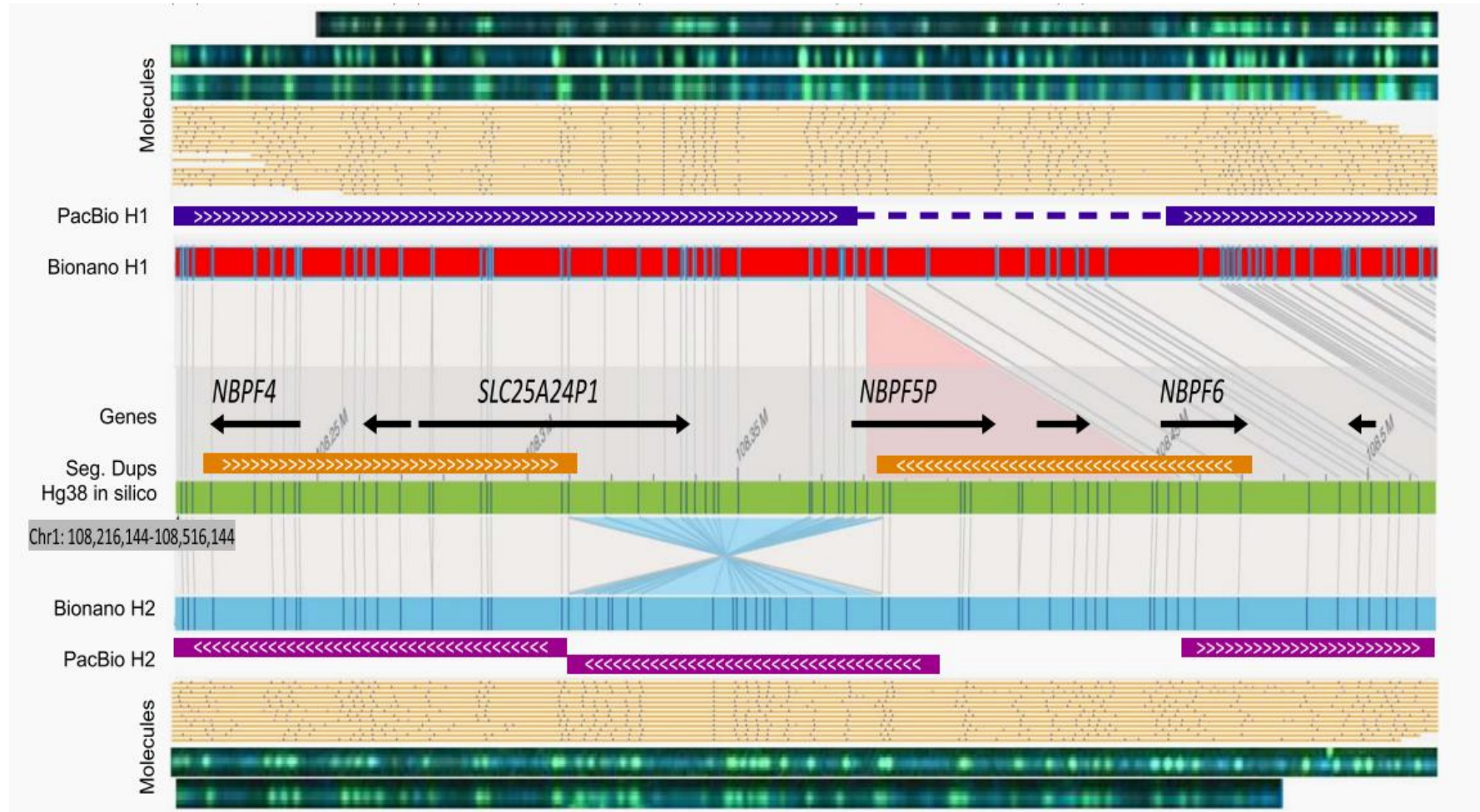
参考元素和非参考SVA元素的VNTR长度的分布



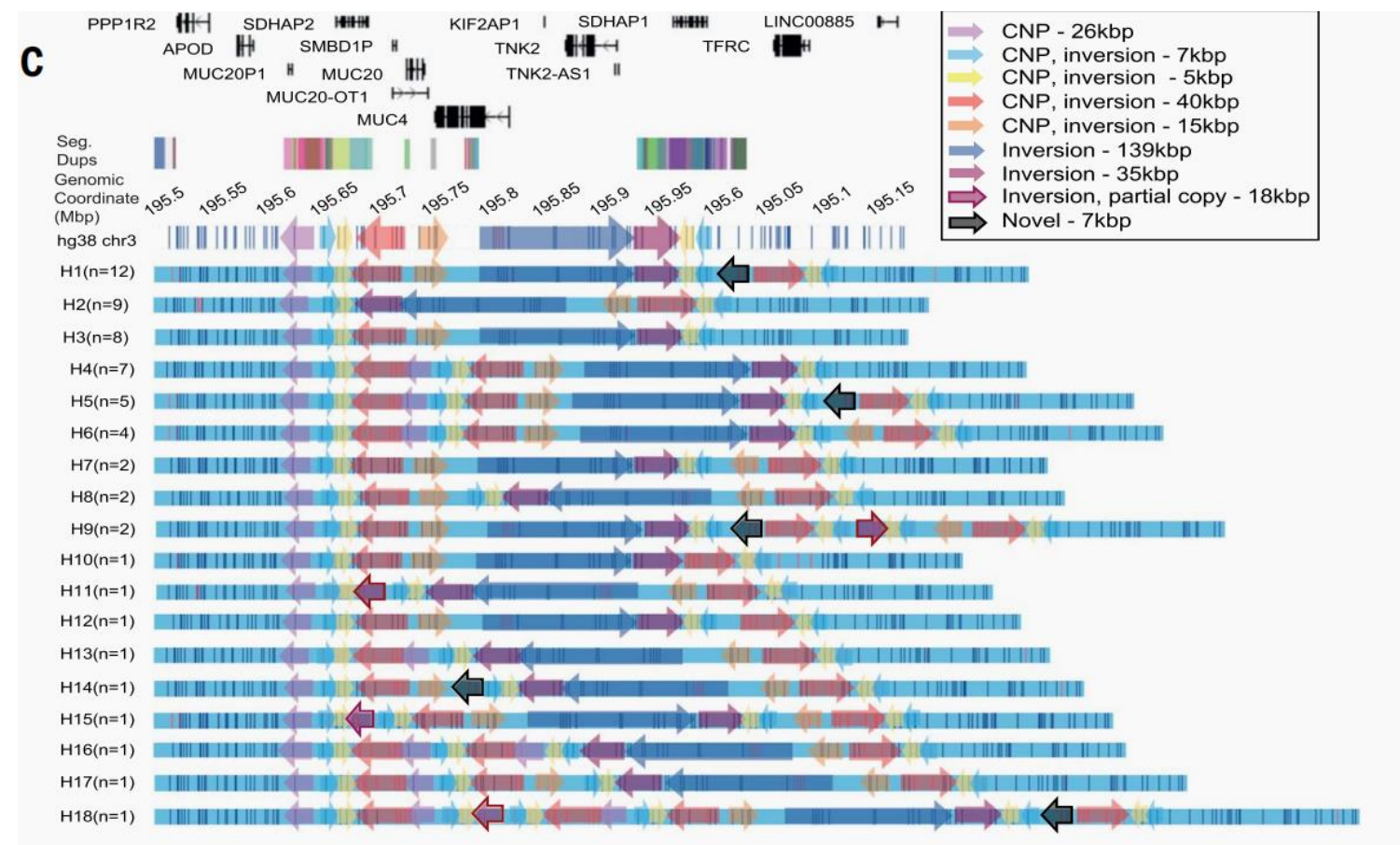
Inversions



Complex structural variation



在染色体3q29处的单倍型结构复杂性



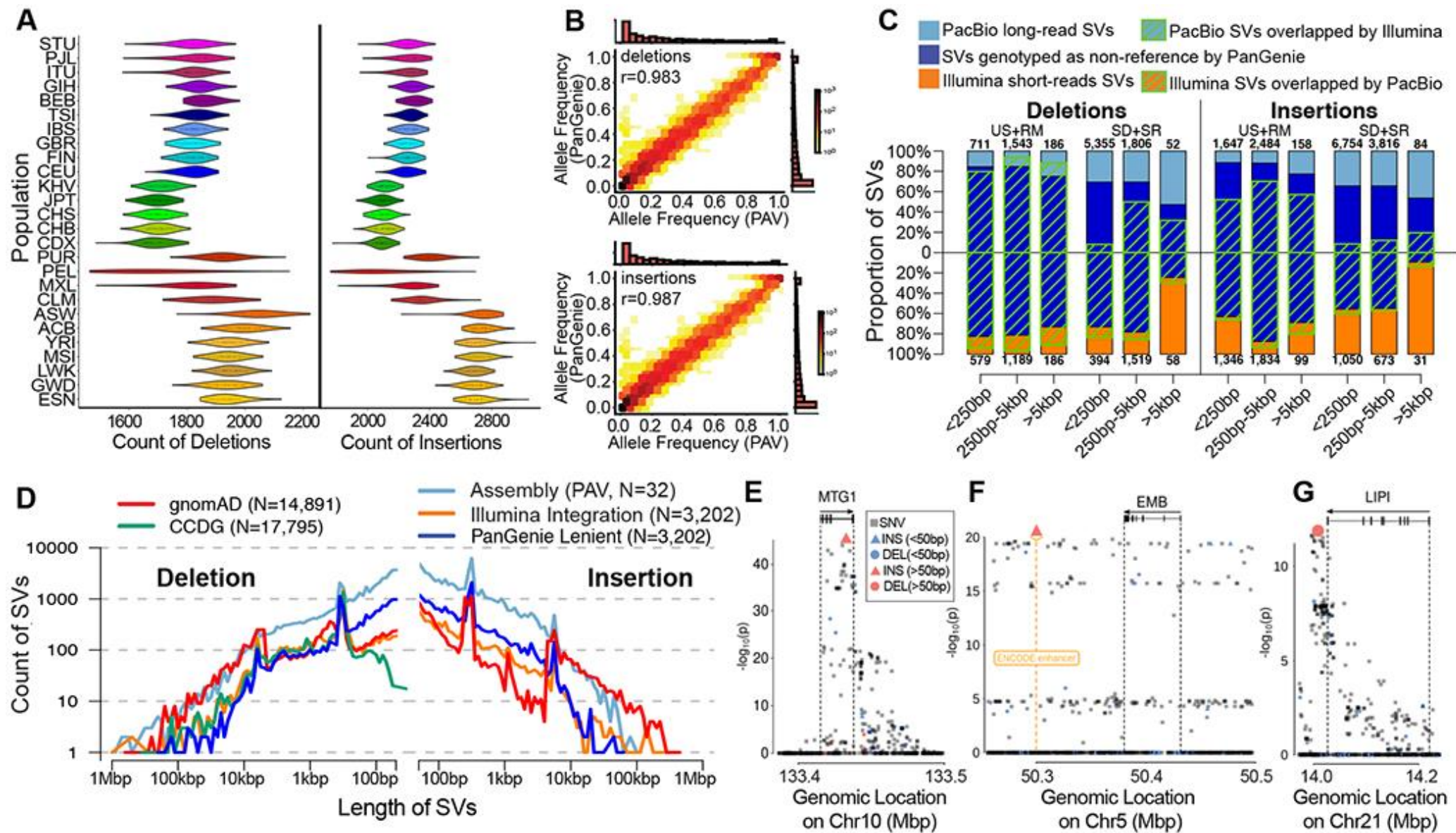
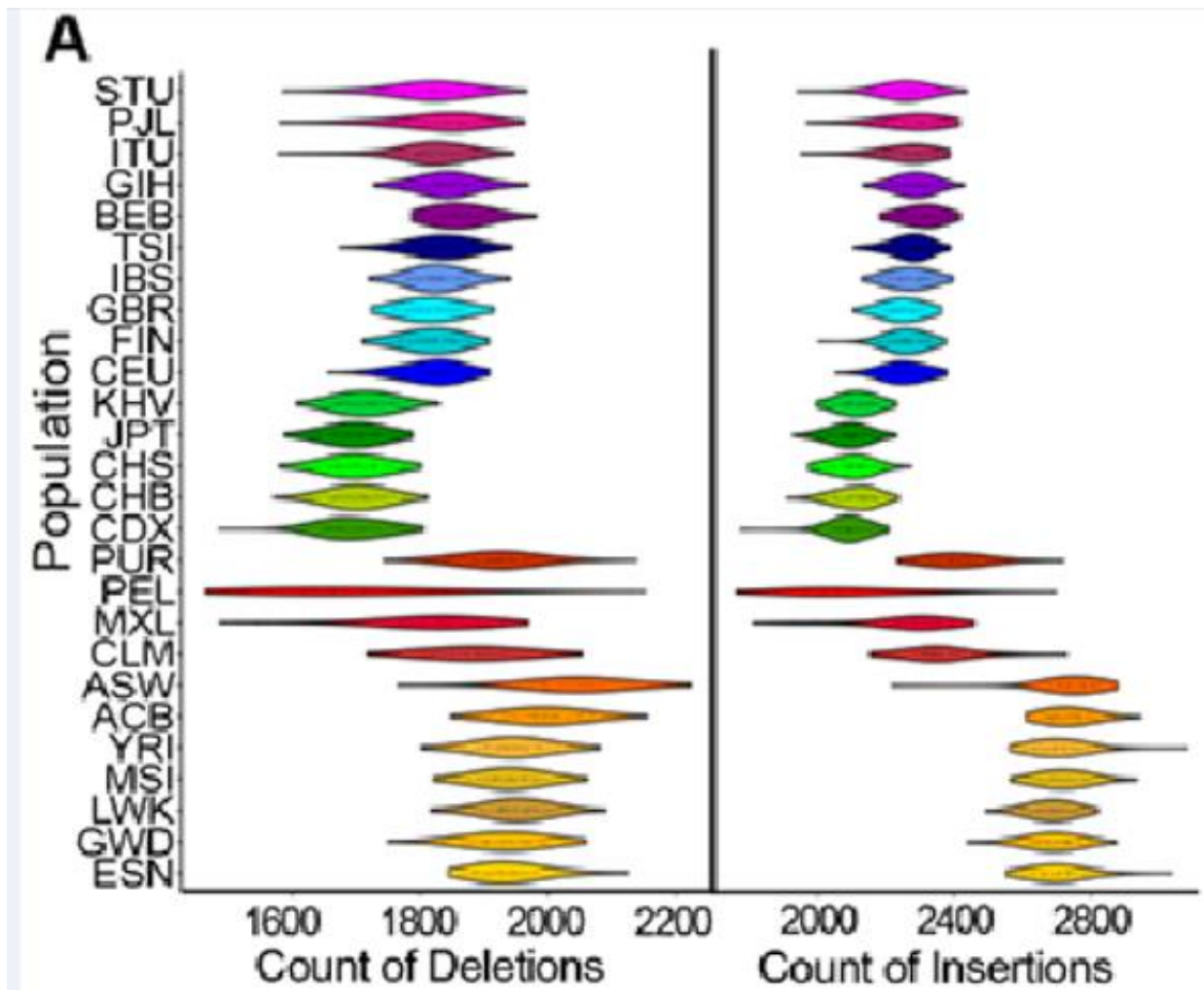


Figure 5: SV genotyping and eQTL analysis

Genotyping



PanGenie: 利用基于单倍型解析的基因组组装构建的泛基因组参考，结合来自原始短读长测序数据的 k-mer 计数信息，对广泛的遗传变异进行基因分型，并提供对短读无法访问的区域的访问。

*使用由15.5M SNVs、1.03M indels (1-49 bp) 和96.1k SVs (其中等位基因脱落率<20%) 组成的参考集进行基因分型步骤，并将这些变体的基因分型纳入1000GP WGS数据集，观察预期的多样性模式。

Figure 5: SV genotyping and eQTL analysis

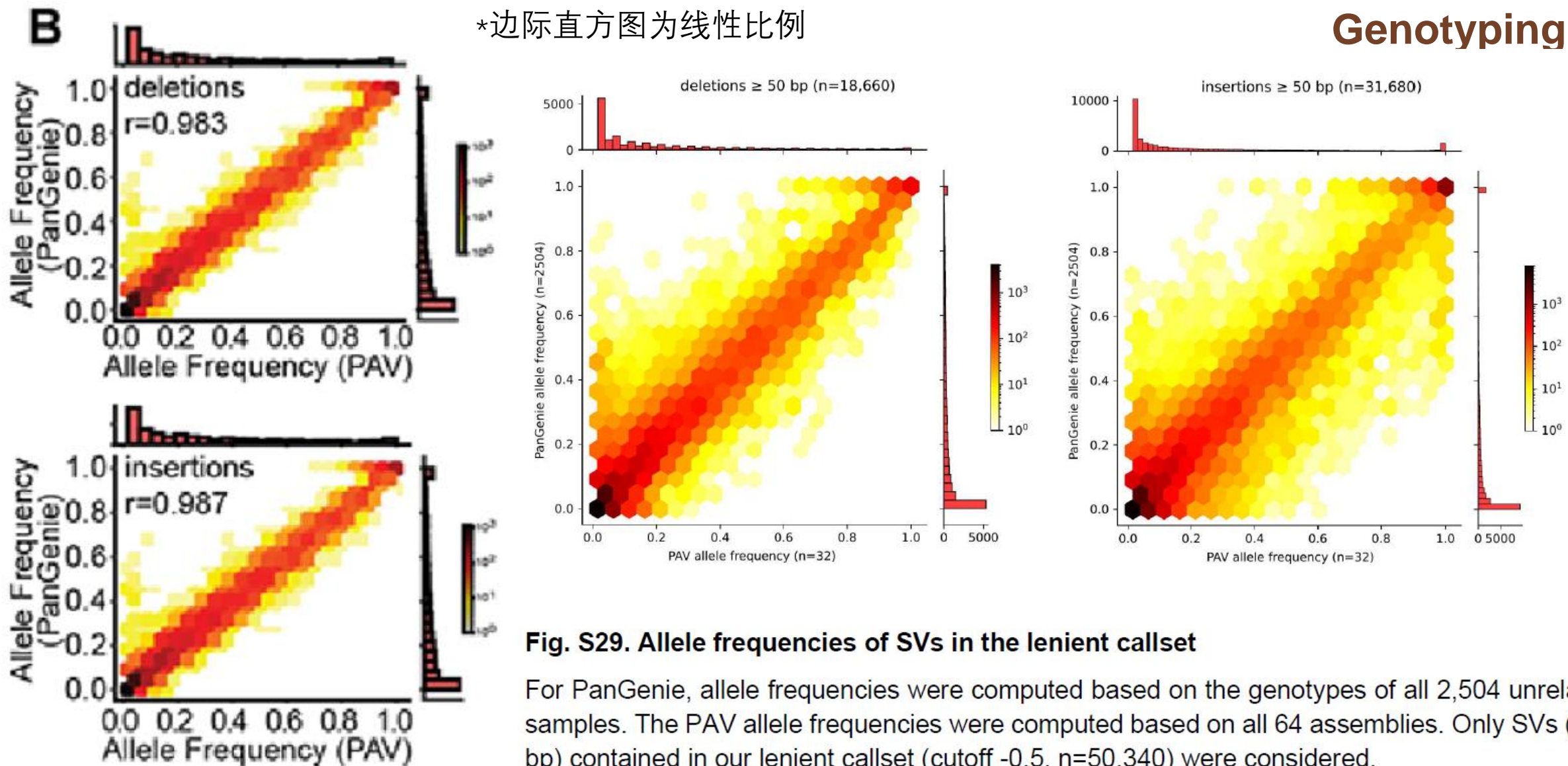


Figure 5: SV genotyping and eQTL analysis

Added value from graph-based genotyping into short read WGS data

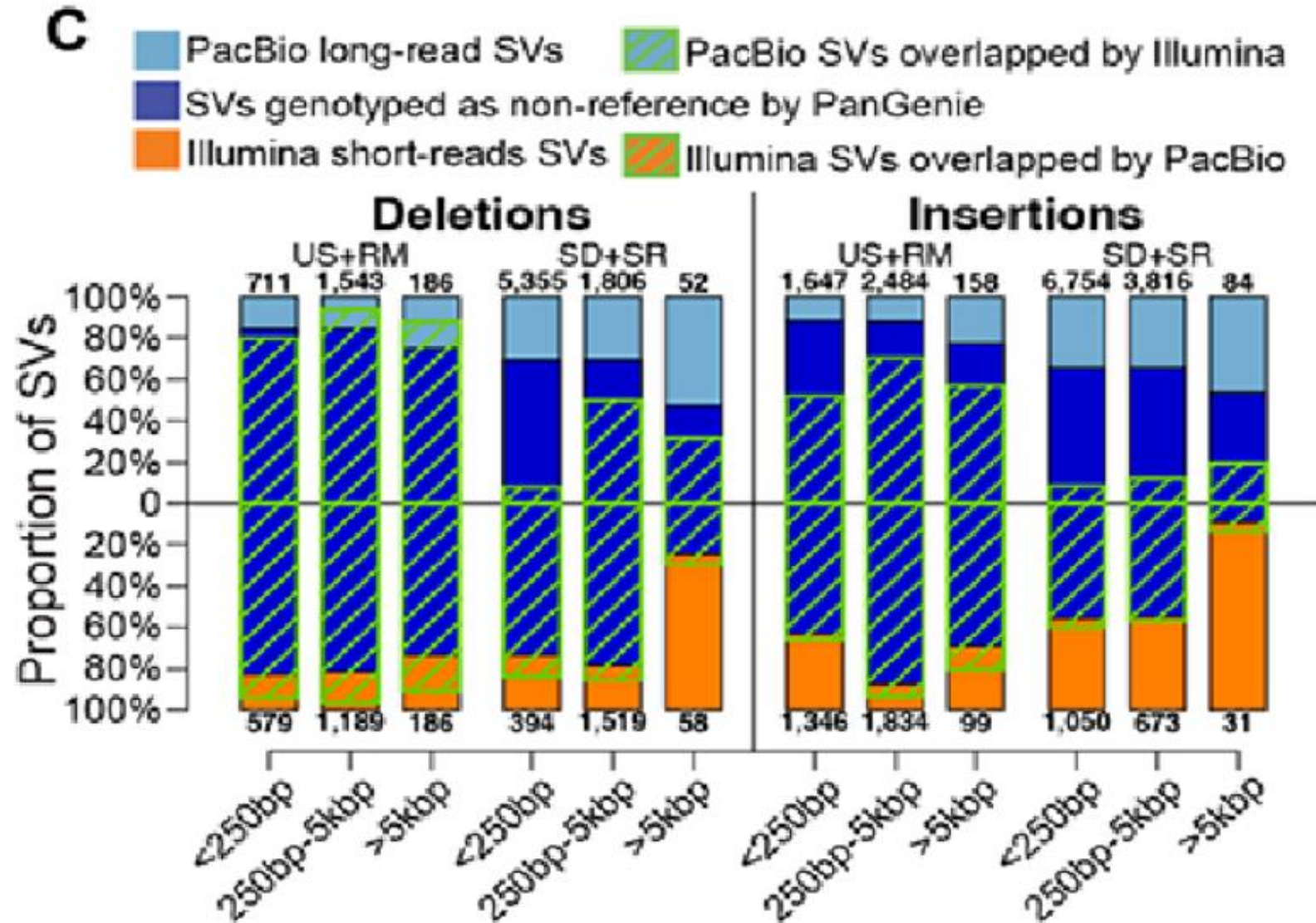
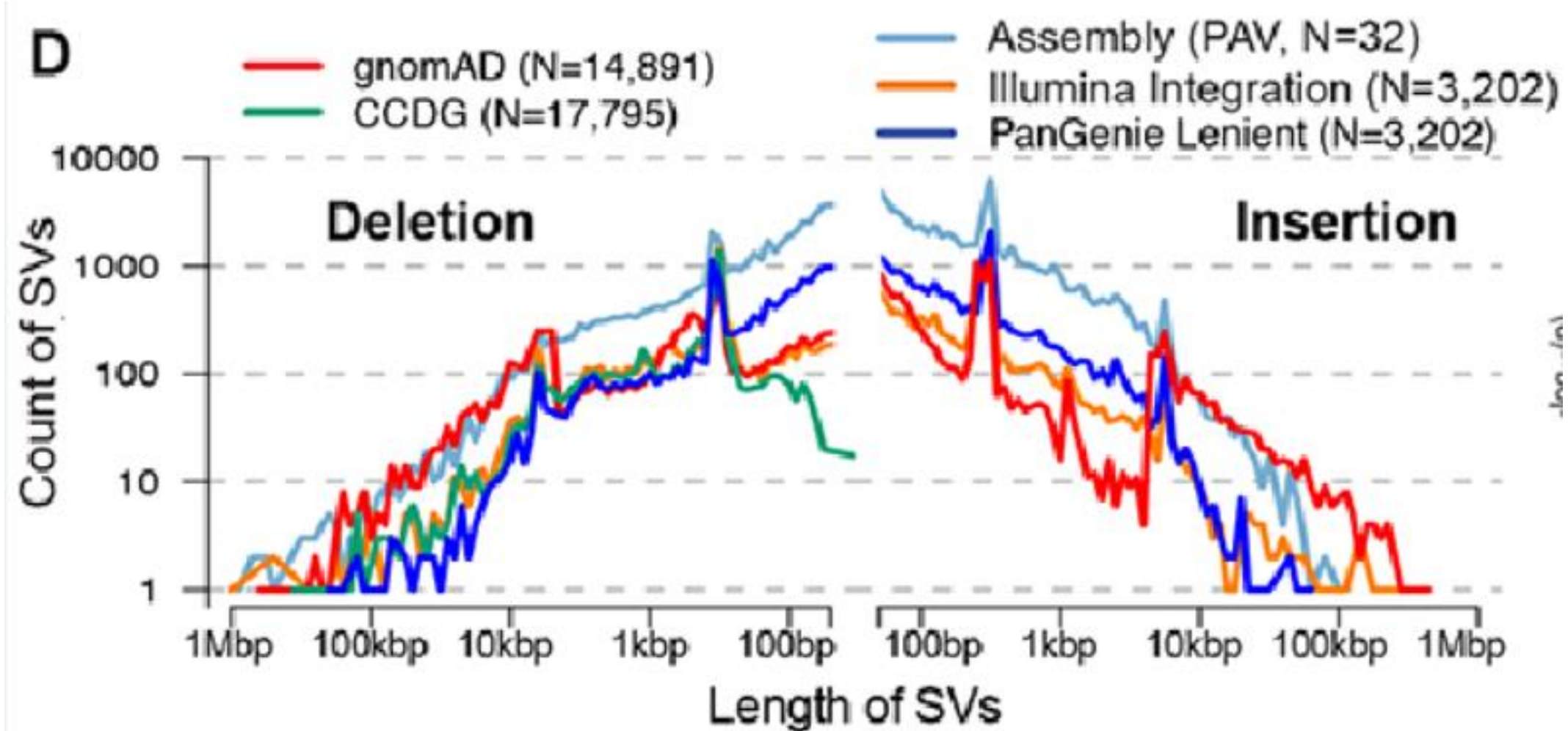


Figure 5: SV genotyping and eQTL analysis

Added value from graph-based genotyping into short read WGS data

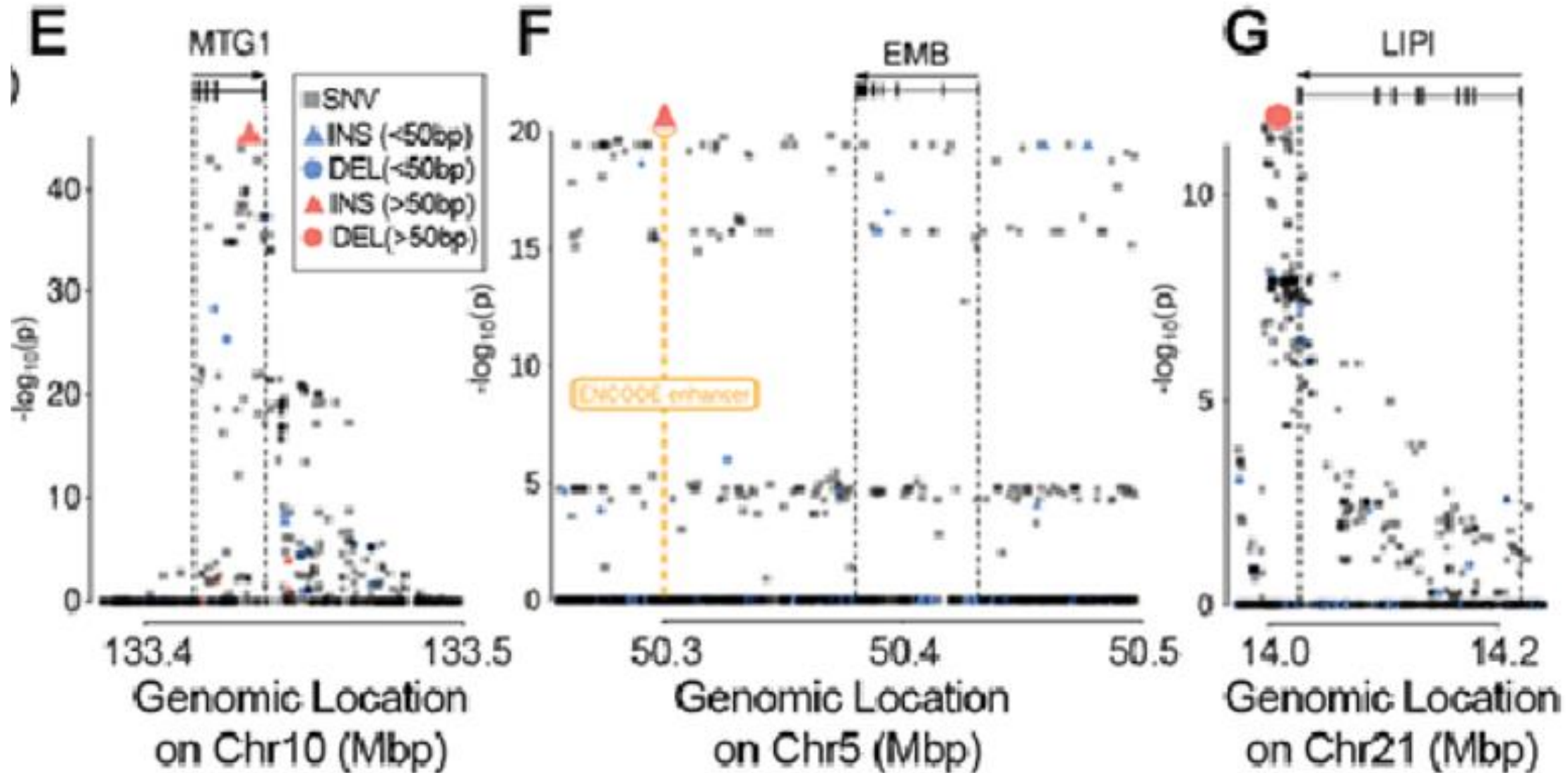


*突变基因数据来自常见疾病基因组学中心 (CCDG) 和基因组聚集数据库 (gnomAD)

Figure 5: SV genotyping and eQTL analysis

*quantitative trait loci (eQTL)

QTL analyses



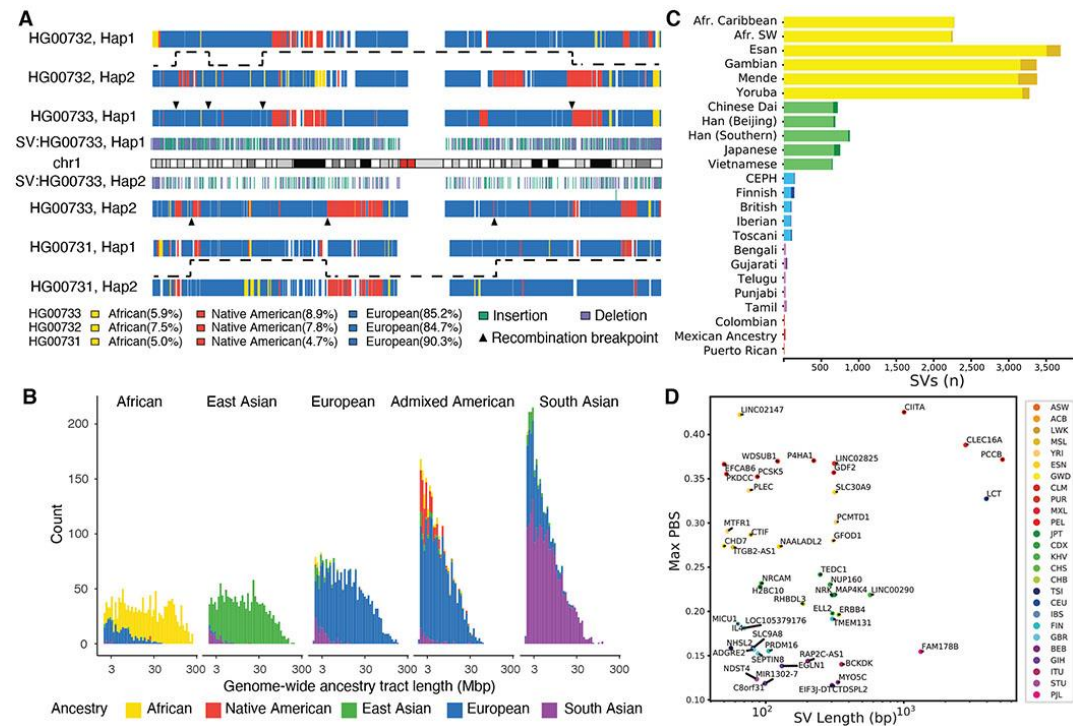


Figure 6

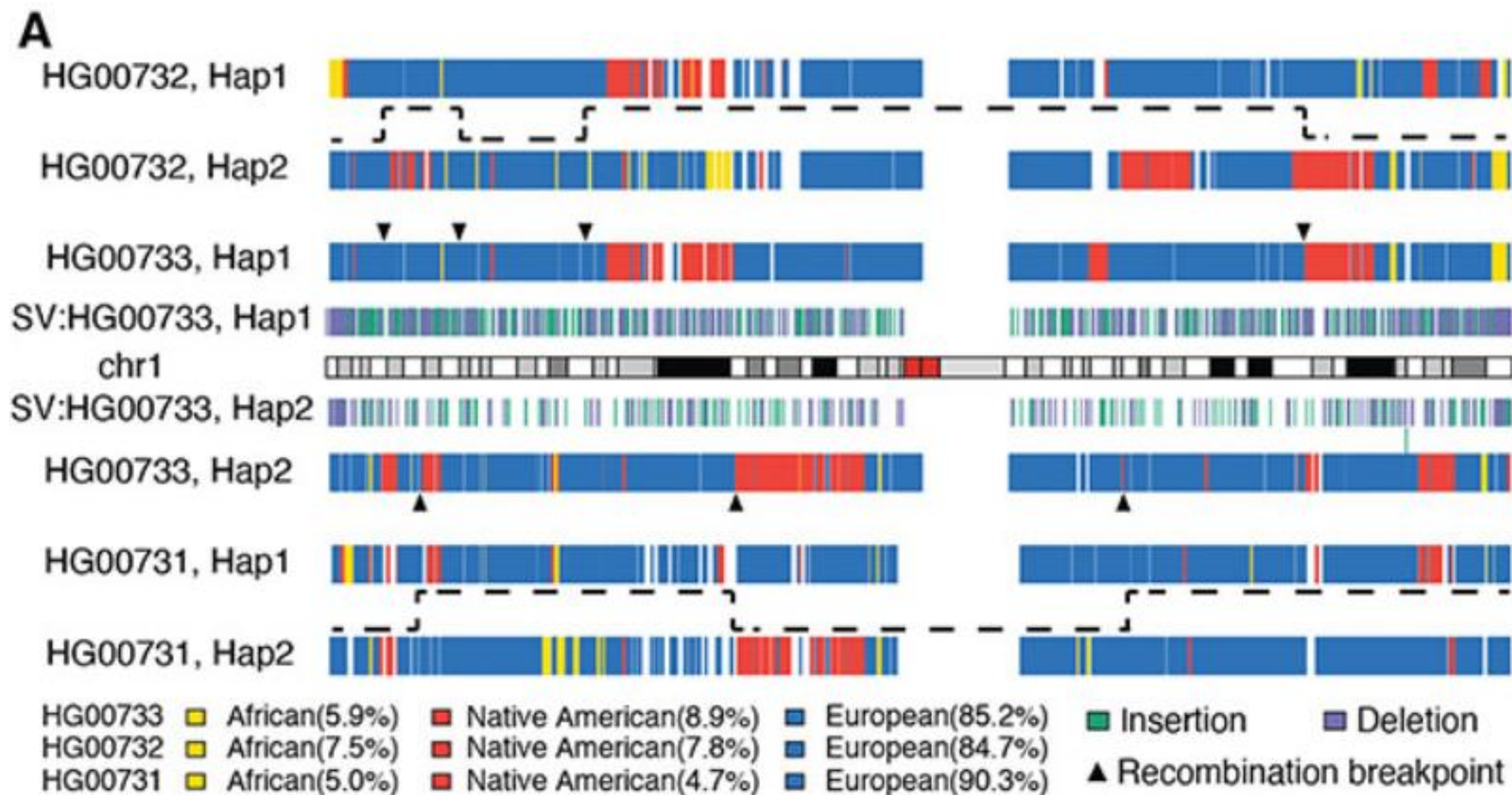


Figure 6

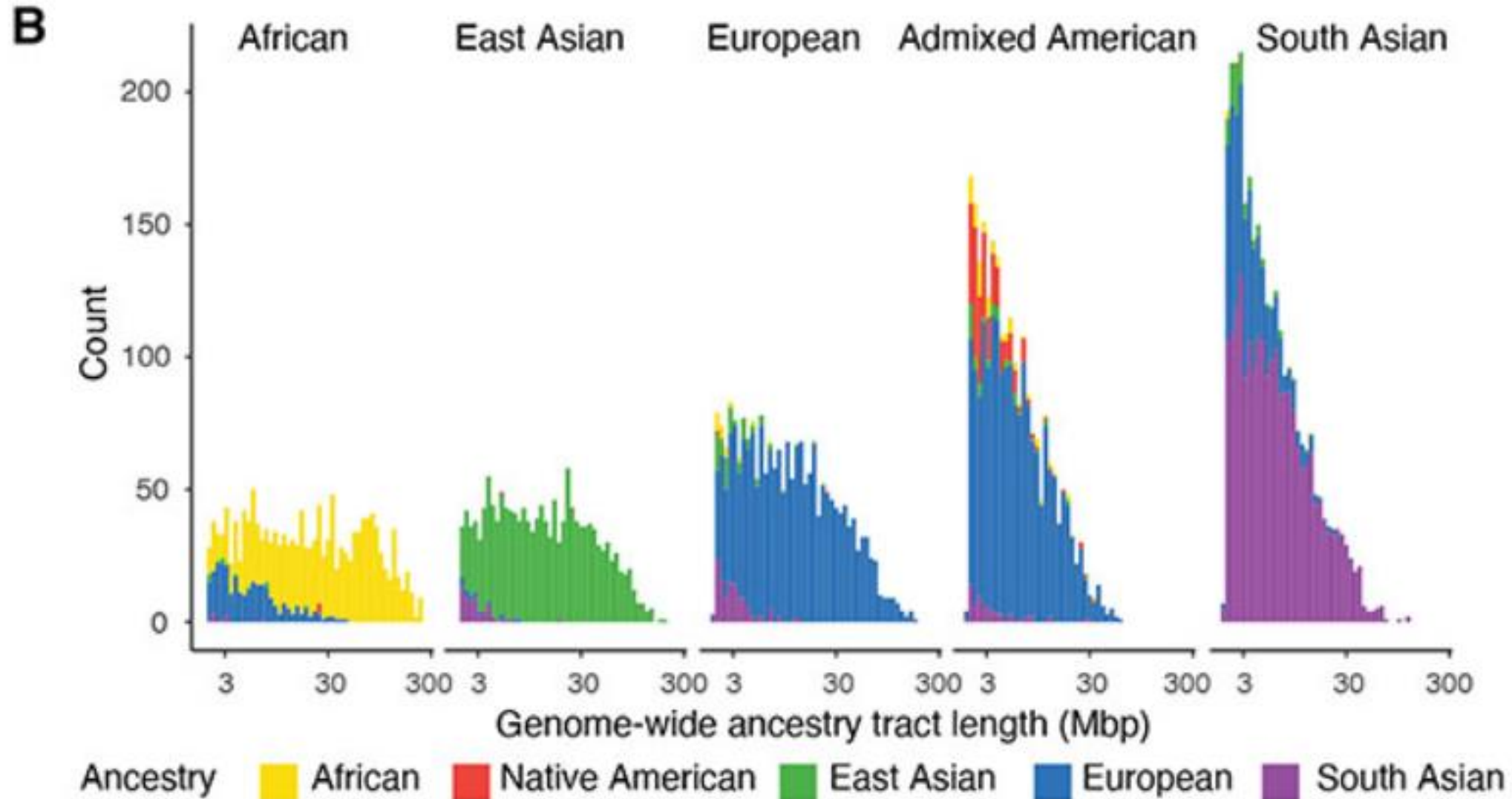


Figure 6

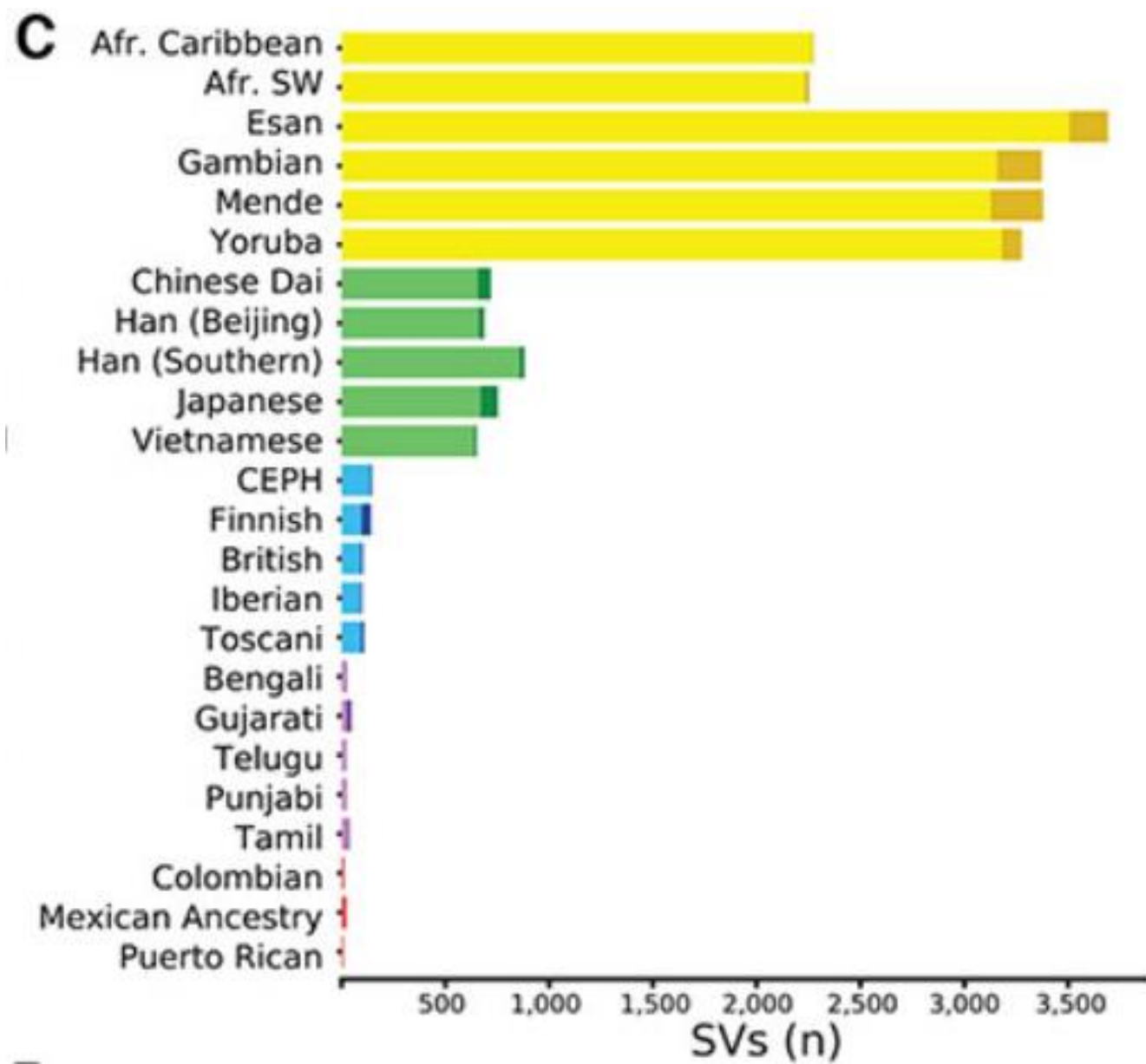
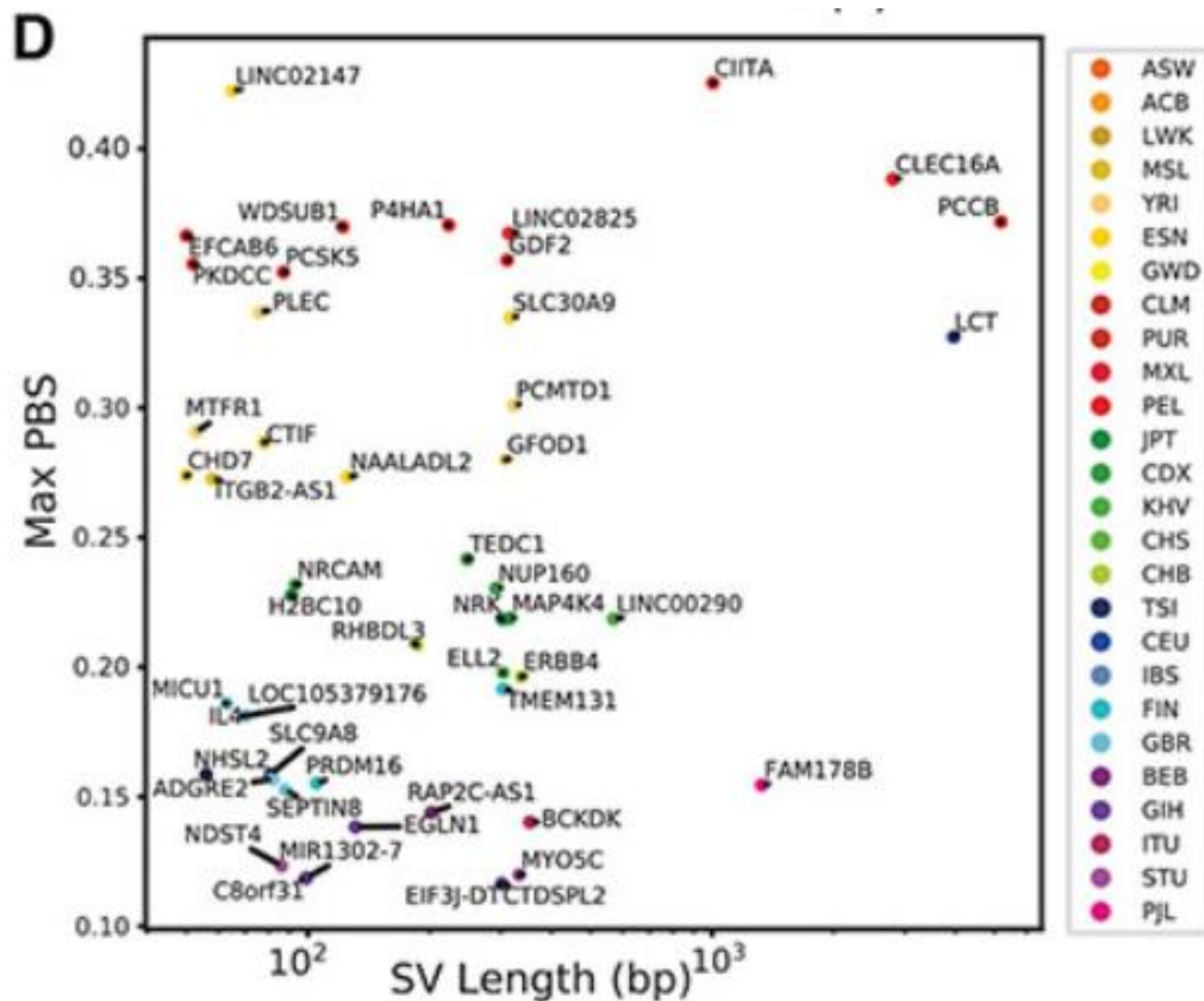


Figure 6



现有成果

- 1.通过比较组装的人类基因组直接发现所有形式的遗传变异（PAV）
- 2.有更多的SVs（63%）被指定为基于同源（>50 bp）的突变机制
3. 全基因组QTL扫描可以弥补分子和临床表型之间的差距，并作为由遗传变异类介导的功能效应的代理。
- 4.具有准确基因型的单倍体解决的SV也将促进SV的进化和种群遗传学研究，包括对反复突变率、种群分层和选择扫荡的估计。

仍需克服的挑战

- 1.长读组装仍未完全解决SD附近和内部更复杂的SV形式
- 2.需要新方法对多平行的VNTRs/STRs以及嵌入在较大重复序列（如SDs和中心粒）中的SVs进行基因分型