000 A SIMULATION-BASED APPROACH TO UNDERSTAND-001 ING GENOMIC VARIANTS IN ADMIXED POPULATIONS 002 003 USING GRAPH-BASED REFERENCE GENOMES 004

Anonymous authors

Paper under double-blind review

1 INTRODUCTION

Admixed populations represent a crucial aspect of global genetic diversity, yet their unique genomic 014 architectures remain underrepresented in genomic studies Stankiewicz & Lupski (2010). Tradi-015 tional genome-wide association studies (GWAS) and variant detection methods often rely on linear 016 reference genomes, which do not fully capture the complexity of genomic variation in admixed indi-017 viduals Gopalan et al. (2022). This study employs a simulation-based approach to generate synthetic 018 genomic data from two ancestral populations, European (EUR) and South American (AMR), and 019 assess the impact of structural variants (SVs) on disease risk in an admixed population. Furthermore, the study evaluates the effectiveness of pangenome Eizenga et al. (2020) graph-based variant 021 detection compared to traditional linear reference models.

022

025

026

006

012 013

2 METHODS

The study is structured into three primary stages: (1) Human Admixture Simulation, (2) Pangenome Generation, and (3) Read Mapping and Variant Detection.

027 029

037

038

041

2.1 HUMAN ADMIXTURE SIMULATION

The study simulates genetic data for 1,000 EUR and 1,000 AMR samples using the chromosomal 031 region chr20:30,000,000-32,000,000. The simulation proceeds over five generations (150 years), after which the minor allele frequency (MAF) of single nucleotide variants (SNVs) is computed. A 033 subset of SNVs unique to AMR (i.e., absent in EUR) is identified, and 20 such SNVs are modified into structural variants:

- 10 SNVs are transformed by adding 1,000 bases of random sequence.
- 10 SNVs are altered by duplicating the 1,000 bases following the SNV site.

039 2.2 PANGENOME GENERATION 040

A representative subset of 20 admixed samples is selected to generate a high-quality reference 042 using T2T (telomere-to-telomere) assemblies, following methodologies employed by the Human 043 Pangenome Project. The GRCh38 and T2T-CHM13 references are incorporated to construct a 044 pangenome graph of the chr20 region. This approach enables the evaluation of graph-based reference genomes in variant detection.

046 047

048

2.3 READ MAPPING AND VARIANT DETECTION

Short-read sequencing data (100bp paired-end reads) is simulated for 100 AMR samples and 100 admixed samples to assess variant detection performance. Reads are mapped to the chr20 region using BWA-MEM with GRCh38 as the reference. The following analyses are conducted:

- 051 052
- The percentage of high-quality mapped reads per sample is computed.
- The percentage of correctly identified variants is determined.

054 055 056	• The detection rates of AMR-specific SVs are assessed by comparing the two types of struc- tural variants (random sequence vs. duplicated sequence).
057	The experiment is repeated using the pangenome reference to compare:
059	• Read mapping quality.
060	Correctly identified variants.
061	Detection rates of AMR-specific SVs in both AMR and admixed samples
062	Differences in detection accuracy between the two structural variant types
063	• Differences in detection accuracy between the two structural variant types.
065	2.4 Repeated Analyses
066 067 068	To account for the stochastic nature of graph construction, the pangenome is reconstructed multiple times using smoothing, varying the sample composition:
069	1. Using a different subset of 20 AMR samples.
070 071	2. Different sample sizes (n=10, n=50) were used to assess the impact of dataset size on pangenome quality.
072 073 074	 Constructing a pangenome with a mix of admixed (n=10) and AMR (n=10) samples to determine whether including ancestral genomes improves variant detection and read map- ning
075	ping.
076 077	3 Results and Expected Impact
078 079	This study is expected to provide insights into the following:
080 081 082	1. The extent to which traditional linear reference genomes fail to detect population-specific structural variations in admixed individuals.
083 084	2. The effectiveness of graph-based reference genomes in improving variant detection, partic- ularly for AMR-specific structural variants.
085 086	3. Pangenome construction parameters (sample composition and size) influence mapping ac- curacy and variant calling performance.
087 088 089	4. By integrating graph-based reference models, potential strategies for improving disease risk variant identification in admixed populations.
090	References
091	Jordan M Fizenga Adam M Novak Jonas A Sibbesen Simon Heymos Ali Ghaffaari Glenn
093	Hickey, Xian Chang, Josiah D Seaman, Robin Rounthwaite, Jana Ebler, Mikko Rautiainen, Shilpa
094	Garg, Benedict Paten, Tobias Marschall, Jouni Sirén, and Erik Garrison. Pangenome graphs.
095	Annu. Rev. Genomics Hum. Genet., 21(1):139–162, 2020.
096	Shyamalika Gopalan, Samuel Pattillo Smith, Katharine Korunes, Iman Hamid, Sohini Ramachan-
097	dran, and Amy Goldberg. Human genetic admixture through the lens of population genomics.
098	Philos. Trans. R. Soc. Lond. B Biol. Sci., 377(1852), 2022.
100	Paweł Stankiewicz and James R Lupski. Structural variation in the human genome and its role in
101	disease. Annu. Rev. Med., 61(1):437-455, 2010.
102	
103	
104	
105	