1  **PopInf: An approach for reproducibly visualizing and assigning population affiliation in**

2  **genomic samples of uncertain origin**

3  Angela M. Taravella Oill[1], Anagha J. Deshpande[1], Heini M. Natri[1], Melissa A. Wilson[1]

4

5  **Author details**

6  1. School of Life Sciences, Center for Evolution and Medicine, The Biodesign Institute,

7  Arizona State University, Tempe, AZ 85282 USA

8

9  **Corresponding author**

10  Melissa A. Wilson

11  School of Life Sciences | Arizona State University | PO Box 874501 | Tempe, AZ 85287-4501

12  mwilsons@asu.edu

13

14

15

16

17

18

19

20

21

22

23

24

25

26

27    **ABSTRACT**

28          Germline genetic variation contributes to cancer etiology, but self-reported race is not

29    always consistent with genetic ancestry, and samples may not have identifying ancestry

30    information. Here we describe a flexible computational pipeline, PopInf, to visualize principal

31    components analysis output and assign ancestry to samples with unknown genetic ancestry,

32    given a reference population panel of known origins. PopInf is implemented as a reproducible

33    workflow in Snakemake with a tutorial on GitHub. We provide a pre-processed reference

34    population panel that can be quickly and efficiently implemented in cancer genetics studies. We

35    ran PopInf on TCGA liver cancer data and identify discrepancies between reported race and

36    inferred genetic ancestry. **Significance.** The PopInf pipeline facilitates visualization and

37    identification of genetic ancestry across samples, so that this ancestry can be accounted for in

38    studies    of    disease    risk.    All    code    and    a    tutorial    are    available    on    Github:

39    https://github.com/SexChrLab/PopInf.

40

41    **Keywords:** population ancestry, principal components analysis, visualization, computational

42    pipeline, cancer GWAS

43

44

45

46

47

48

49  **INTRODUCTION**

50      Cancer is a complex disease with genetic and environmental factors contributing to its risk

51  and progression. The underlying genetic architecture of cancer, like other complex diseases, is

52  influenced by common population-specific genetic variation (1,2). Common genetic variation is

53  shared within populations of shared genetic ancestry. Unaccounted population structure can

54  confound the results of genetic analyses, like in cancer GWAS, by causing spurious associations

55  to disease phenotypes (3). Thus, assessing genetic ancestry and population structure in studies

56  on the effects of genetic loci and genetic background on cancer is crucial.

57      Cancer research has begun to recognize the importance of identifying genetic ancestry

58  across patients in cancer genetic datasets (4) and across cancer cell lines (5). Yuan et al. (4)

59  characterized genetic ancestry across The Cancer Genome Atlas (TCGA) patient cohort to

60  investigate the effect genetic ancestry has on genomic alterations across different cancers and to

61  provide researchers with detailed ancestry information on each patient. Though this publicly

62  accessible resource is of great research value for those using the TCGA data, researchers

63  utilizing other datasets will have to independently infer the ancestry of their samples.

64      Methods and software are currently available to characterize population structure (6,7),

65  estimate local and global ancestry proportions (7,8), or predict ancestry using genomic data (9).

66  These rely on a pre-defined reference panel and may not report admixed samples. Having an

67  easily reproducible and modifiable workflow to visualize PCA and identify ancestry in individuals

68  of unknown ancestral origin would thus be a useful addition to the cancer genetics researchers

69  tool kit.

70      Here we present PopInf v1.0, a pipeline to visualize PCA output and assign ancestry to

71  individuals with unknown ancestry, given a flexible reference population panel of known origins.

72  PopInf v1.0, takes, as input, variants from a sample with unknown or unverified genetic ancestry

73  in variant call format (VCF), compares the variants in the unknown sample to a user defined

74  reference panel, and outputs an inferred ancestry origin report with accompanying PCA plots of

3

75      the unknown samples and the reference panel. We ran PopInf on variants from 148 samples from

76      the Genotype Tissue Expression (GTEx) Project (10) and on 403 samples from the TCGA liver

77      cancer dataset (11) and identify discrepancies between reported race and inferred genetic

78      ancestry. Further, we analyze each sample by chromosome and find cases of chromosome-

79      specific admixture that is not reported in genome-wide analyses.

80

81      **MATERIALS AND METHODS**

82          PopInf v1.0 uses a combination of publicly available software and custom scripts to

83      generate PCA plots and a tab-delimited inferred ancestry report for samples of unknown ancestry

84      or unverified self-reported population ancestry. PopInf v1.0 uses GATK v3.7 (12), VCFtools

85      v.0.1.14 (13), bedtools v.2.27.1 (14), and Plink v.1.9 (15) to prepare the unknown ancestry dataset

86      and reference panel, smartpca - a program within EIGENSOFT v6.0.1 package (6) - for PCA, and

87      a custom R script (16) to infer individuals ancestry and plot the results of PCA of the study samples

88      and reference panel. Our pipeline is incorporated into the reproducible workflow system,

89      Snakemake v5.4.0 (17).

90      **Input**

91          Two sets of variant data are required to use PopInf v1.0: 1) variants from reference

92      populations, and 2) variants from sample(s) of unknown or self-reported race or ancestry. These

93      files need to be mapped to the same reference genome and in VCF file format. Additionally, two

94      sample information text files, one for the reference panel and one for the unknown dataset, are

95      needed for input, each with three tab-delimited columns. For the reference panel sample

96      information text file, column one must contain sample names identical to the naming in the VCF

97      file with one sample per row; column two must specify genetic sex information ("Male" "Female"

98      or "N/A" if unknown, case insensitive); column three must contain population assignment. For the

99      study sample information text file, columns one and two are similar to the reference panel file, but

100     column three is a dummy variable with a single arbitrary value that is the same on every row. For

4

101   example, column three of the sample information text file for the unknown set of samples could

102   be set as "unknown". Finally, the user must provide the FASTA file (.fa) of the reference genome

103   used for read mapping along with a FASTA index file (.fai) and a sequence dictionary file (.dict).

104

105   **Data processing**

106         PopInf v1.0 implements filtering, merging, and file conversion prior to PCA. Single

107   nucleotide polymorphisms (SNPs) are extracted from both the reference panel and study sample

108   VCF files, using GATK v3.7 SelectVariants and merged using GATK v3.7 CombineVariants (12).

109   To ensure PopInf analyzes SNPs that overlap with both the reference and unknown variant sets,

110   missing genotype data is removed using VCFtools v.0.1.14 (vcftools --max-missing flag) (13). If

111   analyzing the X chromosome, the pseudoautosomal regions and X-transposed region (18,19) are

112   masked using bedtools v.2.27.1 (14). Prior to running PCA, the merged VCF file is pruned for

113   linkage disequilibrium (LD) and converted to plink format using Plink v1.9 (15). PCA on a user-

114   defined set of chromosomes (e.g. whole genome, all autosomes, or a single chromosome) is

115   carried out using smartpca (6).

116

117   **Output**

118         PopInf v1.0 generates PCA plots for the first ten PCs for the study samples and the

119   reference panel, and an inferred ancestry report. Genetic ancestry of each study sample is

120   inferred based on the distance between the study sample and the centroid coordinates of PCs 1

121   and 2 of each reference population. A study sample is inferred to originate from a particular

122   population if it falls within N standard deviations (SDs) from the reference population centroid. To

123   provide multiple levels of confidence, the ancestry is inferred using 1, 2, and 3 SDs. If the sample

124   does not fall within three standard deviations of any population, the sample's ancestry will be

125   assigned to the closest population or will be assigned as having admixed ancestry: PopInf

126   calculates the midpoint coordinates between each pairwise combination of reference populations

127     and then compares those distances to the study sample. For a sample to be assigned as admixed,

128     it must be closer to the midpoint of two populations than to the 3rd standard deviation of any

129     population. If the study sample is closer to the 3rd standard deviation of a population than any of

130     the midpoints, it will be assigned to that population.

131

132     **RESULTS**

133     **Usage examples**

134          We ran PopInf v1.0 using variants from two human genetic datasets: one from healthy

135     individuals and one from cancer patients. The GTEx Project (10) dataset consisted of 148

136     individuals and the TCGA liver cancer dataset (11) consisted of 403 individuals (Supplementary

137     Table 1; Supplementary Table 2; Figure 1). Both datasets included self-reported race for most

138     individuals. We inferred the genetic ancestries of these samples based on a reference panel

139     consisting of variants from 986 unrelated individuals from populations across Africa, Europe, East

140     Asian, and South Asia from 1000 Genomes Release 3 (20) (Supplementary Table 3).

141          We find that, using genome-wide genotypes, the genetic ancestry of most study samples

142     does match that which is reported, with notable exceptions, and that we are able to infer ancestry

143     of samples of unreported origin. The inferred ancestry matches closely with the self-reported race

144     information in the GTEx dataset (Supplementary Table 4). One of the GTEx individuals was

145     missing self-reported race. Based on genetic ancestry, this individual was inferred as admixed

146     East Asian and South Asian (Supplementary Table 4). In the TCGA liver cancer dataset, we found

147     11 individuals with discrepancies between self-reported race and inferred ancestry; for all of these

148     individuals, their self-reported race was white and inferred ancestry was South Asian

149     (Supplementary Table 5). We further inferred ancestry for the 10 individuals in the TCGA liver

150     cancer dataset with no self-reported race (Supplementary Table 5).

151          We additionally ran PopInf v1.0 on each autosome and the X chromosome separately,

152     finding that chromosome-specific ancestry does not always match that inferred from the whole

6

153   genome (Figure 2A and C). We identify 16 individuals in the GTEx dataset and 56 individuals in

154   the TCGA dataset (Figure 2 B and D) with variation in chromosome-specific ancestry. All of the

155   admixed individuals had different inferred ancestry results among their chromosomes, as

156   expected. However, there were also 60 (12 from GTEx and 48 from TCGA) individuals inferred

157   as having only one ancestry when analyzing all autosomes together that showed variation in

158   chromosome-specific ancestry (Figure 2 B and D). These ancestry differences across the genome

159   shows that assigning ancestry based only on genome-wide genotypes may result in missing

160   clusters of ancestry across any single chromosome, which may lower our ability to identify risk

161   alleles in datasets consisting of samples of diverse and admixed backgrounds.

162

163   **CONCLUSION**

164       Here, we provide a workflow that will set up and run PCA, summarize the PCA output, and

165   provide the user with plots and an easily searchable inferred ancestry report for samples with

166   unknown or unverified population information. Inferred ancestry results from the GTEx and TCGA

167   datasets revealed heterogeneity in ancestry across the genome, and by chromosome. PopInf can

168   be modified to work with any reference panel, and may be applied to similarly infer chromosomal

169   and genome-wide ancestry in diverse populations.

170

171   **Acknowledgements**

178     Arizona State University for providing high performance computing resources that have

179     contributed to the research results reported within this paper.

180

181     **Authors' contributions**

182     MAW and AMTO conceived the ideas and designed methodology; AMTO and AJD collected the

183     data and analyzed the data. HMN contributed to processing the TCGA data. AMTO and MAW led

184     the writing of the first draft of the manuscript. All authors contributed critically to writing and editing

185     the drafts and gave final approval for publication.

186

187     **Data accessibility**

188     PopInf v1.0, processed 1000 Genomes reference file used in this manuscript, and an

189     accompanying tutorial are available on Github: https://github.com/SexChrLab/PopInf.
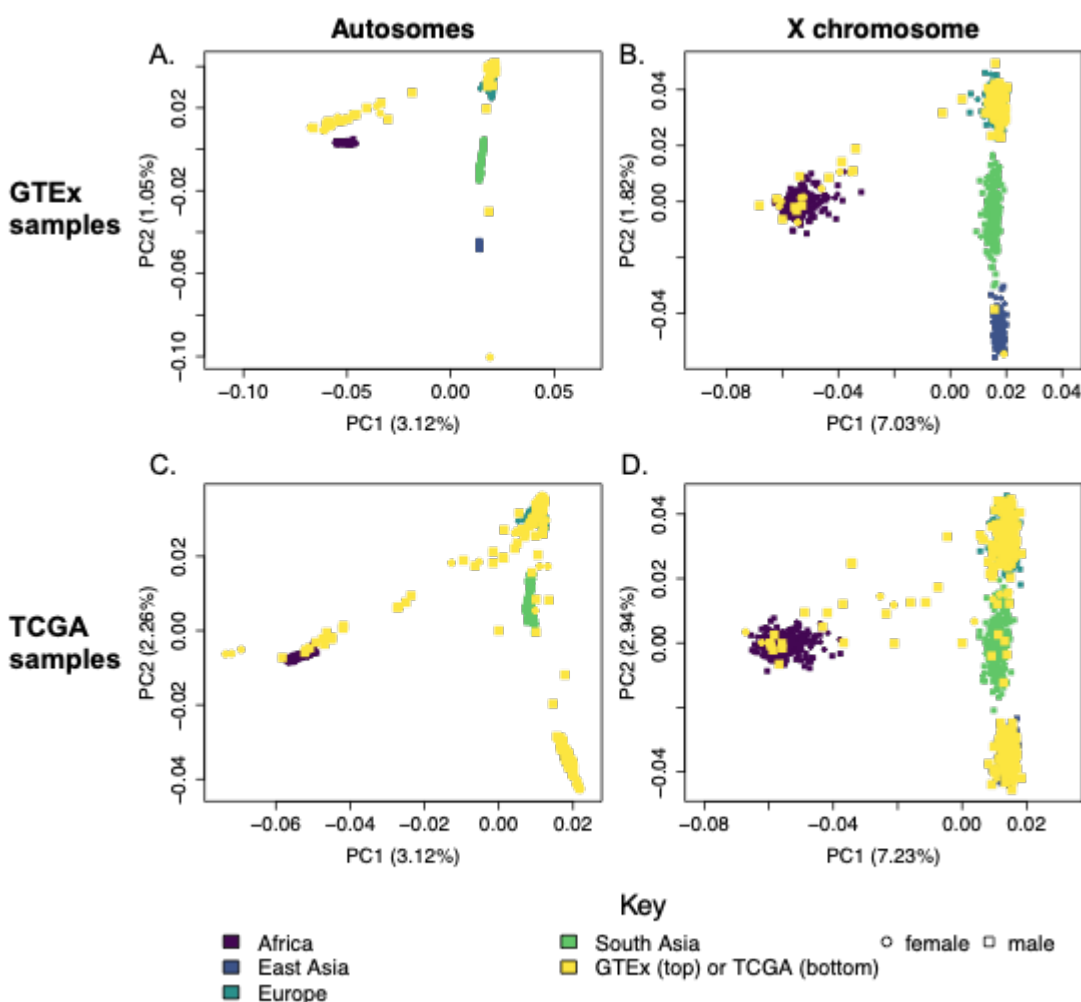
8

190 **FIGURES**

191



**Figure 1. Principal Component Analysis (PCA) output from a sample datasets plotted against the reference dataset.** Principal Components 1 and 2 for all individuals for A) autosomes merged and B) X chromosome for the GTEx dataset, and C) autosomes merged and D) X chromosome for the TCGA dataset. Purple points represent the reference samples of African descent, blue points represent reference samples of East Asian descent, dark green points represent reference samples of European descent, and light green represents reference samples

199    of South Asian descent in the 1000 Genomes reference panel. Yellow points represent samples

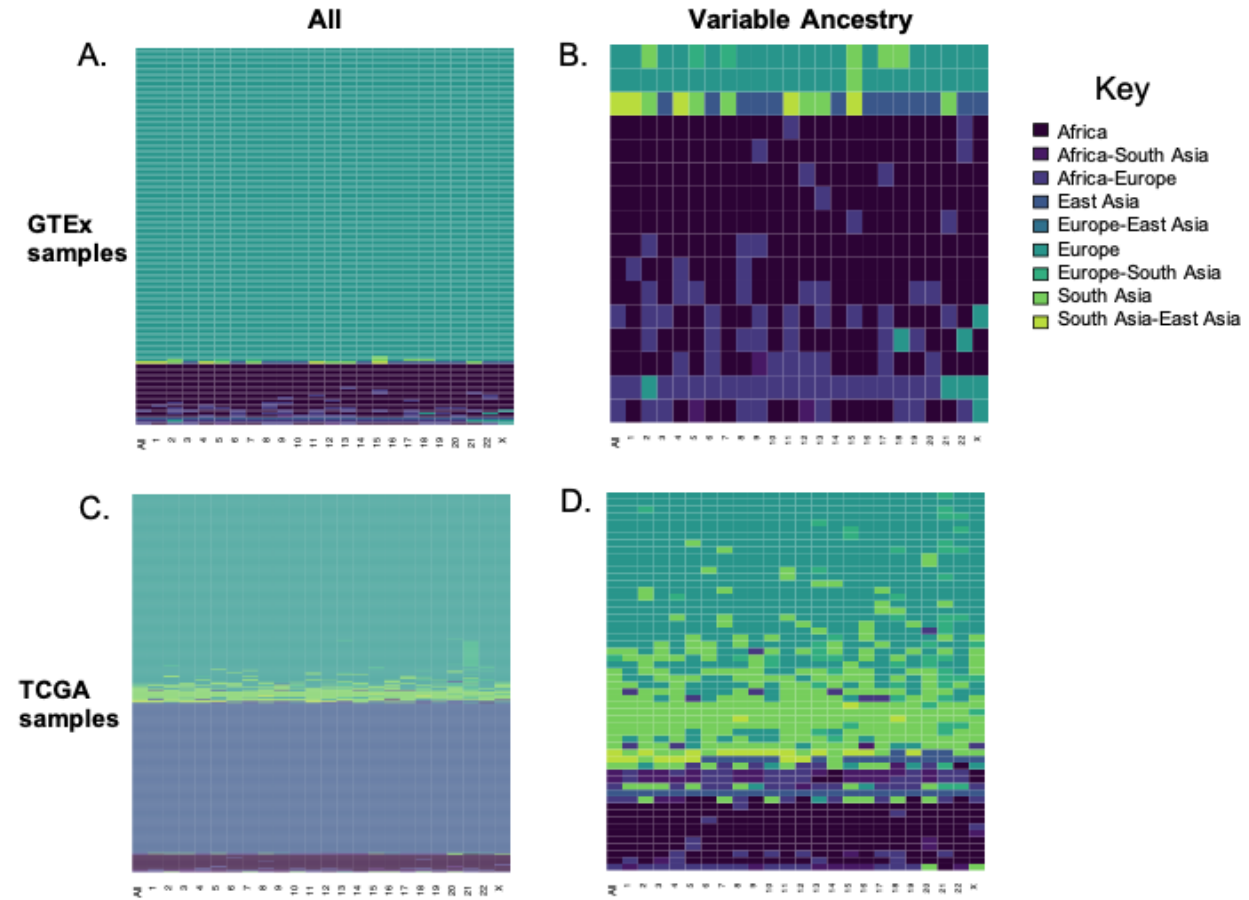200    from the sample datasets (GTEx and TCGA).

201



202

203    **Figure 2. Inferred ancestry for all autosomes combined and each chromosome separately.**

204    A) All 148 GTEx individuals, B) the subset of GTEx individuals with variation in inferred ancestry

205    among their chromosomes. C) All 403 TCGA individuals D) the subset of TCGA individuals with

206    variation in inferred ancestry among their chromosomes. Males and females were run together,

207    and only the autosomes and X chromosome were analyzed. The x-axis represents the

208    chromosome analyzed and the y-axis represents the individual from the GTEx dataset. Colors

209    represent inferred ancestry.

210

211

212     **Supplementary Table 1. GTEx samples used as the unknown sample dataset.** Here we

213     analyzed the population ancestry from whole genome sequence data from 148 samples available

214     from the GTEx dataset. GTEx (release V6p) whole genome sequence data (dbGaP accession

215     #8834) were downloaded from dbGaP.

216

217     **Supplementary Table 2. TCGA samples used as the unknown sample dataset.** Here we

218     analyzed the population ancestry from whole exome sequence data from 403 samples available

219     from the TCGA dataset. TCGA whole exome sequence data (dbGaP accession #11368) were

220     downloaded from NCI Genomic Data Commons (21).

221

222     **Supplementary Table 3. 1000 Genomes samples used for this reference panel.** Here we

223     chose 986 unrelated individuals from 1000 Genomes release 3 data downloaded as VCF mapped

224     to GRCh37 from ftp://ftp.1000genomes.ebi.ac.uk/vol1/ftp/release/20130502/. To include global

225     genetic variation in the reference panel, we chose individuals across populations in Africa, Asia,

226     and Europe.

227

228     **Supplementary Table 4. GTEx inferred ancestry and self-reported race comparison.** We ran

229     PopInf for all autosomes merged and the X chromosome separately on each individual in the

230     GTEx dataset. We compared these results to the self-reported race information for each

231     individual.

232

233     **Supplementary Table 5. TCGA inferred ancestry and self-reported race comparison.** We

234     ran PopInf for all autosomes merged and the X chromosome separately on each individual in the

235     TCGA dataset. We compared these results to the self-reported race information for each

236     individual.

## REFERENCES

1. Timpson NJ, Greenwood CMT, Soranzo N, Lawson DJ, Richards JB. Genetic architecture: the shape of the genetic contribution to human traits and disease. Nature Reviews Genetics. 2018;19:110–24.

2. Hindorff LA, Gillanders EM, Manolio TA. Genetic architecture of cancer and other complex diseases: lessons learned and future directions. Carcinogenesis. 2011;32:945–54.

3. Price AL, Zaitlen NA, Reich D, Patterson N. New approaches to population stratification in genome-wide association studies. Nature Reviews Genetics. 2010;11:459–63.

4. Yuan J, Hu Z, Mahal BA, Zhao SD, Kensler KH, Pi J, et al. Integrated Analysis of Genetic Ancestry and Genomic Alterations across Cancers. Cancer Cell. 2018;34:549-560.e9.

5. Dutil J, Chen Z, Monteiro AN, Teer JK, Eschrich SA. An Interactive Resource to Probe Genetic Diversity and Estimated Ancestry in Cancer Cell Lines. Cancer Res. 2019;79:1263–73.

6. Patterson N, Price AL, Reich D. Population Structure and Eigenanalysis. PLoS Genetics. 2006;2:e190.

7. Alexander DH, Novembre J, Lange K. Fast model-based estimation of ancestry in unrelated individuals. Genome Res [Internet]. 2009 [cited 2018 Jul 12]; Available from: http://genome.cshlp.org/content/early/2009/07/31/gr.094052.109

8. Maples BK, Gravel S, Kenny EE, Bustamante CD. RFMix: A Discriminative Modeling Approach for Rapid and Robust Local-Ancestry Inference. The American Journal of Human Genetics. 2013;93:278–88.

9. Pedersen BS, Quinlan AR. Who's Who? Detecting and Resolving Sample Anomalies in Human DNA Sequencing Studies with Peddy. The American Journal of Human Genetics. 2017;100:406–13.

10. Lonsdale J, Thomas J, Salvatore M, Phillips R, Lo E, Shad S, et al. The Genotype-Tissue Expression (GTEx) project [Internet]. Nature Genetics. 2013 [cited 2018 Jul 5]. Available from: http://www.nature.com/articles/ng.2653

11. Ally A, Balasundaram M, Carlsen R, Chuah E, Clarke A, Dhalla N, et al. Comprehensive and Integrative Genomic Characterization of Hepatocellular Carcinoma. Cell. 2017;169:1327-1341.e23.

12. McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, et al. The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data. Genome Res. 2010;20:1297–303.

13. Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, et al. The variant call format and VCFtools. Bioinformatics. 2011;27:2156–8.

14. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic features. Bioinformatics. 2010;26:841–2.

274    15.  Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ. Second-generation
275          PLINK: rising to the challenge of larger and richer datasets. Gigascience. 2015;4:7.

276    16.  R Development Core Team R. R: A language and environment for statistical computing. R
277          foundation for statistical computing Vienna, Austria; 2011.

278    17.  Koster J, Rahmann S. Snakemake--a scalable bioinformatics workflow engine.
279          Bioinformatics. 2012;28:2520–2.

280    18.  Ross MT, Grafham DV, Coffey AJ, Scherer S, McLay K, Muzny D, et al. The DNA
281          sequence of the human X chromosome. Nature. 2005;434:325–37.

282    19.  Skaletsky H, Kuroda-Kawaguchi T, Minx PJ, Cordum HS, Hillier L, Brown LG, et al. The
283          male-specific region of the human Y chromosome is a mosaic of discrete sequence
284          classes. Nature. 2003;423:825–37.

285    20.  Consortium T 1000 GP. A global reference for human genetic variation. Nature.
286          2015;526:68–74.

287    21.  Grossman RL, Heath AP, Ferretti V, Varmus HE, Lowy DR, Kibbe WA, et al. Toward a
288          Shared Vision for Cancer Genomic Data. New England Journal of Medicine.
289          2016;375:1109–12.

290