# Creating artificial human genomes using generative neural networks

Burak Yelmen, Aurélien Decelle, Linda Ongaro, Davide Marnetto, Corentin Tallec, Francesco Montinaro, Cyril Furtlehner, Luca Pagani, Flora Jay

# Creating Artificial Human Genomes Using Generative Neural Networks

**Authors:** Burak Yelmen[1,2,3*], Aurélien Decelle[3,4], Linda Ongaro[1,2], Davide Marnetto[1], Corentin Tallec[3], Francesco Montinaro[1,5], Cyril Furtlehner[3], Luca Pagani[1,6], Flora Jay[3,*]

**Affiliations:**

*1 Institute of Genomics, University of Tartu, Estonia*

*2 Institute of Molecular and Cell Biology, University of Tartu, Estonia*

*3 Laboratoire de Recherche en Informatique, CNRS UMR 8623, Université Paris-Sud,*

*Université Paris-Saclay, Orsay, France*

*4 Departamento de Física Téorica I, Universidad Complutense, 28040 Madrid, Spain*

*5 Department of Biology-Genetics, University of Bari, Bari, Italy*

*6 APE Lab, Department of Biology, University of Padova, Italy*

*[*]to whom correspondence should be addressed: burakyelmen@gmail.com, flora.jay@lri.fr*

16  **Abstract**

17  Generative models have shown breakthroughs in a wide spectrum of domains due to

18  recent advancements in machine learning algorithms and increased computational

19  power. Despite these impressive achievements, the ability of generative models to

20  create realistic synthetic data is still under-exploited in genetics and absent from

21  population genetics. Yet a known limitation in the field is the reduced access to many

22  genetic databases due to concerns about violations of individual privacy, although they

23  would provide a rich resource for data mining and integration towards advancing

24  genetic studies. In this study, we demonstrated that deep generative adversarial

25  networks (GANs) and restricted Boltzmann machines (RBMs) can be trained to learn

26  the complex distributions of real genomic datasets and generate novel high-quality

27  artificial genomes (AGs) with none to little privacy loss.

28

29  We show that our generated AGs replicate characteristics of the source dataset such

30  as allele frequencies, linkage disequilibrium, pairwise haplotype distances and

31  population structure. Moreover, they can also inherit complex features such as signals

32  of selection. To illustrate the promising outcomes of our method, we showed

33  that  imputation quality for low frequency alleles can be improved by augmenting

34  reference panels with AGs and that the RBM latent space provides a relevant

35  encoding of the data, hence allowing further exploration of the reference dataset and

36  features for solving supervised tasks.

37

38  Generative models and AGs have the potential to become valuable assets in genetic

39  studies by providing a rich yet compact representation of existing genomes and high-

40  quality, easy-access and anonymous alternatives for private databases.

41  **Author Summary**

42  Generative neural networks have been effectively used in many different domains in

43  the last decade, including machine dreamt photo-realistic imagery. In our work, we

44  apply a similar concept to genetic data to automatically learn its structure and, for the

45  first time, produce high quality realistic genomes. These novel genomes are distinct

46  from the original ones used for training the generative networks. We show that artificial

47  genomes, as we name them, retain many complex characteristics of real genomes

48  and the heterogeneous relationships between individuals. They can be used in

49  intricate analyses such as imputation of missing data as we demonstrated. We believe

50  they have a high potential to become alternatives for many genome databases which

51  are not publicly available or require long application procedures or collaborations and

52  remove an important accessibility barrier in genomic research in particular for

53  underrepresented populations.

## Introduction

55  Availability of genetic data has increased tremendously due to advances in

56  sequencing technologies and reduced costs (1). The vast amount of human genetic

57  data is used in a wide range of fields, from medicine to evolution. Despite the

58  advances, cost is still a limiting factor and more data is always welcome, especially in

59  population genetics and genome-wide association studies (GWAS) which usually

60  require substantial amounts of samples. Partially related to the costs but also to the

61  research bias toward studying populations of European ancestry, many

62  autochthonous populations are under-represented in genetic databases, diminishing

63  the extent of the resolution in many studies (2–5). Additionally, the majority of the data

64  held by government institutions and private companies is considered sensitive and not

65  easily accessible due to privacy issues, exhibiting yet another barrier for scientific

66  work. A class of machine learning methods called generative models might provide a

67  suitable solution to these problems.

68

69  Generative models are used in unsupervised machine learning to discover intrinsic

70  properties of data and produce new data points based on those. In the last decade,

71  generative models have been studied and applied in many domains of machine

72  learning (6–8). There have also been a few applications in the genetics field (9–12),

73  one specific study focusing on generating DNA sequences via deep generative models

74  to capture protein binding properties (13).  Among the various generative approaches,

75  we focus on two of them in this study, generative adversarial networks (GANs) and

76  restricted Boltzmann machines (RBMs). GANs are generative neural networks which

77  are capable of learning complex data distributions in a variety of domains (14). A GAN

78  consists of two neural networks, a generator and a discriminator, which compete in a

79  zero-sum game (Supplementary Figure 1). During training, the generator produces

80  new instances while the discriminator evaluates their authenticity. The training

81  objective consists in learning the data distribution in a way such that the new instances

82  created by the generator cannot be distinguished from true data by the discriminator.

83  Since their first introduction, there have been several successful applications of GANs,

84  ranging from generating high quality realistic imagery to gap filling in texts (15,16).

85  GANs are currently the state-of-the-art models for generating realistic images (17).

86

87  A restricted Boltzmann machine, initially called Harmonium, is another generative

88  model which is a type of neural network capable of learning probability distributions

89  through input data (18,19). RBMs are two-layer neural networks consisting of an input

90  (visible) layer and a hidden layer (Supplementary Figure 2). The learning procedure

91  for the RBM consists in maximizing the likelihood function over the visible variables of

92  the model. This procedure is done by adjusting the weights such that the correlations

93  between the visible and hidden variables on both the dataset and sampled

94  configurations from the RBM converge. Then RBM models recreate data in an

95  unsupervised manner through many forward and backward passes between these two

96  layers (Gibbs sampling), corresponding to sampling from the learned distribution. The

97  output of the hidden layer goes through an activation function, which in return becomes

98  the input for the hidden layer. Although mostly overshadowed by recently introduced

99  approaches such as GANs or Variational Autoencoders (20), RBMs have been used

100 effectively for different tasks (such as collaborative filtering for recommender systems,

101 image or document classification) and are the main components of deep belief

102 networks (21–23).

103

104    Here we propose and compare a prototype GAN model along with an RBM model to

105    create Artificial Genomes (AGs) which can mimic real genomes and capture

106    population structure along with other characteristics of real genomes. These

107    prototypes are compared to alternative generative models based on multiple

108    summaries of the data and we explore whether a meaningful encoding of real data

109    has been learned. Finally, we investigate the potential of using AGs as proxies for

110    private datasets that are not accessible in order to address various genomic tasks

111    such as imputation or selection scans.

# Results

**Reconstructing genome wide population structure:**

Initially we created AGs with GAN, RBM, and two simple generative models for comparison: a Bernoulli and a Markov chain model (see Materials & Methods) using 2504 individuals (5008 haplotypes) from 1000 Genomes data (24), spanning 805 SNPs from all chromosomes which reflect a high proportion of the population structure present in the whole dataset (25). Both GAN and RBM models capture a good portion of the population structure present in 1000 Genomes data while the other two models could only produce instances centered around 0 on principal component analysis (PCA) space (Figure 1). All major modes, corresponding to African, European and Asian genomes, are well represented in AGs produced by GAN and RBM models and absent for the Markovian and Bernouilli models.  Wasserstein distances between the 2D PCA representations of real versus generated individuals were closer to 0 for GAN (0.006), RBM (0.006) and the test set (0.01) than for the Markovian (0.124) and Bernoulli (0.240) models. Uniform manifold approximation and projection (UMAP) mapping results (performed on the combined dataset) lead to similar conclusions (Wasserstein 2D distance from real for GAN: 0.021, RBM: 0.091, Markovian: 0.088, Bernoulli: 0.127) although the RBM distribution is slightly shifted compare to the real one (Supplementary Figure 3). We additionally computed the distribution of pairwise differences of haploid genomes within a single dataset or between the real and artificial datasets (Supplementary Figure 4). The real, GAN and RBM distributions present similar multimodal patterns reflecting the underlying population structure (in particular the RBM distribution exhibits three modes corresponding to the average distances between European and Asian, European and African, or African and Asian populations. The overall real pairwise distribution is captured as accurately by the GAN

137   (Wasserstein distance between real and generated distributions: 3.24) and RBM

138   (6.21) models than by a test set (5.06) and those clearly outperform the binomial

139   (42.20) and Markovian (37.92) models. No real genome was copied into the AGs (0

140   identical pair). Since GANs and RBMs showed an excellent performance for this use

141   case, we further explored other characteristics using only these two models.

142

143   **Reconstructing local high-density haplotype structure:**

144   To evaluate if high quality artificial dense genome sequences can also be created by

145   generative models, we applied the GAN and RBM methods to a 10K SNP region using

146   (i) the same individuals from 1000 Genomes data and (ii) 1000 Estonian individuals

147   from the high coverage Estonian Biobank (26) to generate artificial genomes. PCA

148   results of AGs spanning consecutive 10K SNPs show that both GAN and RBM models

149   can still capture the relatively toned-down population structure (Supplementary Figure

150   5; 2D Wasserstein distances for 1000 Genomes and Estonian respectively: 0.004 and

151   0.011 for GAN, 0.011 and 0.006 for RBM, 0.004 and 0.002 for test sets) as well as the

152   overall distribution of pairwise distances (Supplementary Figure 6). Looking at the

153   allele frequency comparison between real and artificial genomes, we see that

154   especially GAN performs poorly for low frequency alleles, due to a lack of

155   representation of these alleles in the AGs whereas RBM seems to have wider

156   distribution over all frequencies (Supplementary Figure 7; correlation between real and

157   generated for 1000 Genomes and Estonian respectively: 0.99 and 0.91 for GAN, 0.94

158   and 0.94 for RBM, 0.99 and 0.99 for test sets). The overall pairwise distributions are

159   fitted better by the RBM than the GAN (Wasserstein distance 117 and 227 for GAN,

160   38 and 26 for RBM, 22 and 16 for test sets). On the other hand, the distribution of the

161   distance of real genomes to the closest AG neighbour shows that GAN model,

162    although slightly underfitting, outperforms RBM model, for which an excess of small

163    distances points towards overfitting, visible by the distribution being closer to the zero

164    point (Supplementary Figure 8).

165

166    Additionally, we performed linkage disequilibrium (LD) analyses comparing artificial

167    and real genomes to assess how successfully the AGs imitate short and long range

168    correlations between SNPs. Pairwise LD matrices for real and artificial genomes all

169    show a similar block pattern demonstrating that GAN and RBM accurately captured

170    the overall structure with SNPs having higher linkage in specific regions (Figure 2a).

171    However, plotting LD as a function of the SNP distance showed that all models capture

172    weaker correlation, with RBM outperforming the GAN model perhaps due to its slightly

173    overfitting characteristic (Figure 2b). However, correlations between real and

174    generated LD across all pairs were similar for GAN and RBM (for 1000 Genomes and

175    Estonian respectively: 0.95 and 0.97 for GAN, 0.94 and 0.98 for RBM) and slightly

176    lower than for test sets (0.99 and 1.0) (Supplementary Figure 9). LD can be seen as

177    a two-point correlation statistic, we also investigated 3-point correlation statistics, that

178    represent the amount of correlation between triplets of SNPs and thus characterize

179    more complex correlation patterns in datasets (Supplementary Figure 10). To further

180    determine the haplotypic integrity of AGs, we performed ChromoPainter (27) and

181    Haplostrips (28) analyses on AGs created using Estonians as the training data. We

182    did not observe separate clustering of real and artificial genomes with Haplostrips

183    (Supplementary Figure 11). However, the majority of the AGs produced via GAN

184    model displayed an excess of short chunks when painted against 1000 Genomes

185    individuals, whereas we do not observe this for RBM AGs (Supplementary Figure 12).

186    Average European chunk length over 100 individuals for GAN AGs was 44.21%, RBM

187    AGs was 54.92%, whereas for real Estonian genomes, it was 62.83%.

188

189    After demonstrating that our models generated realistic AGs according to the

190    described summary statistics, we investigated further whether they respected privacy

191    by measuring the extent of overfitting. We calculated two metrics of resemblance and

192    privacy, the nearest neighbour adversarial accuracy ($AA_{TS}$) and privacy loss presented

193    in a recent study (29).  $AA_{TS}$ score measures whether two datasets were generated by

194    the same distribution based on the distances between all data points and their nearest

195    neighbours in each set. When applied to artificial and real datasets, a score between

196    0.5 and 1 indicates underfitting, between 0 and 0.5 overfitting (and likely privacy loss),

197    and exactly 0.5 indicates that the datasets are indistinguishable. By using an additional

198    real test set, it is also possible to calculate a privacy loss score that is positive in case

199    of information leakage, negative otherwise, and approximately ranges from -0.5 to 0.5.

200    Computed on our generated data, both scores support haplotypic pairwise difference

201    results confirming low privacy loss for GAN AGs with a score similar to the one of an

202    independent Estonian test set never used during training (GAN: 0.027 ; Test: 0.021)

203    and the overfitting nature of RBM AGs with a high risk of privacy leakage (RBM privacy

204    loss: 0.327; Supplementary Figure 13). Using an alternative sampling scheme for the

205    RBM (see Material and Methods) slightly reduced privacy loss (restrained under 0.2

206    for low number of epochs; Supplementary Figure 14). A dataset produced by this

207    alternative scheme had only a slight decrease in quality of the summary statistics while

208    the privacy loss was reduced to 0.1. For this scheme, the correlation between

209    generated and true allele frequencies was 0.92 (instead of 0.95 for the previous RBM

210    and 0.98 for GAN), the correlation for LD values was  0.97 (RBM:0.98, GAN:0.97), the

211    2D-Wasserstein distance for the PCA representations was 0.026 (RBM: 0.006, GAN:

212    0.011, RBM sampling initialized randomly: 0.339), the Wasserstein distance for the

213    pairwise distribution was 97 (RBM: 26, GAN: 227, RBM sampling initialized randomly:

214    689).

215

216    **Selection tests:**

217    We additionally performed cross population extended haplotype homozygosity (XP-

218    EHH) and population branch statistic (PBS) on a 3348 SNP region homogenously

219    dispersed over chromosome 15 to assess if AGs can also capture selection signals.

220    Both XP-EHH and PBS results provided high correlation between the scores of real

221    and artificial genomes (Figure 3). Similar peaks were observed for real and artificial

222    genome datasets (see Discussion).

223

224    **Imputation:**

225    Since it has been shown in previous studies that imputation scores can be improved

226    using additional population specific reference panels (30,31), as a possible future use

227    case, we tried imputing real Estonian genomes using 1000 Genomes reference panel

228    and additional artificial reference panels with Impute2 software (32). Both combined

229    RBM AG and combined GAN AG panels outperformed 1000 Genomes panel for the

230    lowest MAF bin (for MAF < 0.05, 2.5% and 4.4% improvement respectively) which had

231    5926 SNPs out of 9230 total (Figure 4). Also mean info metric over all SNPs were

232    intermediate between the regular imputation scheme (1000 Genomes panel only) and

233    the 'perfect' scheme (panel including private Estonian samples). The scores were

234    1.3%, 2.3%, and 6.9% higher for the combined RBM, GAN and real Estonian panels

235    respectively, compared to the panel with only 1000 Genomes samples. However,

236    aside from the lowest MAF bin, 1000 Genomes panel slightly outperformed both

237    concatenated AG panels for all the higher bins (by 0.05% to 0.6% depending on the

238    frequency bin). This might be a manifestation of haplotypic deformities in AGs that

239    might have disrupted the imputation algorithm.

240

241    **Data encoding and visualization via RBM model:**

242    Furthermore, similarly to tSNE and UMAP, RBMs perform a non-linear dimension

243    reduction of the data and provides a suitable representation of a genomic dataset as

244    a by-product based on the non-linear feature space associated to the hidden layer

245    (Materials & Methods). As Diaz-Papkovich et al (33), we found that the RBM

246    representation differs from the linear PCA ones (Supplementary Figure 15), although

247    the general structure identified by the two lower rank components is highly similar.

248    Like in a PCA, African, East Asian, and to a lesser extent, European populations stand

249    out on the two first components yet the relative positions differ slightly from PCA to

250    RBM. In particular, the Finnish appear slightly more isolated from the other European

251    populations on the first component of the RBM. South Asians are located at the center

252    separated from Europeans, partially overlapping with American populations, and stand

253    out at dimension 5 and higher (versus 3 for the PCA). The third RBM component

254    exhibit a stronger gradient than PCA for Peruvian and Mexican individuals and might

255    reflect their gradient in Native American ancestry. Finally, RBM still exhibits population

256    structure in components higher than 7, contrary to PCA. Interestingly when screening

257    the hidden node activations, we observed that different populations or groups activate

258    different hidden nodes, each one representing a specific combination of SNPs, thereby

259    confirming that the hidden layer provides a meaningful encoding of the data

260    (Supplementary Figure 16).

261

262 **Comparison with alternative generative models:**

263 We additionally performed tests to compare AGs with advanced methods used to

264 generate genomes. One such method is the copying model (34) implemented in

265 HAPGEN2 (35). Although genomes generated via this approach performed very well

266 in terms of SFS, LD and PCA, there was extensive overfitting and privacy loss and

267 multiple individuals (747 identical haplotypes) were directly copied from the original

268 dataset (Supplementary Figure 17).

269

270 Another commonly used approach to generate genomes is coalescent simulations.

271 Although it is inherently difficult to make a fair comparison since coalescent

272 simulations require additional (demographic) parameters and do not provide the

273 desired one-to-one SNP correspondence (see Discussion), we compared SFS and LD

274 decay of AGs with genomes simulated via a previously inferred demographic model

275 (36) using HapMapII genetic map (37) implemented in stdpopsim (38–41). Initially, we

276 performed PCA and checked the allele frequency distribution compared to real

277 genomes (Supplementary Figure 18). The reasoning behind PCA was to demonstrate

278 that coalescent simulation genomes cannot be combined with real genomes since they

279 exist in different planes. Since we had selected SNPs for 1000 Genomes and Estonian

280 datasets to be overlapping, we removed alleles below 0.1 frequency from all datasets

281 to eliminate biases and analyzed LD decay and allele frequency spectrum

282 (Supplementary Figure 19). For both summary statistics, coalescent simulation

283 genomes performed well. Still, direct comparison of frequencies SNPs by SNPs, LD

284 pairs by pairs, PCA, $AA_{TS}$ or distributions of pairwise distances between real and

285 generated data are not feasible for coalescent simulations. Notably, the demographic

286 model we adopted was optimized for another European population (CEU from the

287    1000 Genomes Project), since an in depth study of the demographic properties of

288    Estonians, our target population, required extensive efforts beyond the scope of this

289    paper and in themselves a cost to be considered when adopting coalescent

290    simulations as a generative model.

## Discussion

291

292    In this study, we applied generative models to produce artificial genomes and

293    evaluated their characteristics. To the best of our knowledge, this is the first application

294    of GAN and RBM models in the population genetics context, displaying an overall

295    promising applicability. We showed that population structure and frequency-based

296    features of real populations can successfully be preserved in AGs created using GAN

297    and RBM models. Furthermore, both models can be applied to sparse or dense SNP

298    data given a large enough number of training individuals. Our different trials showed

299    that the minimum required number of individuals for training is highly variable (i.e. to

300    avoid training failures such as mode collapse or incomplete training without converging

301    to an ideal mode) depending on the unknown dimension of the dataset, which is linked

302    to the type of data to be generated and the population(s). Since haplotype data is more

303    informative, we created haplotypes for the analyses but we also demonstrated that the

304    GAN model can be applied to genotype data too, by simply combining two haplotypes

305    if the training data is not phased (see Materials & Methods). In addition, we showed

306    that it might be possible to generate AGs with simple phenotypic traits through

307    genotype data (see Supplementary Table and Supplementary Text). Even though

308    there were only two simple classes, blue and brown eye color phenotypes, generative

309    models can be improved in the future to hold the capability to produce artificial datasets

310    combining AGs with multiple phenotypes

311

312    One major drawback of the proposed models is that, due to computational limitations,

313    they cannot yet be harnessed to create whole artificial genomes but rather snippets or

314    sequential dense chunks. It should be possible to create whole genomes by

315    independently training and generating multiple chunks from different genomic regions

316    using a single uniform population such as Estonians and stitching them together to

317    create a longer, genome-like, sequence for each AG individual. To mitigate possible

318    disruptions in the long-range haplotype structure, these chunks can be defined based

319    on "approximately independent LD blocks" (42). Yet for data with population structure,

320    it would be difficult to decide which combination of chunks can form a single genome.

321    Statistics such as FST or generative models conditioned on group labels might be

322    utilized to overcome this issue. On the other hand, a collection of chunks covering the

323    whole genome can be used for analyses based solely on allele frequencies without

324    any stitching. A technically different approach would be to adapt convolutional GANs

325    for AG generation (43).

326

327    Another problem arose due to rare alleles, especially for the GAN model. We showed

328    that nearly half of the alleles become fixed in the GAN AGs in the 10K SNP dataset,

329    whereas RBM AGs capture more of the rare alleles present in real genomes

330    (Supplementary Figure 20). A known issue in GAN training is mode collapse (44),

331    which occurs when the generator fails to cover the full support of the data distribution.

332    This type of failure could explain the inability of GANs to generate rare alleles. For

333    some applications that depend on rare alleles, GAN models less sensitive to mode

334    dropping may be a promising alternative (45–47).

335

336    An important use case for the AGs in the future might be creating public versions of

337    private genome banks. Through enhancements in scientific and technology

338    knowledge, genetic data becomes more and more sensitive in terms of privacy. AGs

339    might offer a versatile solution to this delicate issue in the future, protecting the

340    anonymity of real individuals. They can be utilized as input for downstream operations

341    such as forward steps of a specific evolutionary process for which they can become

342    variations of the real datasets (similar to bootstrap), or they can be the sole input when

343    the real dataset is not accessible. Initializing a simulation with real data is a procedure

344    that is commonly used in population genetics (48,49). Our results showed that GAN

345    AGs are possibly underfitting while, on the contrary, RBM AGs are overfitting, based

346    on distribution of minimum distance to the closest neighbour (Supplementary Figure

347    8) and $AA_{TS}$ scores (Supplementary Figure 13a), although we showed how overfitting

348    could be restrained by integrating $AA_{TS}$ scores within our models as a criterion for early

349    stopping in training (before the networks start overfitting) and by modifying the RBM

350    sampling scheme. In the context of the privacy issue, GAN AGs have a slight

351    advantage since underfitting and low leakage information is preferable. More distant

352    AGs would hypothetically be harder to be traced back to the original genomes. We

353    also tested the sensitivity of the $AA_{TS}$ score and privacy loss (Supplementary Figure

354    21). It appears that both scores are affected very slightly when we add only a few real

355    genomes to the AG dataset from the training set. Although this case is easily

356    detectable by examining the extreme left tail of the pairwise distribution, it advocates

357    for combining multiple privacy loss criteria and developing other sensitive

358    measurement techniques for better assessment of generated AGs. Additionally, even

359    though we did not detect exact copies of real genomes in AG sets created either by

360    RBM or GAN models, it is a very complicated task to determine if the generated

361    samples can be traced back to the originals. Reliable measurements need to be

362    developed in the future to assure complete anonymity of the source individuals given

363    the released AGs. In particular, we will investigate whether the differential privacy

364    framework is performant in the context of large population genomics datasets (50,51).

365

366    Imputation results demonstrated promising outcomes especially for population specific

367    low frequency alleles. However, imputation with both RBM and GAN AGs integrated

368    reference panels showed slight decrease of info metric for higher frequency alleles

369    compared to only 1000 Genomes panel (Figure 4). Increasing the number of AGs did

370    not affect the results significantly. We initially speculated that this decrease might be

371    related to the disturbance in haplotypic structure and therefore, tried to filter AGs

372    based on chunk counts from ChromoPainter results, preserving only AGs which are

373    below the average chunk count of real genomes. The reasoning behind this was to

374    preserve most realistic AGs with undisturbed chunks. Even with this filtering, slight

375    decrease in higher MAF bins was still present. Yet results of implementation with AGs

376    for low frequency alleles and without AGs for high frequency ones could be combined

377    to achieve the best performance. Although being not very practical in its current form,

378    future improved models can become very useful, largely for GWAS studies in which

379    imputation is a common application to increase resolution. Different generative models

380    such as MaskGAN (16) which demonstrated good results in text gap filling might also

381    be adapted for genetic imputation. RBM is possibly another option to be used as an

382    imputation tool directly by itself, since once the weights have been learned, it is

383    possible to fix a subset of the visible variables and to compute the average values of

384    the unobserved ones by sampling the probability distribution (in fact, it is even easier

385    than sampling entirely new configurations since the fixed subset of variables will

386    accelerate the convergence of the sampling algorithm).

387

388    Scans for detecting selection are another promising use case for AGs as high-fidelity

389    alternatives to real genomes. The XP-EHH and PBS scores computed on AGs were

390    highly correlated with the scores of real genomes. In particular, the highest peak we

391   obtained for Estonian genomes was also present in AGs, although it was the second

392   highest peak in RBM XP-EHH plot (Figure 3). This peak falls within the range of skin

393   color associated *SLC24A5* gene, which is potentially under positive selection in many

394   European populations (52).

395

396   As an additional feature, training an RBM to model the data distribution gives access

397   to a latent encoding of data points, providing a potentially easier to use representation

398   of data (Supplementary Figure 15). Future works could enhance our current GAN

399   model to also provide an encoding mechanism, in the spirit of (53), (54) or (55). It is

400   expected that these interpretable representations of the data will be relevant for

401   downstream tasks (54) and can be used as a starting point for various population

402   genetics analyses.

403

404   We want to highlight that AGs are created without requiring the knowledge of the

405   underlying evolutionary history, or the pre-processing bioinformatic pipelines (SNP

406   ascertainment, data filtering). Unlike coalescent simulations, for which one needs to

407   control parameters, AGs in their current form are solely constructed on raw information

408   of real genomes. Our method offers a direct way to generate artificial genomes for any

409   original dataset. On the other hand, the genomes generated using a coalescent

410   simulator required substantial upstream work (from previous studies) as they were

411   based on an explicitly parameterized model that had been inferred on real data using

412   advanced methods for demographic reconstruction. In particular, this approach is not

413   suitable when we want to generate AGs for highly complex datasets (eg full 1000

414   Genomes) for which it is arduous to infer a full evolutionary model accurately fitting the

415   data and even more so, to mimic all the biases induced by potentially unknown

416    bioinformatic pipelines. Last but not least, this coalescent generated data cannot be

417    merged directly with real public genomes because there is no direct correspondence

418    between the real SNPs and those generated, and coalescent approaches might

419    struggle to match, among other things, real complex patterns of LD (35). To

420    summarize, while the classical coalescent simulator only allows unconditional

421    sampling of a new haplotype $h$ from a predefined distribution $P(h|\theta)$ where the

422    demographic parameters $\theta$ have to be given, our generative models learn how to

423    generate $h$ from the conditional sampling distribution (CSD) $P(h|h_1,...h_n)$, where $(h_i)$

424    are the observed haplotypes. Computing, approximating or sampling from this CSD is

425    known to be a difficult task (34,56,57).

426

427    We believe it will be possible in the future to extend our approach with conditional

428    GAN/RBM methods to allow fine control over the composition of artificial datasets

429    based on (i) additional labels such as population names or any environmental

430    covariate, or (ii) evolutionary parameters. While the former is based only on real

431    datasets, the latter requires training on genetic simulations (coalescent-based or

432    forward) and has a different goal: it may provide an alternative simulator and/or permit

433    inference of evolutionary models.

434

435    We envision three main applications of our generative methods: (i) improve the

436    performance of genomic tasks such as imputation, ancestry reconstruction, GWAS

437    studies, by augmenting public genomic panels with AGs that serve as proxies for

438    private datasets that are not accessible; (ii) enable preliminary genomic analyses and

439    proof-of-concept before committing to long term application protocols and/or to

440    facilitate future collaborations to access private datasets; (iii) use the encoding of real

441    data learned by generative models as a starting input of various tasks, such as

442    recombination, demography or selection inference  or yet unknown tasks.

443

444    Although there are currently some limitations, generative models will most likely

445    become prominent for genetic research in the near future with many promising

446    applications. In this work, we demonstrated the first possible implementations and use

447    of AGs, particularly to be used as realistic surrogates of real genomes which can be

448    accessed publicly without privacy concerns.

## Materials & Methods

**Data:**

We used 2504 individual genomes from 1000 Genomes Project (1000 Genomes Project Consortium 2015) and 1000 individuals from Estonian Biobank (26) to create artificial genomes (AGs). Additional 2000 Estonian genomes were used as a test dataset. Another Estonian dataset consisting of 8678 individuals which were not used in training were used for imputation via Impute2 software (32). Analyses were applied to a highly differentiated 805 SNPs selected as a subset from (25), 3348 SNPs dispersed over the whole chromosome 15 and a dense 10000 SNP range/region from chromosome 15. In the data format we used, rows are individuals/haplotypes (instances) and columns are positions/SNPs (features). Each allele at each position is represented either by 0 or 1. In the case of phased data (haplotypes), each column is one position whereas in the case of unphased data, each two column corresponds to a single position with alleles from two chromosomes. Genomes from Estonian Biobank were accessed with Approval Number 285/T-13 obtained on 17/09/2018 by the University of Tartu Ethics Committee.

**GAN model:**

We implemented the GAN model using python-3.6, Keras 2.2.4 deep learning library with TensorFlow backend (58), pandas 0.23.4 (59) and numpy 1.16.4 (60). We implemented a fully-connected generator network consisting of an input layer with the size of the latent vector size 600, one hidden layer with size proportional to the number of SNPs as SNP_number/1.2 rounded, another hidden layer with size proportional to the number of SNPs as SNP_number/1.1 rounded and an output layer with the size of the number of SNPs. The latent vector is drawn from a Gaussian distribution with zero-

474    mean and unit-variance. The discriminator is also a fully-connected network including

475    an input layer with the size of the number of SNPs, one hidden layer with size

476    proportional to the number of SNPs as SNP_number/2 rounded, another hidden layer

477    with size proportional to the number of SNPs as SNP_number/3 rounded and an

478    output layer of size 1. All layer outputs except for output layers have LeakyReLU

479    activation functions with leaky_alpha parameter 0.01 and L2 regularization parameter

480    0.0001. The generator output layer activation function is tanh and discriminator output

481    layer activation function is sigmoid. Both discriminator and combined GAN were

482    trained thanks to the Adam optimization algorithm with binary cross entropy loss

483    function. We set the discriminator learning rate as 0.0008 and combined GAN learning

484    rate as 0.0001. For 5000 SNP data, the discriminator learning rate was set to 0.00008

485    and the combined GAN learning rate was set to 0.00001. The training to test dataset

486    ratio was 3:1. We used batch size of 32 and trained all datasets up to 20000 epochs.

487    We also investigated stopping the training based on $AA_{TS}$ scores. The score was

488    calculated at 200 epoch intervals. For 805 SNP data, $AA_{TS}$ converged very quickly close

489    to optimum 0.5 score. However, the difference between $AA_{truth}$ and $AA_{syn}$ scores

490    indicates possible overfitting to multiple data points so it was difficult to define a

491    stopping point. For 10K SNP data, convergence was observed after ~30K epochs (to

492    around 0.75) and reduced the number of fixed alleles in AGs but the gain was very

493    minimal (Supplementary Figure 22). Additionally, GAN was prone to mode collapse

494    especially after 20K epochs which resulted in multiple failed training attempts.

495    Therefore, this study presents results for AGs generated at 20k epochs, since the first

496    two PCs of AGs combined with real genomes were visually coherent for all targeted

497    datasets (Figure 1, Supplementary Figure 5). Note that it could be possible to utilize

498    AGs before or after the 20K epoch point. During each batch in the training, when only

499    the discriminator is trained, we applied smoothing to the real labels [1] by vectoral

500    addition of random uniform distribution via numpy.random.uniform with lower bound 0

501    and upper bound 0.1. Elements of the generated outputs were rounded to 0 or 1. After

502    the training is complete, it is possible to generate as many AGs as desired. The code

503    is available at "https://gitlab.inria.fr/ml_genetics/public/artificial_genomes".

504

505    **RBM model:**

506    The RBM model consists of one visible layer of size $N_v$ and one hidden layer of size

507    $N_h$ coupled by a weight matrix W. It is a probabilistic model of the joint distribution of

508    visible $\{v_i, i = 1, \dots N_v\}$ and hidden variables $\{h_j, j = 1, \dots N_h\}$ of the form

509    
$$P(v, h) = e^{-E(v,h)}$$

510    with

511    
$$E(v, h) = \sum_{ij} W_{ij} v_i h_j \; + \; bias\ terms$$

512    Visible variables here are 0,1 as they represent reference/alternative alleles, while the

513    hidden variable type depends on the chosen activation function (sigmoid or RELU).

514    They are there to build dependencies among visible variables which by default are

515    independent, via the interaction strength W. The weight matrix can be used in two

516    different manners to interpret the learned model:

517        1. feature wise: for each hidden variable $j$ the vector $\{W_{ij}, i = 1, \dots N_v\}$

518    represents a certain combination of SNPs which, if activated, will contribute to activate

519    or inhibit this feature $j$. These features are expected to be characteristic of the data

520    structure (such as the population structure) and the vector of feature activations should

521    provide a suitable representation of individuals. If $N_v < N_h$ this corresponds to

522    compressing the input representation.

523      2. direction wise: the SVD decomposition of W provides two sets of singular

524 vectors with one corresponding to the visible space (visible axes) and the other one to

525 the hidden representation (hidden axes). The vectors associated to the largest singular

526 values offer the possibility to project the data in a low dimensional space. Dominant

527 visible axes are expected to be similar to the principal component axes while dominant

528 hidden axes are expected to produce more separable datapoints due to non-linear

529 activation mechanisms. We used the latter (i.e. the projection into the hidden space)

530 to perform our non-linear dimension reduction of the 1000 Genomes data (see

531 Supplementary Figure 13).

532

533 The RBM was coded in Julia (61), and all the algorithm for the training has been done

534 by the authors. The part of the algorithm involving linear algebra used the standard

535 package provided by Julia. Two versions of the RBM were considered. In both

536 versions, the visible nodes were encoded using Bernoulli random variables {0,1}, and

537 the size of the visible layer was the same size as the considered input. Two different

538 types of hidden layers were considered. First with a sigmoid activation function (hence

539 having discrete {0,1} hidden variables), second with ReLu (Rectified Linear unit)

540 activations in which case the hidden variables were positive and continuous (there are

541 distributed according to a truncated gaussian distribution when conditioning on the

542 values of the visible variables). Results with sigmoid activation function were worse

543 compared to ReLu so we used ReLu for all the analyses (Supplementary Figure 23).

544 The number of hidden nodes considered for the experiment was Nh=100 for the 805

545 SNP dataset and Nh=500 for the 10k one. There is no canonical way of fixing the

546 number of hidden nodes, in practice we checked that the number of eigenvalues learnt

547 by the model was smaller than the number of hidden nodes, and that by adding more

548  hidden nodes no improvement were observed during the learning. The learning in

549  general is quite stable, in order to have a smooth learning curve, the learning rate was

550  set between 0.001 and 0.0001 and we used batch size of 32. The negative term of the

551  gradient of the likelihood function was approximated using the PCDk method (62), with

552  k=10 and 100 of persistent chains. As a stopping criterion, we looked at when the $AA_{TS}$

553  score converges to the ideal value of 0.5 when sampling the learned distribution. When

554  dealing with large and sparse datasets for selection tests, RBM model did not manage

555  to provide reasonable $AA_{TS}$ scores because the sampling is intrinsically difficult for

556  large systems with strong correlation. In that case, we used visually coherent PCA

557  results as a stopping criterion. Once the RBM is trained over the dataset, it is possible,

558  in order to avoid running a very long Monte Carlo Markov Chain, to initialize the chain

559  on the training set. However, in the case of the large dataset (Estonian), we observe

560  that the RBM is overfitting the dataset and therefore, starting from the training dataset

561  makes the overfitting even worse. In order to prevent this effect as much as possible,

562  we used another independent dataset of Estonian individuals (denoted sampling set)

563  to start the Monte Carlo Markov Chain. With this trick, we observe that the AATS score

564  exhibits less overfitting than when the Markov Chain was started on the training

565  dataset. We measure the privacy scores for both training and sampling sets compared

566  to a test set. Similar to the GAN, it is possible to generate as many AGs as wanted

567  after       training.       The       relevant       RBM       code       is       available       at

568  "https://gitlab.inria.fr/ml_genetics/public/artificial_genomes".

569

570  **Bernoulli distribution model:**

571   We used python-3.6, pandas 0.23.4 and numpy 1.16.4 for the Bernoulli distribution

572   model code. Each allele at a given position was randomly drawn given the derived

573   allele frequency in the real population.

574

575   **Markov chain model:**

576   We used python-3.6, pandas 0.23.4 and numpy 1.16.4 for the Markov chain model

577   code. For each generated sample alleles were drawn from left (position 0) to right. At

578   the initial position the allele was set by drawing from a Bernoulli distribution

579   parameterized with the real frequency. At a given position $p$ the allele $h_p$ was drawn in

580   *{0,1}* according to its probability given the previous sequence window *of size w, P($h_p$ |*

581   *$h_{p-w}$, …, $h_{p-1}$).* This probability is computed from the observed haplotype frequencies

582   in real data. After the initial position, the sequence window size increased

583   incrementally up to a predefined window size (5 or 10 SNPs). The relevant code is

584   available at "https://gitlab.inria.fr/ml_genetics/public/artificial_genomes".

585

586   **HAPGEN2:**

587   We used HAPGEN2 (35) to generate our targeted region of chromosome 15  for as

588   many individuals as in the original dataset. We provided the training dataset (e.g.

589   either 1000 Genomes or Estonian) and a recombination map (37) of the region as

590   input. We sampled only control individuals and no cases. All other parameters were

591   set to default.

592

593   **Coalescent simulations:**

594   We used stdpopsim (38) with the command line "stdpopsim HomSap -c chr15 -o

595   CEU_chr15.trees -g HapMapII_GRCh37 -d OutOfAfrica_3G09 0 2000 0" to generate

596    2000 CEU haplotypes based on the demographic parameters inferred by Gutenkunst

597    et al. 2009 (36). We then selected the genome region corresponding to one targeted

598    when generating AGs.

599

600    **Summary statistics:**

601    We define here the statistics that are not commonly used in population genetics. The

602    3-point scores measure the correlation patterns for SNP triplets. The 3-point

603    correlation for SNPs $i, j,$ and $k$ is defined as (63):

604        $c_{ijk}(a,b,c) = f_{ijk}(a,b,c) - f_{ij}(a,b) \, f_k(c) - f_{ik}(a,c) \, f_j(b) - f_{jk}(b,c) \, f_i(a) + 2 \, f_i(a) \, f_j(b) \, f_k(c)$ ,

605    where the alleles $(a,b,c) \in \{0,1\}^3$, $f_i(a)$ is the frequency of allele $a$ at SNP $i$, $f_{ij}(a,b)$ is the

606    frequency of the combination of allele $a$ at SNP $i$ and $b$ at SNP $j$, and finally $f_{ijk}(a,b,c)$ is

607    the frequency of the combination $(a,b,c)$ at SNPs $(i,j,k)$. We computed the 3-point

608    correlations for 8,000 randomly-picked triplets under different conditions (SNPs

609    separated by  1, 4, 16, 64, 256, 512 or 1024 SNPs, as well as SNPs chosen at random)

610    in each dataset.

611

612    PCA were computed on all datasets combined (e.g. Figure 1) as well as on "pairs" of

613    datasets (the combination of real and a single type of generated data). 2D-

614    Wasserstein distances for these paired PCA representations were computed based

615    on the entropic regularized optimal transport problem with square euclidean distances

616    computed from PCs 1 and 2 and a regularization parameter set to 0.001 (POT  library,

617    (64)).

618

619    To have reference values regarding the best achievable distances or correlations

620    between real and generated summary statistics we split randomly the 1000Genomes

621    dataset in two and considered half of it as the real dataset and half as a "perfectly

622    generated" dataset (called test).

623

**Chromosome painting:**

625    We compared the haplotype sharing distribution between real and artificial

626    chromosomes through ChromoPainter (27). In detail, we have painted 100 randomly

627    selected "real" and "artificial" Estonians (recipients) against all the 1000 Genome

628    Project phased data (donors). The nuisance parameters -n (348.57) and -M (0.00027),

629    were estimated running 10 iterations of the expectation-maximization algorithm on a

630    subset of 3,800 donor haplotypes.

631

**Haplostrips:**

633    We used Haplostrips (28) to visualize the haplotype structure of real and artificial

634    genomes. We extracted 500 individuals from each sample set (Real, GAN AGs, RBM

635    AGs) and considered them as different populations. Black dots represent derived

636    alleles, white dots represent ancestral alleles. The plotted SNPs were filtered for a

637    population specific minor allele frequency >5%; haplotypes were clustered and sorted

638    for distance against the consensus haplotype from the real set. See the application

639    article for further details about the method.

640

**Nearest Neighbour Adversarial Accuracy (AA$_{TS}$) and privacy loss**

642    We used the following equations for calculating AA$_{TS}$ and privacy loss scores (29) :

$$AA_{truth} = \frac{1}{n}\sum_{i=1}^{n}\mathbf{1}(d_{TS}(i) > d_{TT}(i))$$

644
$$AA_{syn} = \frac{1}{n}\sum_{i=1}^{n} \mathbf{1}(d_{ST}(i) > d_{SS}(i))$$

645
$$AA_{TS} = \frac{1}{2}(AA_{truth} + AA_{syn})$$

646
$$Privacy\ Loss = \ Test\ AA_{TS} - Train\ AA_{TS}$$

647
648   where $n$ is the number of real samples as well as of artificial samples; $\mathbf{1}$ is a function

649   which takes the value 1 if the argument is true and 0 if the argument is false; $d_{TS}(i)$ is

650   the distance between the real genome indexed by i and its nearest neighbour in the

651   artificial genome dataset; $d_{ST}(i)$ is the distance between the artificial genome indexed

652   by i and its nearest neighbour in the real genome dataset; $d_{TT}(i)$ is the distance of the

653   real genome indexed by i to its nearest neighbour in the real genome dataset; $d_{SS}(i)$

654   is the distance of the artificial genome indexed by i to its nearest neighbour in the

655   artificial genome dataset. An $AA_{TS}$ score of 0.5 is optimal whereas lower values indicate

656   overfitting and higher values indicate underfitting. For a better resolution for the

657   detection of overfitting, we also provided $AA_{truth}$ and $AA_{syn}$ metrics identified in the general

658   equation of $AA_{TS}$. If $AA_{TS}$ 0.5 but $AA_{truth}$ 0 and $AA_{syn}$ 1, this means that the model is not

659   overfitting in terms of a single data point but multiple ones. In other words, the model

660   might be focusing on small batches of similar real genomes to create artificial genomes

661   clustered at the center of each batch. Privacy loss is the difference of $AA_{TS}$ score of

662   AGs calculated against the training samples set and a different real test set which was

663   not used in training.

664

665   **Selection tests:**

666   We used scikit-allel package for XP-EHH (65) and PBS (66) tests. We used 1000

667   Estonian individuals (2000 haplotypes) with 3348 SNPs which were homogenously

668    dispersed over chromosome 15 (spanning the whole chromosome with similar

669    distance between SNPs)  for the training of GAN and RBM models. For XP-EHH,

670    Yoruban (YRI, 216 haplotypes) population from 1000 Genomes data was used as the

671    complementary population. For PBS, Yoruban (YRI, 216 haplotypes) and Japanese

672    (JPT, 208 haplotypes) populations from 1000 Genomes data were used as

673    complementary populations. PBS window size was 10 and step size was 5, resulting

674    in 668 windows. 216 real and 216 AG haplotypes were compared for the analyses.

675

## Figure Legends

**Figure 1.** The six first axes of a single PCA applied to real (gray) and artificial genomes (AGs) generated via Bernoulli (green), Markov chain (purple), GAN (blue) and RBM (red) models. There are 5000 haplotypes for each AG dataset and 5008 (2504 genomes) for the real dataset from 1000 Genomes spanning 805 informative SNPs. See Materials & Methods for detailed explanation of the generation procedures.
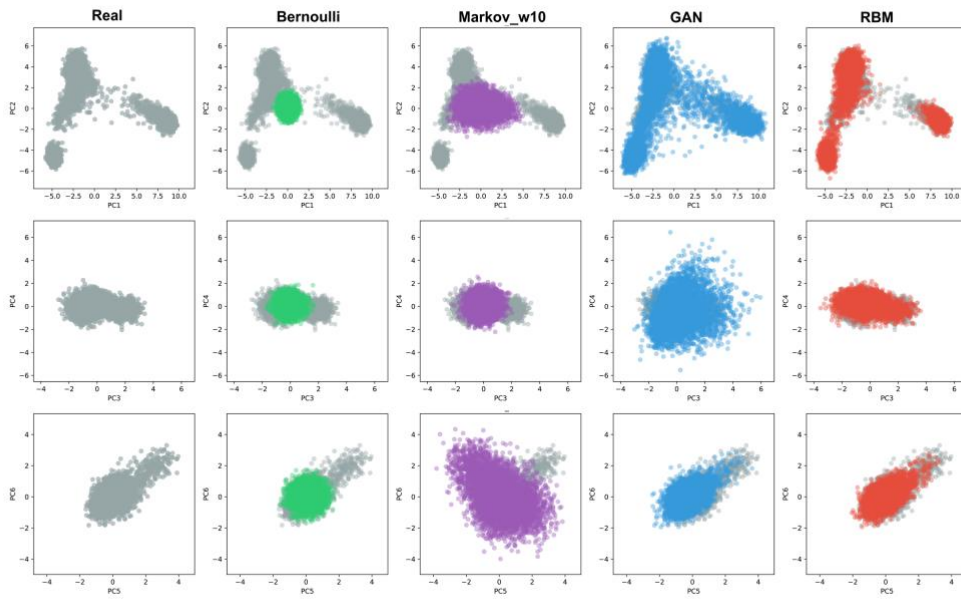
**Figure 2.** Linkage disequilibrium (LD) analysis on real and artificial Estonian genomes. **a)** Correlation ($r^2$) matrices of SNPs. Lower triangular parts are SNP pairwise correlation in real genomes and upper triangular parts are SNP pairwise correlation in artificial genomes. **b)** LD as a function of SNP distance after removing sites that are fixed in at least in one dataset. Pairwise SNP distances were stratified into 50 bins and for each distance bin, the correlation was averaged over all pairs of SNPs belonging to the bin.

**Figure 3.** Selection tests on chromosome 15. **a)** Standardized XP-EHH scores of real and artificial Estonian genomes using 1000 Genomes Yoruba population (YRI) as the complementary population. Correlation coefficient between real and GAN XP-EHH score is 0.902, between real and RBM XP-EHH score is 0.887. **b)** PBS scores of real and artificial Estonian genomes using 1000 Genomes Yoruba (YRI) and Japanese (JPT) populations as the complementary populations. PBS window size is 10 and step size is 5. Dotted black line corresponds to the 99$^{th}$ percentile. Correlation coefficient between real and GAN PBS score is 0.923, between real and RBM PBS score is 0.755. Highest peaks are marked by an asterisk.

**Figure 4.** Imputation evaluation of three different reference panels based on Impute2 software's info metric. Imputation was performed on 8678 Estonian individuals (which were not used in training of GAN and RBM models) using only 1000 Genomes panel

701    (gray), combined 1000 Genomes and Estonian genomes used in training (green),

702    combined 1000 Genomes and GAN artificial genomes panel (blue) and combined

703    1000 Genomes and RBM artificial genomes panel (red). SNPs were divided into 10

704    MAF bins, from 0.05 to 0.5, after which mean info metric values were calculated. Grey

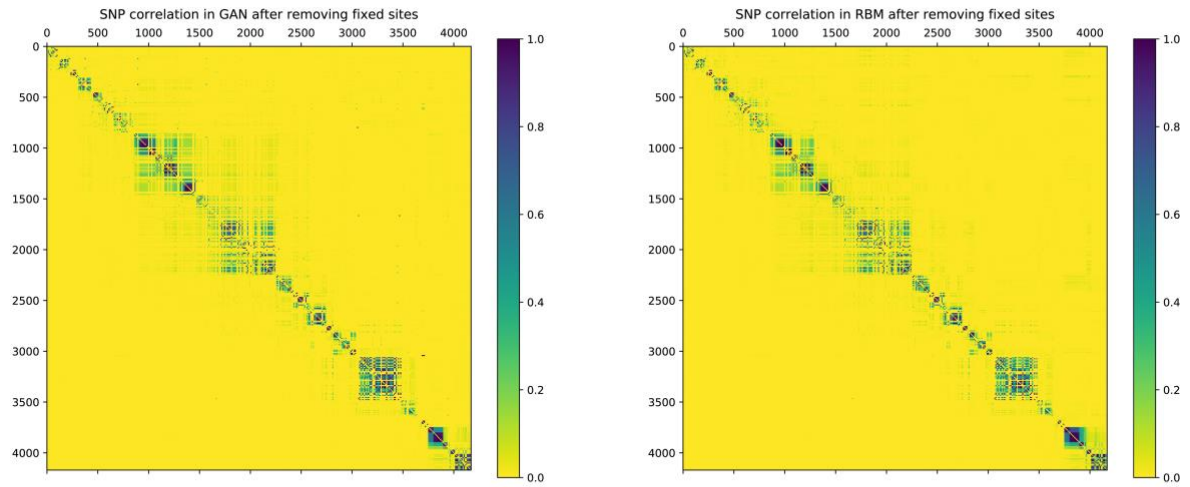705    bars show the percentage of SNPs which belong to each bin.
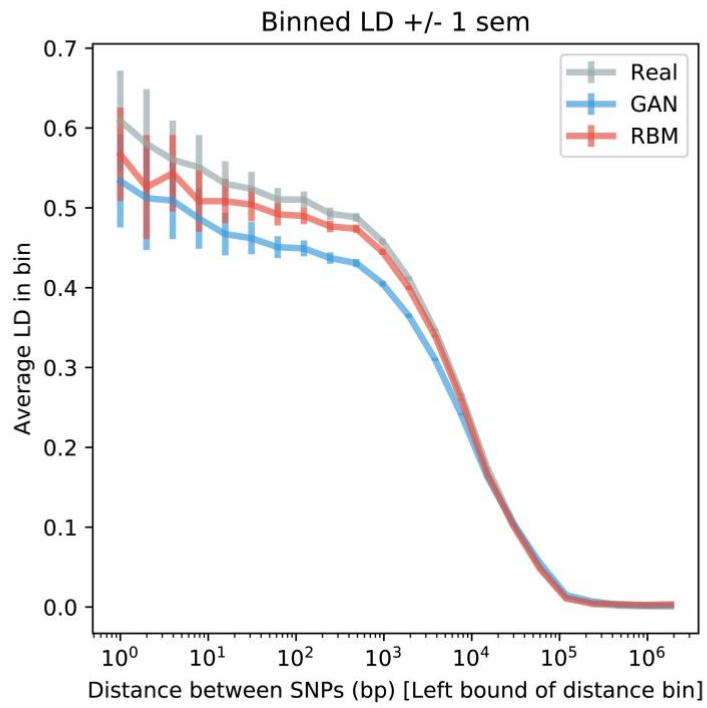
706

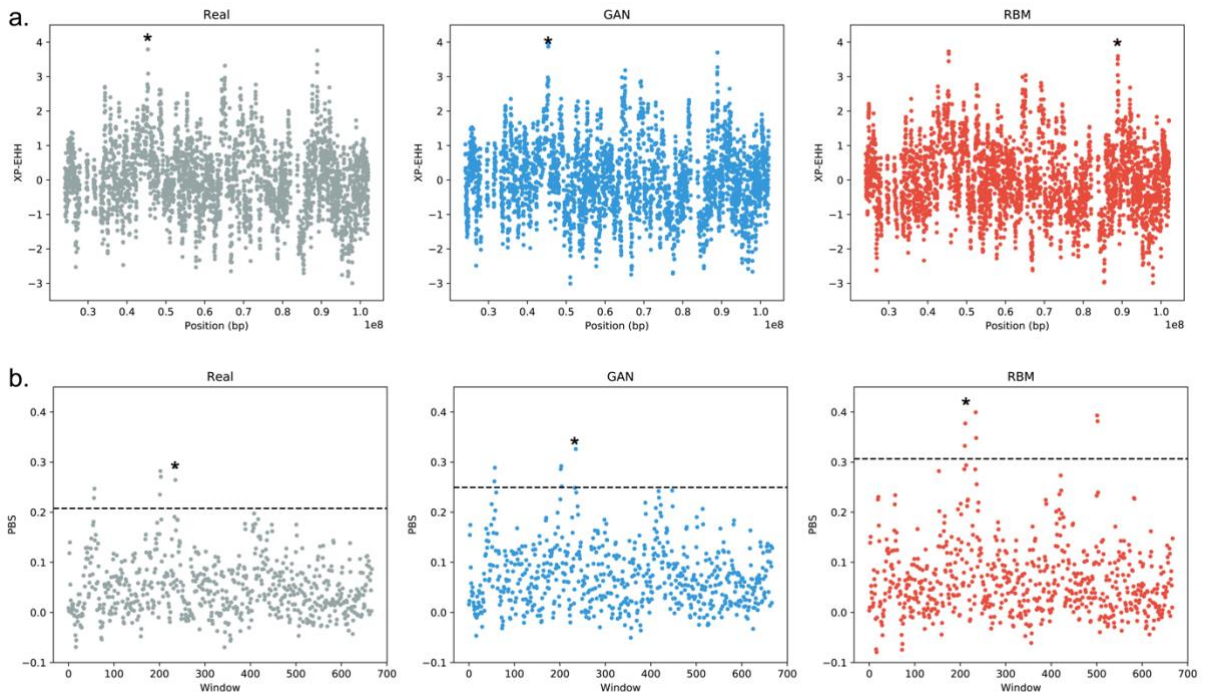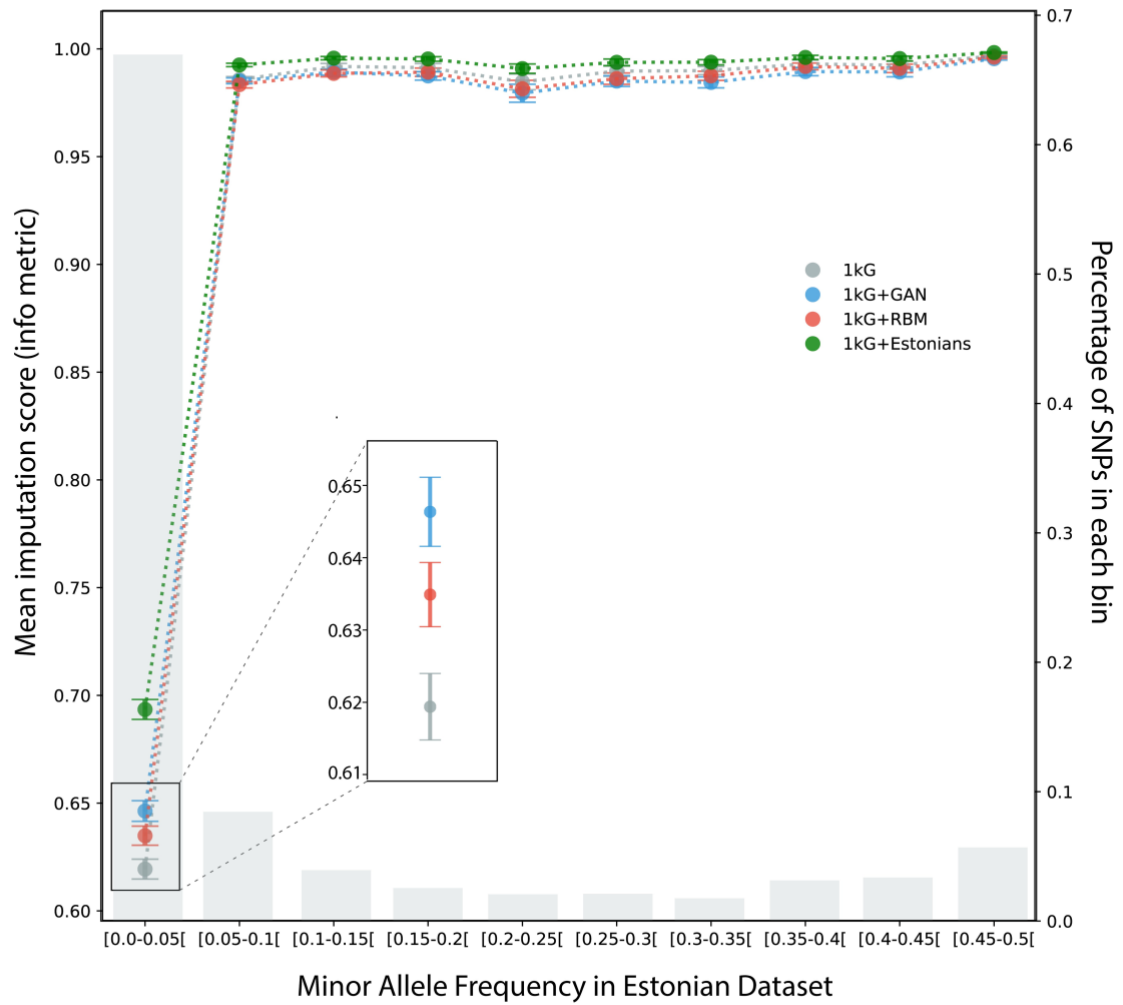707    **Figure 1.**



708

709     **Figure 2.**

a.



b.



710

711  **Figure 3.**

712



713

714    **Figure 4.**



715

## Supporting Information Legends

716

717    **Supplementary Figure 1.** Generative adversarial network (GAN) scheme.

718    **Supplementary Figure 2.** Restricted Boltzmann machine (RBM) scheme.

719    **Supplementary Figure 3.** Uniform manifold approximation and projection (UMAP) of real

720    genomes from 1000 Genomes data spanning 805 SNPs along with artificial genome

721    counterparts created via **a)** Bernoulli, **b)** Markov chain (with 10 window length), **c)** GAN and

722    **d)** RBM models.

723    **Supplementary Figure 4.** Distribution of haplotypic pairwise difference within (left) and

724    between (right) datasets of real genomes from 1000 Genomes data spanning 805 SNPs and

725    artificial genome counterparts generated using different models.

726    **Supplementary Figure 5.** PCA of real genomes (gray) from **a)** 1000 Genomes data and **b)**

727    Estonian Biobank spanning 10K SNPs along with artificial genome counterparts generated

728    using GAN (blue) and RBM (red) models.

729    **Supplementary Figure 6.** Distribution of haplotypic pairwise difference within (left) and

730    between (right) datasets of real genomes from **a)** 1000 Genomes data and **b)** Estonian

731    Biobank spanning 10K SNPs and artificial genome counterparts generated using GAN and

732    RBM models.

733    **Supplementary Figure 7.** Allele frequency comparison of corresponding SNPs between

734    real genomes from Estonian Biobank spanning 10K SNPs and artificial genome counterparts

735    generated using GAN and RBM models as **a)** the whole range and **b)** zoomed to low

736    frequencies. Clustering below the diagonal in the low frequency section for the GAN plot

737    indicates insufficient representation of rare alleles in artificial genomes.

738    **Supplementary Figure 8.** Distribution of minimum distance to the closest neighbour for real

739    genomes from **a)** 1000 Genomes data and **b)** Estonian Biobank spanning 10K SNPs along

740    with artificial genome counterparts generated via GAN and RBM models.

741    **Supplementary Figure 9.** LD comparison of real (Estonian) vs generated datasets.

742    **Supplementary Figure 10.** 3-point correlation statistics for SNPs separated by different

743    distances.

744    **Supplementary Figure 11.** Haplostrips showing the mixed nature of haplotype structures for

745    real Estonian (gray rows) along with GAN (blue rows) and RBM (red rows) haplotypes.

746    **Supplementary Figure 12.** Chromosome painting of two **a)** real Estonian genomes, **b)** GAN

747    and **c)** RBM artificial Estonian genomes with 1000 Genomes donors colored based on super

748    population codes. EUR – European, EAS – East Asian, AMR – Admixed American, SAS –

749    South Asian, AFR – African.

750    **Supplementary Figure 13. a)** Nearest neighbour adversarial accuracy ($AA_{TS}$) scores of

751    artificial genomes generated from Estonian Biobank. Black line indicates the

752    optimum value whereas values below the line indicate overfitting and values above

753    the line indicate underfitting. **b)** Privacy loss. Test1 is a separate set of real Estonian

754    genomes. Positive values indicate information leakage, hence overfitting.

755    **Supplementary Figure 14.** $AA_{TS}$ and privacy loss change of RBM AGs over epochs.

756    **Supplementary Table.** Genotype/phenotype contingency table for real and artificial

757    Estonian genomes (AG). Ancestral allele "A" is associated with brown eye color and derived

758    allele "G" is associated with blue eye color phenotype.

759    **Supplementary Text.** Preliminary analysis on generating artificial genomes with

760    corresponding phenotypes.

761    **Supplementary Figure 15.** Comparison of PCA (right column) and non-linear dimension

762    reduction via RBM (left column) for real genomes from 1000 Genomes data spanning 805

763    SNPs. The RBM reduction was obtained by projecting the real data into the hidden space of

764    the RBM (see Materials & Methods). Population codes are as defined by the 1000 Genomes

765    Project.

766    **Supplementary Figure 16.** Activations of each of the 100 nodes belonging to the RBM

767    hidden layer when applied to the real genomes from 1000 Genomes data spanning 805

768    SNPs. For each hidden node the X-axis corresponds to the real haplotypes and Y-axis to the

769    activation of the node by a single haplotype. On the X-axis, haplotypes are ordered by region

770    (Africa, America, East Asia, European, East Asia) and colored by population. Because this

771    RBM activation function is a ReLU with threshold 0 (by design), all values are positive and a

772    zero-value indicates that the node is not activated by a given haplotype. The ordering of

773    nodes has no specific meaning.

774    **Supplementary Figure 17.** Analyses of artificial genomes generated by HAPGEN2 showing

775    **a)** PCA of generated (green) performed with real Estonian genomes (grey) and **b)**

776    distribution of minimum distance to the closest neighbour displaying real Estonian genomes

777    (grey), HAPGEN2 (green), GAN (blue) and RBM (red) artificial genomes.

778    **Supplementary Figure 18**. **a)** PCA of real (Estonian) and artificial genomes simulated via

779    coalescent approach using stdpopsim (CEU). **b)** Allele frequency quantiles of real (Estonian)

780    vs artificial genomes simulated via coalescent approach using stdpopsim (CEU).

781    **Supplementary Figure 19. a)** LD as a function of SNP distance after removing sites that are

782    fixed in at least one dataset and removing alleles below 0.1 frequency from all datasets.

783    Pairwise SNP distances were stratified into 50 bins and for each distance bin, the correlation

784    was averaged over all pairs of SNPs belonging to the bin. Allele frequency quantiles of real

785    (Estonian) vs **b)** GAN Estonian artificial genomes, **c)** RBM Estonian artificial genomes and

786    **d)** artificial genomes simulated via coalescent approach using stdpopsim (CEU).

787    **Supplementary Figure 20.** Comparison of sites which are polymorphic in real genomes

788    from Estonian Biobank but fixed in artificial genome counterparts generated via GAN and

789    RBM models.

790    **Supplementary Figure 21.** Sensitivity tests for **a)** $AA_{TS}$ (scores over 0.5 indicate underfitting

791    and below 0.5 indicate overfitting) and **b)** privacy scores (orange and red lines to mark the

792    difference between RBM trained up to 350 and 690 epochs). All datasets consist of 2000

793    samples. Test1 and Test2 are real Estonian individuals who were not used in training.

794    Mixed1 dataset has 1 real individual from the training dataset, Mixed2 has 10, Mixed3 has

795    50, Mixed4 has 100, Mixed5 has 500 and Mixed6 has 1000 individuals.

796    **Supplementary Figure 22.** Evaluation of $AA_{TS}$ scores of the GAN model for artificial

797    Estonian genomes spanning **a)** 805 highly informative SNPs and **b)** dense 10K SNPs along

798    with the total fixed sites for the outputs of epochs at 200 intervals.

799    **Supplementary Figure 23.** Comparison of **a)** $AA_{TS}$ score and **b)** linkage disequilibrium of

800    artificial genomes created via RBM model with sigmoid and ReLu activation functions.

# Acknowledgements

# References

812

813    1.    Mardis ER. DNA sequencing technologies: 2006-2016. Nature Protocols. 2017.

814    2.    Cann HM. A Human Genome Diversity Cell Line Panel. Science (80- ). 2002;

815    3.    Mallick S, Li H, Lipson M, Mathieson I, Gymrek M, Racimo F, et al. The Simons

816          Genome Diversity Project: 300 genomes from 142 diverse populations. Nature. 2016;

817    4.    Popejoy AB, Fullerton SM. Genomics is failing on diversity. Nature. 2016.

818    5.    Sirugo G, Williams SM, Tishkoff SA. The Missing Diversity in Human Genetic Studies.

819          Cell. 2019.

820    6.    Libbrecht MW, Noble WS. Machine learning applications in genetics and genomics.

821          Nature Reviews Genetics. 2015.

822    7.    Zhang H, Xu T, Li H, Zhang S, Wang X, Huang X, et al. StackGAN: Text to Photo-

823          Realistic Image Synthesis with Stacked Generative Adversarial Networks. In:

824          Proceedings of the IEEE International Conference on Computer Vision. 2017.

825    8.    Rolnick D, Dyer EL. Generative models and abstractions for large-scale

826          neuroanatomy datasets. Current Opinion in Neurobiology. 2019.

827    9.    Davidsen K, Olson BJ, DeWitt WS, Feng J, Harkins E, Bradley P, et al. Deep

828          generative models for T cell receptor protein sequences. Elife [Internet]. 2019 Sep 5

829          [cited 2019 Sep 12];8. Available from: https://elifesciences.org/articles/46935

830    10.   Liu Q, Lv H, Jiang R. HicGAN infers super resolution Hi-C data with generative

831          adversarial networks. In: Bioinformatics. Oxford University Press; 2019. p. i99–107.

832    11.   Tubiana J, Cocco S, Monasson R. Learning protein constitutive motifs from sequence

833          data. Elife. 2019;

834    12.   Shimagaki K, Weigt M. Selection of sequence motifs and generative Hopfield-Potts

835          models for protein families. bioRxiv. 2019 Sep 5;652784.

836    13.   Killoran N, Lee LJ, Delong A, Duvenaud D, Frey BJ. Generating and designing DNA

837          with deep generative models. 2017 Dec 17 [cited 2019 Sep 15]; Available from:

838          http://arxiv.org/abs/1712.06148

839    14.    Goodfellow I, Pouget-Abadie J, Mirza M. Generative Adversarial Networks (GANs) -

840          Tutorial. Neural Inf Process Syst. 2014;

841    15.    Ledig C, Theis L, Huszár F, Caballero J, Cunningham A, Acosta A, et al. Photo-

842          realistic single image super-resolution using a generative adversarial network. In:

843          Proceedings - 30th IEEE Conference on Computer Vision and Pattern Recognition,

844          CVPR 2017. 2017.

845    16.    Fedus W, Goodfellow I, Dai AM. MaskGAN: Better Text Generation via Filling in

846          the_____. 2018 Jan 23 [cited 2019 Aug 26]; Available from:

847          http://arxiv.org/abs/1801.07736

848    17.    Brock A, Donahue J, Simonyan K. Large Scale GAN Training for High Fidelity Natural

849          Image Synthesis. 2018 Sep 28 [cited 2019 Sep 13]; Available from:

850          http://arxiv.org/abs/1809.11096

851    18.    Smolensky P. Information processing in dynamical systems: Foundations of harmony

852          theory. In: Parallel Distributed Processing Explorations in the Microstructure of

853          Cognition. 1986.

854    19.    Teh YW, Hinton GE. Rate-coded restricted boltzmann machines for face recognition.

855          In: Advances in Neural Information Processing Systems. 2001.

856    20.    Kingma DP, Welling M. Auto-Encoding Variational Bayes. 2013 Dec 20 [cited 2019

857          Aug 26]; Available from: http://arxiv.org/abs/1312.6114

858    21.    Hinton GE, Salakhutdinov RR. Reducing the dimensionality of data with neural

859          networks. Science (80- ). 2006;

860    22.    Hinton GE. Learning multiple layers of representation. Trends in Cognitive Sciences.

861          2007.

862    23.    Larochelle H, Bengio Y. Classification using discriminative restricted boltzmann

863          machines. In: Proceedings of the 25th International Conference on Machine Learning.

864          2008.

865    24.    1000 Genomes Project Consortium, Auton A, Brooks LD, Durbin RM, Garrison EP,

866          Kang HM, et al. A global reference for human genetic variation. Nature [Internet].

867          2015;526(7571):68–74. Available from:

868          http://www.ncbi.nlm.nih.gov/pubmed/26432245%0Ahttp://www.pubmedcentral.nih.gov

869          /articlerender.fcgi?artid=PMC4750478

870    25.    Colonna V, Ayub Q, Chen Y, Pagani L, Luisi P, Pybus M, et al. Human genomic

871          regions with exceptionally high levels of population differentiation identified from 911

872          whole-genome sequences. Genome Biol. 2014;

873    26.    Leitsalu L, Haller T, Esko T, Tammesoo ML, Alavere H, Snieder H, et al. Cohort

874          profile: Estonian biobank of the Estonian genome center, university of Tartu. Int J

875          Epidemiol. 2015;

876    27.    Lawson DJ, Hellenthal G, Myers S, Falush D. Inference of population structure using

877          dense haplotype data. PLoS Genet. 2012;

878    28.    Marnetto D, Huerta-Sánchez E. Haplostrips: revealing population structure through

879          haplotype visualization. Methods Ecol Evol. 2017;8(10):1389–92.

880    29.    Yale A, Dash S, Dutta R, Guyon I, Pavao A, Bennett K. Privacy Preserving Synthetic

881          Health Data. 2019 Apr 24 [cited 2019 Aug 27]; Available from: https://hal.inria.fr/hal-

882          02160496/

883    30.    Gurdasani D, Carstensen T, Tekola-Ayele F, Pagani L, Tachmazidou I, Hatzikotoulas

884          K, et al. The African Genome Variation Project shapes medical genetics in Africa.

885          Nature. 2015;

886    31.    Mitt M, Kals M, Pärn K, Gabriel SB, Lander ES, Palotie A, et al. Improved imputation

887          accuracy of rare and low-frequency variants using population-specific high-coverage

888          WGS-based imputation reference panel. Eur J Hum Genet. 2017;

889    32.    Howie B, Marchini J, Stephens M. Genotype imputation with thousands of genomes.

890          G3 Genes, Genomes, Genet. 2011;

891    33.    Diaz-Papkovich A, Anderson-Trocme L, Gravel S. Revealing multi-scale population

892          structure in large cohorts. bioRxiv. 2019;

893    34.    Li N, Stephens M. Modeling Linkage Disequilibrium and Identifying Recombination

894          Hotspots Using Single-Nucleotide Polymorphism Data. Genetics. 2003;

895    35.    Su Z, Marchini J, Donnelly P. HAPGEN2: Simulation of multiple disease SNPs.

896            Bioinformatics. 2011;

897    36.    Gutenkunst RN, Hernandez RD, Williamson SH, Bustamante CD. Inferring the joint

898            demographic history of multiple populations from multidimensional SNP frequency

899            data. PLoS Genet. 2009;

900    37.    Frazer KA, Ballinger DG, Cox DR, Hinds DA, Stuve LL, Gibbs RA, et al. A second

901            generation human haplotype map of over 3.1 million SNPs. Nature. 2007;

902    38.    Adrion JR, Cole CB, Dukler N, Galloway JG, Gladstein AL, Gower G, et al. A

903            community-maintained standard library of population genetic models. Elife. 2020;

904    39.    Kelleher J, Etheridge AM, McVean G. Efficient Coalescent Simulation and

905            Genealogical Analysis for Large Sample Sizes. PLoS Comput Biol. 2016;

906    40.    Tian X, Browning BL, Browning SR. Estimating the Genome-wide Mutation Rate with

907            Three-Way Identity by Descent. Am J Hum Genet. 2019;

908    41.    Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, et al. Initial

909            sequencing and analysis of the human genome. Nature. 2001;

910    42.    Berisa T, Pickrell JK. Approximately independent linkage disequilibrium blocks in

911            human populations. Bioinformatics. 2016;

912    43.    Radford A, Metz L, Chintala S. Unsupervised Representation learning with Deep

913            Convolutional GANs. Int Conf Learn Represent. 2016;

914    44.    Salimans T, Goodfellow I, Zaremba W, Cheung V, Radford A, Chen X. Improved

915            techniques for training GANs. In: Advances in Neural Information Processing

916            Systems. 2016.

917    45.    Arjovsky M, Chintala S, Bottou L. Wasserstein generative adversarial networks. In:

918            34th International Conference on Machine Learning, ICML 2017. 2017.

919    46.    Lucas T, Tallec C, Verbeek J, Ollivier Y. Mixed batches and symmetric discriminators

920            for GAN training. In: 35th International Conference on Machine Learning, ICML 2018.

921            2018.

922    47.    Dieng AB, Ruiz FJR, Blei DM, Titsias MK. Prescribed generative adversarial

923        networks. arXiv Prepr arXiv191004302. 2019;

924    48.    Martin MD, Jay F, Castellano S, Slatkin M. Determination of genetic relatedness from

925        low-coverage human genome sequences using pedigree simulations. Mol Ecol. 2017;

926    49.    Fortes-Lima C, Laurent R, Thouzeau V, Toupance B, Verdu P. Complex genetic

927        admixture histories reconstructed with Approximate Bayesian Computations. bioRxiv.

928        2019;761452.

929    50.    Dwork C, McSherry F, Nissim K, Smith A. Calibrating noise to sensitivity in private

930        data analysis. In: Lecture Notes in Computer Science (including subseries Lecture

931        Notes in Artificial Intelligence and Lecture Notes in Bioinformatics). 2006.

932    51.    Torkzadehmahani R, Kairouz P, Ai G, Paten B. DP-CGAN : Differentially Private

933        Synthetic Data and Label Generation [Internet]. [cited 2019 Oct 7]. Available from:

934        https://github.com/tensorflow/privacy

935    52.    Basu Mallick C, Iliescu FM, Möls M, Hill S, Tamang R, Chaubey G, et al. The Light

936        Skin Allele of SLC24A5 in South Asians and Europeans Shares Identity by Descent.

937        PLoS Genet. 2013;

938    53.    Dumoulin V, Belghazi I, Poole B, Mastropietro O, Lamb A, Arjovsky M, et al.

939        Adversarially Learned Inference. 2016 Jun 2 [cited 2019 Aug 28]; Available from:

940        http://arxiv.org/abs/1606.00704

941    54.    Chen X, Duan Y, Houthooft R, Schulman J, Sutskever I, Abbeel P. InfoGAN:

942        Interpretable representation learning by information maximizing generative adversarial

943        nets. In: Advances in Neural Information Processing Systems. 2016.

944    55.    Donahue J, Krähenbühl P, Darrell T. Adversarial Feature Learning. 2016 May 31

945        [cited 2019 Aug 28]; Available from: http://arxiv.org/abs/1605.09782

946    56.    Paul JS, Steinrücken M, Song YS. An accurate sequentially markov conditional

947        sampling distribution for the coalescent with recombination. Genetics. 2011;

948    57.    Paul JS, Song YS. A principled approach to deriving approximate conditional

949        sampling distributions in population genetics models with recombination. Genetics.

950        2010;

951   58.   Chollet F. Keras: Deep Learning library for Theano and TensorFlow. GitHub Repos.
952         2015;

953   59.   McKinney W. Data Structures for Statistical Computing in Python. In: Proceedings of
954         the 9th Python in Science Conference (SCIPY 2010). 2010.

955   60.   Oliphant TE. SciPy: Open source scientific tools for Python. Computing in Science
956         and Engineering. 2007.

957   61.   Bezanson J, Edelman A, Karpinski S, Shah VB. Julia: A fresh approach to numerical
958         computing. SIAM Rev. 2017;

959   62.   Brügge K, Fischer A, Igel C. The flip-the-state transition operator for restricted
960         Boltzmann machines. Mach Learn. 2013;

961   63.   Shimagaki K, Weigt M. Selection of sequence motifs and generative Hopfield-Potts
962         models for protein families. Phys Rev E. 2019;

963   64.   Flamary R, Courty N. POT Python Optimal Transport library [Internet]. 2017. Available
964         from: https://pythonot.github.io/

965   65.   Sabeti PC, Varilly P, Fry B, Lohmueller J, Hostetter E, Cotsapas C, et al. Genome-
966         wide detection and characterization of positive selection in human populations.
967         Nature. 2007;

968   66.   Yi X, Liang Y, Huerta-Sanchez E, Jin X, Cuo ZXP, Pool JE, et al. Sequencing of 50
969         human exomes reveals adaptation to high altitude. Science (80- ). 2010;

970

971