

1 Article

2 Genotype fingerprints enable fast and private 3 comparison of genetic testing results for research and 4 direct-to-consumer applications

5 Max Robinson¹ and Gustavo Glusman^{2,*}

6 ¹ Institute for Systems Biology, 401 Terry Ave N, Seattle, WA 98109; Max.Robinson@SystemsBiology.org

7 ² Institute for Systems Biology, 401 Terry Ave N, Seattle, WA 98109; Gustavo@SystemsBiology.org

8 * Correspondence: Gustavo@SystemsBiology.org; Tel.: +1-206-732-1273

9 Received: date; Accepted: date; Published: date

10 **Abstract:** Genetic testing has expanded out of the research laboratory into medical practice and the
11 direct-to-consumer market, and rapid analysis of the resulting genotype data can now have
12 significant impact. We present a method for summarizing personal genotypes as ‘genotype
13 fingerprints’ that meet these needs. Genotype fingerprints can be derived from any single
14 nucleotide polymorphism (SNP)-based assay, and remain comparable as chip designs evolve to
15 higher marker densities. We demonstrate that they support distinguishing types of relationships
16 among closely related individuals and closely related individuals from individuals from the same
17 background population, as well as high-throughput identification of identical genotypes,
18 individuals in known background populations, and de novo separation of subpopulations within
19 a large cohort through extremely rapid comparisons. While fingerprints do not preserve
20 anonymity, they provide a useful degree of privacy by summarizing a genotype in a way that
21 prevents reconstruction of individual marker states. Genotype fingerprints are therefore well-
22 suited as a format for public aggregation of genetic information to support ancestry and
23 relatedness determination without revealing personal health risk status.

24 **Keywords:** computational genomics, genome comparison, algorithms, genetic testing, privacy,
25 direct-to-consumer, study design, population genetics

27 1. Introduction

28 A very large number of genotypes have been produced by DNA hybridization, employing a
29 variety of array designs [1]. The low cost of hybridization assays relative to sequencing, including
30 whole-genome sequencing (WGS), exome sequencing, and other forms of targeted sequencing, has
31 led to the commoditization of array-based genotyping and has enabled commercial companies
32 (including 23andMe, AncestryDNA, Family Tree DNA, and others [2]) to offer this service directly
33 to consumers (DTC). DTC services typically yield results with high concordance [3] and low no-call
34 rates [4]. Nevertheless, genotyping the same individual using different array designs can yield
35 slightly different results, as each technology has its own biases. Even when using the same
36 technology, genotype reference version and variant encoding format, genotyping the same
37 individual repeatedly can give slightly different results due to the stochastic nature of genome
38 processing and analysis, batch effects, or differences in the computational pipelines used. Some
39 companies regularly reanalyze the raw data for all customers, refining the results over time; as a
40 result, customers who download their genotype data repeatedly over the years may have slightly
41 differing results even from the same sample. In addition to relatedness applications, array-based
42 genotyping is also used as a quality-control step prior to WGS when comprehensive variant
43 information is desired.

44 Many methods exist for comparing genome-wide genotypes in order to infer relatedness, with
45 varying degrees of accuracy [5]. Most methods are computationally demanding and require full
46 access to the genotype data of the individuals to be compared, potentially precluding their
47 application to the study of samples with restricted access, or direct use of these methods by non-
48 specialists interested in exploring their ancestry and genealogy.

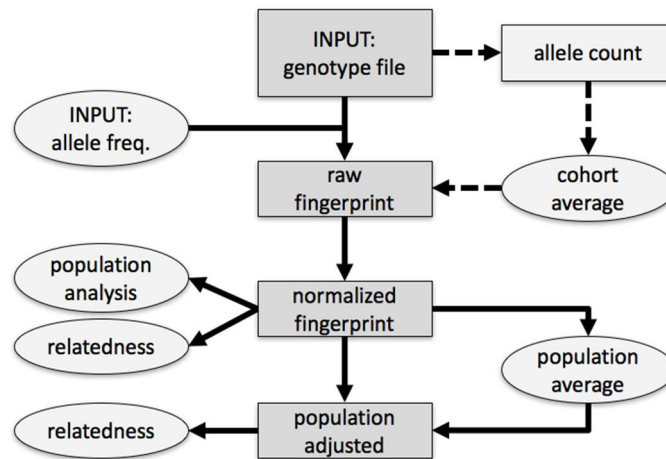
49 A genotype determined by DNA hybridization enumerates all observed alleles for a
50 predefined set of variants of interest, typically common single-nucleotide polymorphisms (SNPs).
51 In this format, each SNP is identified by its identifier ('rsid', reference SNP identifier) in the dbSNP
52 database [6]. For each rsid, the observed genotype of the individual is stated, including those for
53 which the individual is homozygous for the reference allele. The chromosome and coordinate of the
54 SNP, relative to a version of the reference that is (hopefully) stated in the 'header' of the genotype
55 file, is implied by the rsid and not recorded in the genotype file.

56 We have recently published a method for converting personal whole genome sequence data
57 into 'genome fingerprints' that facilitate (and greatly accelerate) their comparison [7]. Our method
58 encodes the characteristics of pairs of consecutive single nucleotide variants (SNVs) relative to a
59 reference, as represented in variant call format (VCF) files or structurally equivalent formats. In
60 contrast to the format of genotyping results, WGS results in VCF format typically encode only
61 differences from the reference; genomic locations in which the individual is homozygous for the
62 reference allele are typically not stated, achieving a more compact representation.

63 We present here an analogous method for summarizing personal genotypes, yielding
64 'genotype fingerprints' that can be readily compared to estimate relatedness. The genotype
65 fingerprints can be computed starting from any of several chip array designs, with genome
66 coordinates expressed relative to any reference version; the resulting fingerprints are directly
67 comparable without further conversion. Computation on the genotype fingerprints is fast and
68 requires little memory, enabling comparison of large sets of genomes. No individual variants or
69 other detailed features of the personal genome can be reconstructed from the fingerprint, thereby
70 allowing private information to be more closely guarded and protected and decoupling genome
71 comparison from genome interpretation. Fingerprints of different sizes allow balancing the speed
72 and accuracy of the comparisons. Due to the high value of estimating relatedness, the potential
73 applications of genotype fingerprinting range from basic science (study design, population studies)
74 to personalized medicine, forensics, and data privacy.

75 2. Materials and Methods

76 *Methodology overview.* We fingerprint genotypes in four stages. First we summarize a genotype
77 as a tally of biallelic SNPs, stratified by observed alleles, by variant identifiers, and accounting for
78 allele frequencies ('raw' fingerprint, Figure 1). We then normalize the raw fingerprint to account for
79 systematic methodological patterns. The resulting 'normalized' fingerprint preserves differences
80 between individuals from different groups (populations), and are appropriate for clustering
81 individuals by population. Whether individuals are assigned to populations *a priori* or via
82 clustering, we average the normalized fingerprints of the individuals in a population to produce a
83 'population' fingerprint, which characterizes the population rather than an individual. To improve
84 detection of relationships within a population, we derive a 'population-adjusted' fingerprint from
85 an individual's normalized fingerprint by subtracting the associated population fingerprint.
86 Documentation, code, and sample datasets are available at [8].



87
88
89
90
91
92
93
94

Figure 1. Overview of method. SNPs in the input genotype file are encoded into a table (raw) by observed alleles and rsid numerical value, taking allele frequencies into account; this can optionally be approximated by subtracting allele counts estimated from a simple model of an observed cohort (dashed arrows). The raw fingerprint is then normalized and may be adjusted to represent deviation from the center of the closest population. Rectangles and ellipses pertain to individual genotypes or to multiple genotypes, respectively; darker gray denotes the flow of information for one genotype, from the input file to the normalized and adjusted fingerprints.

95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119

Raw fingerprints. A ‘raw’ fingerprint is a 4 row by L column table of SNP allele counts, where L is the main parameter of the method and determines the information content of the fingerprint; by default $L=1000$. The four rows correspond to the permitted alleles of binary SNPs (A, C, G and T); variants with other possible alleles (including insertions, deletions, and multi-nucleotide variants) are ignored, partly because binary SNPs are so abundant but also because other variant classes are subject to more variation in how each observable allele is reported. We tally all observed alleles for each SNP (reference and alternate). Each SNP is tallied in a column determined from its rsid, its reference number in the dbSNP database [6]. The full process is as follows:

1. Filter out all variants that are not autosomal, biallelic SNPs with reference and alternate alleles limited to A, C, G, and T. Optionally, include only SNPs in a preselected set, e.g., 23andMe V2 and V3.

The remaining steps are applied to each retained SNP.

2. Determine a fingerprint column as the rsid modulo L (when $L=1000$, SNP rs1801133 will be recorded in column 133).
3. For each SNP, count observations of each nucleotide: n_A is the count of A alleles (0, 1, or 2), n_C is the count of C alleles, etc., with $n_N = n_A + n_C + n_G + n_T = 2$.
4. Determine the expected count $ex = n_N f_X$ for each nucleotide X from a set of known allele frequencies $f_A + f_C + f_G + f_T = 1$. Depending on the context, these frequencies may be specific to each SNP (e.g. for human data, extracted from dbSNP), or without reference to external data, may be computed per column from all SNPs contributing to the column among an observed cohort of genotypes (detailed below).
5. Tally differences $n_X - ex$ from expectation in each row and column: add $n_A - ex_A$ to the row for A, $n_C - ex_C$ to the row for C, etc., in the column corresponding to the rsid.

120
121
122
123

Retrieving the allele frequencies for each observed rsid requires prior knowledge and can incur significant computational costs. A more efficient variant involves computing expected frequencies for each row and column directly from a cohort of genotypes of a common type (determined with the same assay, array design, etc.) as follows:

- 124 1. Tally allele counts separately for each individual as above, except in step 5 increment
125 the value in the $[4 \times L]$ matrix by one for each observed allele. (In case of
126 homozygosity, incrementing twice results in an increment by two.) These tallies result
127 in summed N_A , N_C , N_G , and N_T values in each column, with column total $N_N =$
128 $N_A + N_C + N_G + N_T = 2k$, twice the number of SNPs assigned to the column. The steps
129 below do not require reprocessing the full genotypes, only these $4 \times L$ tallies per
130 individual.
- 131 2. Compute cohort average frequencies by summing the tallied allele counts in each
132 column across all individuals and dividing by the column total. This produces four
133 allele frequencies F_x for each column.
- 134 3. Finish computing each entry in each individual's raw fingerprint by subtracting the
135 expected count $e_x = N_N F_x$.

136

137 *Fingerprint normalization.* In this step, we account for unequal assortment of genotype
138 information within the fingerprint. This imbalance is due to a methodological aspect of the
139 fingerprinting process (grouping of variants by rsid), not a source of information about the
140 fingerprinted individuals.

- 141 1. We subtract the mean and divide by the standard deviation of each column, which
142 mitigates differences between columns in the types of nucleotide substitutions
143 (transitions and transversions), which derive from the set of rsids assigned to each
144 column.
- 145 2. We then subtract the mean and divide by the standard deviation of each row. This
146 further mitigates methodological differences between values within the fingerprint,
147 which primarily reflect the genotyping methodology rather than variation between
148 individuals.

149

150 *Population fingerprints.* We compute a population fingerprint as the average of the normalized
151 fingerprints from the individuals in the population. Note that genotype fingerprints are only
152 directly comparable when computed using the same format parameter L ; different values of L cause
153 rsids to be grouped into columns differently. However, different versions of a genotype array
154 design contain substantial overlaps in the set of SNPs the array contains, and rsids will be grouped
155 in the same manner for a given value of L regardless of array design. Thus, genotypes from the
156 same population on slightly different variants of the same array design may be mixed in computing
157 the population fingerprint.

158

159 *Adjusting fingerprints for population.* We then compute a population-adjusted fingerprint for an
160 individual by subtracting a population fingerprint from the normalized fingerprint of that
161 individual. These individual fingerprints and the population fingerprint must have been computed
162 using the same parameters L .

163

164 *Fingerprint comparison.* To compare two fingerprints, concatenate the rows of each fingerprint
165 matrix into a vector and compute the Spearman correlation between the two vectors. This same
166 procedure is appropriate for comparing two normalized fingerprints or two population-adjusted
167 fingerprints, whether adjusted to the same or different populations.

168

169 *Family analysis.* We obtained 23andMe SNP chip genotype data for a family of five [9],
170 including Mother, Father, Son, Daughter and Aunt. Son is 23andMe V2 data and the rest of the
171 family are 23andMe V3 data. We computed normalized genotype fingerprints ($L=5000$) for the five
172 individuals and performed all pairwise comparisons. We also extracted from these samples the lists
173 of rsids observed in V2 and V3, for use in further analyses below.

174

175 *Population structure analysis.* Principal Components Analysis is a standard method for
176 characterizing population structure prior to genome-wide association studies (GWAS). We
177 therefore compared well-characterized population structures within data from the 1000 Genomes
178 Project (release 20130502,
179 ALL.chrNN.phase3_shapeit2_mvncall_integrated_v5.20130502.genotypes.vcf.gz). As a genomic
180 method, we identified SNPs with a minor allele frequency of 5% or more, removed SNPs in
181 complete linkage disequilibrium with a SNP to the left (i.e., a smaller chromosomal position),
182 retained 5% at random (298,454 SNPs) and counted occurrences of the minor allele (0, 1, or 2) in
183 each genome to form a 2504 × 298,454 genotype matrix M.

184 As a fingerprint-based method, we extracted observed genotypes for each of 2504 genomes
185 twice, once for each rsid in the 23andMe V2 SNP list and once using the V3 SNP list. We then
186 computed genotype fingerprints from these extracted genotypes using $L=500$, 1000, and 5000,
187 resulting in fingerprint data matrices of 2504 × 2000, 2504 × 4000, and 2504 × 20,000 entries,
188 respectively. We performed Standard PCA separately on the six resulting matrices (V2 or V3 and
189 $L=500$, 1000, or 5000) using the R function call `prcomp(M,center=TRUE,scale.=TRUE)`.

190
191 *Evaluation of “nearest population fingerprint” for population assignment.* We computed a
192 population fingerprint for each of the 26 populations selected for study in the 1000 Genomes
193 Project. We then re-classified each individual via fingerprint comparison against the 26 population
194 fingerprints, as described above for comparing individual fingerprints. Each individual was
195 considered classified as belonging to the population with the closest population fingerprint. To
196 avoid overfitting, we excluded each individual from the computation of their own population
197 fingerprint in leave-one-out fashion.

198 3. Results

199 3.1. A method for encoding genotyping data

200 We developed a locality sensitive hashing [10] algorithm for computing ‘fingerprints’ from
201 genotype data, including data produced by direct-to-consumer (DTC) genetics companies (e.g.,
202 23andMe, AncestryDNA). These genotype fingerprints meet the characteristics of genotype data:
203 they can be rapidly computed starting from any of several chip array designs, with genome
204 coordinates expressed relative to any reference version, and the resulting fingerprints are directly
205 comparable so long as the same fingerprint length L is used. We describe fingerprints generated
206 using SNP lists derived from two array designs used by 23andMe: V2, based on Illumina
207 HumanHap550 Genotyping BeadChip (~550,000 SNPs) and V3, based on Illumina OmniExpress
208 Genotyping BeadChip (~960,000 SNPs). The fingerprints are a reduced representation of the
209 genotype data computed once per individual, and can be efficiently databased and compared to
210 determine whether two genotypes represent the same individual, closely related individuals, or
211 unrelated individuals. As with our previously reported genome fingerprints [7], individual alleles
212 cannot be reconstructed from the genotype fingerprint beyond what is predictable from detectable
213 population and family relationships, enabling sharing of fingerprints for comparison when privacy
214 concerns prevent sharing the full genotype file itself.

215 The main parameter of our algorithm, L , determines the size and SNP groupings of the
216 fingerprint. Smaller fingerprints (e.g., $L=100$) average variants over fewer bins and are useful for
217 extremely fast, low-resolution comparisons e.g. to determine identity, while larger fingerprints
218 ($L=1000$ or 5000) are higher resolution and better support detailed analyses, including population
219 reconstruction.

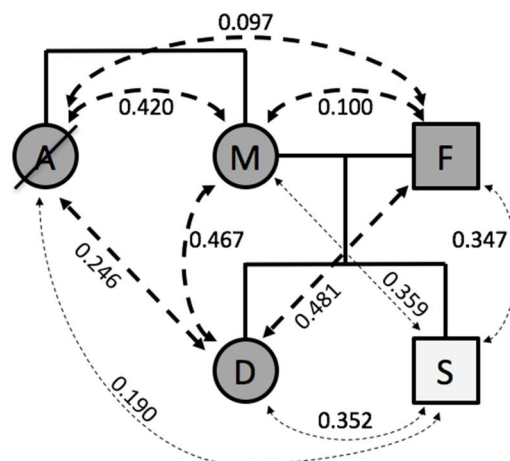
220 3.2. Computation on genotype fingerprints is fast

221 Generation of raw genotype fingerprints is quite simple, and its speed is limited by the time
222 taken to read the file rather than the fingerprint size L ; on our workstations, it requires 10-15 s per

223 500,000-960,000 SNP genotype. For higher speeds, each raw fingerprint is independent and the
 224 process is readily parallelized. In our population studies, population fingerprints for all 26
 225 populations of the 1000 Genomes cohort required less than 60 s, fingerprint normalization averaged
 226 0.13 s per genotype, and serializing the 2504 fingerprints for efficient comparison took only 37 s.
 227 The 3,133,756 “all-against-all” comparisons for this data set took 15 CPU seconds for $L=1000$ (4.8
 228 microseconds per comparison) and 79 CPU seconds for $L=5000$ (25.2 microseconds per comparison).
 229 Although parallelization was not required for this cohort, it is straightforward and may be helpful
 230 for all-against-all comparisons of cohorts of tens to hundreds of thousands of genotypes.

231 3.3. Rapid relationship detection

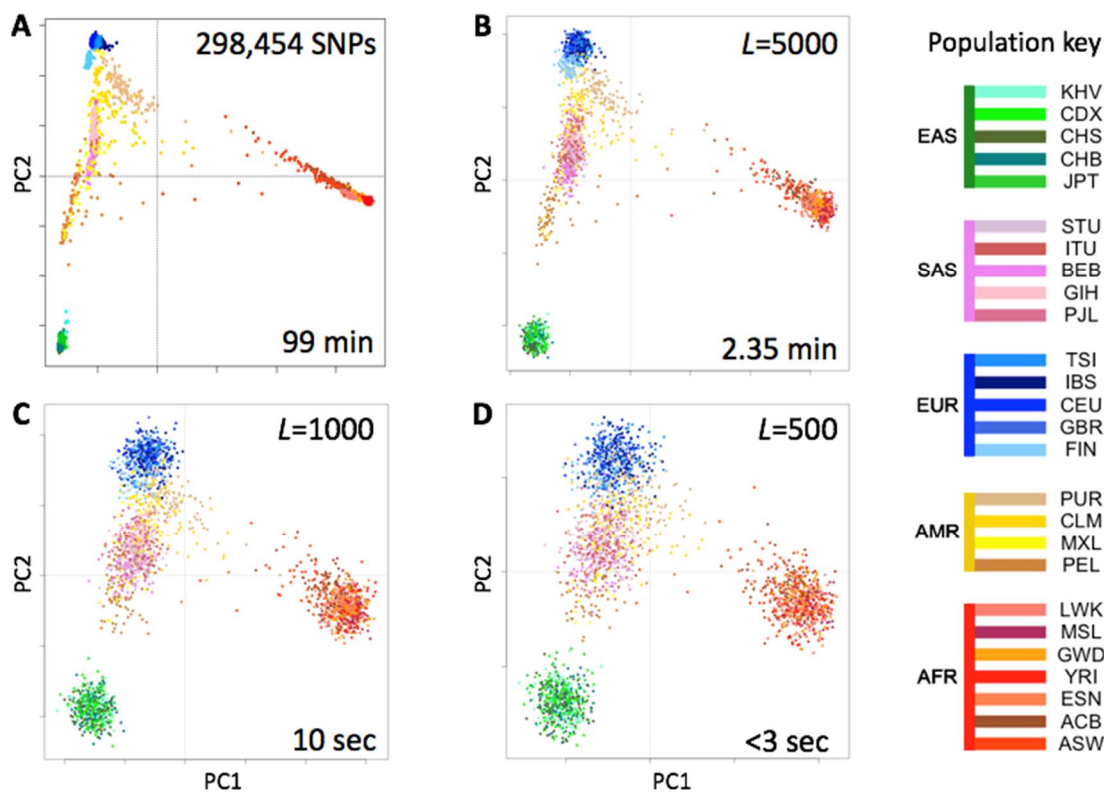
232 Here we illustrate the use of genotype fingerprints for characterizing family relationships
 233 within a family of five [4], who had previously made their 23andMe genotype results publicly
 234 available. Comparisons of these fingerprints resulted in similarity scores (Spearman’s rho values)
 235 that are consistent with the known family relationship types (Figure 2). Rho values for full sibling
 236 pairs (Aunt and Mother .420, Daughter and Son .352) and parent-offspring pairs (.481, .467, .359,
 237 .347) are higher than for avuncular relationships (Aunt and Daughter .246, Aunt and Son .190),
 238 which in turn are higher than unrelated pairs (Aunt and Father .097, Mother and Father .100). The
 239 correlations between the Son and the other family members was reduced, as expected for
 240 comparisons across SNP lists (V2 for Son, V3 for all others) with both substantial overlap and
 241 differences.



242 **Figure 2.** Comparison within a family of five. A: Aunt (deceased); M: Mother; F: Father; D:
 243 Daughter; S: Son. Dashed lines represent family relationships; thin lines denote comparison
 244 between individuals assayed on different versions of the genotyping platform.
 245

246 3.4. Rapid analysis of population structure

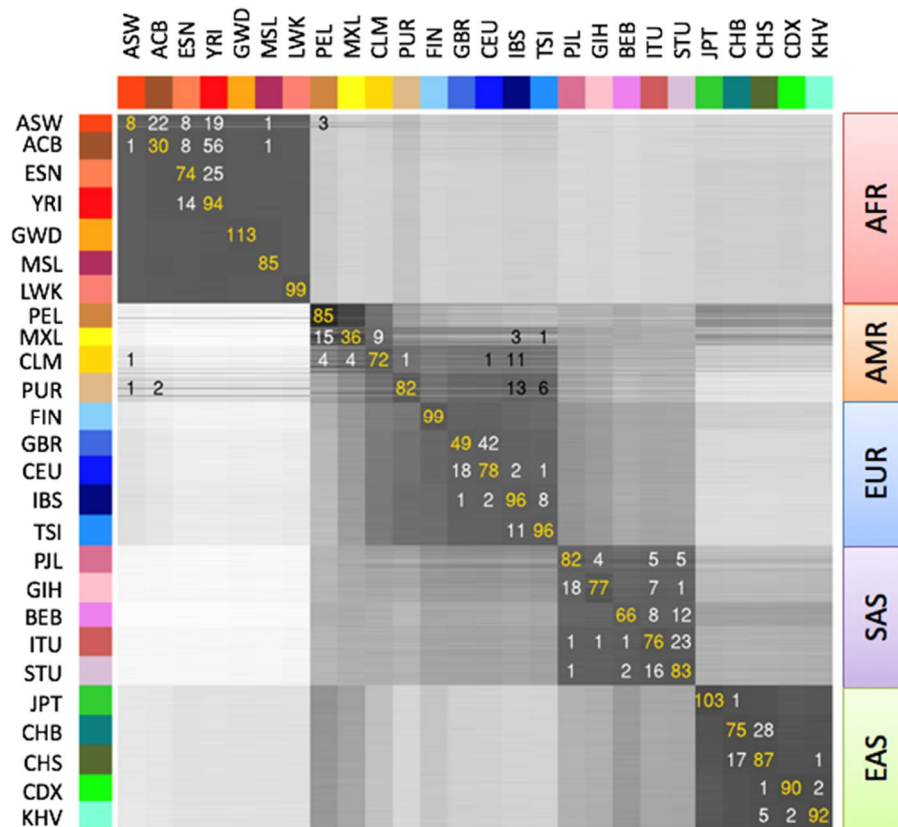
247 We tested the utility of genotype fingerprints for population studies. We extracted “genotypes”
 248 for all 2504 individuals of the 1000 Genomes Project cohort from their VCF-format genomes using
 249 the 23andMe V3 SNP list, fingerprinted each extracted genotype, and used PCA to reconstruct the
 250 known population structure (Figure 3). This entire process took a fraction of the time and memory
 251 required to perform the same task using a more standard approach (see Methods). As expected, the
 252 quality of the reconstruction depended on fingerprint resolution (L): fingerprints with $L=5000$
 253 yielded excellent population structure reconstruction, comparable to the results of population
 254 reconstruction using high-resolution genome fingerprints (compare Fig. 3, Fig. 5 from [7]).
 255 Genotype fingerprints with smaller values of L yielded progressively lower-resolution results at
 256 proportionally higher speeds (seconds rather than more than an hour, Figure 3). Thus, genotype
 257 fingerprints allow a suitable balance between resolution and speed to be achieved, and supports
 258 scaling of population structure studies to whole population-scale genotype data, without requiring
 259 analysis of linkage disequilibrium and other complications prior to analysis.



260
 261 **Figure 3.** Estimates of population structure in the 1000 Genomes Project data set at different
 262 resolutions. Individuals are color coded according to their population as per the key to the right.
 263 EAS, SAS, EUR, AMR and AFR: East Asian, South Asian, European, Admixed American, and
 264 African, respectively. (A) Principal components analysis (PCA) of the 2504 individuals using
 265 ~300,000 SNPs. (B) PCA on genotype fingerprints with $L=5000$. (C) PCA on genotype fingerprints
 266 with $L=1000$. (D) PCA on genotype fingerprints with $L=500$.

267 3.5. Rapid population assignment

268 We computed “population fingerprints” in the 1000 Genomes data set by averaging genotype
 269 fingerprints (V3 set, $L=5000$) of the individuals in each population. To determine each individual’s
 270 population of origin, we then compute the correlation between the fingerprint of a query genome
 271 and the fingerprint of each population (Figure 4), and classified each individual as belonging to the
 272 population with the strongest fingerprint correlation. We evaluated this classification method by
 273 “leave one out” cross-validation. The annotated population had the highest correlation for 2027 of
 274 2504 samples (81%), or among the top 2 (92.9%) or top 3 (96.1%) most correlated. The only
 275 misclassifications to a population from another continent involved the Admixed American
 276 populations (AMR); excluding these populations increased correct classifications to 85.7% (best
 277 match), 97.8% (top 2) and 99.2% (top 3). In general, misclassifications both between and within
 278 continental groups were between historically or geographically associated population pairs (ASW,
 279 ACB with Nigerian populations ESN and YRI; Latin American populations MXL, CLM, and PUR
 280 with Mediterranean populations IBS, TSI; northern and southern Han Chinese populations CHB,
 281 CHS; south Indian populations ITU, STU), suggesting that admixture was the principal source of
 282 misclassification. Even for very closely related population pairs (e.g. CEU, GBR), individuals were
 283 more often correctly classified than misclassified.



284

285

286

287

288

289

Figure 4. Correlations between the genotype fingerprints of the 2504 individuals (rows) and the average fingerprints of the 26 populations (columns) in the 1000 Genomes Project. Population codes and colors as in Figure 3. Numbers in gold, white and black denote population assignments: to the same annotated population, to the same continent but different population, or to a different continent, respectively.

290

3.6. Robustness to SNP list

291

292

293

294

295

296

297

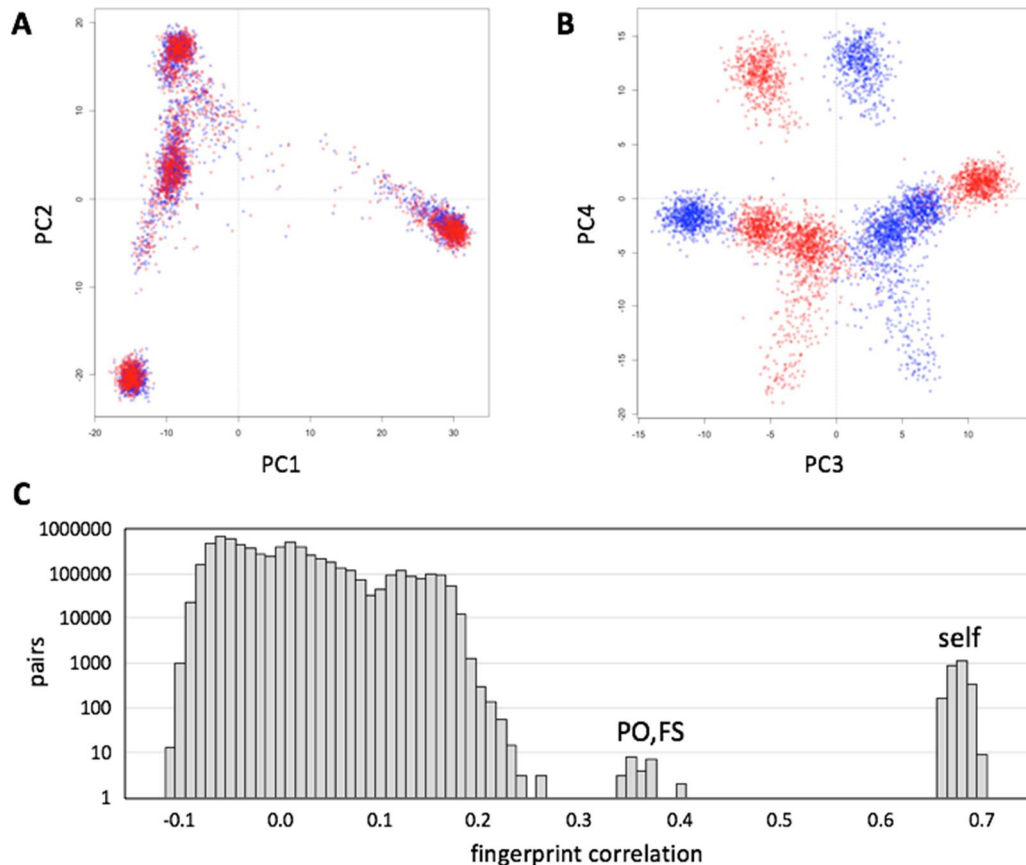
298

299

300

301

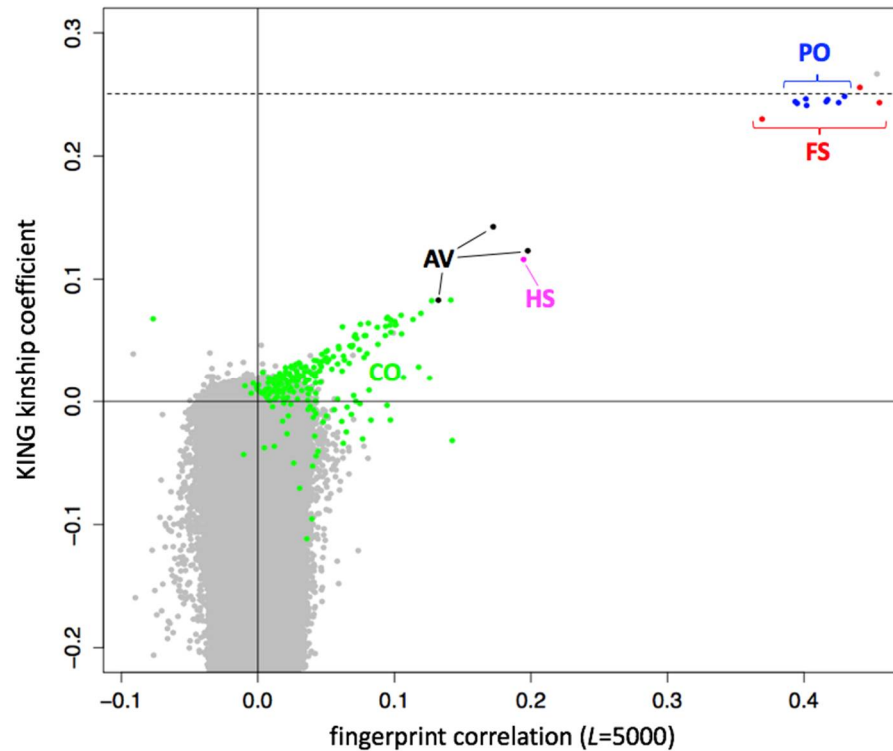
We evaluated whether genotype fingerprints can be compared across chip array designs. We fingerprinted genotypes for each individual in the 1000 Genomes Project data set extracted using the 23andMe V2 and V3 SNP lists (548,911 and 902,448 SNPs, respectively), yielding a mixed set of 5008 fingerprints (V2 and V3 fingerprints for each of 2504 individuals). We studied this joint set using PCA (Figure 5A) and observed that the first two principal components reconstruct the known population structure (Figure 3). PC3 separates between fingerprints computed on V2 and V3 versions (Figure 5B). Furthermore, the correlations between the two versions of each individual (self, Figure 5C) are always higher than those between related individuals (PO, FS), which are in turn higher than those between unrelated individuals. Thus, while fingerprints derived using different versions of the same chip design are distinguishable, comparisons between them are still useful for detecting identical individuals, family analysis, and population analysis.



302 **Figure 5.** Comparison of genotype fingerprints relative to different SNP lists. We deduced
 303 normalized genotype fingerprints ($L=5000$) for the 1000 Genomes Project cohort using the 23andMe
 304 V2 (red) and V3 (blue) SNP lists. (A) First two principal components, showing population structure.
 305 (B) Third and fourth principal components, showing separation between the two SNP lists. (C)
 306 Distribution of cross-correlations between the two sets of genotype fingerprints (all possible pairs of
 307 V2 vs. V3). Comparisons between the two genotype fingerprints for the same individual (self) and
 308 comparisons between parent/offspring and full-sibling pairs (PO, FS) formed distinct, high-
 309 correlation subsets.
 310

311 3.7. Fast detection of close relationships

312 Based on our work with genome fingerprints, we reasoned that using population fingerprints
 313 to cancel out correlations due to information shared among a population, therefore allowing close
 314 relationships to be distinguished from shared population background. We therefore adjusted the
 315 $L=5000$ V3-type fingerprints (see 3.6 above) for their annotated population of origin and performed
 316 all pairwise comparisons of these adjusted fingerprints. We compared the fingerprint correlations
 317 with kinship coefficients computed using KING [11] and with previously reported relationships
 318 [12] computed using RELPAIR [13] (Figure 6). As expected, the correlation between individuals
 319 from the same annotated population (Figure 4), but not related within a few generations, is
 320 removed by adjustment to the population average (Figure 6), and population-adjusted fingerprints
 321 for unrelated individuals were essentially uncorrelated. Comparison of population-adjusted
 322 genotype fingerprints therefore supports the detection of individuals in the 1000 Genomes cohort
 323 previously reported as closely related [12]. The highly-correlated pairs correspond to relationships
 324 of degrees varying from full siblings to cousins; for these pairs, fingerprint correlations show a
 325 linear relationship with kinship coefficients computed by KING (Figure 6, colored points).
 326 Parent/offspring and full sibling relationships, which have the same expected KING kinship
 327 coefficient (0.25) but different variance from that expected value, produced equivalent high
 328 fingerprint correlations (around 0.4).



329
330
331
332
333
334
335

Figure 6. Identification of close relationships in the 1000 Genomes Project. Comparison between the correlations of population-adjusted genotype fingerprints (V3 set, $L=5000$) and the kinship coefficient as computed using KING, highlighting close relationships identified using RELPAIR. FS: full siblings (red). PO: parent/offspring (blue). HS: half siblings (magenta). AV: avuncular (black). CO: cousins (green). All other pairs in gray. One FS pair (HG03873 and HG03998, with maximal kinship, in gray) was not identified by RELPAIR.

336

4. Discussion

337
338
339
340
341
342
343

We have presented a method for computing ‘fingerprints’ of genomewide SNP array genotypes as reported by DTC genetics companies, using 23andMe data as an example. Like our previously reported fingerprints from whole-genome resequencing data, genotype fingerprints retain sufficient information to enable ultrafast comparison of genotypes, without retaining the sensitive, individual SNP data necessary to predict phenotypes; genotype fingerprints are therefore suitable for databasing and sharing for ancestry and close relationship determination without exposing more sensitive, health-related information.

344
345
346
347
348
349

We demonstrated the utility of genotype fingerprints for rapid versions of common tasks: identifying genotypes from the same individual, from closely related individuals, or from a known population, and *de novo* clustering of individuals into subpopulations. Comparing fingerprints derived from two different SNP lists (23andMe V2 and V3), our genotype fingerprints were robust to differences in the number of SNPs assayed for detecting identity, detecting close relationships, and for delineating populations.

350
351
352
353
354
355
356
357
358

Conceptually, genotype fingerprints are an adaptation of our genome fingerprinting method [7] to more widely available, more standardized, but lower-resolution genotype data. While genome fingerprints facilitate comparison of data across different reference sequence versions by encoding consecutive SNV pairs, genotype fingerprints achieve a similar interoperability by encoding individual SNPs using annotated rsids, alleles, and allele frequencies. SNPs are simply SNVs with high population frequency, but this frequency difference has practical consequences. While whole genome sequencing is expected to reveal an increasing number of rare SNVs, the vast majority of SNPs have already been identified, evaluated for linkage, assigned stable identifiers (rsids), and incorporated into high-throughput assays. In contrast, many SNVs either lack

359 identifiers or have been assigned preliminary identifiers still subject to change (e.g., by merging
360 with a different identifier representing the same variant). Stable identifiers facilitate matching
361 variants across genome reference versions and assays, enabling the desired robustness to a
362 changing reference genome using the simpler encoding method presented here.

363 Insertions and deletions have also been assigned standard representations in genotype files
364 (symbols I and D, respectively), but are much less abundant in the genome than SNPs, are not as
365 widely assayed, and require normalization prior to extraction as genotypes from WGS or exome
366 data [14]. For simplicity and consistency, we therefore chose to exclude them from analysis, as we
367 did for computing genome fingerprints from VCF files. We also chose to exclude SNPs on the sex
368 chromosomes, which vary in count between males and females and may lead to distorted similarity
369 values.

370 Sharing genetic information raises privacy considerations of several kinds. Much attention has
371 been paid to the risk of re-identification of de-identified samples [15], even when querying genetic
372 data sets via bandwidth-limiting interfaces like the GA4GH beacons. These concerns have given rise
373 to privacy preservation strategies such as obscuring rare variants and budgeting queries [16]. While
374 enabling an important and powerful query - namely, "has this allele been seen before?" [17] - these
375 strategies for preventing re-identification preclude multiple other potential applications, thus
376 limiting the utility of genome data sharing. There are however genetic data sharing scenarios in
377 which anonymity is not an issue, but phenotype prediction is. For example, an individual may wish
378 to compare their genotype (obtained via a DTC genetic testing company) to the genotypes of other
379 individuals for ancestry and relationship determination, but without revealing whether their
380 genome harbors alleles associated with a specific phenotype, e.g., Alzheimer's disease - both
381 currently known alleles, and ones whose significance may be discovered in the future. Like genome
382 fingerprints, genotype fingerprints decouple genotype comparison from genotype interpretation,
383 supporting the identification of closely related individuals without exposing individual variant
384 states.

385 At present, the number of private individuals who have used DTC genetics services to
386 ascertain their own genotype vastly exceeds the number of individuals with full genome data. We
387 expect genotype fingerprints to have immediate applicability for facilitating genotype comparisons,
388 empowering citizen science without concomitantly revealing sensitive private genetic information.

389 **Supplementary Materials:** Documentation, code, and sample datasets are available at [8].

390 **Author Contributions:** GG designed the study. MR and GG performed analyses, wrote the manuscript and
391 approved its final version.

392 **Funding:** This research was funded by NIH grant number U54 EB020406.

393 **Acknowledgments:** We wish to thank the Corpas family for releasing their individual genotypes as CC0.

394 **Conflicts of Interest:** GG and MR hold a provisional patent application on the method described in this
395 manuscript. GG holds stock options in Arivale, Inc. Arivale, Inc. did not fund the study and was not involved
396 in its design, implementation, or reporting. The funders had no role in the design of the study; in the collection,
397 analyses, or interpretation of data; in the writing of the manuscript, and in the decision to publish the results.

398 References

- 399 1. Canada, R. Exploring Microarray Chips Available online: [http://haplogroup.org/exploring-microarray-](http://haplogroup.org/exploring-microarray-chips/)
400 [chips/](http://haplogroup.org/exploring-microarray-chips/) (accessed on Aug 17, 2018).
- 401 2. List of DNA testing companies - ISOGG Wiki Available online:
402 https://isogg.org/wiki/List_of_DNA_testing_companies (accessed on Aug 17, 2018).
- 403 3. Imai, K.; Kricka, L. J.; Fortina, P. Concordance study of 3 direct-to-consumer genetic-testing services. *Clin.*
404 *Chem.* **2011**, *57*, 518–521.
- 405 4. Glusman, G.; Cariaso, M.; Jimenez, R.; Swan, D.; Greshake, B.; Bhak, J.; Logan, D. W.; Corpas, M. Low
406 budget analysis of Direct-To-Consumer genomic testing familial data. *F1000 Research* **2012**, *1*,
407 doi:10.3410/f1000research.1-3.v1.

- 408 5. Ramstetter, M. D.; Dyer, T. D.; Lehman, D. M.; Curran, J. E.; Duggirala, R.; Blangero, J.; Mezey, J. G.;
409 Williams, A. L. Benchmarking Relatedness Inference Methods with Genome-Wide Data from Thousands
410 of Relatives. *Genetics* **2017**, *207*, 75–82.
- 411 6. Sherry, S. T.; Ward, M. H.; Kholodov, M.; Baker, J.; Phan, L.; Smigielski, E. M.; Sirotkin, K. dbSNP: the
412 NCBI database of genetic variation. *Nucleic Acids Res.* **2001**, *29*, 308–311.
- 413 7. Glusman, G.; Mauldin, D. E.; Hood, L. E.; Robinson, M. Ultrafast Comparison of Personal Genomes via
414 Precomputed Genome Fingerprints. *Front. Genet.* **2017**, *8*, 136.
- 415 8. Genotype fingerprints' homepage Available online:
416 http://db.systemsbio.org/gestalt/genotype_fingerprints (accessed on Aug 17, 2018).
- 417 9. Glusman, G.; Cariaso, M.; Jimenez, R.; Swan, D.; Greshake, B.; Bhak, J.; Logan, D. W.; Corpas, M. 23andMe
418 SNP chip genotype data 2012.
- 419 10. Indyk, P.; Motwani, R. Approximate nearest neighbors. In *Proceedings of the thirtieth annual ACM*
420 *symposium on Theory of computing - STOC '98*; ACM Press: New York, New York, USA, 1998; pp. 604–613.
- 421 11. Manichaikul, A.; Mychaleckyj, J. C.; Rich, S. S.; Daly, K.; Sale, M.; Chen, W.-M. Robust relationship
422 inference in genome-wide association studies. *Bioinformatics* **2010**, *26*, 2867–2873.
- 423 12. Gazal, S.; Sahbatou, M.; Babron, M.-C.; Génin, E.; Leutenegger, A.-L. High level of inbreeding in final
424 phase of 1000 Genomes Project. *Sci. Rep.* **2015**, *5*, 17453.
- 425 13. Epstein, M. P.; Duren, W. L.; Boehnke, M. Improved inference of relationship for pairs of individuals. *Am.*
426 *J. Hum. Genet.* **2000**, *67*, 1219–1231.
- 427 14. Tan, A.; Abecasis, G. R.; Kang, H. M. Unified representation of genetic variants. *Bioinformatics* **2015**, *31*,
428 2202–2204.
- 429 15. Erlich, Y.; Williams, J. B.; Glazer, D.; Yocum, K.; Farahany, N.; Olson, M.; Narayanan, A.; Stein, L. D.;
430 Witkowski, J. A.; Kain, R. C. Redefining genomic privacy: trust and empowerment. *PLoS Biol.* **2014**, *12*,
431 e1001983.
- 432 16. Raisaro, J. L.; Tramèr, F.; Ji, Z.; Bu, D.; Zhao, Y.; Carey, K.; Lloyd, D.; Sofia, H.; Baker, D.; Flicek, P.;
433 Shringarpure, S.; Bustamante, C.; Wang, S.; Jiang, X.; Ohno-Machado, L.; Tang, H.; Wang, X.; Hubaux, J.-P.
434 Addressing Beacon re-identification attacks: quantification and mitigation of privacy risks. *J. Am. Med.*
435 *Inform. Assoc.* **2017**, *24*, 799–805.
- 436 17. Glusman, G.; Caballero, J.; Mauldin, D. E.; Hood, L.; Roach, J. C. Kaviar: an accessible system for testing
437 SNV novelty. *Bioinformatics* **2011**, *27*, 3216–3217.

438