

Structural variant discovery and characterization from *de novo* assembly of Khoë-Sān genomes

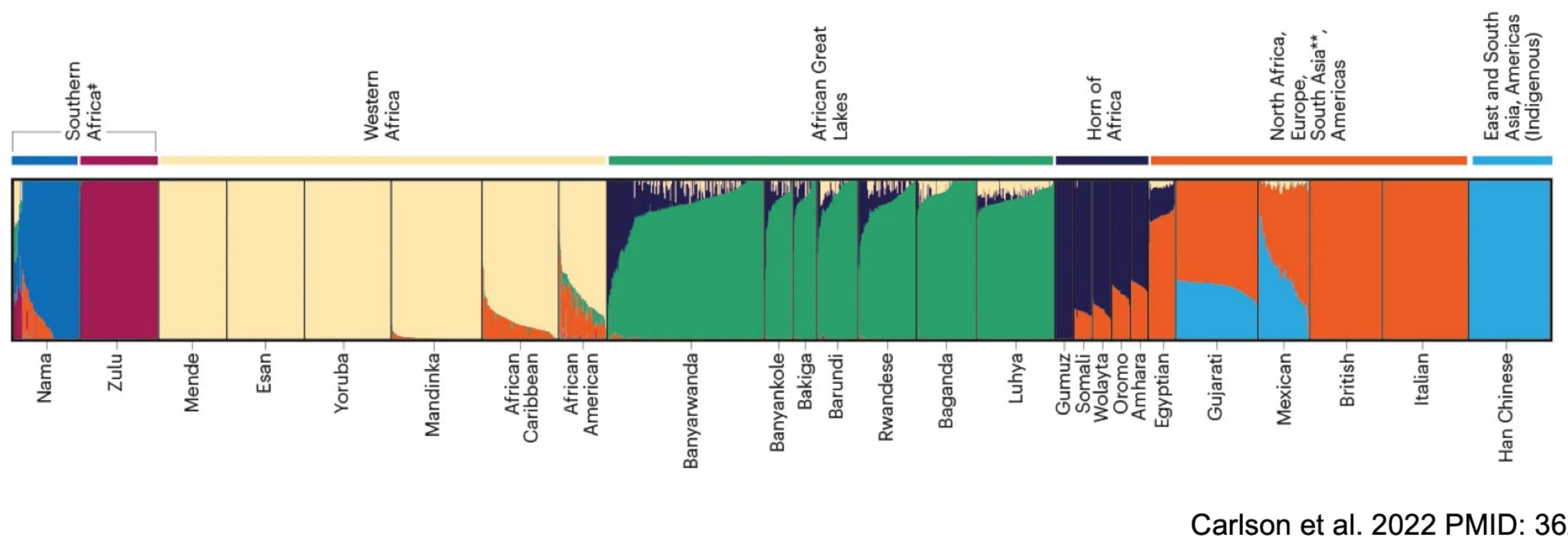
Zoeb N. Jamal^{1*}, Daniela C. Soto^{1*}, Kristin Hardy^{1*}, Mohamed Abuelanin¹, William Palmer¹, Javier Prado-Martinez², Paul Norman³, Marlo Moller⁴, Brenna M. Henn¹, Megan Y. Dennis¹

¹Genome Center, University of California, Davis, CA. ²Institute of Evolutionary Biology, PRBB, Barcelona, Spain.

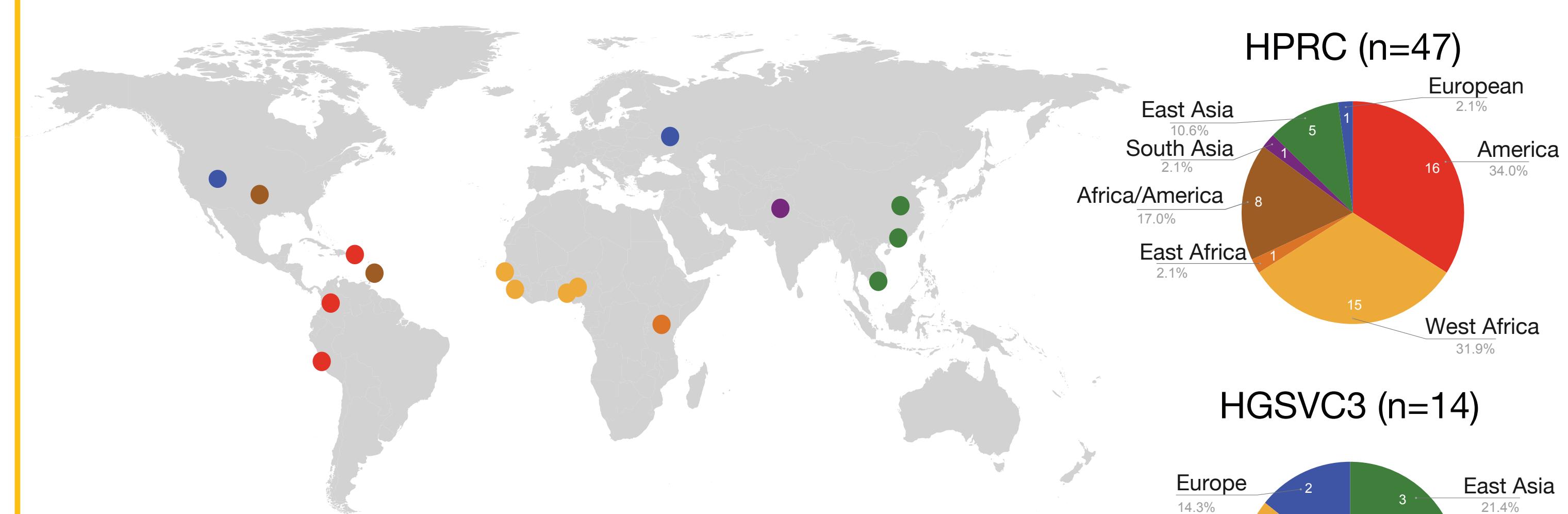
³University of Colorado, Biomedical Informatics, Aurora, CO. ⁴Stellenbosch University, Molecular Biology and Human Genetics, Stellenbosch, South Africa.

Background

- Long read sequencing (LRS) of large cohorts is enhancing our understanding of the genomic variation landscape
- Genetic diversity within African populations can exceed that between continental groups, yet many remain underrepresented



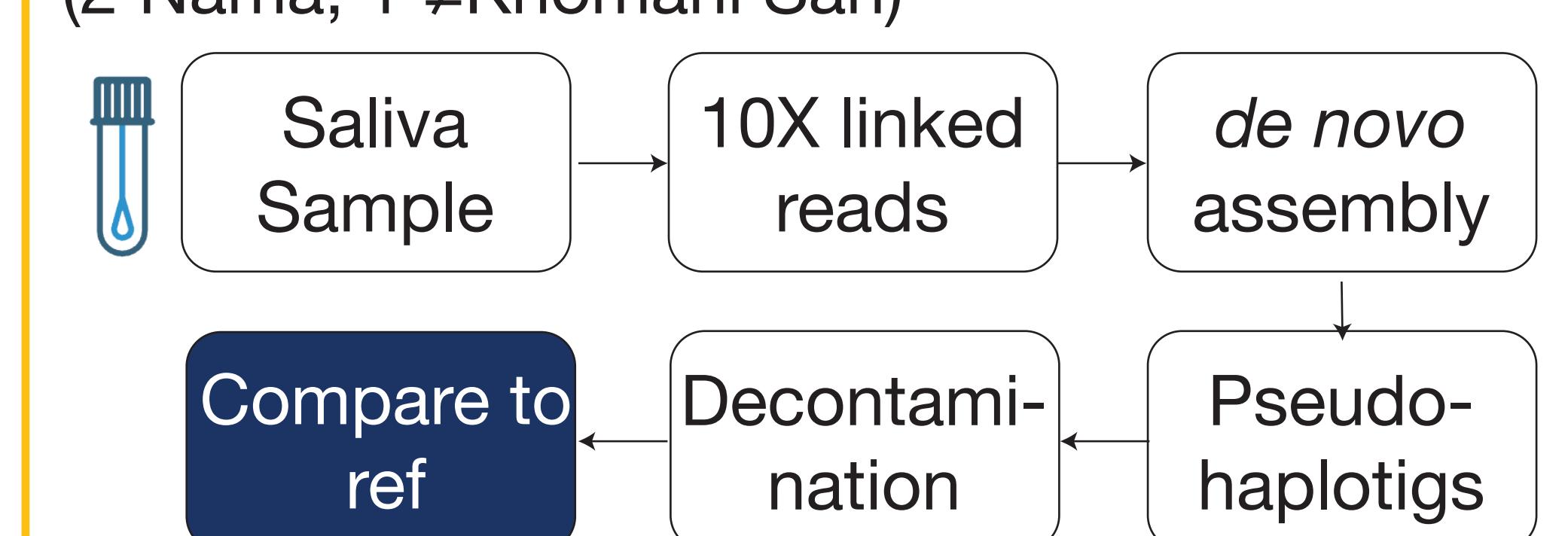
- This gap suggests ample genetic diversity is missing in large datasets



- In particular, the Khoë-Sān indigenous peoples of Southern Africa carry some of the most divergent haplotypes in extant human lineages

Methods

- We performed 10X genomics long range sequencing and *de novo* assembly for 3 Khoë-Sān individuals (2 Nama, 1 ≠Khomani San)



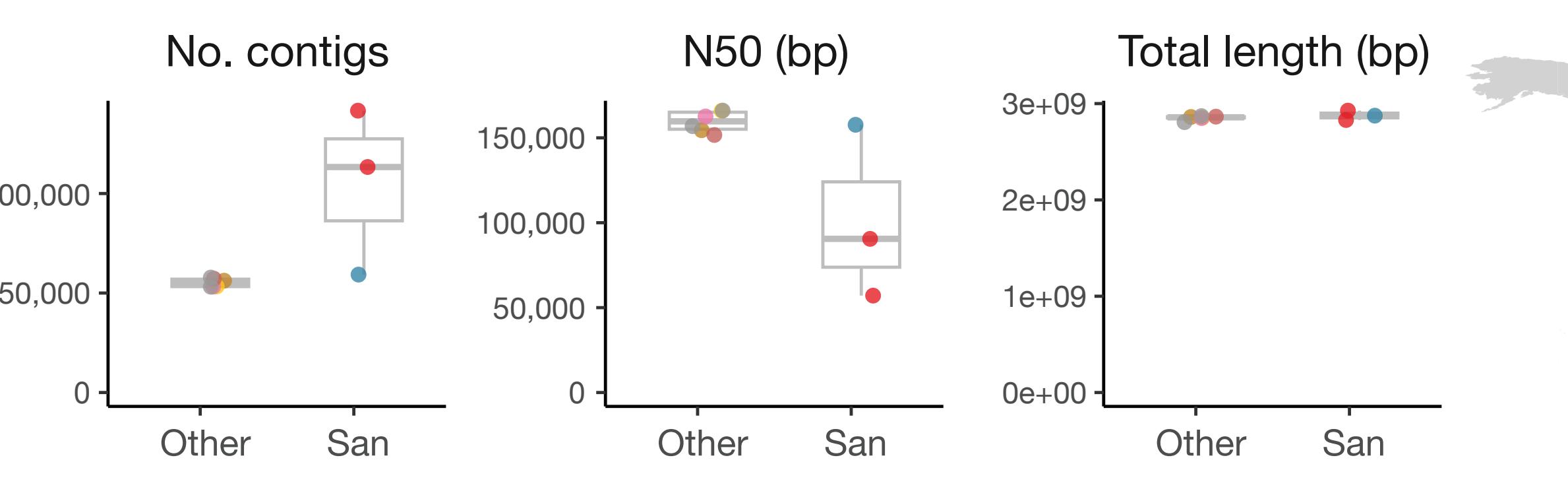
- We called structural variants (SVs) against T2T-CHM13 for these Khoë-Sān genomes and other 10X linked-read assemblies
- Used whole genome short read sequencing (SRS) data from a cohort of 93 Khoë-Sān individuals to genotype discovered SVs and investigate copy number variability
- Used VEP to query the impact of genotyped SVs

References

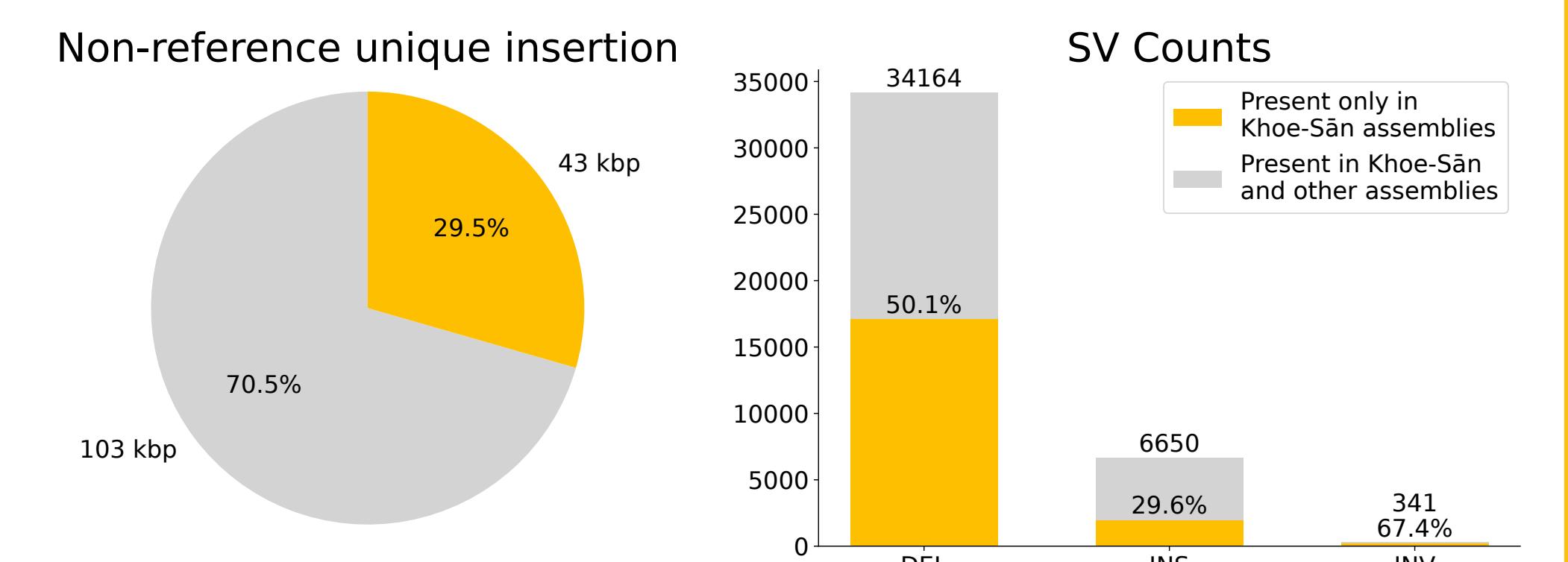
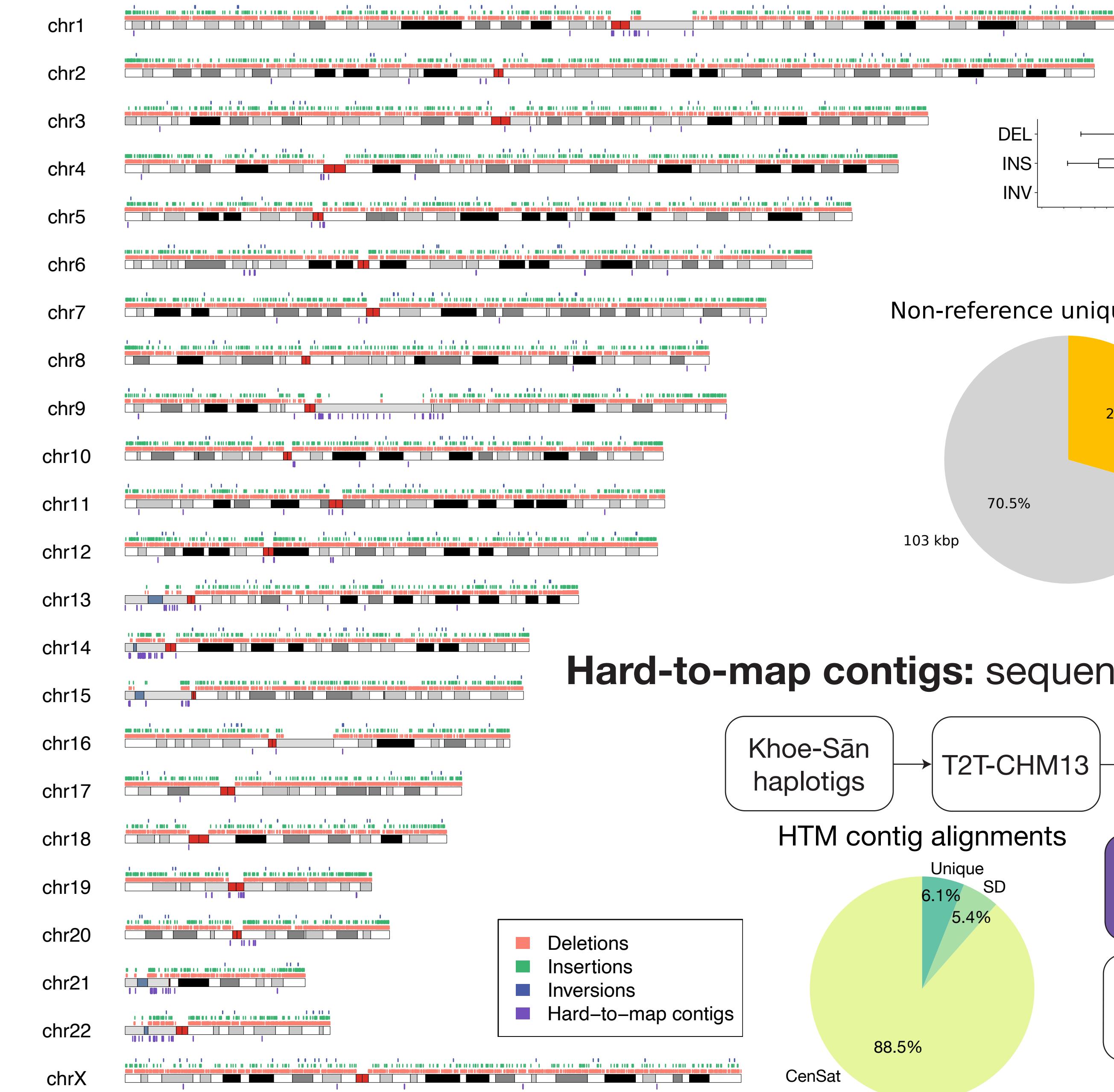
Carlson et al. 2022 PMID: 36261568, Chen et al. 2019 PMID: 31856913
Shen and Kidd. 2020 PMID: 32013076

Results

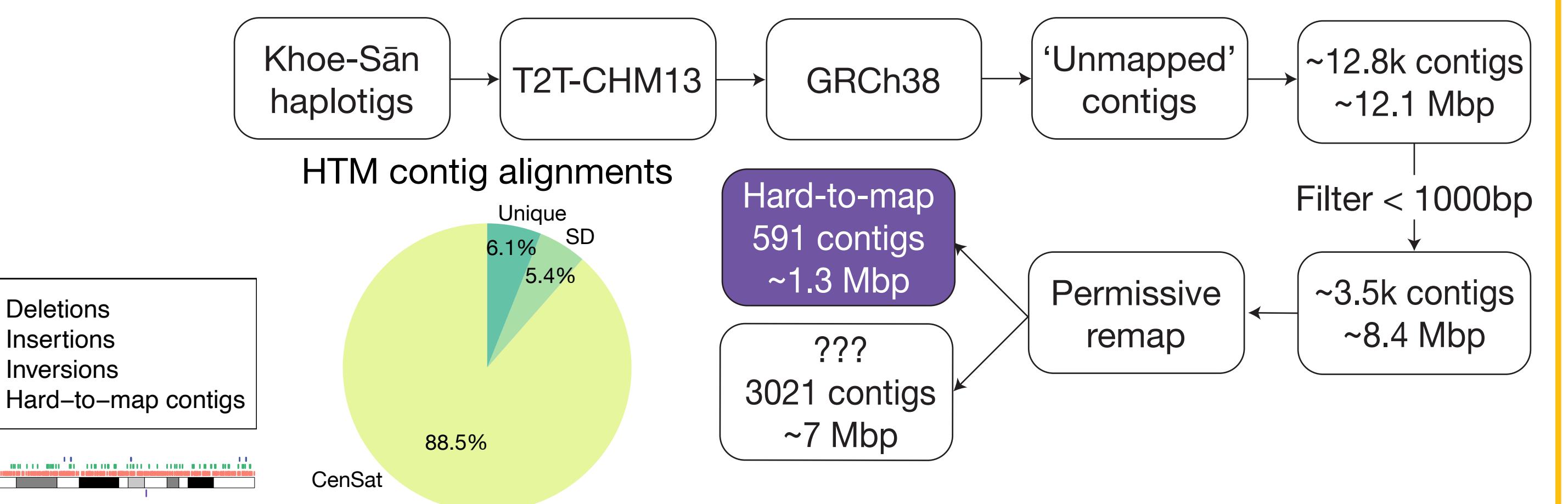
Assembly statistics: *de novo* assembly compared to other 10X linked read assemblies



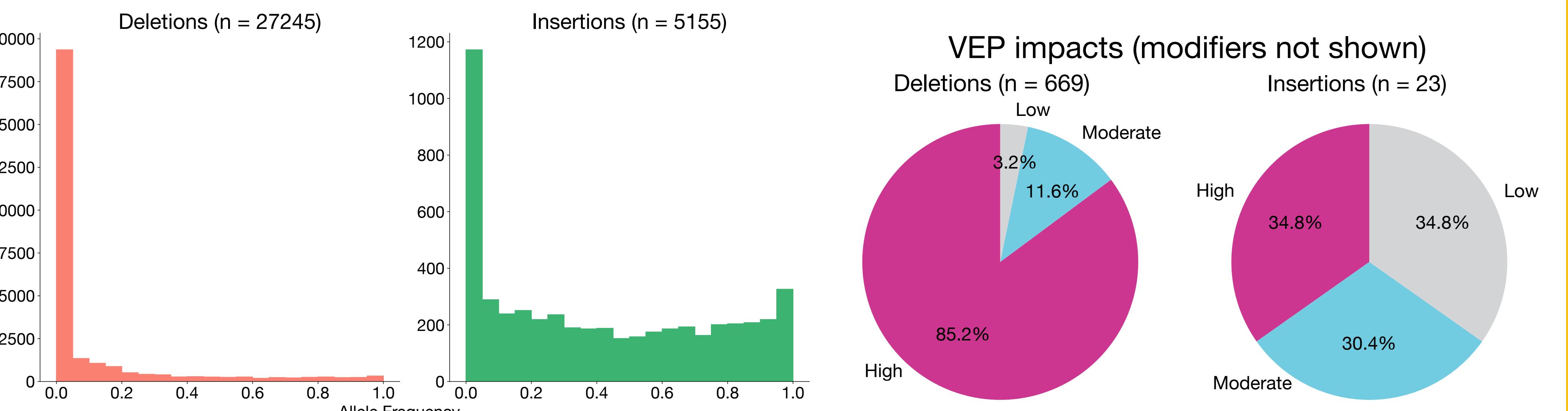
Structural variation: assembly to reference comparison



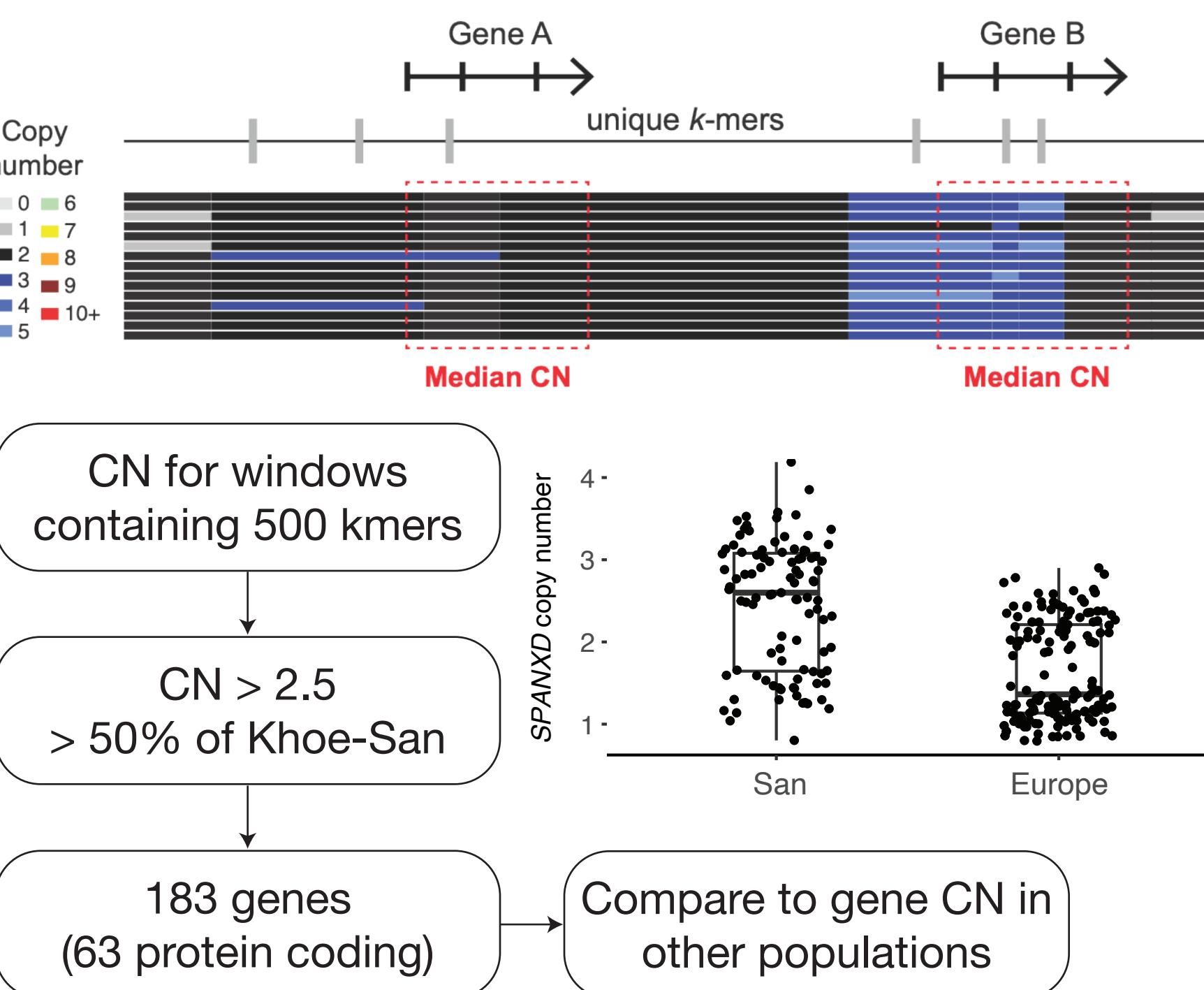
Hard-to-map contigs: sequence not found in reference genomes



SV genotyping: using short read sequencing data from a cohort of 93 Khoë-Sān



Copy number variation: using short reads from 93 Khoë-Sān individuals



Conclusions

- A lack of diversity in African genomes in LRS projects limits us from exploring human genetic variation
- We generated 10X linked-read and *de novo* assembly for 3 Khoë-Sān individuals
- We found ~7Mbp not mapping to human reference genomes, even after relaxing parameters
- We ascertained allele frequencies of many novel structural variants through genotyping in SRS Khoë-Sān samples
- We detected copy number variations that may differ from other human populations

Future directions

- Robust investigation of non-mapping contigs, namely aligning short read Khoë-Sān datasets to them
- Genotype discovered SVs in other populations from the HGDP, and identify Khoë-Sān specific variants
- Calculate F_{ST} and V_{ST} between Khoë-Sān and HGDP to study stratification
- Further explore high impact variants highlighted by VEP
- Leverage long haplotypes to investigate structurally complex loci, including HLA
- Assess improved mappability for Khoë-Sān individuals using new assemblies

Acknowledgements

Thank you to funding sources (NIMH: R01MH132818, NIGMS: R35GM133531) and collaborators for their help facilitating this project, as well as the HGDP and HPRC for access to their data.