

The American Journal of Human Genetics, Volume 111

Supplemental information

**Structural and genetic diversity in the secreted
mucins *MUC5AC* and *MUC5B***

Elizabeth G. Plender, Timofey Prodanov, PingHsun Hsieh, Evangelos Nizamis, William T. Harvey, Arvis Sulovari, Katherine M. Munson, Eli J. Kaufman, Wanda K. O'Neal, Paul N. Valdmanis, Tobias Marschall, Jesse D. Bloom, and Evan E. Eichler

Note S1. Using tandem repeats finder, we discovered an 8-mer VNTR in intron 15 of *MUC5AC* with a canonical motif of TCACCCAC in all human haplotypes. Haplogroup 3 (H3) alleles harbor more copies of the repeat and a lower percent motif identity compared to haplogroup 1 and 2 alleles in humans. STREME (Sensitive, Thorough, Rapid, Enriched Motif Elicitation) analysis reveals five 24-mers that are exclusive to H3 alleles for genotyping this haplogroup (ACCATTCACCTCACCCATTACCCATTACCC, ACTCACCCACTCACCCATTACCCATTAC, ACTCACCCACTCACTCACTCACCTACTCAA, CAGTGGGTGATTGAGTGGGTGAATGGGTGA, and TAAGTTGAGTGAGTGGTGAGTGAGTGGA). The canonical 8-mer motif is exclusive to human haplotypes, with all NHPs featuring motif lengths of 12-28 nucleotides. Total VNTR copy number is also highly variable among the NHPs, leading to differences in total sequence length for intron 15 across haplotypes (Fig. S5).

H1 H2 H3

MUC5AC

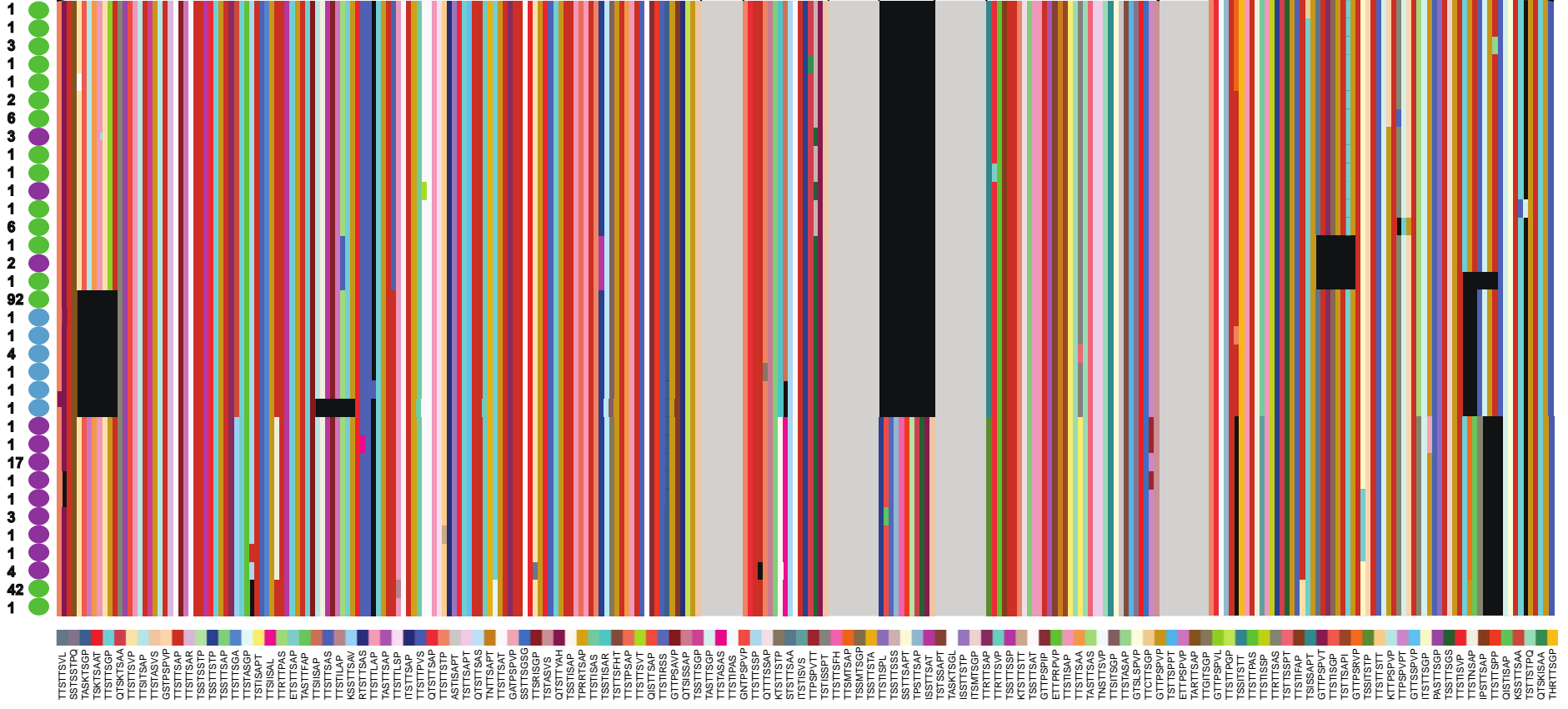


Figure S1. MUC5AC protein VNTR motif utilization across 206 human haplotypes. MUC5AC protein 8-mer variation plots across the five possible VNTR domains. Domain plot height corresponds to the total number of unique alleles across the large central exon. Numbers per unique allele correspond to the number of haplotypes that are matching in motif composition across all domains in linear sequence space. Black boxes indicate an absence of sequence. Gray represents cys domain sequence between distinct VNTR domains. Alleles are ordered vertically by hierarchical clustering of motif composition. Motifs are shifted in sequence space based on the same identity (indicating linear shifts of individual motifs).

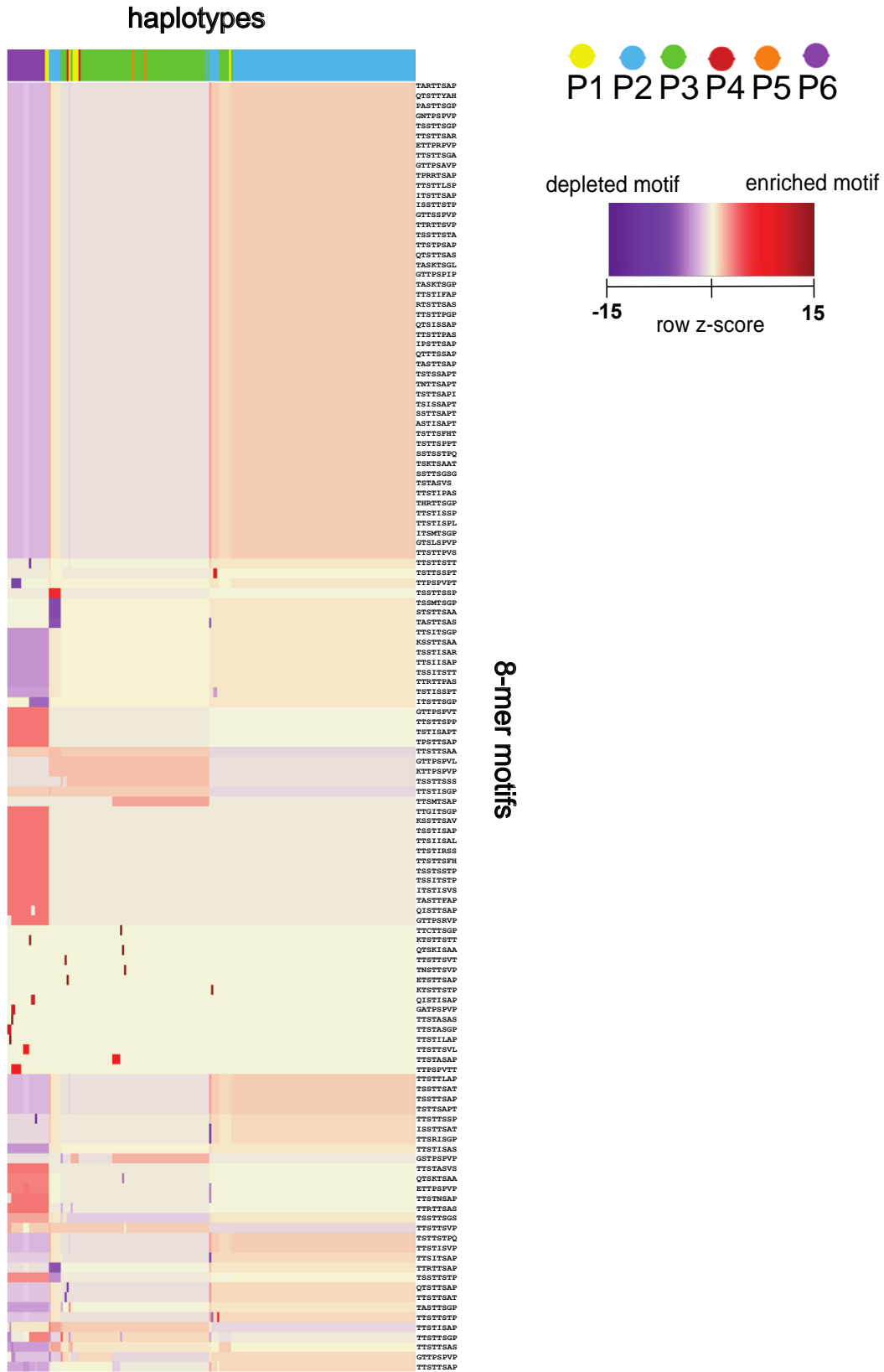


Figure S2. Extended version of Figure 1G.

Figure S3. MUC5B protein VNTR motif utilization across 206 human haplotypes.

MUC5B protein 29-mer variation plots across the five possible VNTR domains. Domain plot height corresponds to the total number of unique alleles across the large central exon. The numbers per unique allele correspond to the number of haplotypes that are matching in motif composition across all domains in linear sequence space. Black boxes indicate an absence of sequence. Gray represents cys domain sequence between distinct VNTR domains. Alleles are ordered vertically by hierarchical clustering of motif composition. Motifs are shifted in sequence space based on the same identity (indicating linear shifts of individual motifs).

haplotypes

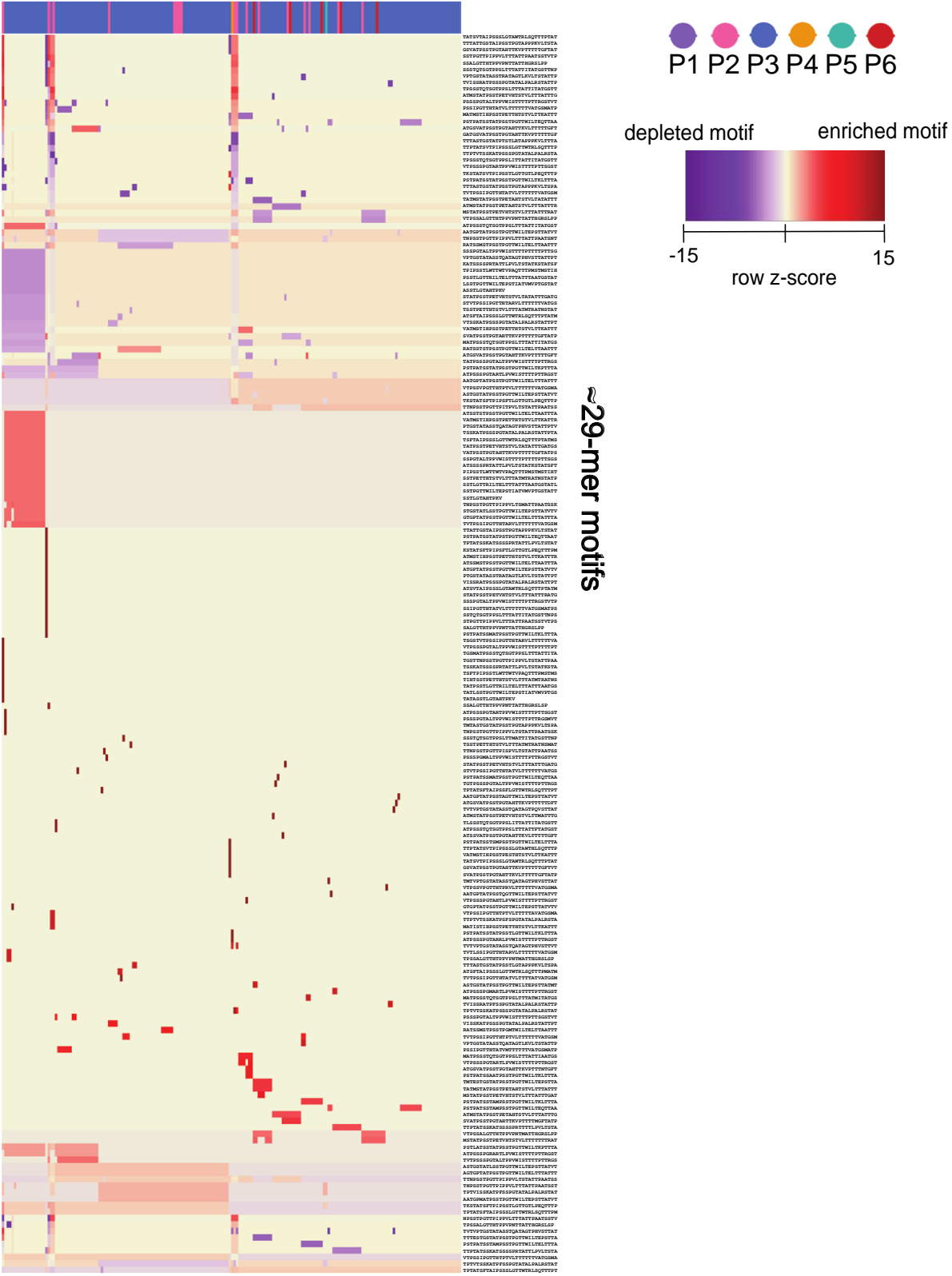


Figure S4. Extended version of Figure 2G.

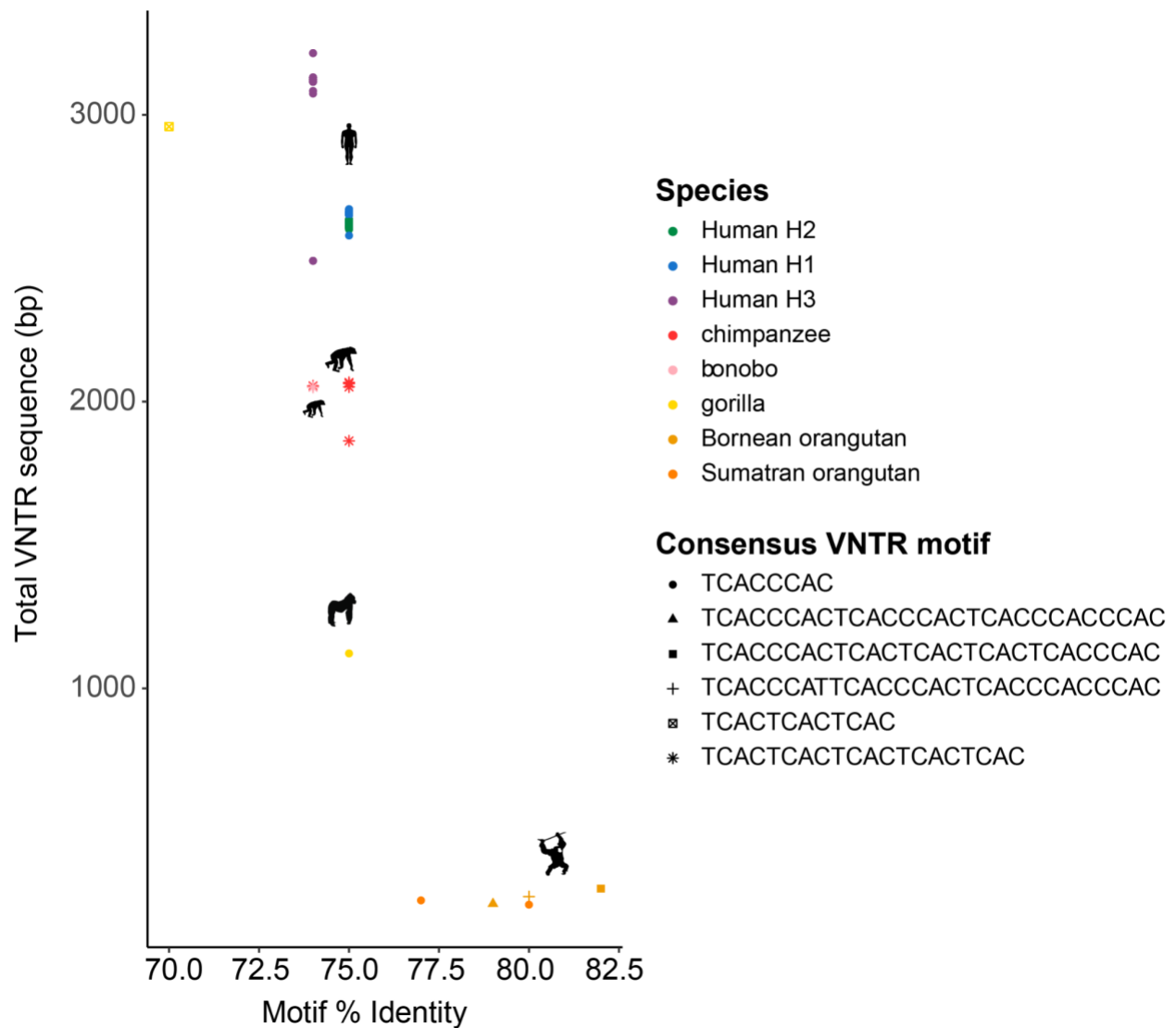


Figure S5. Intron 15 VNTR sequence characterization in *MUC5AC* for humans and nonhuman apes. Comparison of total VNTR sequence length (in base pairs) and canonical motif percent identity across 206 human haplotypes and at least two haplotypes per nonhuman ape species. Species/haplogroup identity coded by color and shape coded by motif.

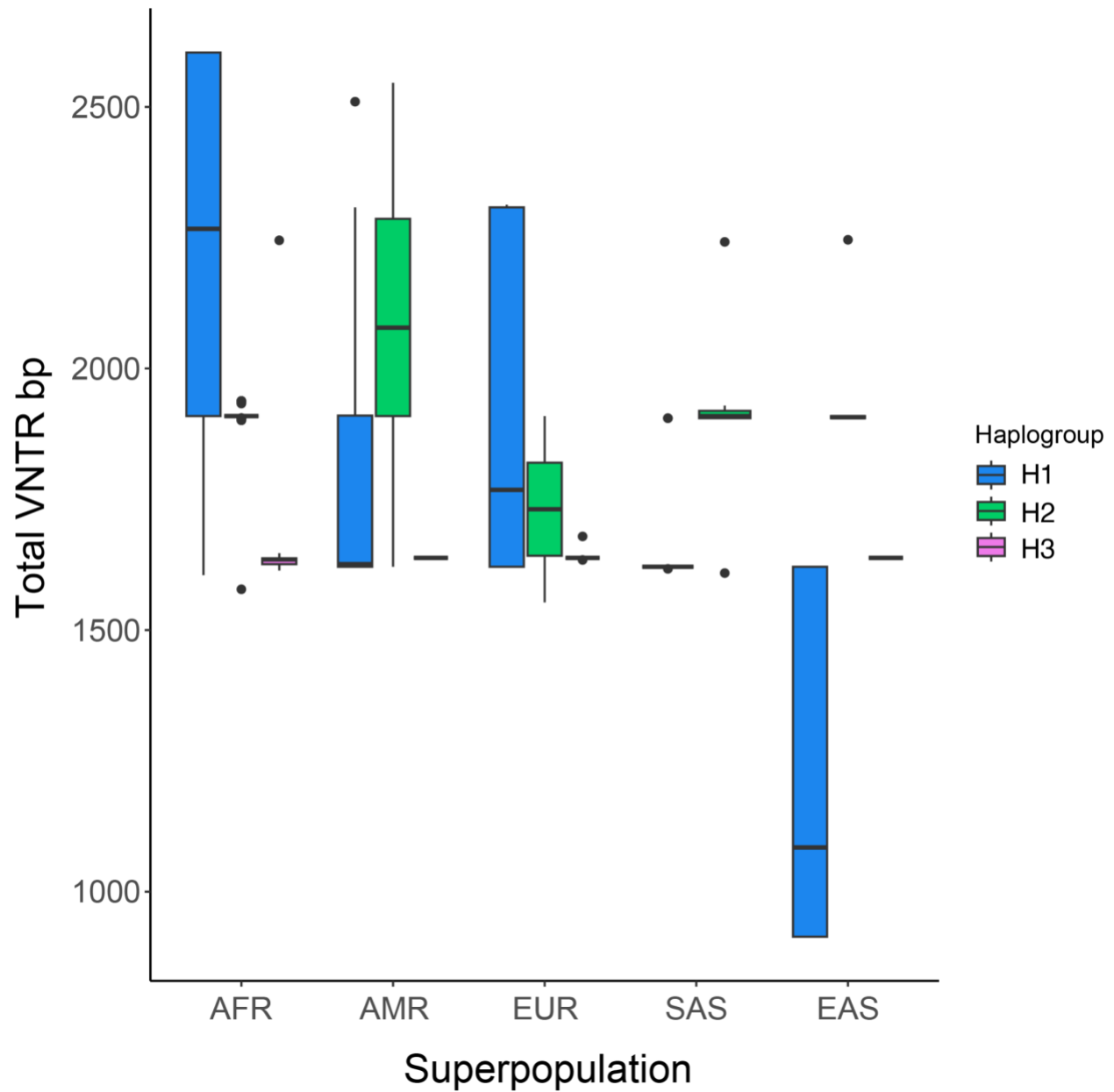


Figure S6. *MUC5AC* VNTR enhancer polymorphism across human super populations in the HGSCV/HPRC sample set, stratified by haplogroup identity.

Total VNTR bp indicates the longest continuous track of degenerate 8-mer VNTR sequence in the enhancer region of *MUC5AC*. Haplogroups correspond to phylogenetic clades of the *MUC5AC* locus.

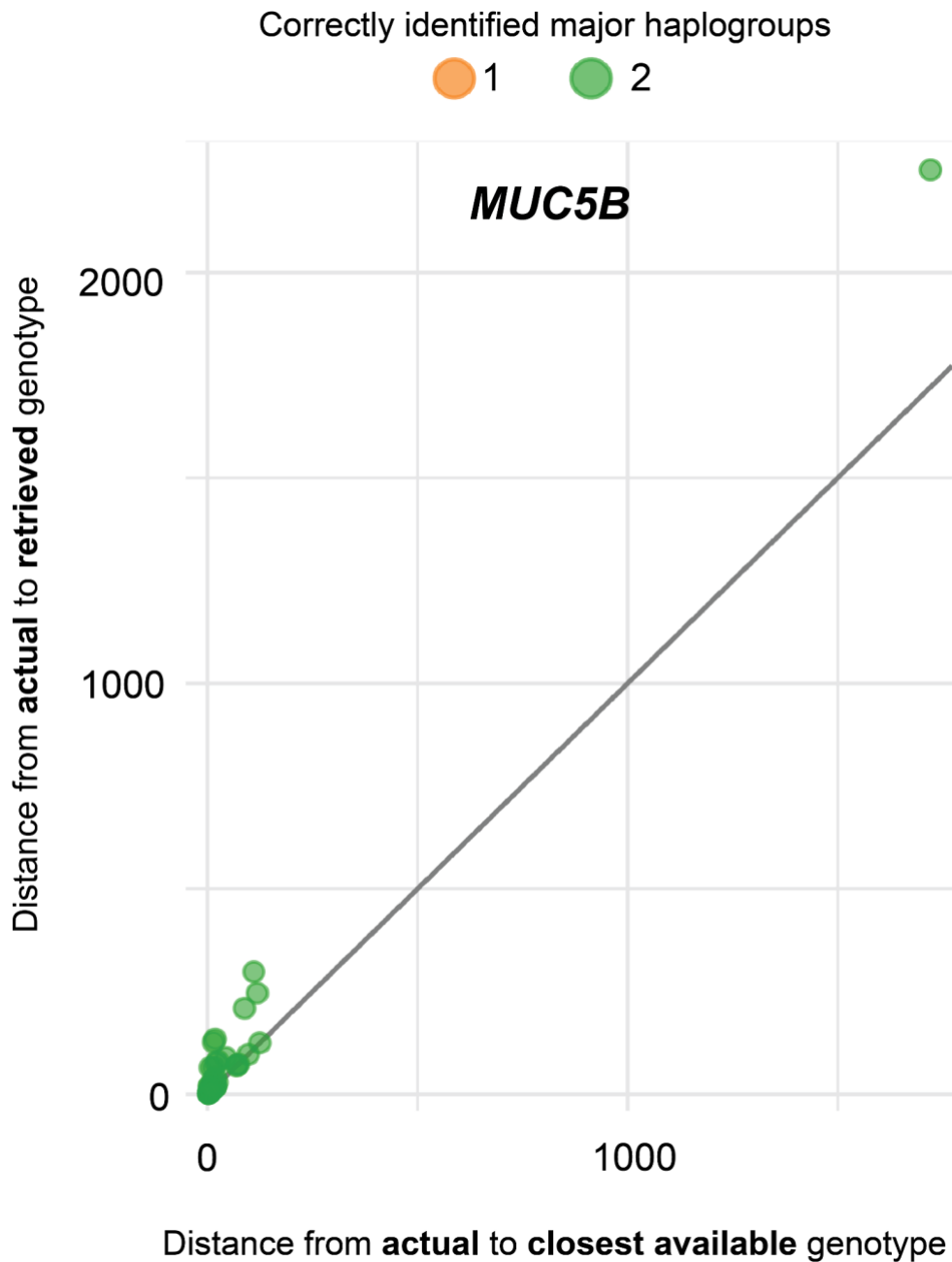


Figure S7. Genotyping accuracy of *MUC5B* haplogroups with Locityper. Locityper leave-one-out results comparing edit distances between actual and retrieved genotype (predicted from genotyper) versus edit distances between actual and closest possible genotype (best possible reference genotype from multiple sequence alignment with true genotype) for *MUC5B*. Dot color corresponds to the number of haplotypes in diploid sample sets that were correctly genotyped.

A**MUC5B promoter polymorphism, all samples: rs35705950**

▲ Risk
▼ Protective
— Bonferroni correction

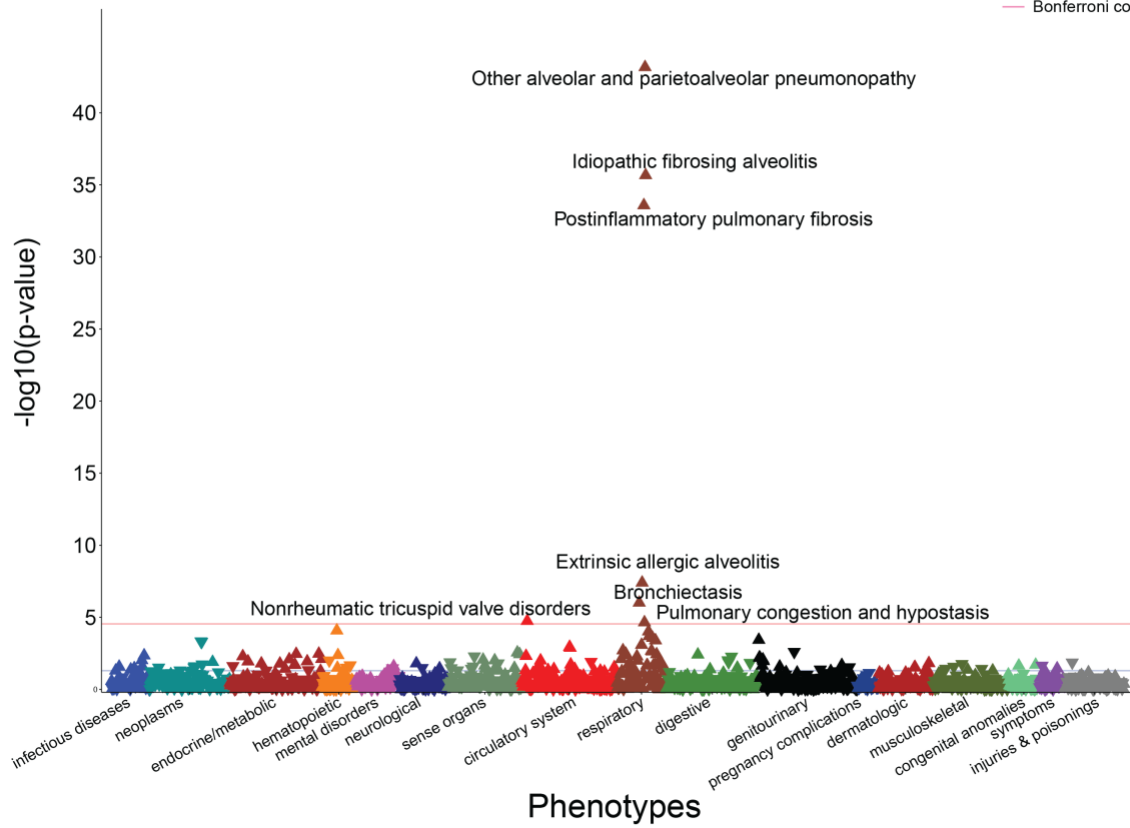
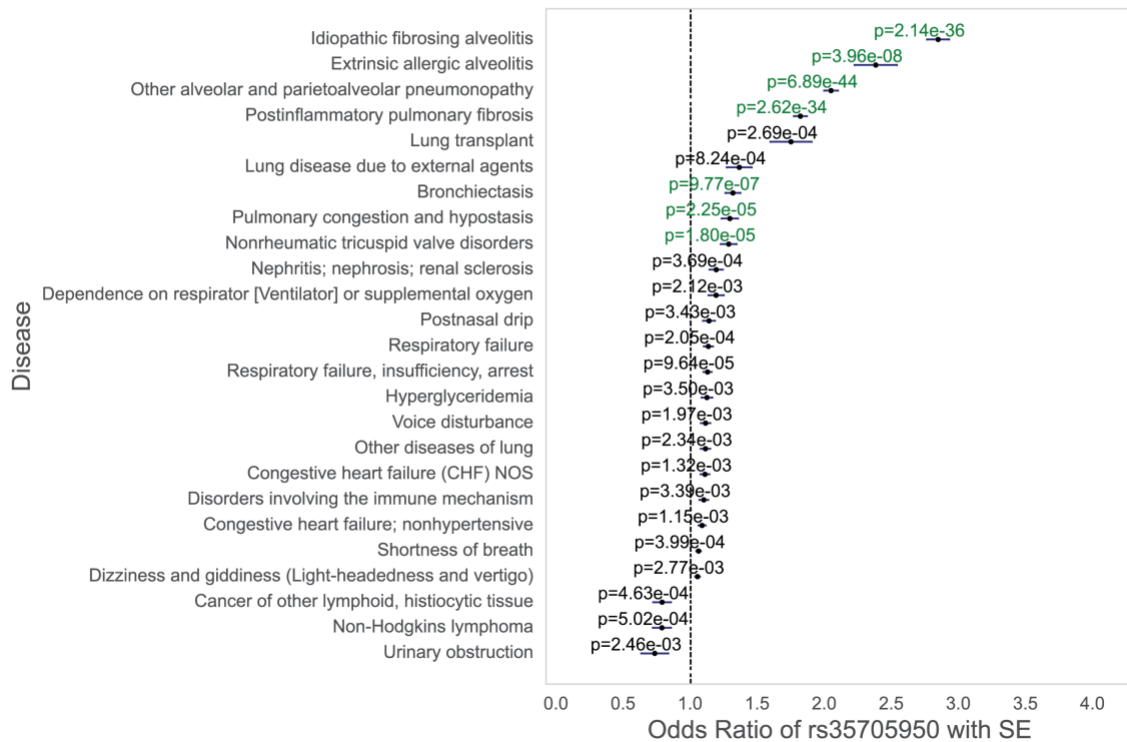
**B**

Figure S8. PheWAS of *MUC5B* promoter polymorphism rs35705950 in *All of Us*.

- A. Phenome-wide association study (PheWAS) Manhattan Plot of phenotype groupings and associations with the *MUC5B* promoter polymorphism rs35705950 in all populations (i.e., all samples). Only phenotypes that passed Bonferroni correction are noted by name.
- B. Odds ratios of the top 25 diseases (based on ordered p-values) for associations with rs35705950. Green indicates phenotypes that passed Bonferroni correction. SE = standard error (horizontal bars corresponding to each p-value).