**SNVstory: A dockerized algorithm for rapid and accurate inference of sub-continental ancestry**

Audrey E. Bollas[1], Andrei Rajkovic[1], Defne Ceyhan[1], Jeffrey B. Gaither[1], Elaine R. Mardis[1,2], Peter White[1,2,#]

[1]The Steve and Cindy Rasmussen Institute for Genomic Medicine, The Abigail Wexner Research Institute, Nationwide Children's Hospital, Columbus, OH USA

[2]Department of Pediatrics, The Ohio State University College of Medicine, Columbus, OH USA

[#]Corresponding Author:

      Peter White, Ph.D.

      The Institute for Genomic Medicine

      The Abigail Wexner Research Institute

      Nationwide Children's Hospital

      575 Children's Crossroad

      Columbus, OH 43215

      Phone: +1 (614) 355-2671

      Email: peter.white@nationwidechildrens.org

24    **Abstract**

25    Knowing a patient's genetic ancestry is crucial in clinical settings, providing benefits such as tailored

26    genetic testing, targeted health screening based on ancestral disease-predisposition rates, and

27    personalized medication dosages. However, self-reported ancestry can be subjective, making it

28    difficult to apply consistently. Moreover, existing approaches utilize genome sequencing data to infer

29    ancestry at the continental level, creating the need for methods optimized for individual ancestry

30    assignment. We present SNVstory, a method built upon three independent machine learning models

31    for accurately inferring the sub-continental ancestry of individuals. SNVstory includes a feature-

32    importance scheme, unique among open-source ancestral tools, which allows the user to track the

33    ancestral signal broadcast by a given gene or locus. We apply SNVstory to a clinical dataset,

34    comparing self-reported ethnicity and race to our inferred genetic ancestry. SNVstory represents a

35    significant advance in methods to assign genetic ancestry, predicting ancestry across 36 different

36    populations with high accuracy.

37

38    **Introduction**

39    Ancestry derived from genomic data, referred to as genetic ancestry, is a measurable and biologically

40    defined parameter. Although much of the human genome is identical across all populations, it is

41    estimated that depending on an individual's ancestry, 0.1% to 0.4% may differ from the human

42    reference genome. While this genetic variation includes structural variants (SVs), copy number

43    variants (CNVs), and small insertions or deletions (indels), by far the largest and easiest to detect

44    category occurs in the form of single nucleotide variants (SNVs), many of which are unique to

45    genetically distinct populations[1].

46

47    Knowledge of a patient's genetic ancestry has clinical implications, ranging from genetic testing to

48    health screening based on ancestral disease-predisposition rates, and in some cases, may inform

2

49    what medicine dosage to prescribe a patient[2–4]. However, self-reported race is frequently used in the

50    research and clinical setting and is often inconsistent with genetic ancestry, potentially driving health

51    disparities[5–8]. Genome sequencing-based diagnostic testing in patients suspected of having a rare

52    genetic disorder requires accurate data filtering to remove variants common to a given population.

53    Precise identification of the patient's ancestry improves the identification of rare disease-causal

54    variants. Therefore, developing methods to report ancestry accurately and consistently is essential.

55

56    In addition to clinical importance, knowing the ancestral composition of an individual or a population

57    is essential in the genetic research setting. For example, signals from genome-wide association

58    studies (GWAS) or whole genome sequencing cohorts can be reassessed based on population

59    stratification, whereby loci associated with disease may be more accurately identified by discarding

60    rare variants associated with an individual's ancestry rather than with the disease in question[9,10].

61

62    Given the importance of ancestry, several ancestry inference algorithms that operate on genomic

63    data have been developed that can be divided into two broad types: parametric and non-parametric.

64    Parametric learning algorithms estimate a finite set of parameters from the data to establish a

65    relationship between the independent and dependent variables. Two widely-used parametric tools

66    are STRUCTURE[11] and ADMIXTURE[12], which estimate the proportions of different ancestries (or

67    ancestral populations) for each individual, known as admixture. Recently, Archetypal Analysis was

68    shown to be more computationally efficient and provide more interpretable results than

69    ADMIXTURE[13]. In contrast, non-parametric methods do not have a finite set of parameters and

70    instead rely on the intrinsic structure of the data to determine which data points best resemble each

71    other.

72

73 The emergence of population-scale genome sequencing datasets with a form of self-reported

74 ancestry allows models to be built with prior knowledge of represented ancestries. In place of

75 individualized genetic data, large databases house genomic summary results, such as aggregate

76 variant allele frequencies stratified by population. For example, the Single Nucleotide Polymorphism

77 database (dbSNP) is the largest genomic aggregate database with 11 different populations from over

78 one million samples[14]. However, the 11 distinct populations contain a high degree of overlap and

79 primarily represent continental groupings[15]. The Genome Aggregation Database (gnomAD) is

80 another aggregate database with allele frequencies from 140,000 subjects from 26 populations[16]. In

81 addition to these large-scale repositories of aggregate allele frequencies, there exist a few datasets at

82 the level of the individual, such as the 1000 Genomes Project (1kGP)[1] and the Simons Genome

83 Diversity Project (SGDP)[17], which are much smaller in sample size, with 2,504 and 279 samples,

84 respectively. Nevertheless, the 1kGP and SGDP have been critical in characterizing ancestry and

85 human history as they contain the most granular population labels.

86

87 Taken together, these curated variant datasets enable an alternative class of models to be used to

88 predict ancestry based upon samples labeled with known ancestry[18–28]. However, many methods

89 suffer shortcomings, including not having discrete ancestry labels beyond the main continental

90 groups or, for those methods using the 1kGP, not considering that many subjects are within the same

91 families and, therefore, fail to satisfy the principle of independent and identically distributed data. As

92 such, there is a critical need for methods to accurately predict an individual's genetic ancestry from

93 genome sequencing data by implementing supervised models.

94

95 Here, we address some limitations surrounding supervised learning of ancestry by developing three

96 independent models from gnomAD, 1kGP, and SGDP. Our models estimate ancestry from 36 different

97 populations with high accuracy. Furthermore, we provide software that enables users to run our

98    models on their data, taking the widely accepted variant call format (VCF) files as input and

99    outputting predictions and a graphical representation of the likelihood of a given genetic ancestry.

100   As a form of validation, we apply these models to our in-house clinical research dataset and correlate

101   the estimates with those of self-reported ancestry.

102

103   **Materials and Methods**

104   **Training Datasets**

105   Genomic datasets from gnomAD, 1kGP, and SGDP were processed separately (**Figure 1**), as described

106   below.  The gnomAD variants are provided on reference genome GRCh37, and the 1kGP and SGDP

107   were called on reference genome GRCh38.

108

109   **The Genome Aggregation Database (gnomAD)**

110   The gnomAD v2.1 exome and genome sequencing variant dataset provides aggregated data from 17

111   populations, meaning allele frequencies of each population for 17 million exome variants. We

112   reduced the number of input features for machine learning by following a similar protocol to the one

113   described by the MacArthur lab by filtering for high call rates, biallelic-only sites, and a frequency

114   greater than 0.1% (https://macarthurlab.org/2018/10/17/gnomad-v2-1/). After this filtering,

115   81,398 SNVs remained, formatted as a matrix of ancestries and corresponding SNV frequencies.

116

117   To obtain SNV calls for individuals, as is provided in standard VCF format, we simulated individuals

118   from each ancestry by effectively flipping a weighted coin for each individual and their respective

119   variant (**Figure 1**). This resulted in a synthetic-based matrix of samples spanning the ancestry

120   classifications in gnomAD v2.1 and SNVs, coded as reference, heterozygous, or homozygous for each

121   SNV position. Although this approach does not capture haplotypes, the simulated samples are

122   genetically typical examples of the chosen ancestry to a first approximation.

123

**The 1000 Genomes Project (1kGP)**

125 The New York Genome Center performed genome sequencing (GS) on 3,202 samples, including 602

126 trios, from the 1kGP cohort at 30x coverage, released in 2020[29]. The data were aligned to GRCh38

127 using BWA-MEM[30], and variants were called by GATK *HaplotypeCaller* (GATK version 3.5.0) using

128 default settings. The dataset contains 126,659,422 SNVs from 26 populations spanning East and

129 South Asia, North and South America, Africa, and Europe. Sample sizes were not uniformly

130 represented across the different populations, i.e., the dataset was imbalanced. Due to the high genetic

131 similarity between individuals from Utah and the United Kingdom, the Utah population was removed

132 from the analysis.

133

**The Simons Genome Diversity Project (SGDP)**

135 The SGDP consists of GS of 300 individuals from seven major population groups, 75 countries, and

136 142 diverse populations. GS FASTQ files from 279 samples were downloaded from the European

137 Nucleotide Archive (PRJEB9586). Sequencing reads were aligned to genome assembly GRCh38 using

138 BWA-MEM. SNV and INDEL calling was performed with GATK version 4.1.9, described below. GATK

139 *HaplotypeCaller* was run on each sample using the GVCF workflow to generate a per-sample

140 intermediate GVCF. The GATK *GenotypeGVCFs* function was used to perform base calling across all

141 samples jointly to obtain genotypes for each sample in VCF format. We then performed variant

142 recalibration and filtering in the two-stage process using the GATK functions *VariantRecalibration*

143 and *ApplyVQSR*. The final combined data set contained a total of 48,815,712 SNVs.

144

**Quality Control**

146 Quality control of the gnomAD (https://macarthurlab.org/2018/10/17/gnomad-v2-1/) and 1kGP[29]

147 were as previously described. For the SGDP dataset, we ran several quality-control tools to detect

148    any issues with sequencing quality and sample contamination. We ran Picard *CollectMultipleMetrics*

149    on the aligned bam files to collect alignment summary, quality score, and GC bias metrics (**Table S1**).

150    Sequencing read allocation was calculated using samtools. Coverage information was collected using

151    mosdepth[31]. The average coverage for all realigned samples was 40X (ranging from 31X to 77X).

152    Sample contamination level was determined by the number of reads inconsistent with the genotype

153    in dbSNP[14] sites. One sample was flagged for possible sample contamination (**Supplemental**

154    **Materials and Methods**).

155

156    **Removal of Related Samples**

157    Related samples of the third degree (e.g., first cousins, great grandparents, or great-grandchildren)

158    or closer were identified by the relationship inference tool, KING[32]. Data from the 1kGP and SDGP

159    were preprocessed using PLINK2 with the following parameters: *"--new-id-max-allele-len 10000 --*

160    *max-alleles 2*"[33]. KING recommends performing as little filtering as possible. However, an additional

161    filtering step was performed to prevent the computation from running out of memory. Therefore, the

162    analysis was restricted to variants shared by at least two individuals: "*--maf 0.0007*" in the case of the

163    1kGP and "*--maf 0.007*" for SDGP. After removing the variants present in only one sample, KING was

164    executed on the resulting bed file, with the "*--kinship*" option set to report pairwise relatedness

165    inference. Samples from the analysis were flagged that had a third-degree kinship coefficient cutoff

166    >= 0.0442, a value previously established by the authors of KING[32]. Four samples were removed from

167    further analysis in the SGDP dataset based on the KING relatedness results (**Supplemental Materials**

168    **and Methods**).

169

170    Because some samples from the 1kGP are related to more than one other individual in the cohort, the

171    following procedure was implemented to remove the fewest number of samples. Considering only

172    the relationships with coefficients exceeding the third-degree cutoff, a graph-based method was

173     implemented to recursively identify nodes (samples) with the largest number of edges

174     (relationships) and remove those nodes until all subgraphs had, at most, a single connection. For

175     subgraphs with a single connection, one sample was randomly selected from the pair, while all

176     singletons were included in the list of samples to keep. From 167 samples with at least one close

177     relationship, 117 were flagged for inclusion in downstream analysis. The remaining samples were

178     removed with PLINK2.

179

180     **Variant Selection and Preprocessing**

181     Variants from 1kGP and SGDP underwent a final filtering step by taking the intersection of targeted

182     exonic regions of the exome capture reagent used routinely in our clinical lab (IDT xGen Exome Hyb

183     Panel v2 targets hg38 BED file) with the set of genetic variants from the unrelated individuals using

184     BEDTools *intersect* (v2.30.0)[34]. The resulting VCF was converted into a numerical encoding

185     homozygous alternative = 2, heterozygous = 1, reference or missing = 0. The vectors of genotypes

186     were combined to form a matrix of variants by genotypes. For variant selection from gnomAD, see

187     the following gnomAD section in Model training and cross-validation below.

188

189     **Model Training and Cross-Validation**

190     The models were trained on each dataset separately, as required by their differing labeling strategies

191     (**Figure 1**).

192

193     **gnomAD:** Because our gnomAD algorithm uses synthetic data, we must consider two parameters: a

194     population size that balances the model's accuracy with training time and resources and a p-value

195     from a Chi-Square test that removes uninformative SNVs. This was accomplished using a nested for

196     loop to iterate over all combinations of population sizes and p-values for SNV removal (**Figure S1**).

197     For each combination, we generated a set of 80/20 training/validation splits of the data. A Chi-Square

198   test was applied to each SNV (feature) in the training data to determine whether it was informative

199   for distinguishing ancestry in the population. SNVs were removed that did not meet the p-value

200   threshold. We used a gradient-boosted decision tree from XGBoost to train the model on the training

201   set and then test on the validation set[35]. Fold generation and training were performed five times for

202   each p-value, and the accuracy was averaged to represent the accuracy for each p-value. Once all the

203   p-values were tested, the p-value with the highest accuracy was selected (**Figure S2**). Then, the

204   model was retrained on all the data for that specific population size and tested on a synthetic hold-

205   out set. The accuracy for the hold-out set is representative of that population. A continental model

206   (population size of 4,084 individuals; SNV p-value threshold of 7.5e-49) was built to predict six

207   groups: Africa, South Asia, Europe, East Asia, America, and Ashkenazi Jewish. Two sub-continental

208   classifiers were built to predict ancestry within the East Asian (**Figure S2A.**; population size of

209   13,593 individuals; SNV p-value threshold of 1.78e-09) and European groups (**Figure S2B.**;

210   population size of 45,243 individuals; SNV p-value threshold of 1.78e-24).

211

212   **1kGP:** For the 1kGP dataset, the support vector machine (SVM) library from scikit-learn[36] was used

213   to train a classifier to predict the continental groups: Africa, Europe, South Asia, East Asia, and

214   America. In addition, multiple classifiers were trained independently for each sub-continental group,

215   i.e., Kenya or African Caribbean in Barbados. All SVMs were trained using the radial basis function

216   (RBF) kernel and with the gamma parameter fixed as the default. Hyperparameter tuning of the C

217   penalty term was accomplished by performing cross-validation using the scikit-learn stratified k-fold

218   library. The default five splits were chosen, and the shuffle variable was set to true. The F1 macro

219   average was selected to represent a model's performance.

220

221   **SGDP:** The SVM library from scikit-learn was used to train the model for the SGDP dataset. Stratified

222   k-fold cross-validation was performed using the standard scikit-learn library. Seven continental

9

223    groups were predicted from this cohort (Africa, West Eurasia, East Asia, South Asia, Oceania, Central

224    Asia Siberia, and America), as the subcontinental groups needed more samples per group to train an

225    accurate model. The F1 macro average was chosen as a representation of a model's performance to

226    account for the imbalanced data.

227

228    **Results**

229    **Model Performance**

230    We report the performance of the gnomAD, 1kGP, and SGDP continental models using external

231    validation sets (**Figure 2A-F**), and cross-validation results on the subcontinental models (**Figures S2**

232    **and S3**) were performed because additional datasets with the same subcontinental labels were not

233    available.

234

235    Confusion matrices are shown in **Figures 2A-D**, providing the ancestry prediction for each sample in

236    the validation data. In the 1kGP and SGDP models, we see some discrepancies between the European

237    and American groups. In the case of the 1kGP model (**Figure 2A**), some SGDP samples labeled as

238    European are predicted to be American. Similarly, in the SGDP model, some 1kGP samples labeled as

239    American are predicted as European. This may be due to a higher similarity of the feature space

240    between European and American samples than other groups (**Figure S3**). The gnomAD model is

241    validated with 1kGP (**Figure 2C**) and SGDP (**Figure 2D**) samples. Overall, all continental models have

242    a high area under the curve in both ROC (**Figure 2E**) and Precision-Recall (**Figure 2F**) curves,

243    described in the figure legend.

244

245    The gnomAD East Asian and European subcontinental models have accuracies of 99.90% and

246    80.92%, respectively (**Figure S2A, B**). The results for the 1kGP subcontinental model are obtained

247    by averaging the probabilities for each sample across cross-validation folds and then computing the

10

248  confusion matrix (**Figure S4**). The accuracies for the 1kGP subcontinental models are as follows:

249  Africa, 90.26%; America, 93.06%; East Asia, 87.23%; Europe, 94.29%; South Asia, 85.86%.

250

251  **Feature Interpretation**

252  Feature importance for the gnomAD continental model was calculated using SHAP[37] values to

253  provide insight into which SNVs and their corresponding genes have the most impact on the model

254  predictions. SHAP values for the 1kGP and SGDP models were not calculated because the memory

255  requirement for the kernel explainer was too high due to the number of features in the models.

256

257  Global feature importance for the gnomAD continental model is reported by aggregating SHAP values

258  across each gene and taking the mean absolute value of each gene across 2,800 of the training

259  samples (**Figure S5**). The 'knownCanonical' genes table was downloaded from the UCSC Table

260  Browser using assembly GRCh37 to get the genomic interval for each gene. If a region contains

261  multiple genes, we combine the genes to form a non-overlapping genomic interval (e.g., ANKRD45,

262  TEX50). Of the 77,402 variants used to train the model, 3,231 were not located in gene regions and

263  were removed from further analysis. The most significant gene impacting the model is Keratin

264  Associated Protein 19-8. Samples with a variant in this gene are more likely to be predicted as

265  American.

266

267  We also aggregated SHAP values across larger cytolocations to visualize which regions across the

268  genome  are  most  impactful  in  the  model  predictions  (accessed  using  this  file:

269  (https://hgdownload.soe.ucsc.edu/goldenPath/hg19/database/cytoBand.txt.gz). **Figure S6** shows

270  the feature importance for an individual from the training data labeled as African. Regions are

271  colored by population label with the maximum absolute SHAP value. Regions that have the most

11

272    impact on predicting the sample African are 'chromosome 1: 172,900,000-176,000,000' and

273    'chromosome 5: 63,200,000_66,700,000'.

274

275    **Comparison of Genetic vs. Self-Reported Ancestry in Clinical Samples**

276    SNVstory was implemented on an in-house dataset of clinical exome sequencing testing from 293

277    individuals generated by the Institute for Genomic Medicine Clinical Laboratory to demonstrate the

278    application of our models. We compare the model predictions to the self-reported ancestry of the

279    proband (**Table S2**). Self-reported race is derived from the paternal/maternal ethnic background.

280    Ethnicity is categorized into one of three groups: Non-Hispanic or Latino, Hispanic or Latino, and

281    Unknown/Not Reported Ethnicity. Race is classified into one of five groups: White, Asian, Bi-

282    racial/Multi-racial, Black or African American, and Unknown/Unspecified. Due to the broadness of

283    these categories, we report the comparison between predicted genetic ancestry for the continental

284    models only (**Table 1**).

285

286    Most of the individuals share agreement between genetic ancestry and ethnicity/race, e.g., for those

287    predicted to be European, a match of White / Non-Hispanic or Latino for race /ethnicity occurs in

288    92.5%, 96.7%, and 89.1% of individuals by the gnomAD (**Table 1A**), 1kGP (**Table 1B**), and SGDP

289    (**Table 1C**) models, respectively. However, several cases exist where individuals are self-reported as

290    White while having a different genetic ancestry across multiple models, and vice versa. Additionally,

291    13 of our cases have either Unknown/Not Reported Ethnicity or Unknown/Unspecified Race. As

292    discussed in the Introduction, the ability to refine or add genetic ancestry information in these cases

293    is helpful for added diagnostic precision in variant filtering/prioritization.

294

295    **Model Interpretation for Indeterminant Samples**

296    Most of our in-house dataset has agreement across all three continental models (81.9% of samples)

297    and even more across at least two continental models (98.0%). A disproportionate number of

298    individuals share disagreement across all three models between those that are self-reported as Bi-

299    racial or Multi-racial vs. those that are White, Asian, Black or African American (50% vs. 9%

300    disagreement, respectively). Those individuals with Unknown/Unspecified Race are not included in

301    this calculation. These results suggest our models have worse performance on admixed samples,

302    where two or more populations may be present. In reporting results, we use the label with the highest

303    probability. Some discrepancies between model results may be mitigated by adding a minimum

304    threshold on the probability required to obtain a result.

305

306    **Individualized Ancestry Report**

307    Here, we illustrate the ability of SNVstory to provide ancestry predictions in an easily visualized

308    format for individual samples (**Figure 3**). The probabilities for the gnomAD and the 1kGP continental

309    models were 100% European, while the SGDP continental model was 95% West Eurasia. The

310    gnomAD subcontinental model has the highest probability (48%) for North-Western European

311    (nfe_nwe), and the 1kGP subcontinental model has the highest probability (100%) for British From

312    England and Scotland (eur_gbr). The subcontinental model probabilities are weighted by the

313    continental probabilities, which are returned as 0% probability for the remaining models. These

314    predictions agree with the true sample ancestry taken from the 1kGP validation set.

315

316    **Discussion**

317    We have described a method to predict ancestry from genomic data that provides multiple

318    improvements over existing ancestry inference tools. Firstly, SNVstory incorporates

319    samples/variants from three different curated datasets, expanding the number of labels and the

320    granularity of the model classification beyond the main continental divisions. Secondly, drawing

13

321 upon the gnomAD database produces a much larger number of variants on which our models were

322 trained, providing the opportunity to classify ancestry on a wider (or more diverse) range of features.

323 Thirdly, SNVstory excludes consanguineous samples from training, ensuring that the

324 overrepresentation of closely related individuals does not bias the model. Finally, our novel

325 implementation is optimized for individualized results rather than clustering large cohorts of

326 samples into shared ancestral groups.

327

328 In our gnomAD model, we introduce a method to simulate individual samples from aggregate allele

329 frequencies of a known population. This is potentially useful for any study requiring access to

330 reference variants from a population where data from individual samples is obfuscated. One

331 limitation in our approach is that we did not account for linkage disequilibrium between variants

332 when simulating individual samples. This could result in some samples with patterns of variants that

333 do not exist in actual samples. An improvement in future models would be to remove variants with

334 high levels of linkage disequilibrium between them. If high recognizability to actual samples is

335 required, established metrics of linkage disequilibrium, such as the correlation coefficient $r^2$, could

336 be used to measure the 'realness' of a simulated sample based on existing variant patterns, and

337 simulated VCFs could be validated based on this quality. However, in practice, the larger pool of

338 variants provided by gnomAD more than compensates for the lost dependence among proximal

339 groups of variants. We have demonstrated that the performance of the gnomAD models with

340 simulated individuals is comparable to that of models trained with actual samples.

341

342 With the growing number of reference datasets containing individuals from diverse ancestral

343 backgrounds, it is possible to build ancestry prediction models that reflect these populations.

344 However, there is room for improvement, as our most diverse dataset (SGDP) includes the fewest

345 samples. We could not build subcontinental models as granular as the labels provided because there

346  were as few as two samples per label for many instances. Additionally, our model cannot accurately

347  predict ancestry proportions in samples with admixed ancestry. Most admixture prediction software

348  depends on a priori knowledge of the number of non-admixed populations and requires

349  representation from such populations. There is limited availability of reference samples from

350  admixed individuals, so our training data lacked representation from any admixed samples. Efforts

351  to expand the number of reference sequences for diverse and admixed populations will provide

352  opportunities to fill this gap.

353

354  SNVstory's feature-importance capacity is unique among ancestral tools and could have significant

355  clinical utility. The clinical application of most ancestral prediction tools is limited to simply

356  predicting the patient's ancestry. However, SNVstory's unique capability to describe a given locus as

357  characteristic, or atypical, of a given ancestry could lead to improved prioritization of variants. For

358  example, SNVstory finds the most ancestrally informative gene on average to be KRTAP19-8, which

359  is greatly enriched for SNVs predictive of Native American/Latino ancestry (**Figure S5**). This gene is

360  a known driver of thyroid lymphoma[38], a disorder that is the second-most-common type of cancer

361  among Hispanic women[39] but not even among the top five cancer types among women worldwide[40].

362  The inferred distinctiveness of Latino copies of KRTAP19-8 suggests that rare founder mutations in

363  this gene may contribute to increased rates of thyroid cancer among women of Hispanic ancestry.

364  The ability to target variants in genes inherited from specific populations adds a new tool to the

365  diagnostician's toolkit and could lead to improved patient outcomes.

366

367  Finally, our approach allows users to reliably execute our models given a single-sample or multi-

368  sample VCF, with results tailored toward ancestry assignment for an individual sample. This provides

369  immediately useful ancestry information in the clinical setting, where ancestry can be used to inform

370  diagnostic or therapeutic decisions. Specifically, a subject's ancestry can be used to help prioritize

371     variants that may be rare in one population but not another. In the clinical setting, it may be essential

372     to recognize the difference between ethnicity, race, and genetic ancestry in determining the optimal

373     therapy or drug dosage.

374

375     Given the widespread availability of genome sequencing data and models like SNVstory that can

376     accurately predict ancestry, we advocate for genetic ancestry to become the standard classification

377     reported for genetic studies and clinical applications, where appropriate. Genetic ancestry offers

378     enormous advantages over other self-reported information, such as ethnicity or race, because it

379     supplies biological characteristics of a population and is consistently measurable. This advantage will

380     only increase as more populations are sequenced and ancestry prediction becomes more reliable,

381     and we improve our ability to contextualize the impact of genetic ancestry on clinical decision-

382     making.

383

384     **Acknowledgments**

387

388     **Author Contributions**

389     AB and AR processed data and trained models for gnomAD, 1kGP, and SGDP. AR designed the

390     methods to simulate data from gnomAD allele frequencies and cross-validation architecture. AB and

391     DC prepared figures and tables. AB wrote the first draft of the paper. JG, DC, AR, and PW assisted in

392     preparing or revising the paper. AR and AB wrote the SNVstory software package. PW and EM

393     supervised the project.

394

395     **Data Availability**

396 The training data for our model are available as follows. gnomAD v2.1 data is available from

397 https://gnomad.broadinstitute.org/downloads/. 1000 Genomes Project data is shared via the

398 International Genome Sample Resource and can be accessed from

399 https://www.internationalgenome.org/data-portal/data-collection/30x-grch38. Simons Genome

400 Diversity Project data is available from the European Nucleotide Archive under project PRJEB9586.

401 SNVstory is an open-source model and is available from https://github.com/nch-igm/snvstory.

402

403 **Funding**

404 This work was supported by the Nationwide Children's Foundation and The Abigail Wexner Research

405 Institute at Nationwide Children's. The funders had no role in study design, data collection, data

406 analysis, the decision to publish, or manuscript preparation.

407

408 **Ethics Approval and Consent to Participate**

409 This study was reviewed and approved by the Institutional Review Board (IRB) of The Abigail

410 Wexner Research Institute at Nationwide Children's Hospital (Office for Human Research Protections

411 (OHRP) IORG0000326; IRB00000568) as IRB17-00206 ("Institute for Genomic Medicine

412 Comprehensive Profiling for Cancer, Blood, and Somatic Disorders"). The participant's legal

413 guardian/next of kin provided written informed consent to participate in this study.

414

415 **Competing Interests**

416 No Competing interests: Audrey Bollas, Andrei Rajkovic, Defne Ceyhan, Jeffrey Gaither, and Peter

417 White. Elaine Mardis: Qiagen N.V., supervisory board member, honorarium, and stock-based

418 compensation. Singular Genomics Systems, Inc., board of directors, honorarium, and stock-based

419 compensation.

420

17

**References**

1. Auton, A., Abecasis, G.R., Altshuler, D.M., Durbin, R.M., Abecasis, G.R., Bentley, D.R., Chakravarti, A., Clark, A.G., Donnelly, P., Eichler, E.E., et al. (2015). A global reference for human genetic variation. Nature *526*, 68–74. 10.1038/nature15393.

2. Hauser, D., Obeng, A.O., Fei, K., Ramos, M.A., and Horowitz, C.R. (2018). Views Of Primary Care Providers On Testing Patients For Genetic Risks For Common Chronic Diseases. Health Aff. Proj. Hope *37*, 793–800. 10.1377/hlthaff.2017.1548.

3. Jorde, L.B., and Bamshad, M.J. (2020). Genetic Ancestry Testing What Is It and Why Is It Important? JAMA *323*, 1089–1090. 10.1001/jama.2020.0517.

4. Ramamoorthy, A., Pacanowski, M.A., Bull, J., and Zhang, L. (2015). Racial/ethnic differences in drug disposition and response: review of recently approved drugs. Clin. Pharmacol. Ther. *97*, 263–273. 10.1002/cpt.61.

5. Fujimura, J.H., and Rajagopalan, R. (2011). Different differences: The use of 'genetic ancestry' versus race in biomedical human genetic research. Soc. Stud. Sci. *41*, 5–30.

6. Shraga, R., Yarnall, S., Elango, S., Manoharan, A., Rodriguez, S.A., Bristow, S.L., Kumar, N., Niknazar, M., Hoffman, D., Ghadir, S., et al. (2017). Evaluating genetic ancestry and self-reported ethnicity in the context of carrier screening. BMC Genet. *18*, 99. 10.1186/s12863-017-0570-y.

7. Mersha, T.B., and Abebe, T. (2015). Self-reported race/ethnicity in the age of genomic research: its potential impact on understanding health disparities. Hum. Genomics *9*, 1. 10.1186/s40246-014-0023-x.

8. Gomes, M.B., Gabrielli, A.B., Santos, D.C., Pizarro, M.H., Barros, B.S.V., Negrato, C.A., Dib, S.A., Porto, L.C., and Silva, D.A. (2018). Self-reported color-race and genomic ancestry in an admixed population: A contribution of a nationwide survey in patients with type 1 diabetes in Brazil. Diabetes Res. Clin. Pract. *140*, 245–252. 10.1016/j.diabres.2018.03.021.

9. Brown, R., Lee, H., Eskin, A., Kichaev, G., Lohmueller, K.E., Reversade, B., Nelson, S.F., and Pasaniuc, B. (2016). Leveraging ancestry to improve causal variant identification in exome sequencing for monogenic disorders. Eur. J. Hum. Genet. *24*, 113–119. 10.1038/ejhg.2015.68.

10. Purcell, S.M., Wray, N.R., Stone, J.L., Visscher, P.M., O'Donovan, M.C., Sullivan, P.F., Sklar, P., Purcell (Leader), S.M., Stone, J.L., Sullivan, P.F., et al. (2009). Common polygenic variation contributes to risk of schizophrenia and bipolar disorder. Nature *460*, 748–752. 10.1038/nature08185.

11. Pritchard, J.K., Stephens, M., and Donnelly, P. (2000). Inference of population structure using multilocus genotype data. Genetics *155*, 945–959. 10.1093/genetics/155.2.945.

12. Alexander, D.H., Novembre, J., and Lange, K. (2009). Fast model-based estimation of ancestry in unrelated individuals. Genome Res. *19*, 1655–1664. 10.1101/gr.094052.109.

13. Gimbernat-Mayol, J., Mantes, A.D., Bustamante, C.D., Montserrat, D.M., and Ioannidis, A.G. (2022). Archetypal Analysis for population genetics. PLOS Comput. Biol. *18*, e1010301. 10.1371/journal.pcbi.1010301.

18

459   14. Sherry, S.T., Ward, M.-H., Kholodov, M., Baker, J., Phan, L., Smigielski, E.M., and Sirotkin, K.
460       (2001). dbSNP: the NCBI database of genetic variation. Nucleic Acids Res. *29*, 308–311.

461   15. Jin, Y., Schaffer, A.A., Feolo, M., Holmes, J.B., and Kattman, B.L. (2019). GRAF-pop: A Fast
462       Distance-Based Method To Infer Subject Ancestry from Multiple Genotype Datasets Without
463       Principal Components Analysis. G3 Bethesda Md *9*, 2447–2461. 10.1534/g3.118.200925.

464   16. Karczewski, K.J., Francioli, L.C., Tiao, G., Cummings, B.B., Alföldi, J., Wang, Q., Collins, R.L.,
465       Laricchia, K.M., Ganna, A., Birnbaum, D.P., et al. (2020). The mutational constraint spectrum
466       quantified from variation in 141,456 humans. Nature *581*, 434–443. 10.1038/s41586-020-
467       2308-7.

468   17. Mallick, S., Li, H., Lipson, M., Mathieson, I., Gymrek, M., Racimo, F., Zhao, M., Chennagiri, N.,
469       Nordenfelt, S., Tandon, A., et al. (2016). The Simons Genome Diversity Project: 300 genomes
470       from 142 diverse populations. Nature *538*, 201–206. 10.1038/nature18964.

471   18. Kumar, A., Montserrat, D.M., Bustamante, C., and Ioannidis, A. (2020). XGMix: Local-Ancestry
472       Inference with Stacked XGBoost (Genomics) 10.1101/2020.04.21.053876.

473   19. Maples, B.K., Gravel, S., Kenny, E.E., and Bustamante, C.D. (2013). RFMix: A Discriminative
474       Modeling Approach for Rapid and Robust Local-Ancestry Inference. Am. J. Hum. Genet. *93*, 278–
475       288. 10.1016/j.ajhg.2013.06.020.

476   20. Sheehan, S., and Song, Y.S. (2016). Deep Learning for Population Genetic Inference. PLOS
477       Comput. Biol. *12*, e1004845. 10.1371/journal.pcbi.1004845.

478   21. Hwa, H.-L., Wu, M.-Y., Lin, C.-P., Hsieh, W.H., Yin, H.-I., Lee, T.-T., and Lee, J.C.-I. (2019). A single
479       nucleotide polymorphism panel for individual identification and ancestry assignment in
480       Caucasians and four East and Southeast Asian populations using a machine learning classifier.
481       Forensic Sci. Med. Pathol. *15*, 67–74. 10.1007/s12024-018-0071-y.

482   22. Durand, E.Y., Do, C.B., Mountain, J.L., and Macpherson, J.M. (2014). Ancestry Composition: A
483       Novel, Efficient Pipeline for Ancestry Deconvolution (Bioinformatics) 10.1101/010512.

484   23. Chu, B.B., Sobel, E.M., Wasiolek, R., Ko, S., Sinsheimer, J.S., Zhou, H., and Lange, K. (2021). A fast
485       Data-Driven method for genotype imputation, phasing, and local ancestry inference:
486       MendelImpute.jl. Bioinforma. Oxf. Engl., btab489. 10.1093/bioinformatics/btab489.

487   24. Shi, G., and Kuang, Q. (2021). Ancestral Spectrum Analysis With Population-Specific Variants.
488       Front. Genet. *12*.

489   25. Wang, Y., Song, S., Schraiber, J.G., Sedghifar, A., Byrnes, J.K., Turissini, D.A., Hong, E.L., Ball, C.A.,
490       and Noto, K. (2021). Ancestry inference using reference labeled clusters of haplotypes. BMC
491       Bioinformatics *22*, 459. 10.1186/s12859-021-04350-x.

492   26. Soumare, H., Rezgui, S., Gmati, N., and Benkahla, A. (2021). New neural network classification
493       method for individuals ancestry prediction from SNPs data. BioData Min. *14*, 30.
494       10.1186/s13040-021-00258-7.

495   27. Dalfovo, D., and Romanel, A. (2023). Analysis of Genetic Ancestry from NGS Data Using EthSEQ.
496       Curr. Protoc. *3*, e663. 10.1002/cpz1.663.

497    28. Karim, M.R., Cochez, M., Zappa, A., Sahay, R., Beyan, O., Schuhmann, D.-R., and Decker, S. (2020).
498        Convolutional Embedded Networks for Population Scale Clustering and Bio-ancestry
499        Inferencing.

500    29. Byrska-Bishop, M., Evani, U.S., Zhao, X., Basile, A.O., Regier, A.A., Corvelo, A., Clarke, W.E.,
501        Musunuri, R., Fairley, S., Runnels, A., et al. High coverage whole genome sequencing of the
502        expanded 1000 Genomes Project cohort including 602 trios. 45.

503    30. Li, H. (2013). Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
504        10.48550/arXiv.1303.3997.

505    31. Pedersen, B.S., and Quinlan, A.R. (2018). Mosdepth: quick coverage calculation for genomes and
506        exomes. Bioinformatics *34*, 867–868. 10.1093/bioinformatics/btx699.

507    32. Manichaikul, A., Mychaleckyj, J.C., Rich, S.S., Daly, K., Sale, M., and Chen, W.-M. (2010). Robust
508        relationship inference in genome-wide association studies. Bioinformatics *26*, 2867–2873.
509        10.1093/bioinformatics/btq559.

510    33. Chang, C.C., Chow, C.C., Tellier, L.C., Vattikuti, S., Purcell, S.M., and Lee, J.J. (2015). Second-
511        generation PLINK: rising to the challenge of larger and richer datasets. GigaScience *4*, s13742-
512        015-0047–0048. 10.1186/s13742-015-0047-8.

513    34. Quinlan, A.R., and Hall, I.M. (2010). BEDTools: a flexible suite of utilities for comparing genomic
514        features. Bioinformatics *26*, 841–842. 10.1093/bioinformatics/btq033.

515    35. Chen, T., and Guestrin, C. (2016). XGBoost: A Scalable Tree Boosting System. In Proceedings of
516        the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining, pp.
517        785–794. 10.1145/2939672.2939785.

518    36. Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,
519        Prettenhofer, P., Weiss, R., Dubourg, V., et al. (2011). Scikit-learn: Machine Learning in Python. J.
520        Mach. Learn. Res. *12*, 2825–2830.

521    37. Lundberg, S., and Lee, S.-I. (2017). A Unified Approach to Interpreting Model Predictions.
522        10.48550/arXiv.1705.07874.

523    38. Qiu, K., Li, K., Zeng, T., Liao, Y., Min, J., Zhang, N., Peng, M., Kong, W., and Chen, L. (2021).
524        Integrative Analyses of Genes Associated with Hashimoto's Thyroiditis. J. Immunol. Res. *2021*,
525        8263829. 10.1155/2021/8263829.

526    39. Estrada-Florez, A.P., Bohórquez, M.E., Sahasrabudhe, R., Prieto, R., Lott, P., Duque, C.S., Donado,
527        J., Mateus, G., Bolaños, F., Vélez, A., et al. (2016). Clinical features of Hispanic thyroid cancer
528        cases and the role of known genetic variants on disease risk. Medicine (Baltimore) *95*, e4148.
529        10.1097/MD.0000000000004148.

530    40. Ferlay, J., Colombet, M., Soerjomataram, I., Parkin, D.M., Piñeros, M., Znaor, A., and Bray, F.
531        (2021). Cancer statistics for the year 2020: An overview. Int. J. Cancer. 10.1002/ijc.33588.

532

533

534   **Figure Titles and Legends**

535   **Figure 1. Schematic of ancestry inference model strategy.** The workflow visualizes each dataset

536   separately with colored boxes and arrows: gnomAD (blue), 1kGP (yellow), and SGDP (red). For the

537   gnomAD synthetic-based matrix, allele frequencies for each variant for each population given in

538   gnomAD are used to create a distribution of reference, heterozygous and homozygous alleles for each

539   population. A matrix format is created by converting the distributions into 0's, 1's, and 2's for each

540   locus for samples in each population. For 1kGP and SGDP, a matrix format is built directly from

541   variants in the VCF. For the model architecture, continental model labels are shown in white boxes,

542   and the number of labels in the corresponding subcontinental models is below in brackets.

543

544   **Figure 2. Continental ancestry inference model performance.** A-D. Confusion matrices of the

545   1kGP model using SGDP as validation (A), SGDP model using 1kGP as validation (B), gnomAD model

546   using 1kGP as validation (C), and gnomAD model using SGDP as validation (D). E. Macro-averaged

547   ROC curves. F. Macro-averaged precision-recall curves.

548

549   **Figure 3. SNVstory ancestry report.** The representative output of model results from SNVstory for

550   a European sample taken from the 1kGP dataset.

551 **Tables**

552 **Table 1. Genetic ancestry versus self-reported ethnicity and race.** Value counts of genetic

553 ancestry model predictions trained using gnomad **(A)**, 1kGP **(B)**, and SGDP **(C)** compared to self-

554 reported ethnicity and race.

555     **A. gnomAD**

| Model Labels | Ethnicity | Race | Counts |
|---|---|---|---|
| afr | Non-Hispanic or Latino | Black or African American | 20 |
| | | Bi-racial/Multi-racial | 10 |
| | Unknown/Not Reported Ethnicity | Bi-racial/Multi-racial | 3 |
| | Hispanic or Latino | Bi-racial/Multi-racial | 2 |
| | | White | 1 |
| | Non-Hispanic or Latino | White | 1 |
| amr | Hispanic or Latino | White | 8 |
| | | Unknown/Unspecified | 5 |
| | | Black or African American | 2 |
| | Non-Hispanic or Latino | White | 1 |
| | Hispanic or Latino | Bi-racial/Multi-racial | 1 |
| asj | Non-Hispanic or Latino | White | 1 |
| eas | Non-Hispanic or Latino | Asian | 3 |
| | | White | 2 |
| | Hispanic or Latino | Bi-racial/Multi-racial | 1 |
| eur | Non-Hispanic or Latino | White | 210 |
| | | Bi-racial/Multi-racial | 5 |
| | Hispanic or Latino | Bi-racial/Multi-racial | 5 |
| | | White | 3 |
| | Unknown/Not Reported Ethnicity | White | 3 |
| | | Bi-racial/Multi-racial | 1 |
| | Hispanic or Latino | Unknown/Unspecified | 1 |
| sas | Non-Hispanic or Latino | Asian | 3 |
| | | White | 1 |

556

22

557     **B.  1kGP**

| Model Labels | Ethnicity | Race | Counts |
|---|---|---|---|
| afr | Non-Hispanic or Latino | Black or African American | 19 |
| | | Bi-racial/Multi-racial | 2 |
| | Hispanic or Latino | Bi-racial/Multi-racial | 1 |
| | Unknown/Not Reported Ethnicity | Bi-racial/Multi-racial | 1 |
| amr | Hispanic or Latino | White | 12 |
| | Non-Hispanic or Latino | Bi-racial/Multi-racial | 10 |
| | Hispanic or Latino | Bi-racial/Multi-racial | 8 |
| | Non-Hispanic or Latino | White | 8 |
| | Hispanic or Latino | Unknown/Unspecified | 6 |
| | Hispanic or Latino | Black or African American | 2 |
| | Unknown/Not Reported Ethnicity | Bi-racial/Multi-racial | 2 |
| | Non-Hispanic or Latino | Black or African American | 1 |
| eas | Non-Hispanic or Latino | Asian | 3 |
| eur | Non-Hispanic or Latino | White | 207 |
| | Unknown/Not Reported Ethnicity | White | 3 |
| | Non-Hispanic or Latino | Bi-racial/Multi-racial | 3 |
| | Unknown/Not Reported Ethnicity | Bi-racial/Multi-racial | 1 |
| sas | Non-Hispanic or Latino | Asian | 3 |
| | | White | 1 |

558

23

559    **C.  SGDP**

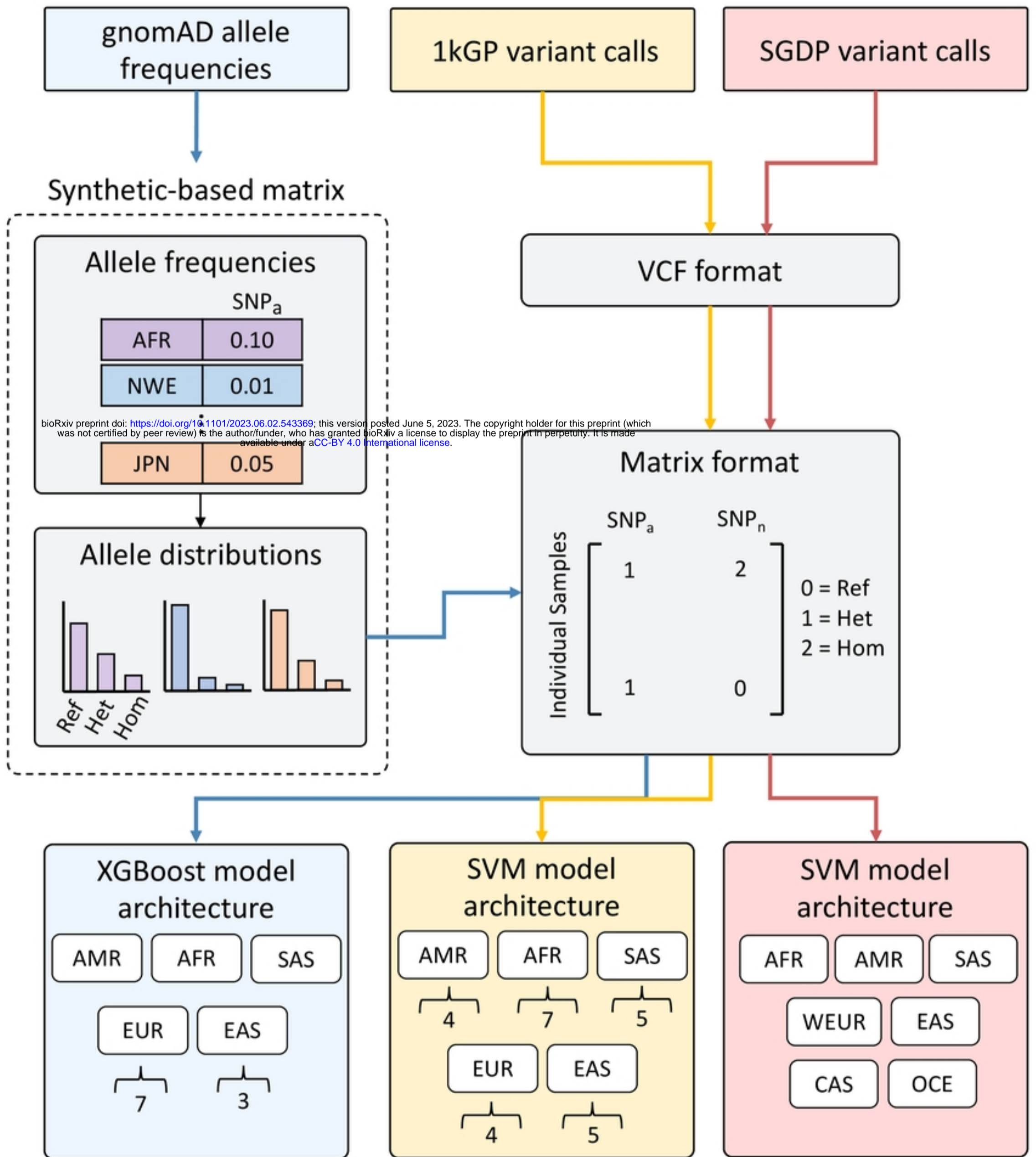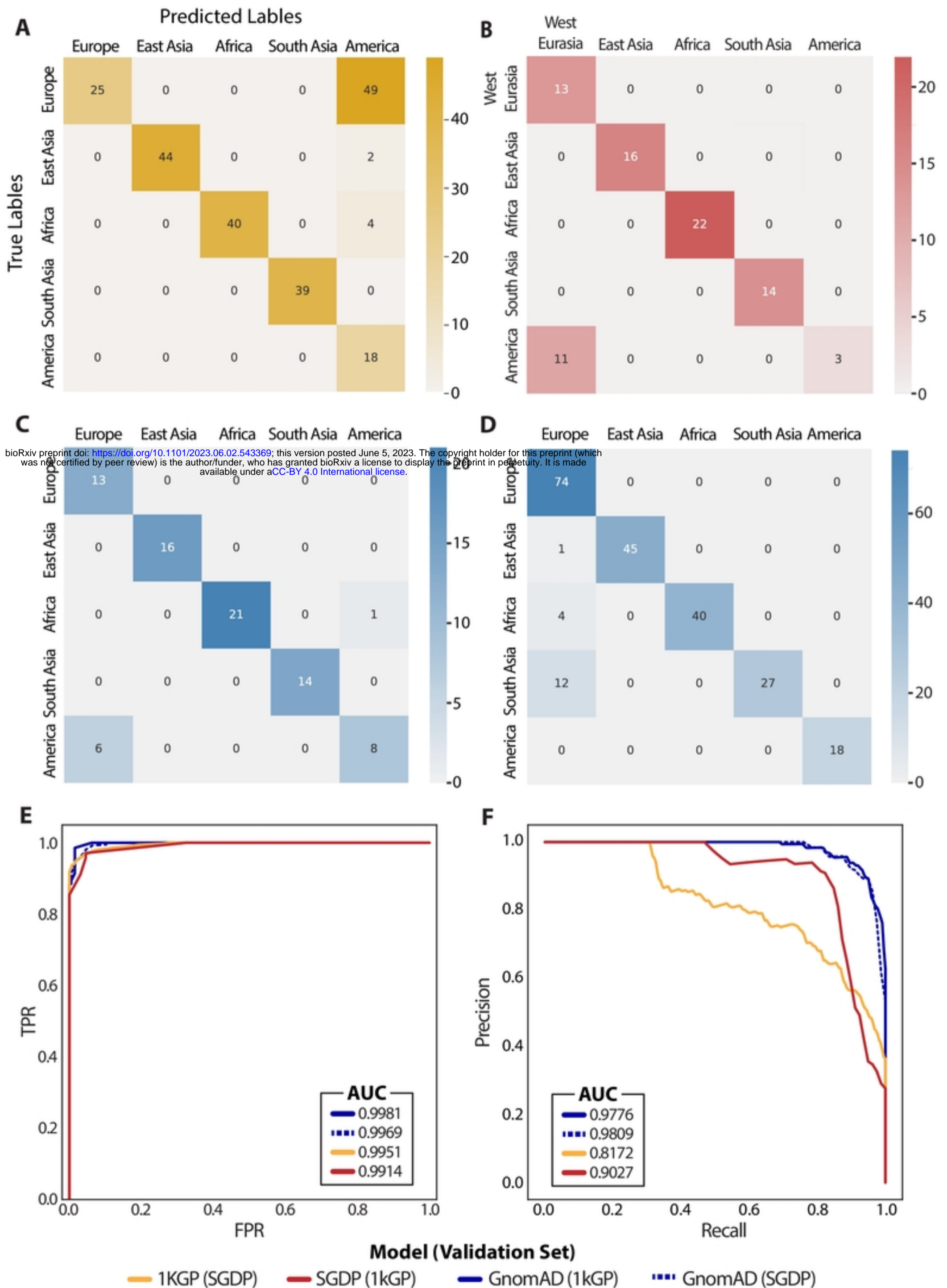| Model Labels | Ethnicity | Race | Counts |
|---|---|---|---|
| Africa | Non-Hispanic or Latino | Black or African American | 20 |
| | | Bi-racial/Multi-racial | 9 |
| | Hispanic or Latino | Bi-racial/Multi-racial | 3 |
| | Unknown/Not Reported Ethnicity | Bi-racial/Multi-racial | 3 |
| | Hispanic or Latino | Black or African American | 2 |
| | Non-Hispanic or Latino | White | 1 |
| CentralAsiaSiberia | Hispanic or Latino | Unknown/Unspecified | 3 |
| | | White | 1 |
| EastAsia | Non-Hispanic or Latino | Asian | 3 |
| SouthAsia | Hispanic or Latino | White | 4 |
| | Non-Hispanic or Latino | Asian | 3 |
| | | White | 3 |
| WestEurasia | Non-Hispanic or Latino | White | 212 |
| | Hispanic or Latino | White | 7 |
| | Non-Hispanic or Latino | Bi-racial/Multi-racial | 6 |
| | Hispanic or Latino | Bi-racial/Multi-racial | 6 |
| | | Unknown/Unspecified | 3 |
| | Unknown/Not Reported Ethnicity | White | 3 |
| | | Bi-racial/Multi-racial | 1 |

560

Figure 1

Figure 2

# HG00096

Figure 3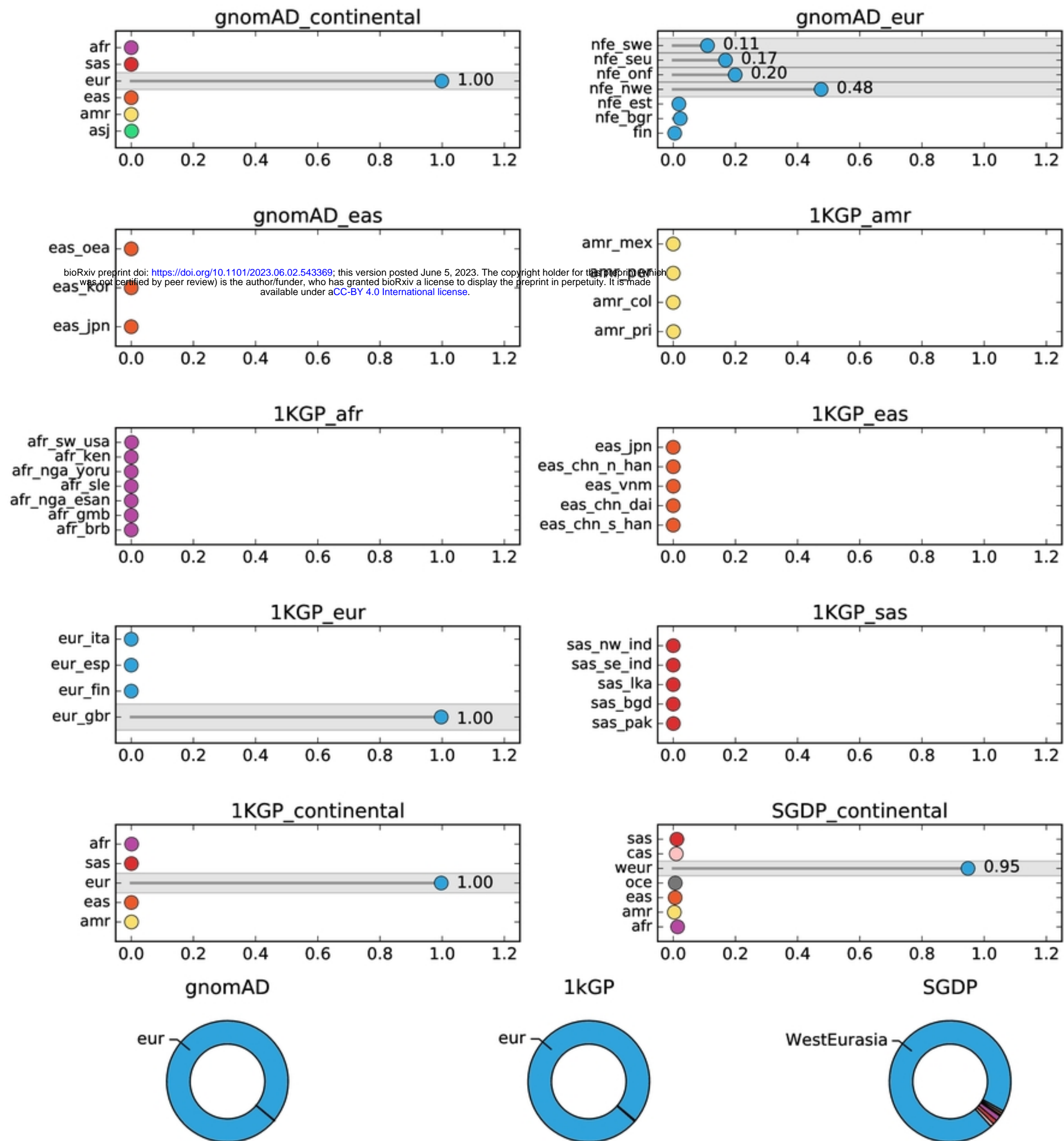