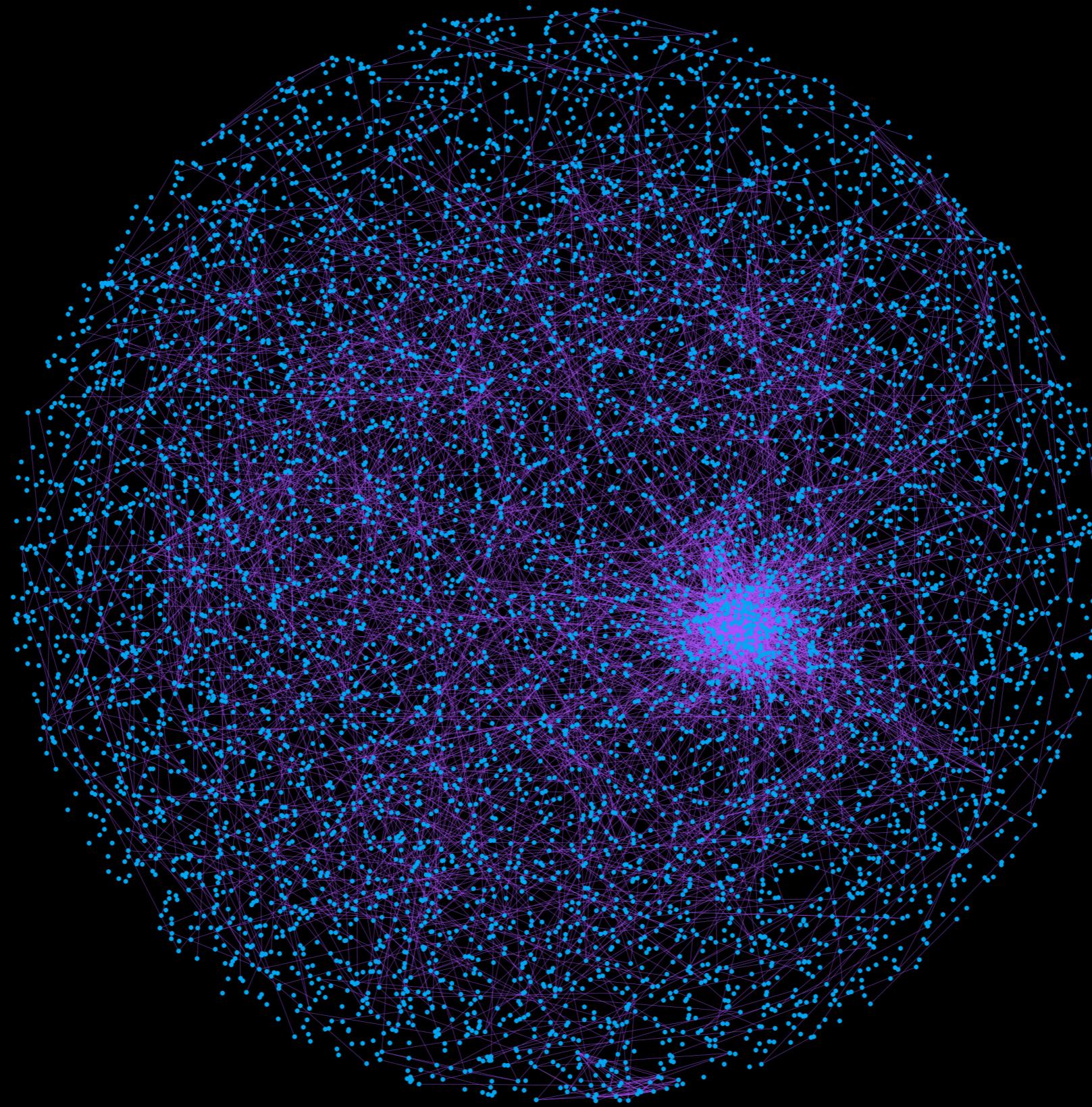


Machine Learning in population genetics

Alexander Rakitko, MSU

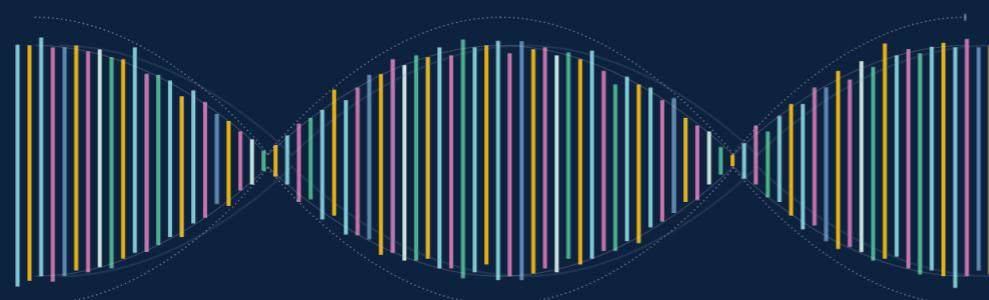
Network of genetic relatives



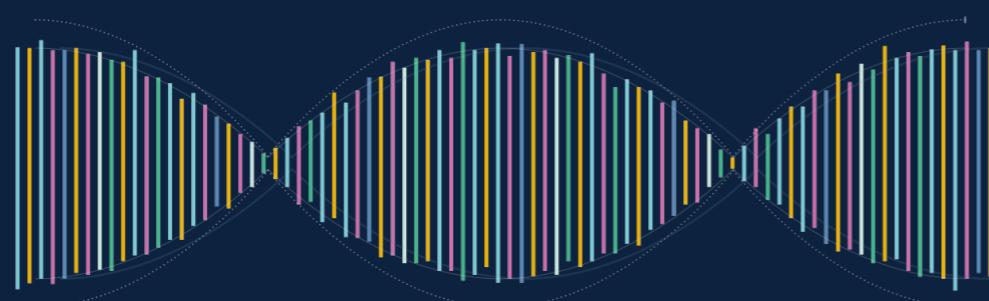
We work with SNPs



... A A T G C A A T G C G A ...

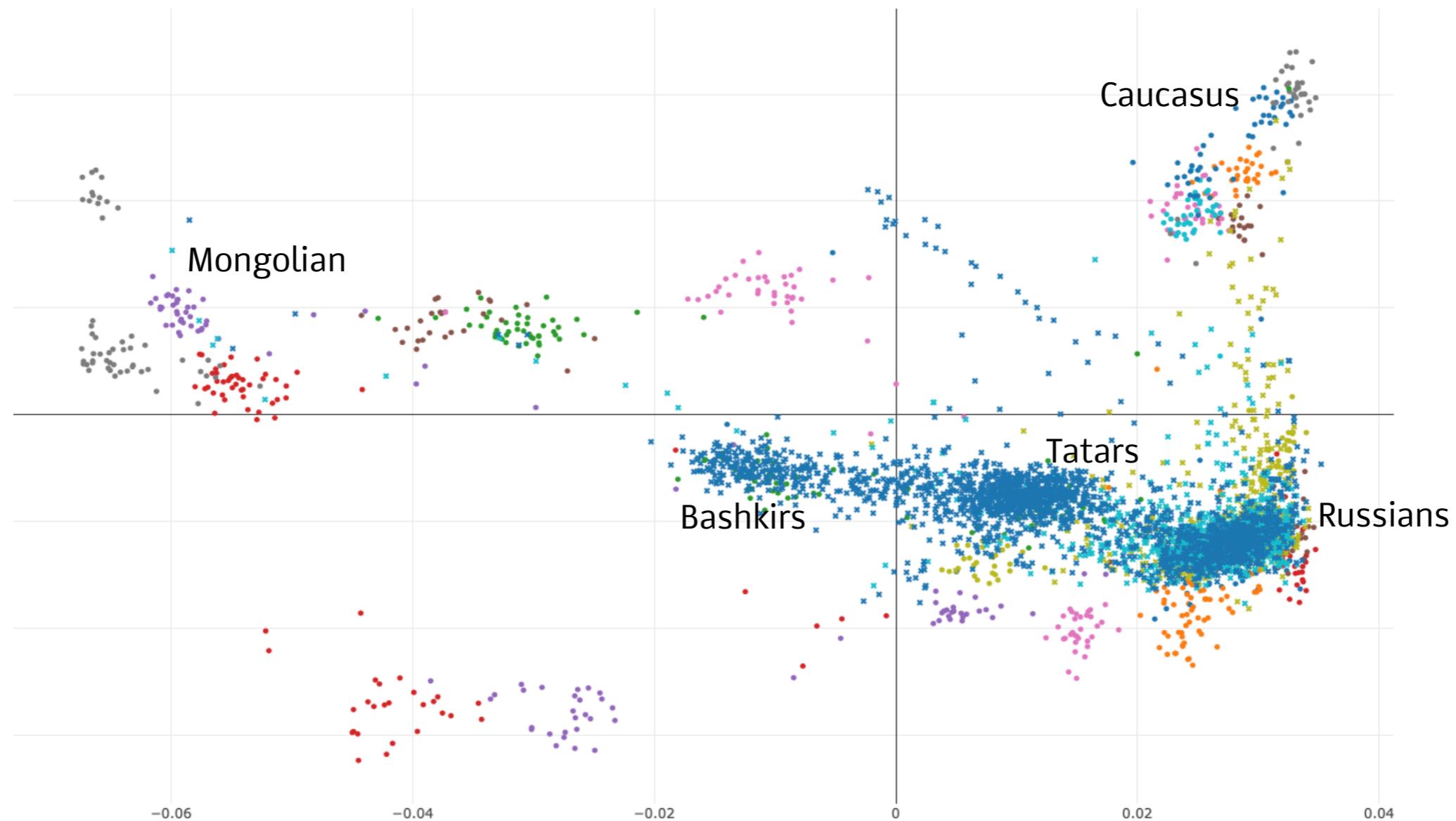


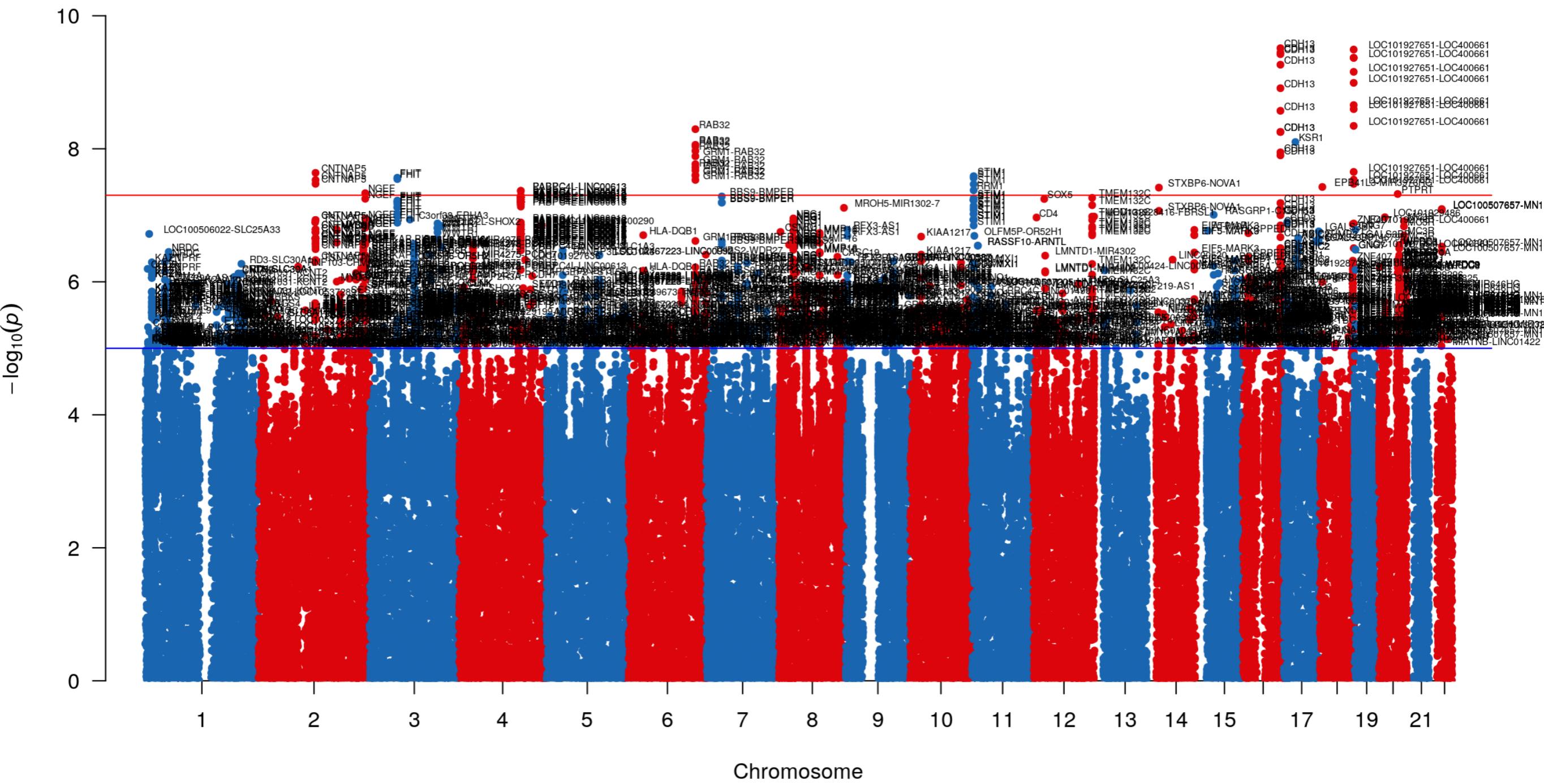
... A A T G C A G T G C G A ...

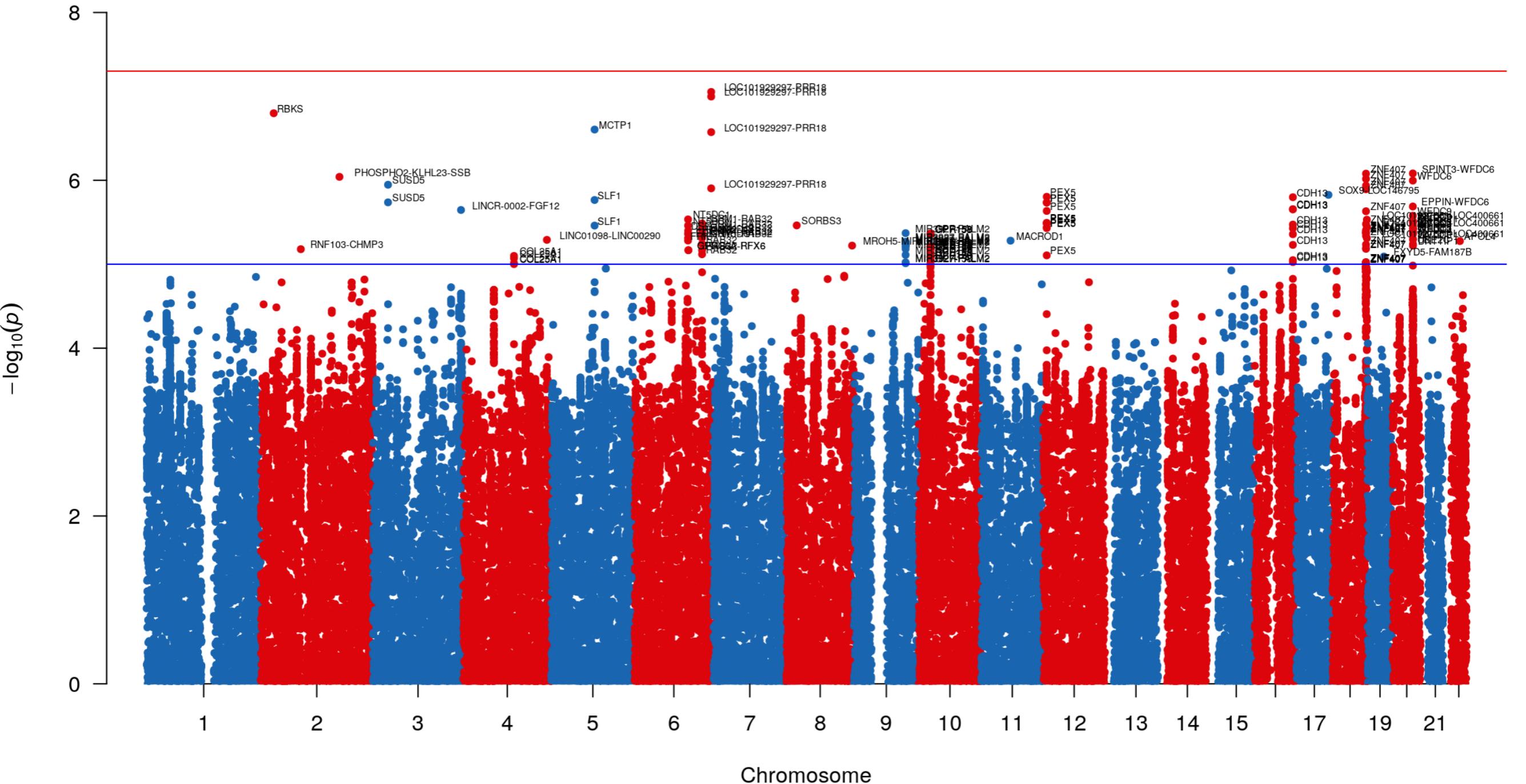


... A A T G C A C T G C G A ...

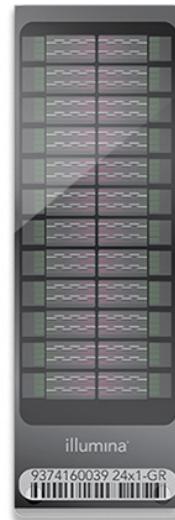
GWAS of schizophrenia with Russian National Consortium for Psychiatric Genetics (RNCPG)



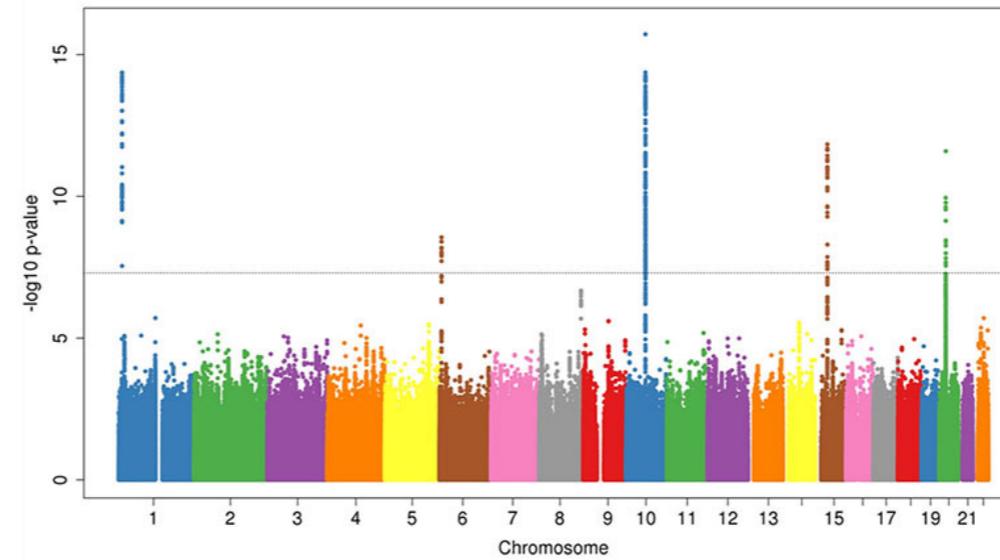




Polygenic Risk Scores



Microarray Genotyping



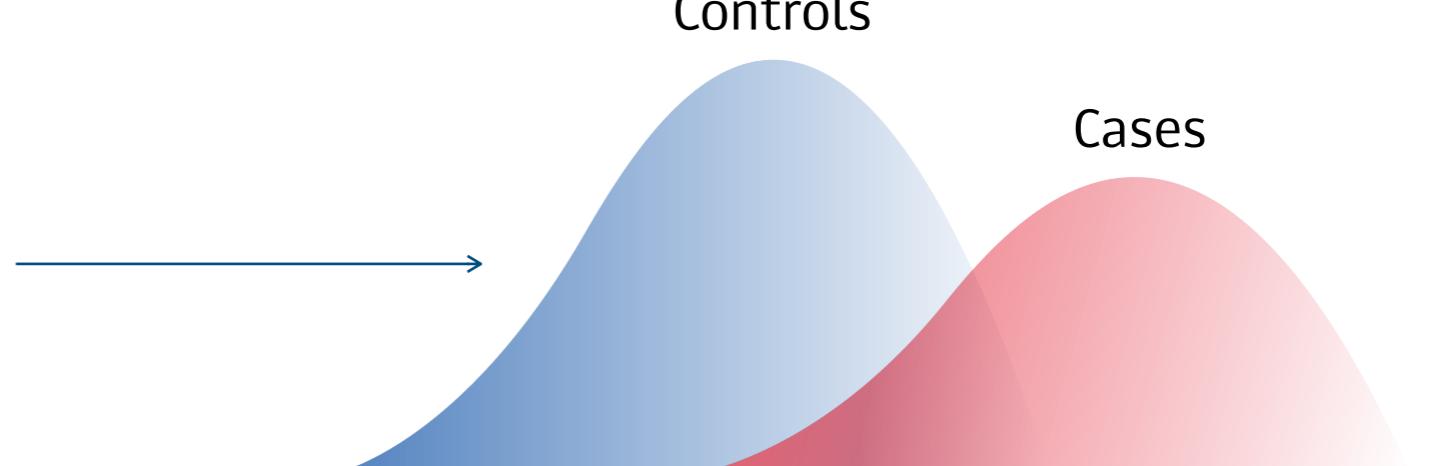
GWAS analysis

$$\text{PRS} = \beta_1 X_1 + \dots + \beta_n X_n$$

Genome-Wide Polygenic Risk Score

Controls

Cases



PRS for coronary artery disease (CAD)

Polygenic Risk Score

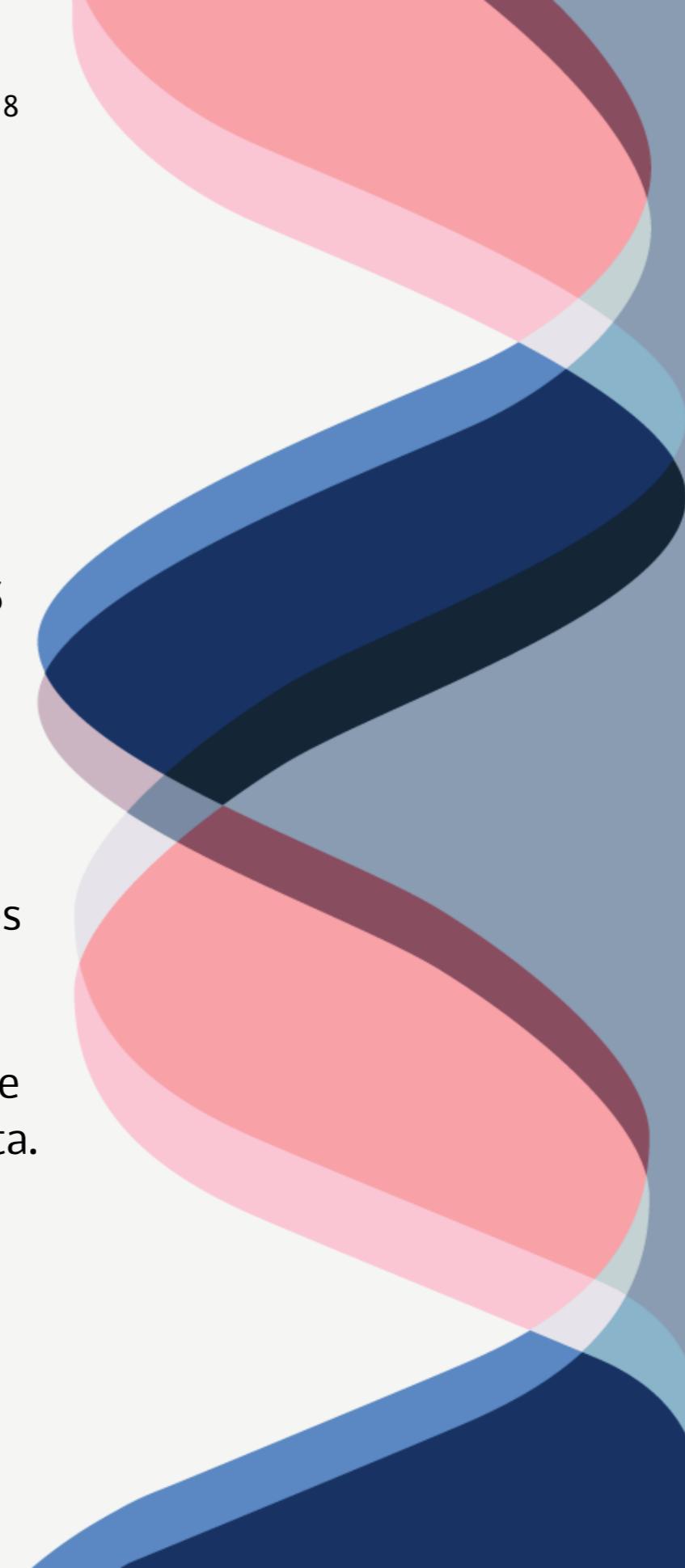
The PRS previously described in Khera et al. 2018 was used as a predictive model. The PRS was based on analysis of 6,630,150 SNPs.

Test dataset

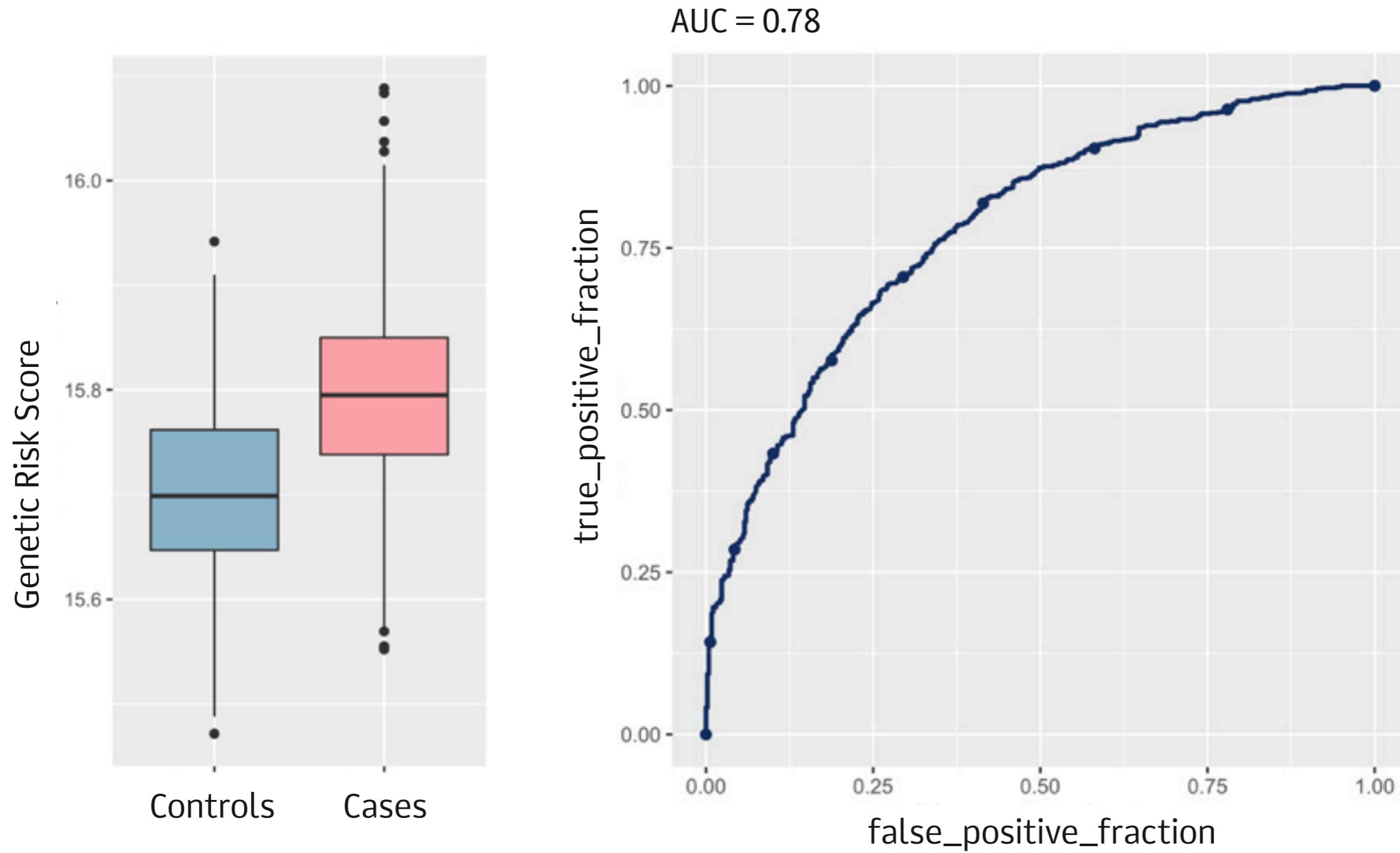
A PennCath study subset was used as a test sample (Reilly et al. 2011). The total sample size was 1,401 individuals, including 933 cases with CAD and 468 controls. The genotyping array contains 500,000 SNPs. Age, sex, triglyceride, LDL and HDL levels were available for every individual in addition to genetic data.

Khera, Amit V., et al. "Genome-wide polygenic scores for common diseases identify individuals with risk equivalent to monogenic mutations." *Nature genetics* 50.9 (2018): 1219.

Reilly, Muredach P., et al. "Identification of ADAMTS7 as a novel locus for coronary atherosclerosis and association of ABO with myocardial infarction in the presence of coronary atherosclerosis: two genome-wide association studies." *The Lancet* 377.9763 (2011): 383-392.

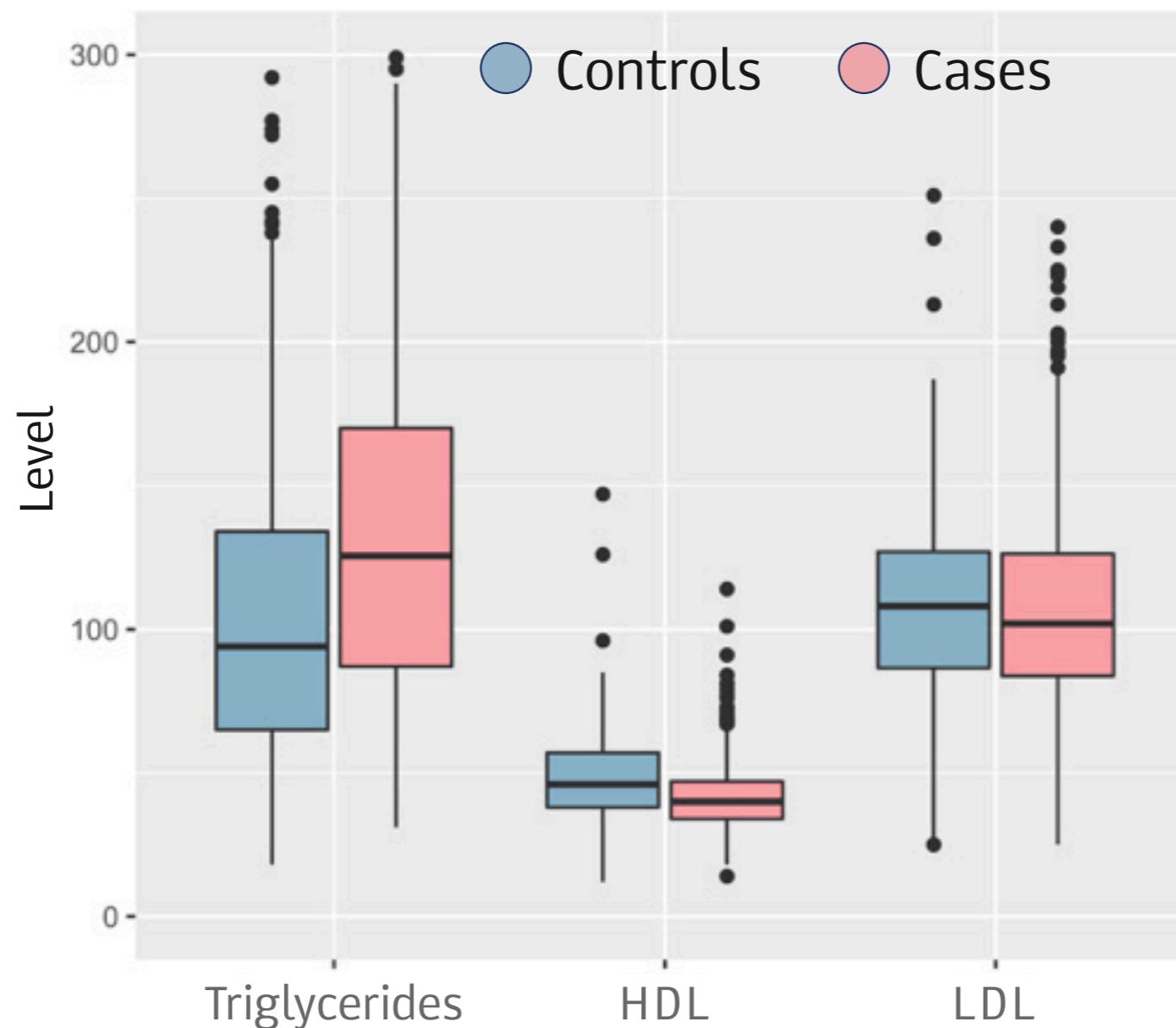


Validation of Polygenic Risk Score



Non-genetic risk factors

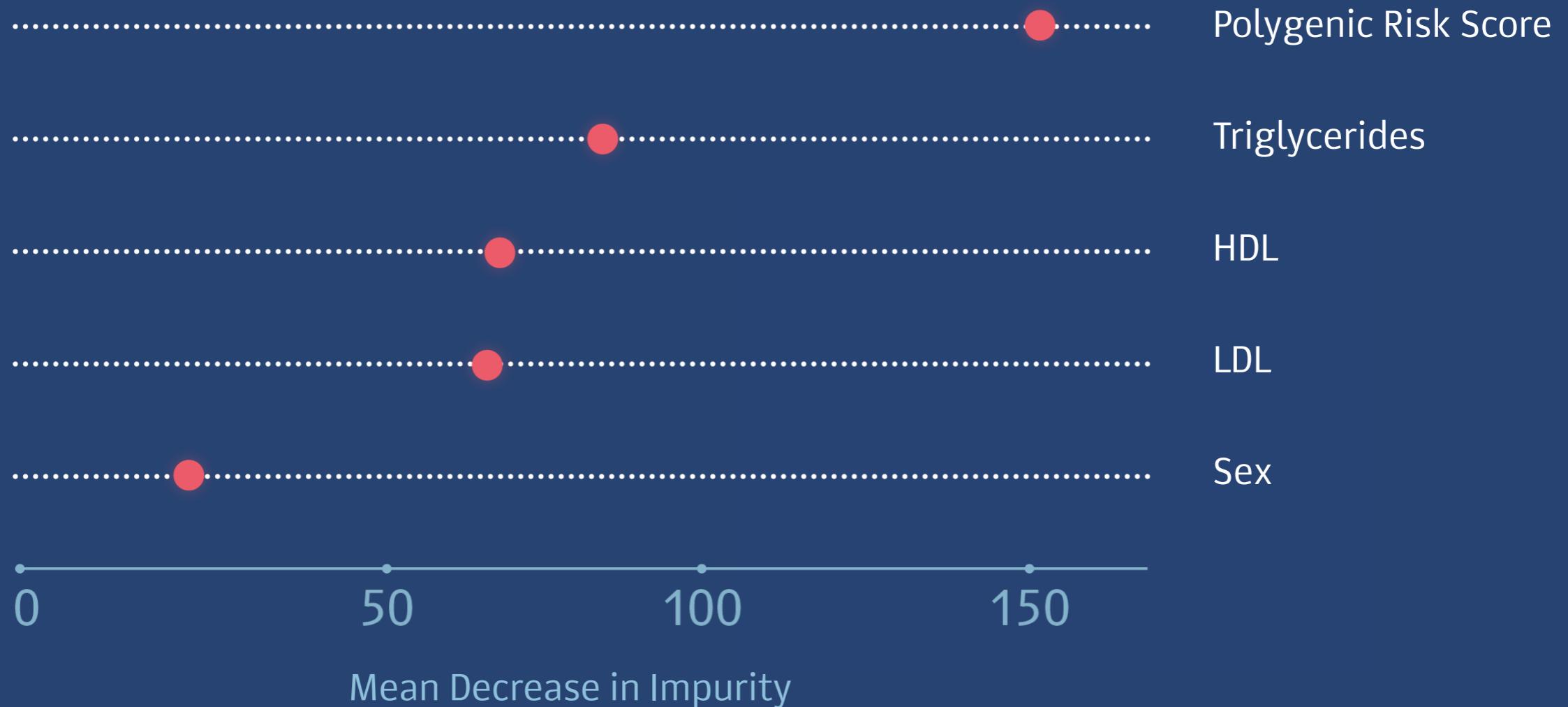
Coronary artery disease



Comparison of predictive models

Model	Non-genetic factors	Genetic + non-genetic factors
Logistic regression	69,6 % (3%)	75,3 % (3,5%)
Random forests	65,3% (3%)	74,5 % (2%)
Gradient boosting	67,2 % (2%)	75,8 % (3%)

Feature importance (RF model)



Population adjustment for Polygenic Risk Scores

Regarding stratification, most PRS methods do not explicitly address recent admixture, and none consider recently admixed individuals' unique local mosaics of ancestry; thus, further methodological development is needed.

Martin, Alicia R., et al. "Clinical use of current polygenic risk scores may exacerbate health disparities." *Nature genetics* 51.4 (2019): 584.

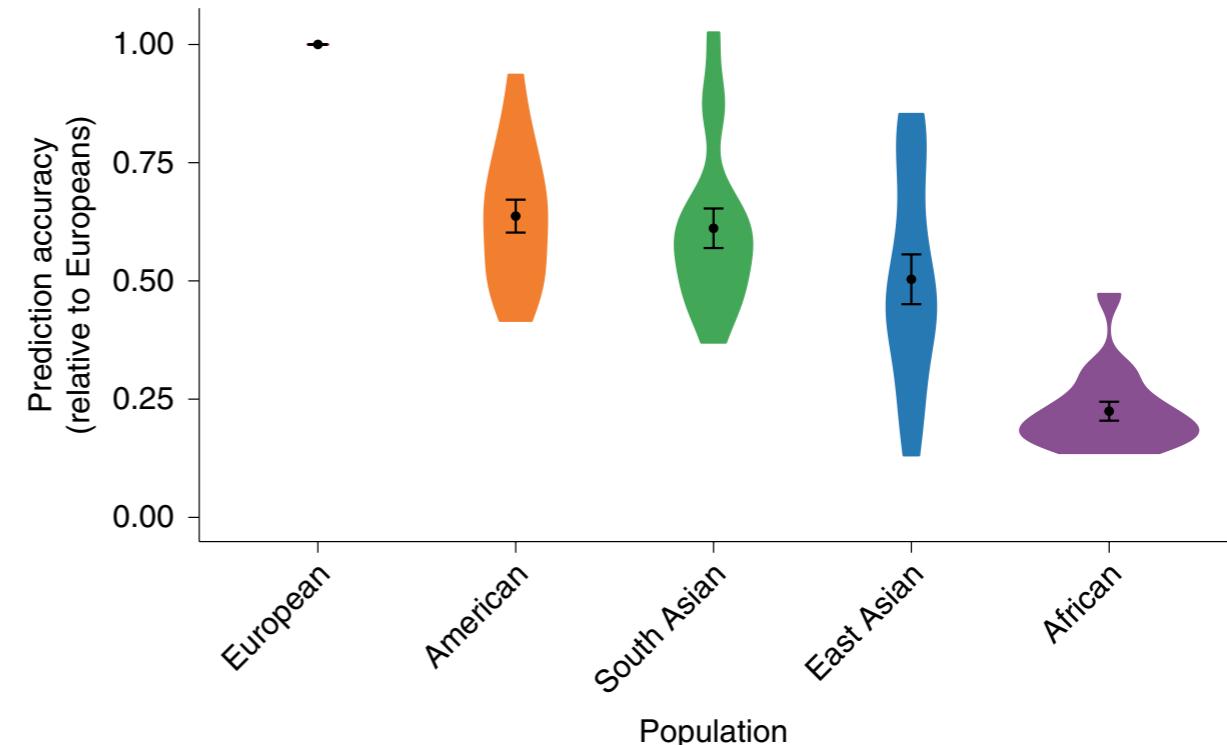
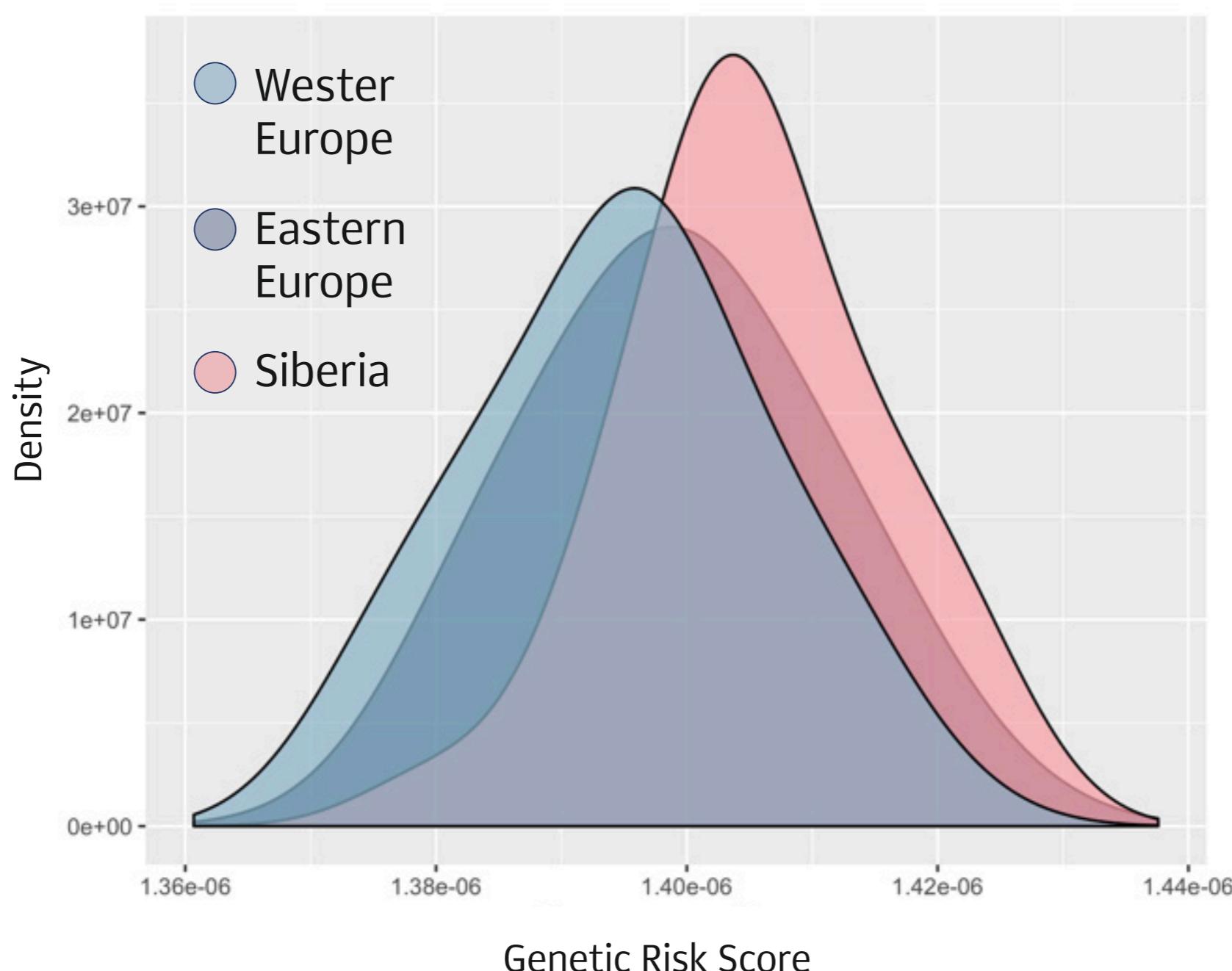
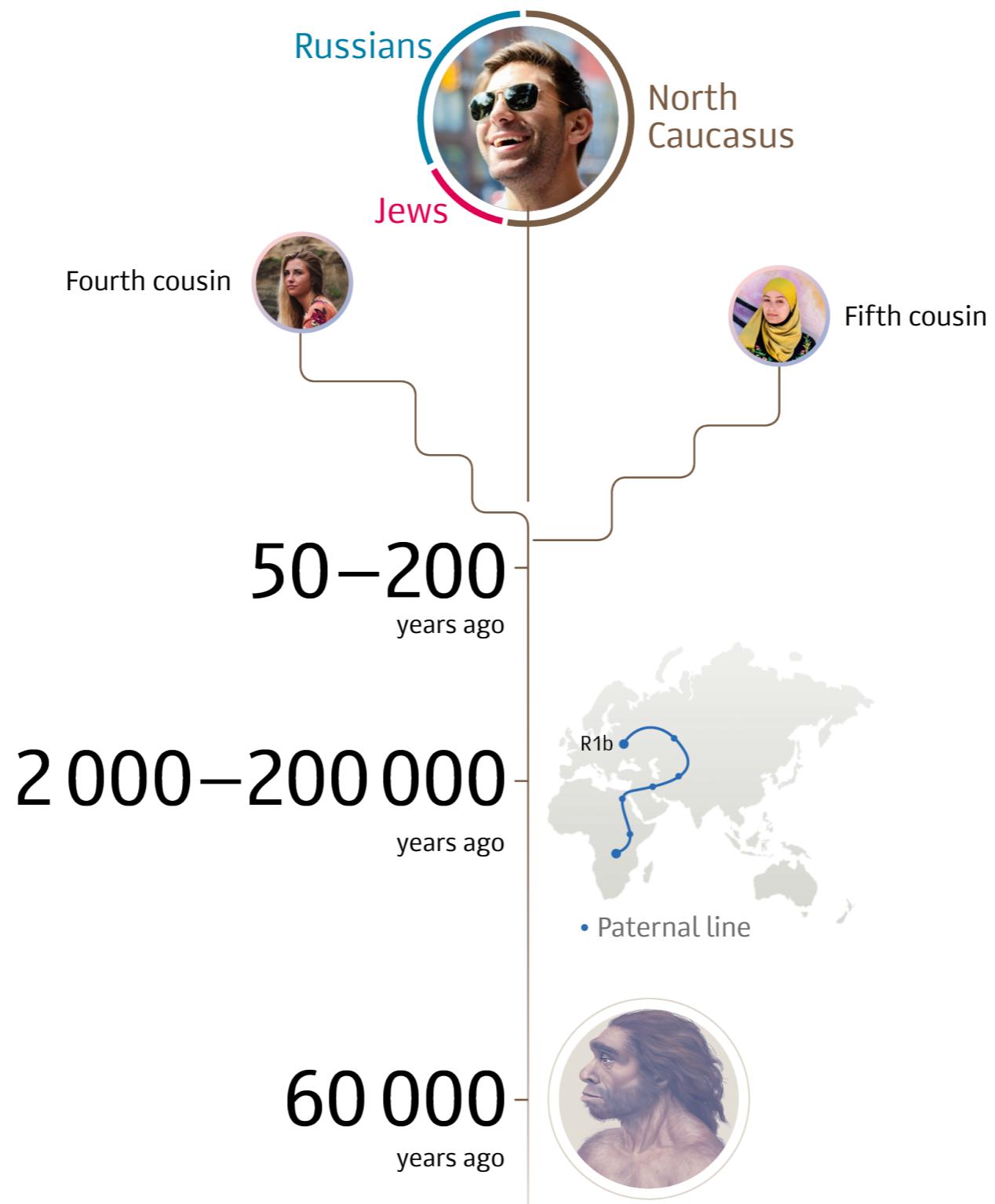


Fig. 3 | Prediction accuracy relative to European-ancestry individuals across 17 quantitative traits and 5 continental populations in the UKBB. All phenotypes shown here are quantitative anthropometric and blood-panel traits, as described in Supplementary Table 6, which includes discovery-cohort sample sizes. Prediction target individuals do not overlap with the discovery cohort and are unrelated; sample sizes are shown in Supplementary Table 7. Violin plots show distributions of relative prediction accuracies, points show mean values, and error bars show s.e.m. values. Prediction R^2 for each trait and population are shown in Supplementary Fig. 12.

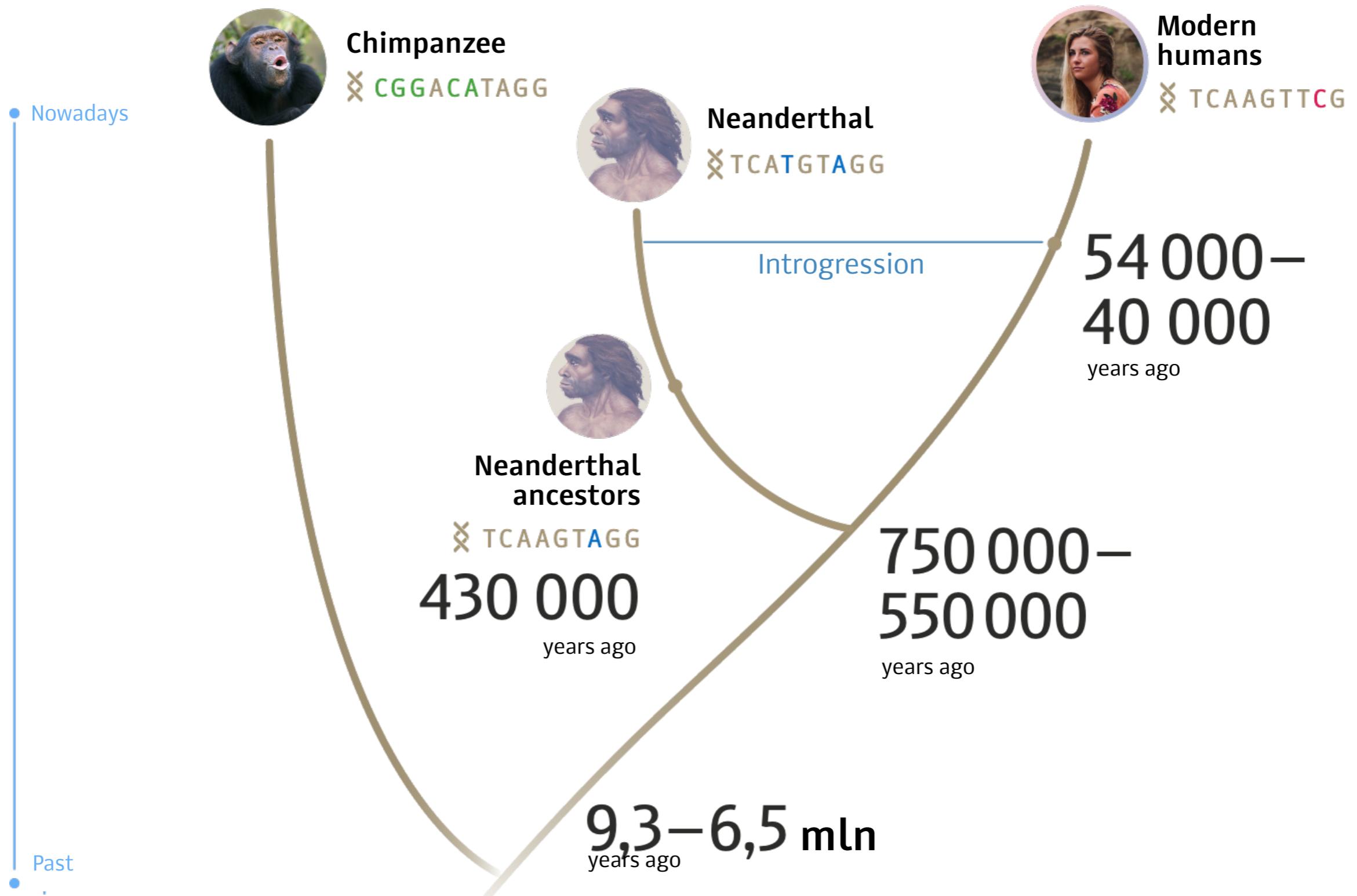
Population differences



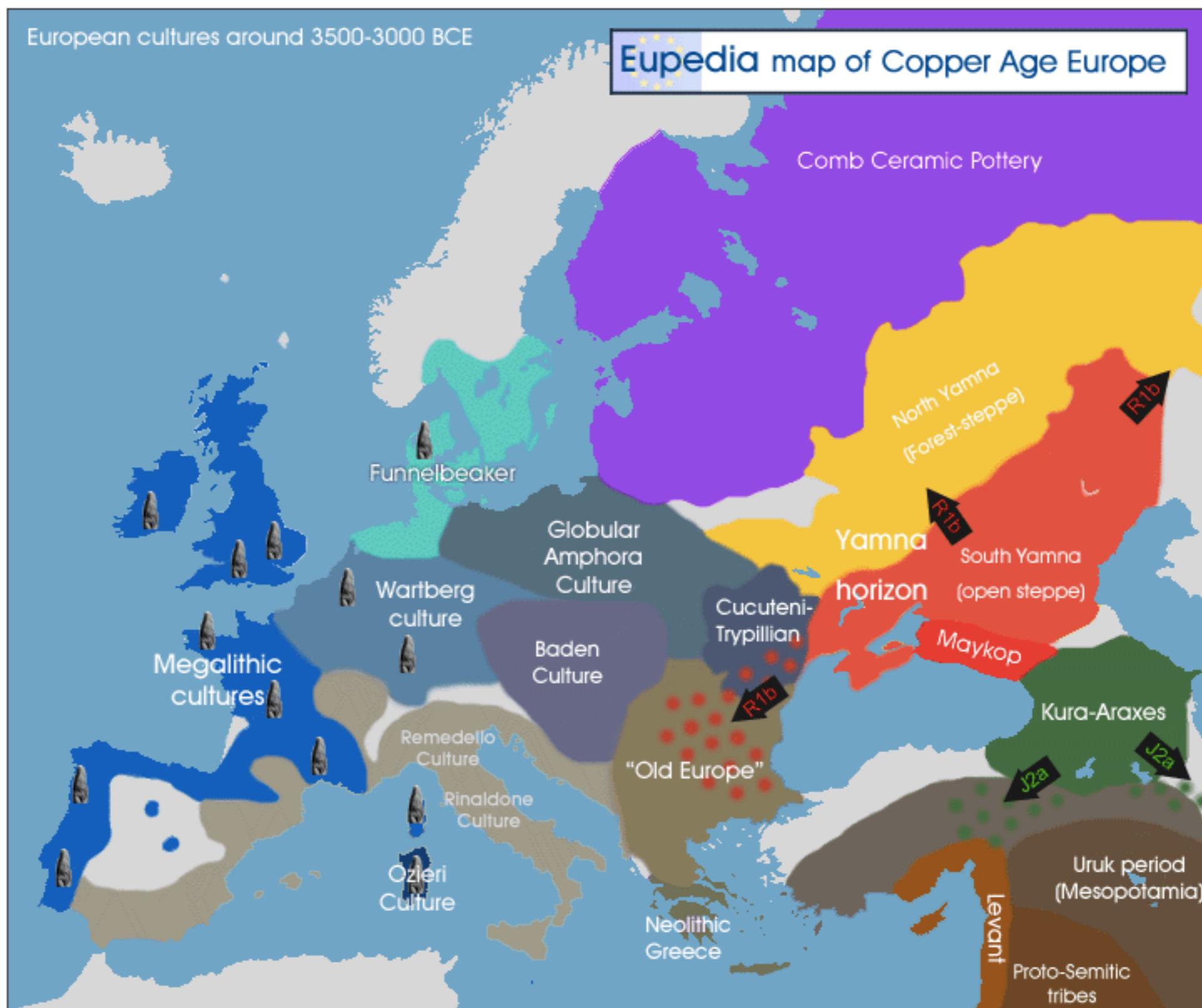
The goal of hackathon: discover your ancestry



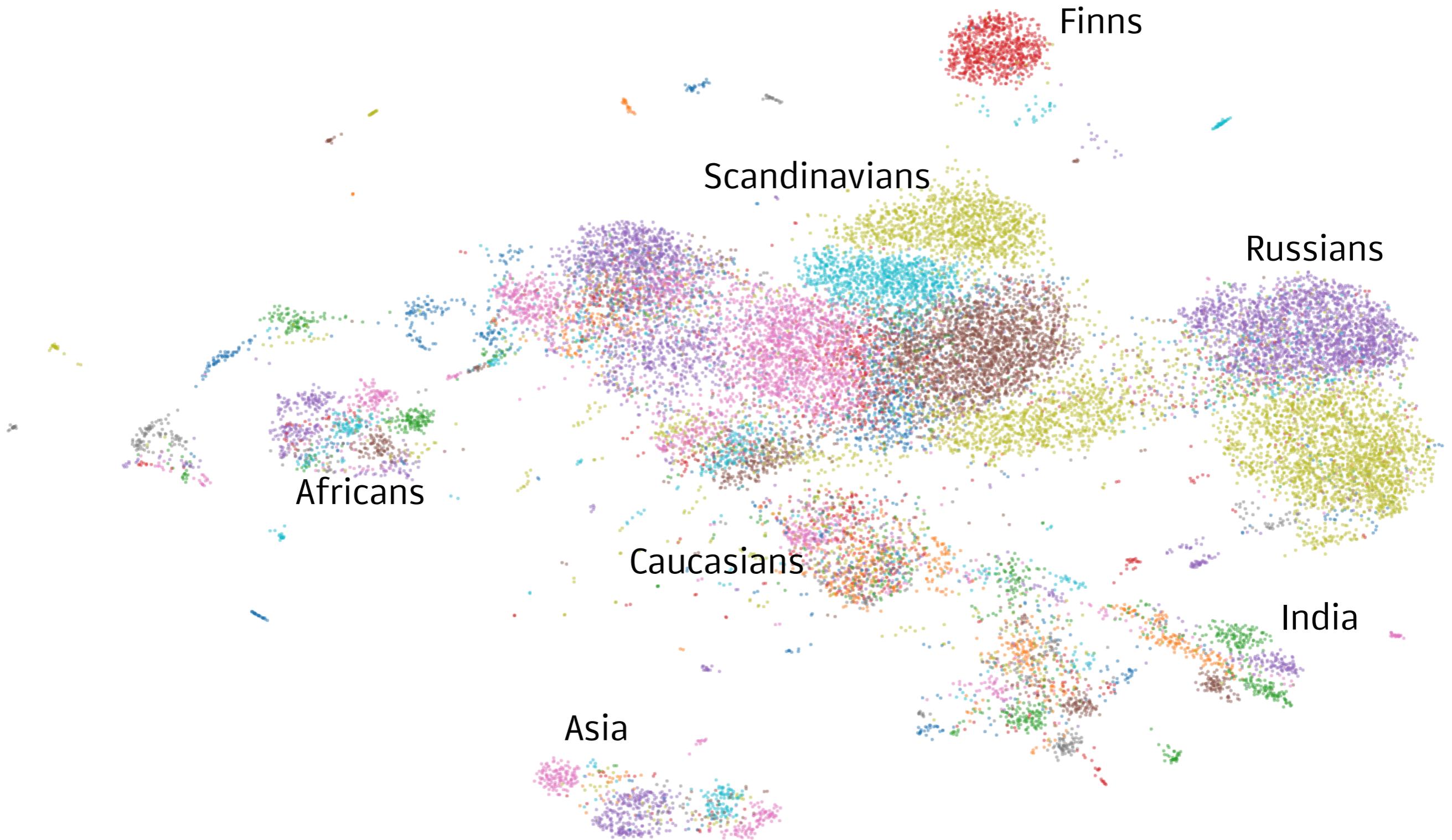
Find population with the highest Neandethal introgression



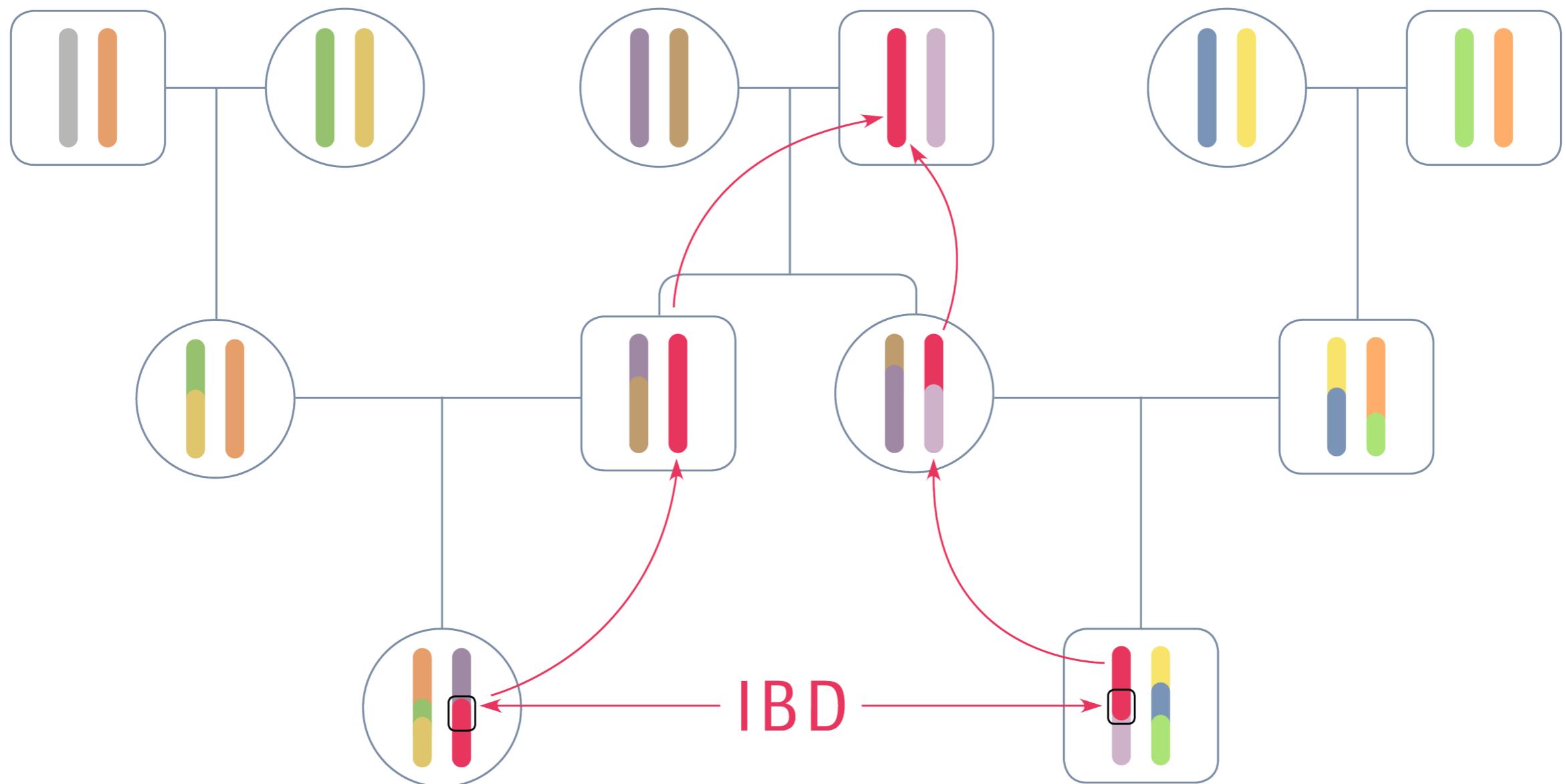
Proportion of Yamna culture in modern Polish genomes



Plot the genetic map of Europe



Paint chromosomes into modern populations



Hackathon summary

Algorithms

- Neanderthal genes (phylogenetic metrics)
- Ancient cultures (Max Likelihood Est)
- Genetic map of Europe (PCA, tSNE, UMAP)
- Find populations (k-means, HDBSCAN)
- Paint chromosomes (SVM / k-NN / NN / RF + HMM)

Dataset

- 2K+ ancient genomes
- 5K+ modern genomes

Languages

- R
- Python
- Bash