

Visualizations of Genotype Datasets

This repository contains the visualizations of 1000 Genomes dataset ([1000 Genomes Project Consortium](#)) and the modern Eurasian population dataset ([Lamnidis et al.](#)). The dimensional reduction techniques used include PCA, t-SNE, UMAP, and Adversarial Autoencoders ([Makhzani et al.](#))

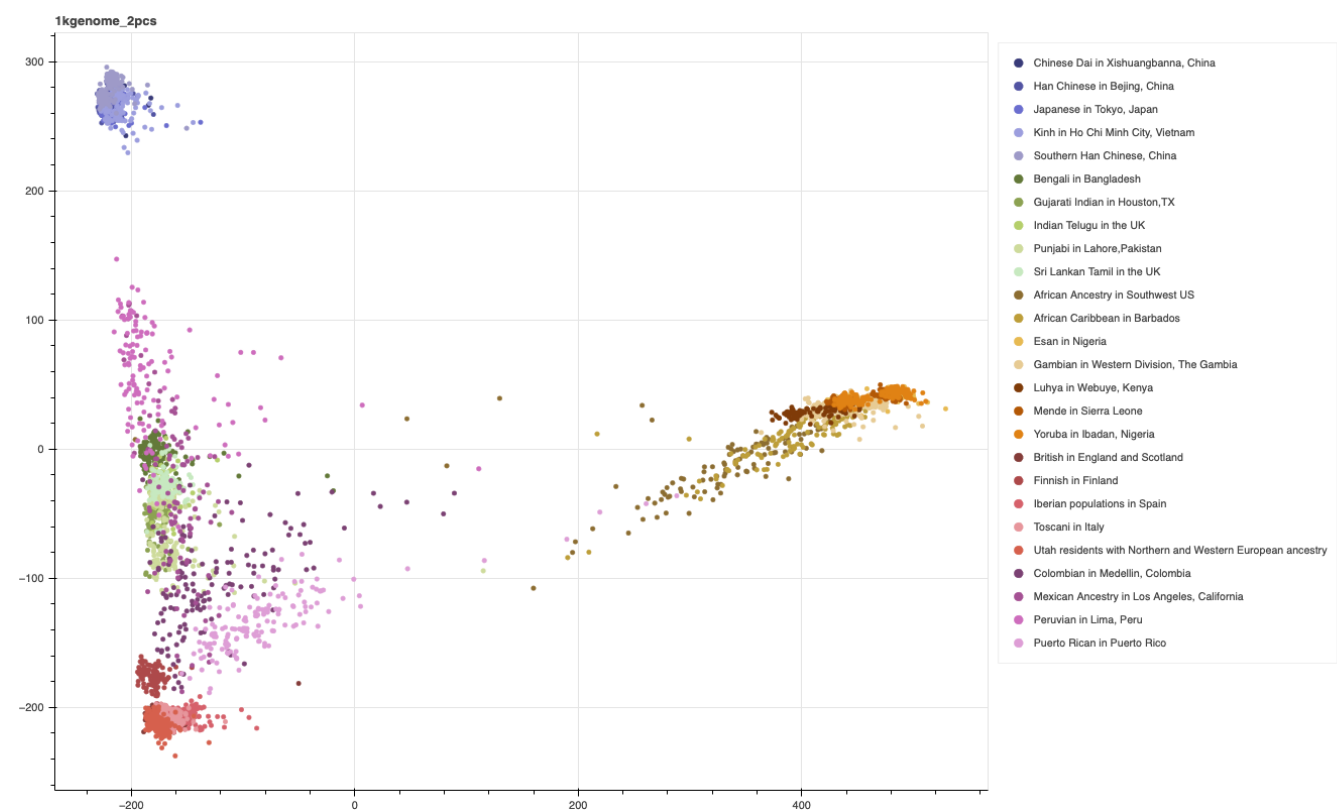
All visualizations can be founded in the form of interactive htmls in the [htmls](#) folder. The interactive htmls let you pan, zoom, hover on data points, and also hide data of particular labels by clicking on the legend.

We also hand pick some visualizations and display them as static images below.

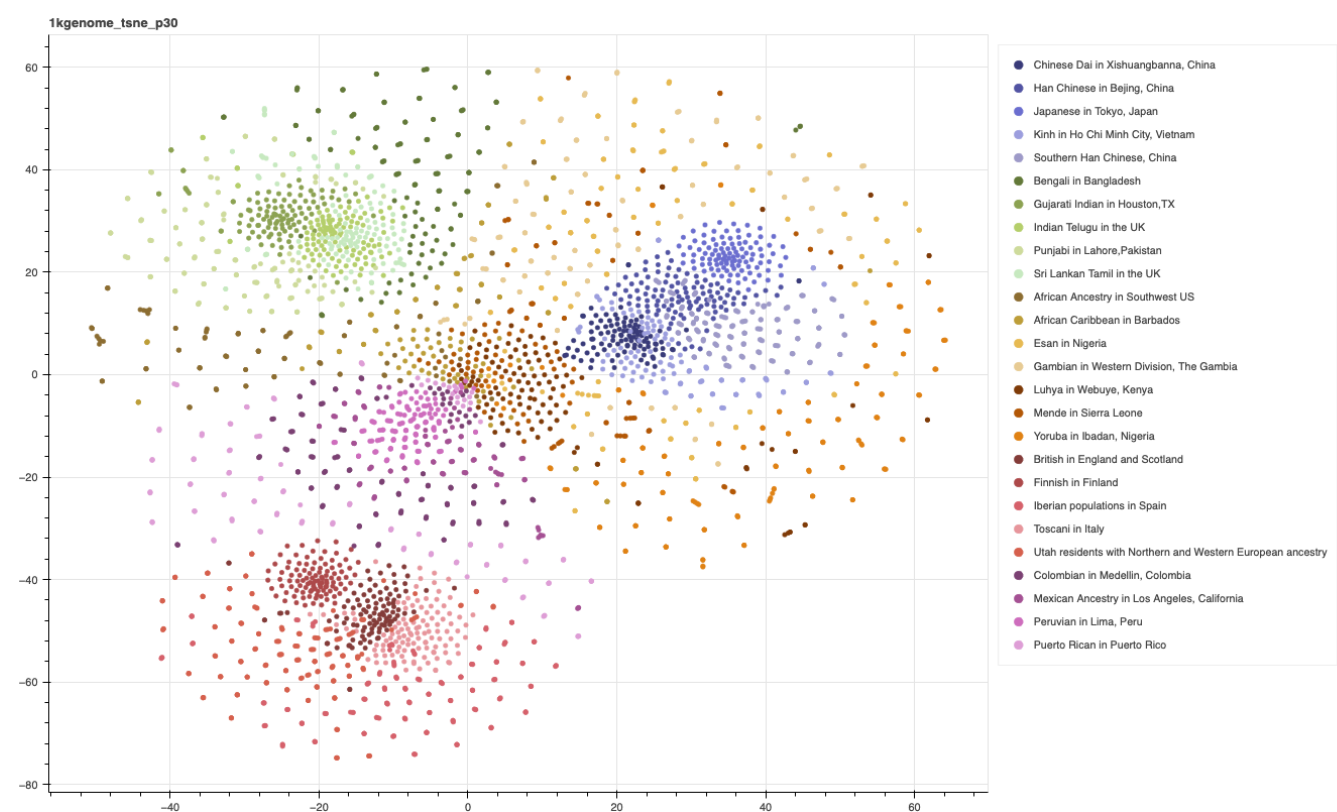
1000 Genomes

For reference, [Diaz-Papkovich et al.](#) have applied PCA, t-SNE, and UMAP to the 1000 Genomes dataset. The following figure is the first figure in the paper

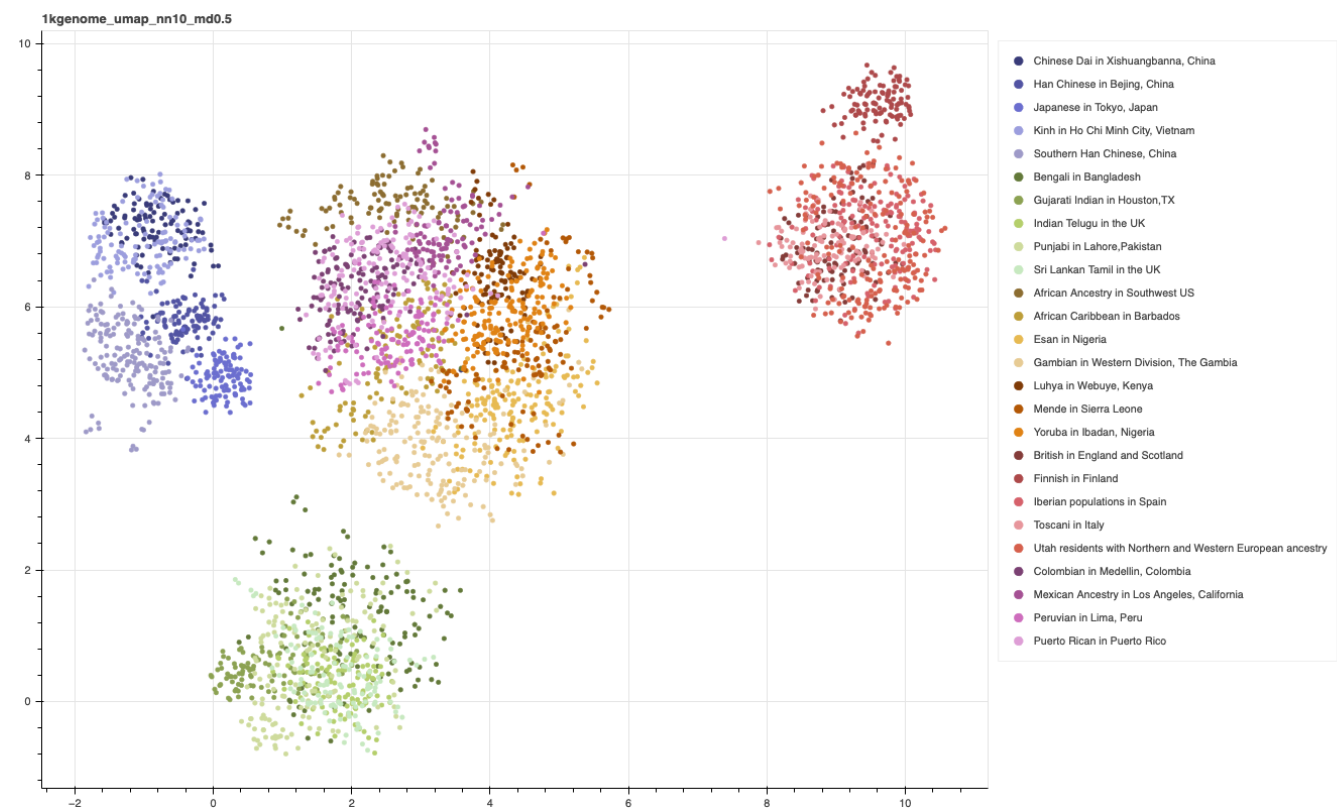




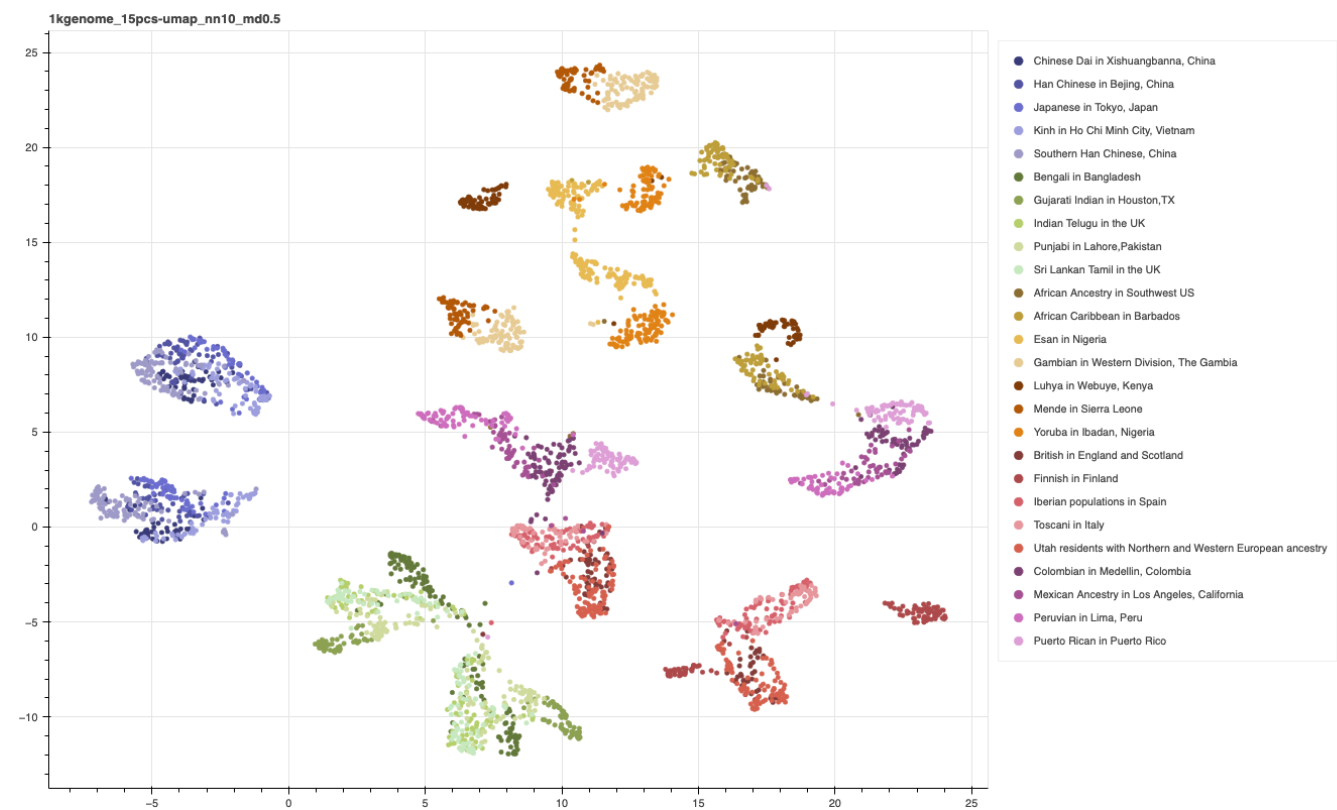
• t-SNE



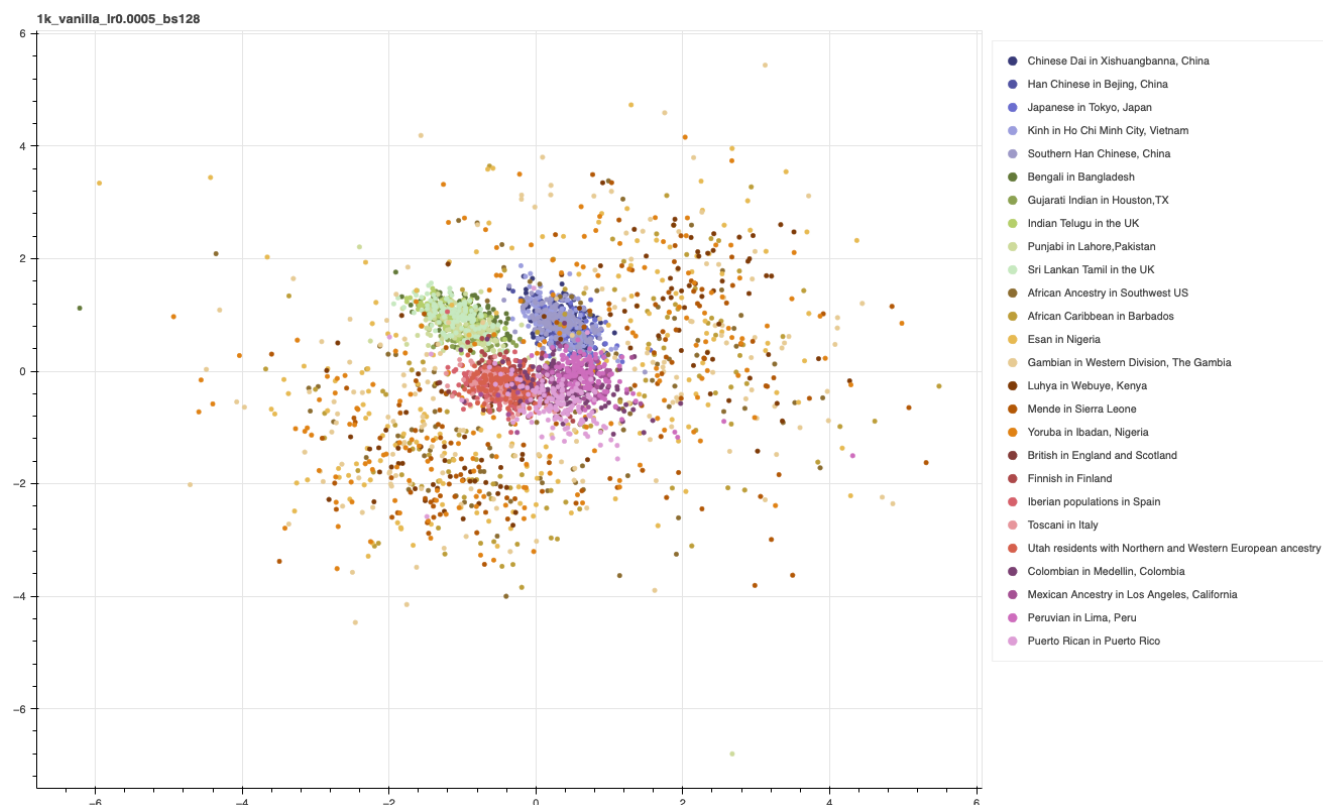
• UMAP



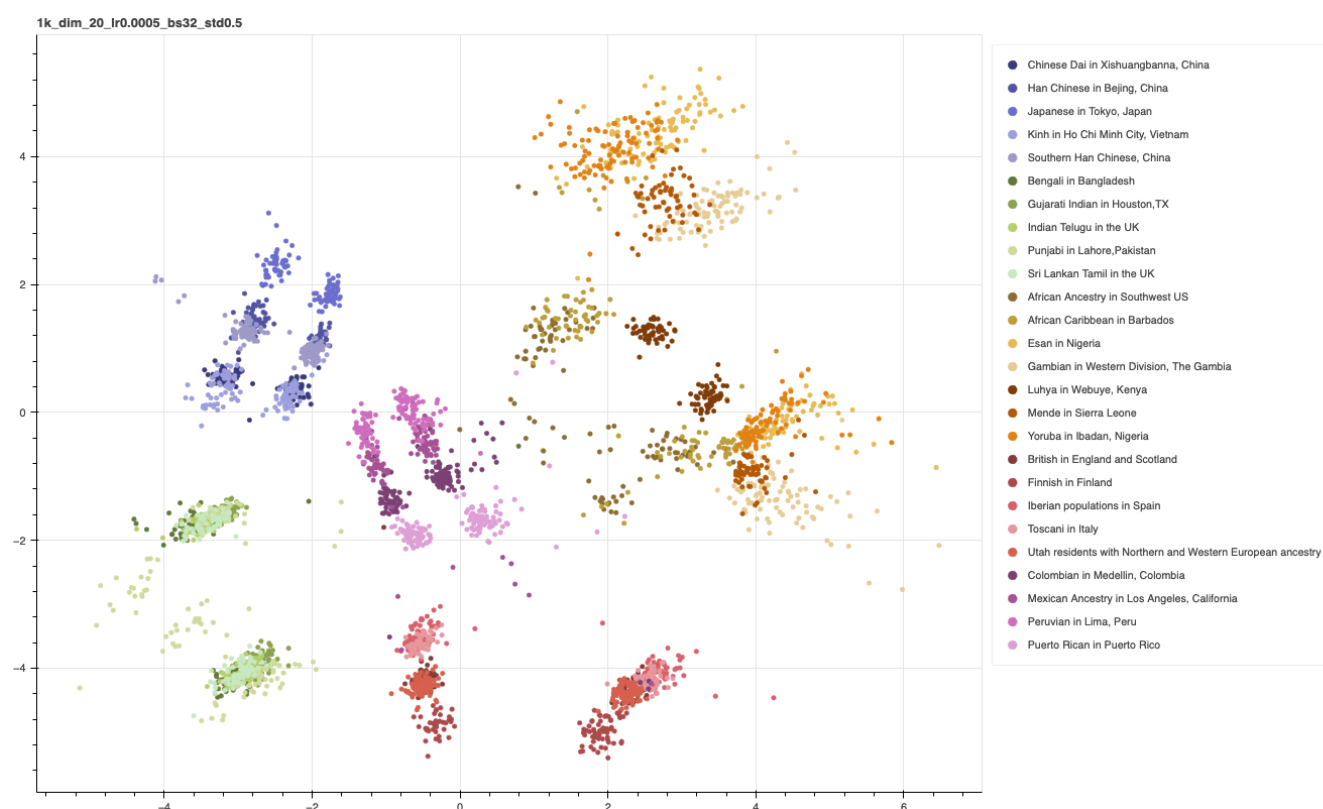
• UMAP on 15 PCs



• Vanilla Adversarial Autoencoders



- Adversarial Autoencoders with Cluster Heads

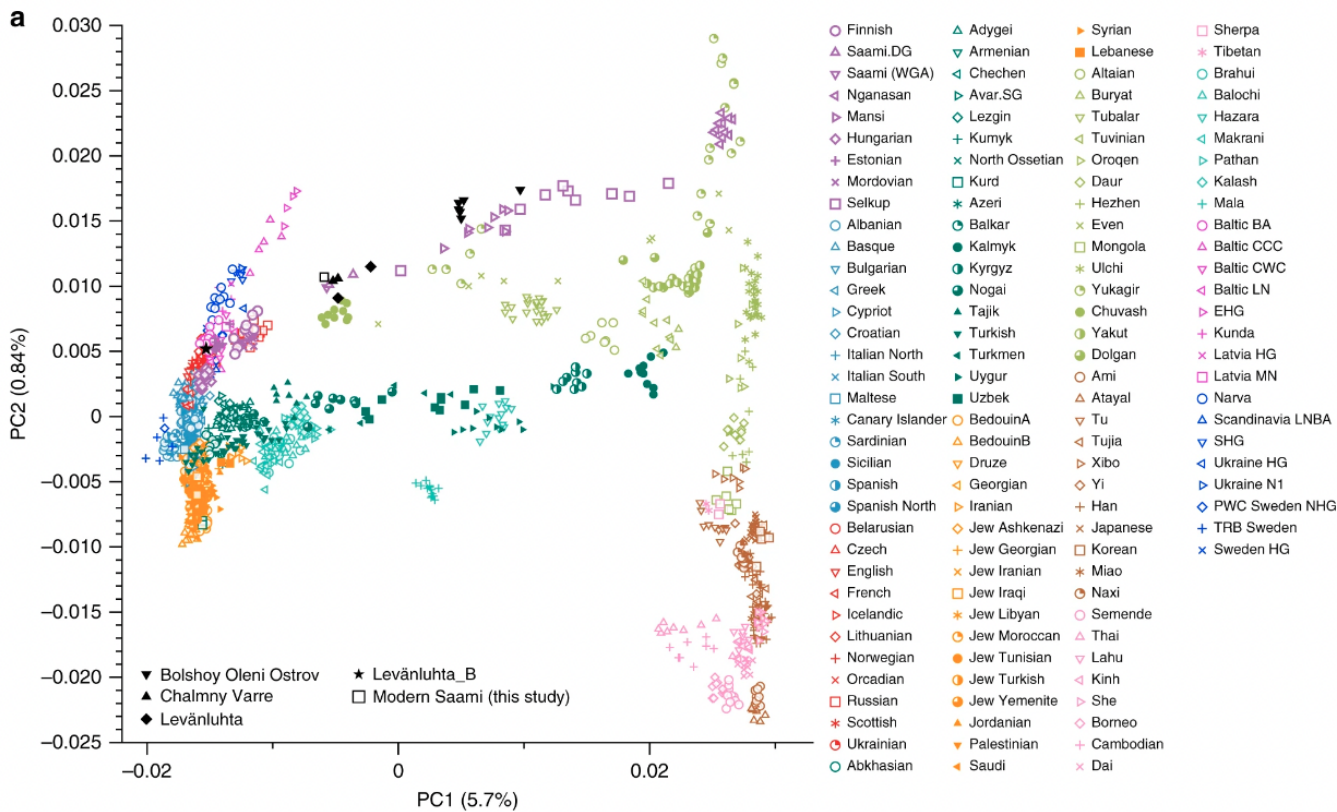


We haven't managed to exactly reproduce the results from Diaz-Papkovich et al; we're working on it.

And feel free to explore the visualizations further using the interactive htmls that can be found in the [htmls](#) folder

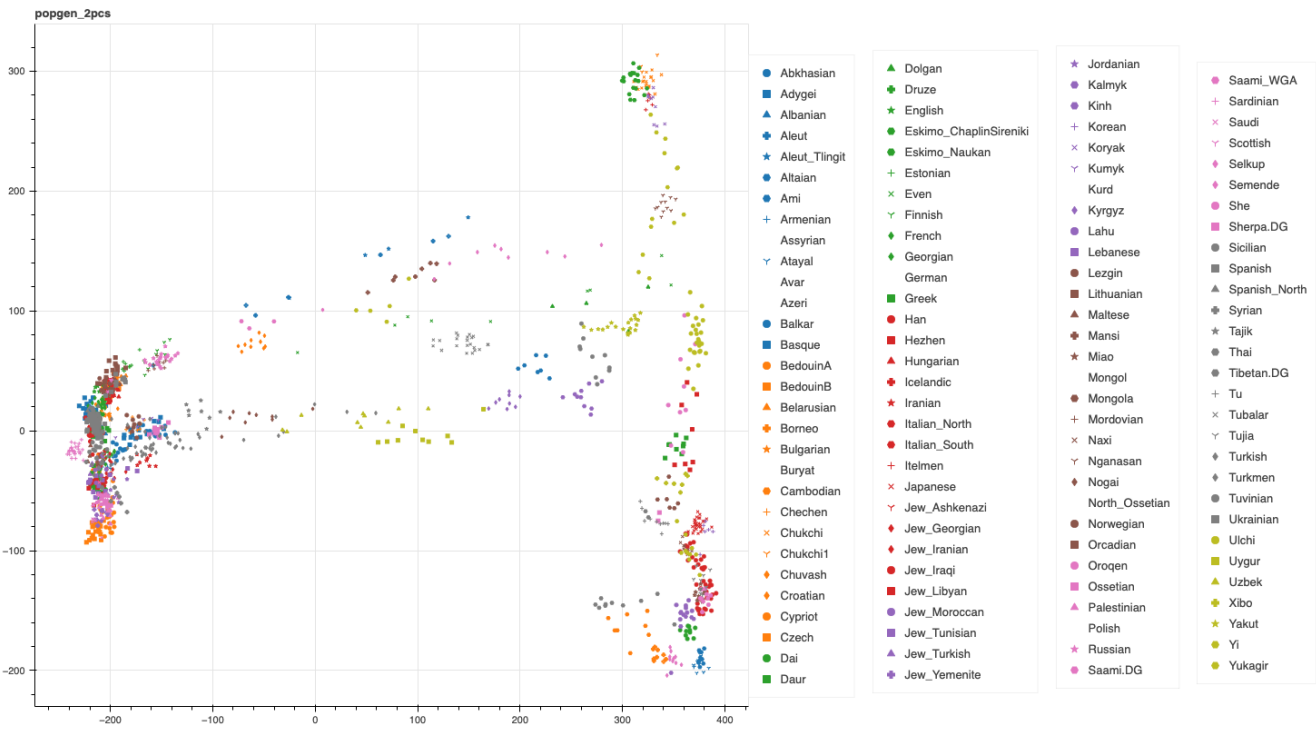
Modern Eurasian population

For reference, [Lamnidis et al.](#) have applied PCA to this dataset:

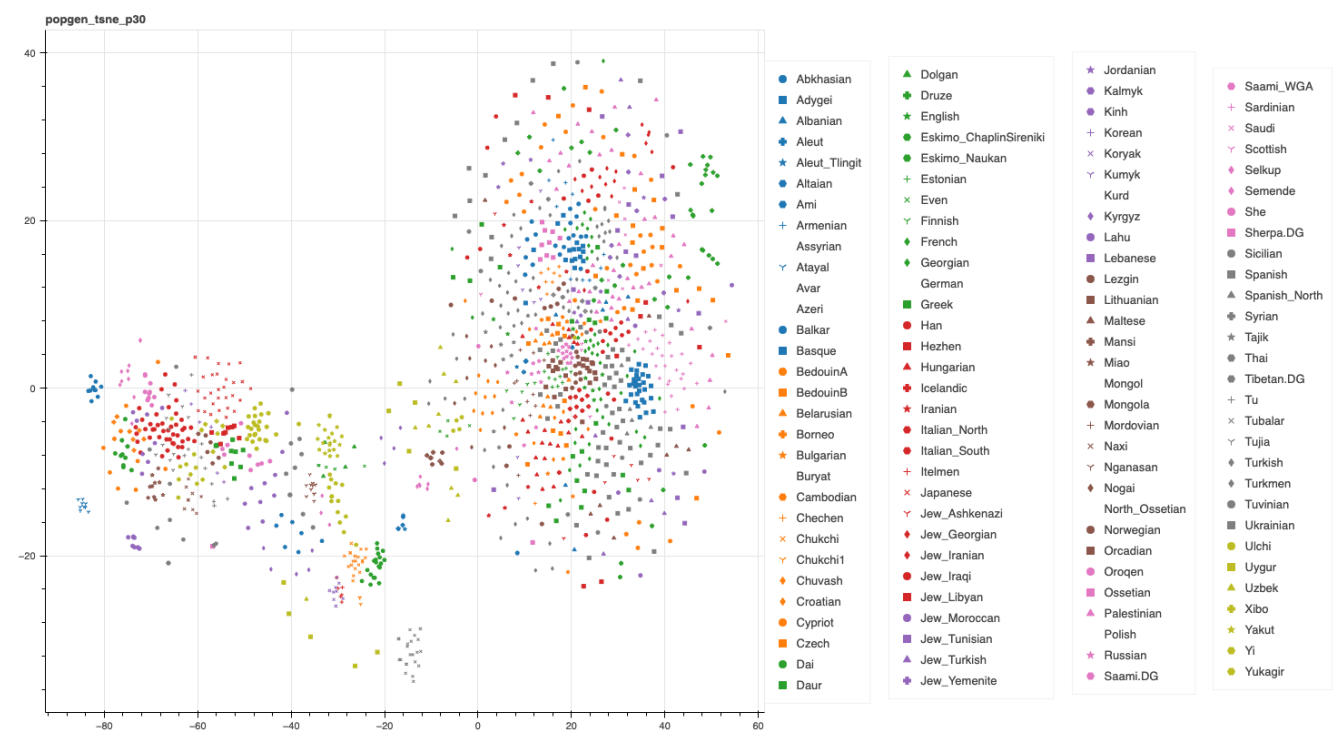


Our results are below:

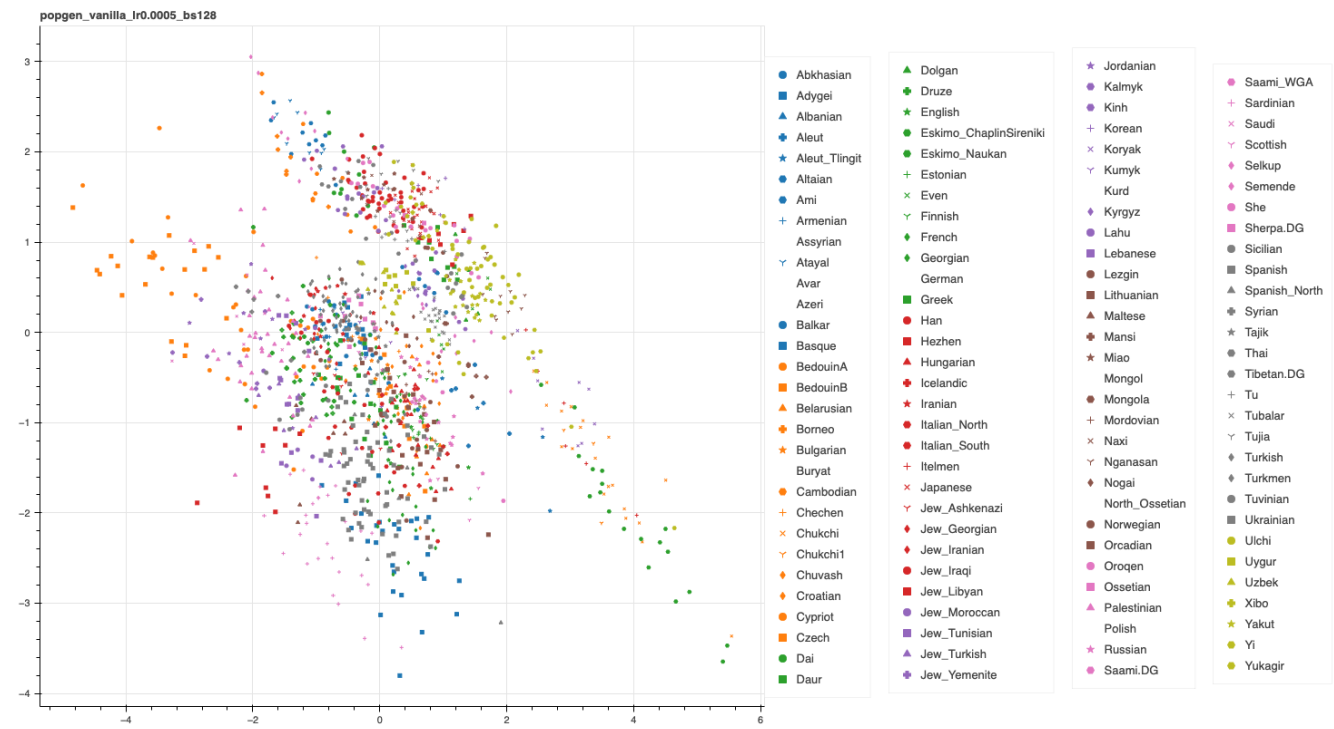
• PCA



• t-SNE



• Vanilla Adversarial Autoencoders



• Adversarial Autoencoders with Cluster Heads

