

Manuscript Number: FSIGEN-D-15-00279R1

Title: The Global AIMS Nano set: A 31-plex SNaPshot assay of ancestry-informative SNPs

Article Type: Research Paper

Keywords: SNPs; AIMS; Biogeographical ancestry; SNaPshot; Population-specific Divergence

Corresponding Author: Dr. Chris Phillips, PhD

Corresponding Author's Institution: Forensic Genetics Unit

First Author: Maria de la Puente

Order of Authors: Maria de la Puente; Carla Santos; Manuel Fondevila; Laura Manzo; Ángel Carracedo; María Victoria Lareu; Chris Phillips, PhD

Abstract: A 31-plex SNaPshot assay, named 'Global AIMS Nano', has been developed by reassembling the most differentiated markers of the EUROFORGEN Global AIM-SNP set. The SNPs include three tri-allelic loci and were selected with the goal of maintaining a balanced differentiation of: Africans, Europeans, East Asians, Oceanians and Native Americans. The Global AIMS Nano SNP set provides higher divergence between each of the five continental population groups than previous small-scale AIM sets developed for forensic ancestry analysis with SNaPshot. Both of these characteristics minimise potential bias when estimating co-ancestry proportions in individuals with admixed ancestry; more likely to be observed when using markers disproportionately informative for only certain population group comparisons. The optimised multiplex is designed to be easily implemented using standard capillary electrophoresis regimes and has been used to successfully genotype challenging forensic samples from highly degraded material with low level DNA. The ancestry predictive performance of the Global AIMS Nano set has been evaluated by the analysis of samples previously characterised with larger AIM sets.

Suggested Reviewers:

The Global AIMs Nano set: a 31-plex SNaPshot assay of ancestry-informative SNPs

M. de la Puente¹, C. Santos¹, M. Fondevila¹, L. Manzo¹; The EUROFORGEN-NoE Consortium; Á. Carracedo^{1,2}, M.V. Lareu¹, C. Phillips^{1*}.

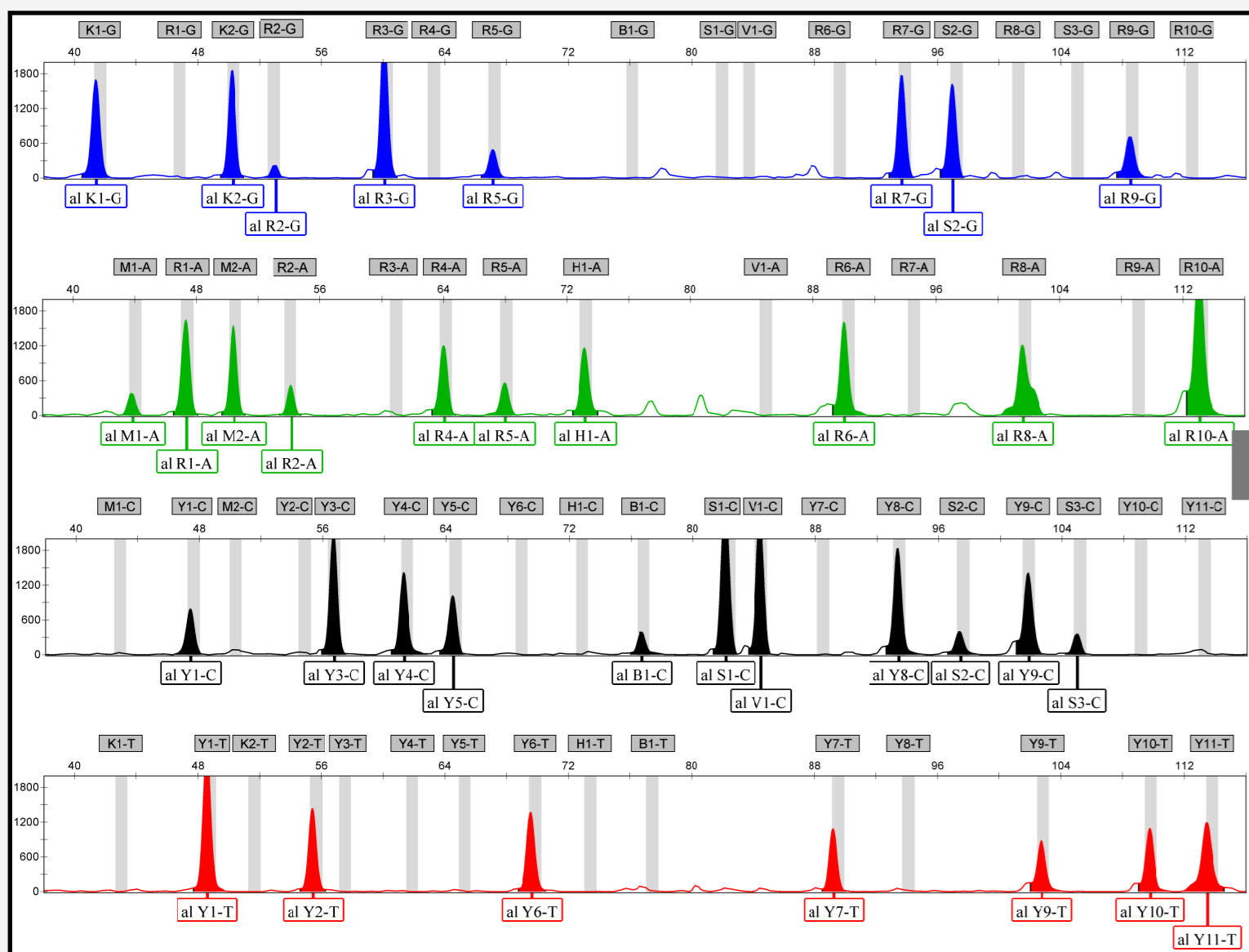
¹ Forensic Genetics Unit, Institute of Forensic Sciences, University of Santiago de Compostela, Spain

² Center of Excellence in Genomic Medicine Research, King Abdulaziz University, Jeddah, Saudi Arabia

* Corresponding author.

E-mail address: c.phillips@mac.com (C. Phillips).

31 SNPs



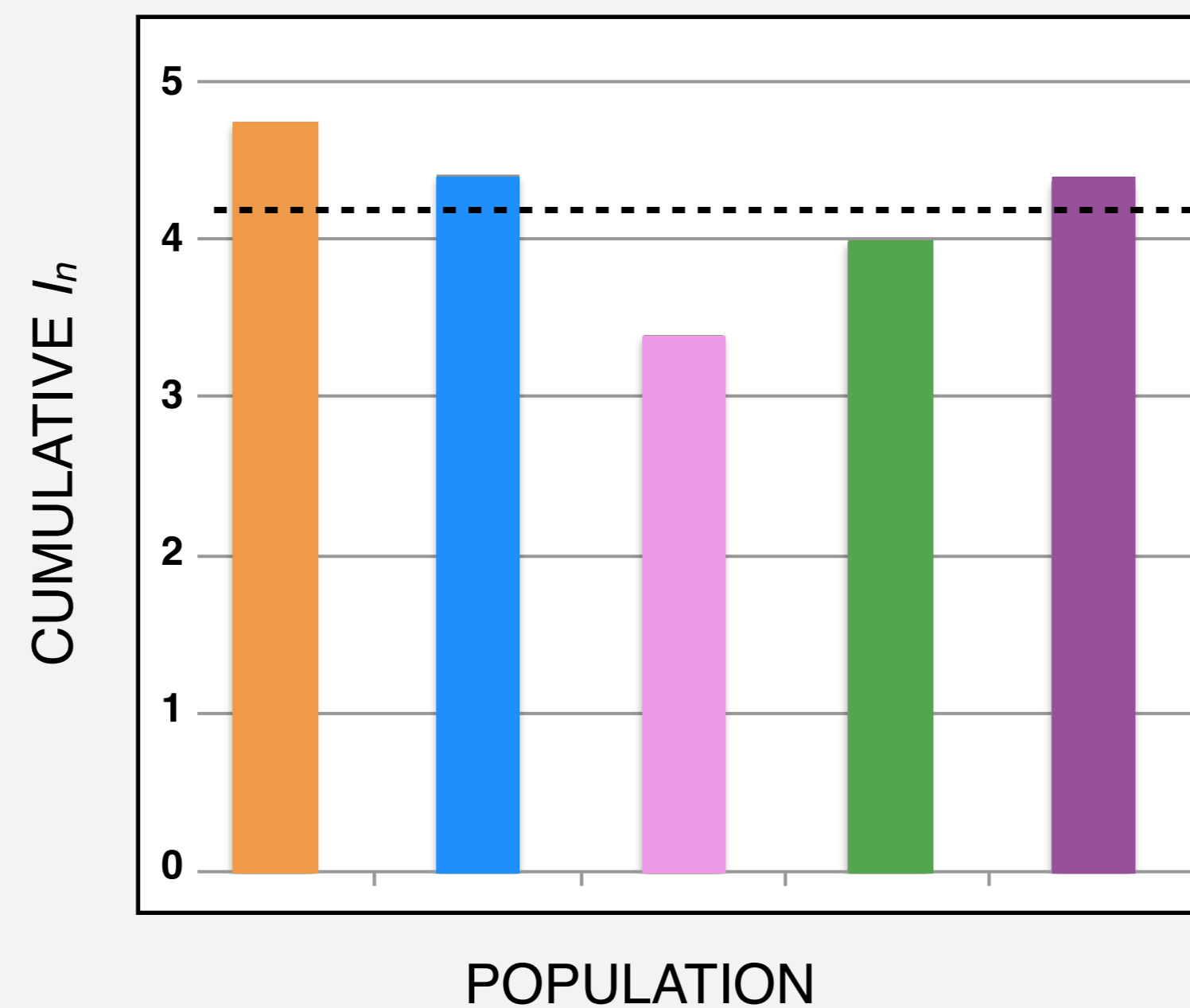
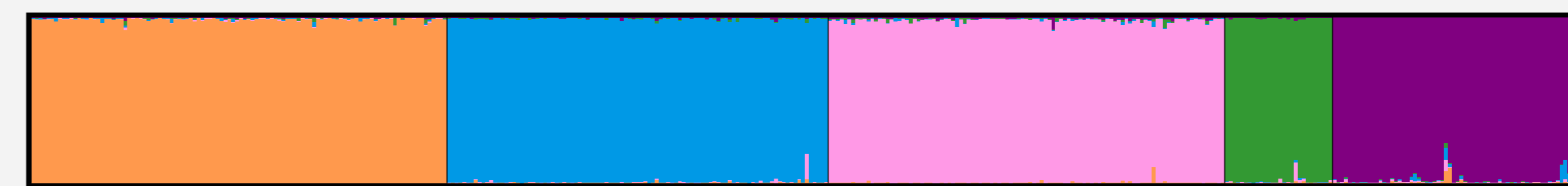
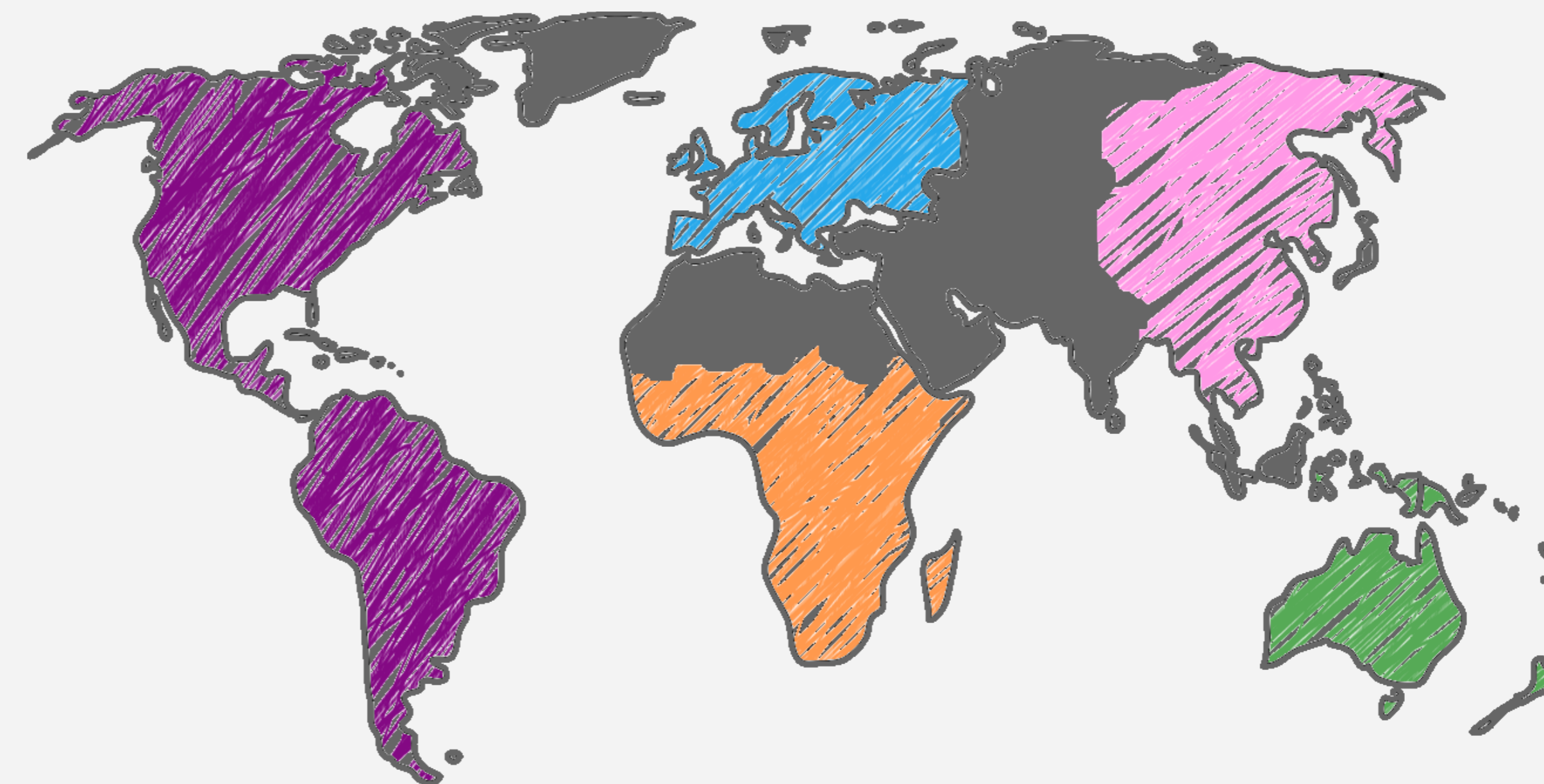
Single CE assay

Balanced population-specific Divergence

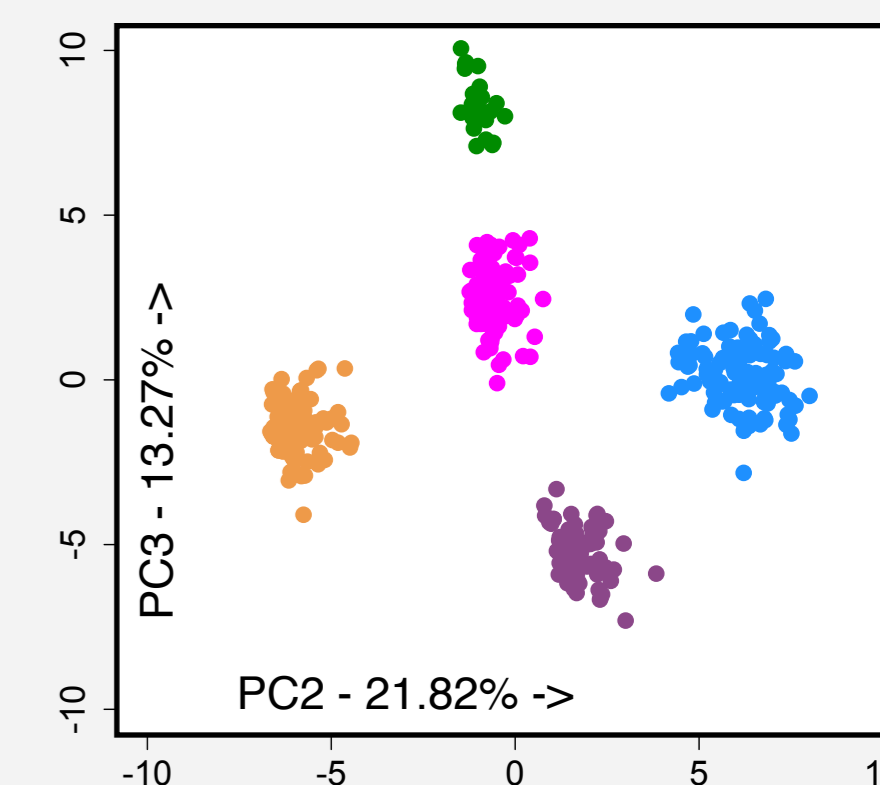
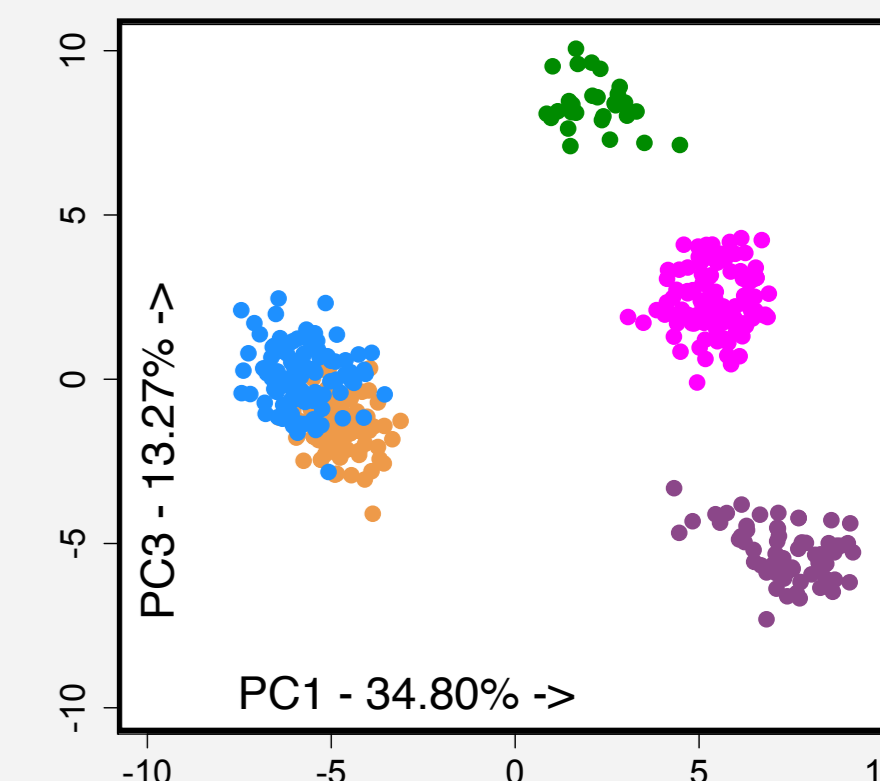
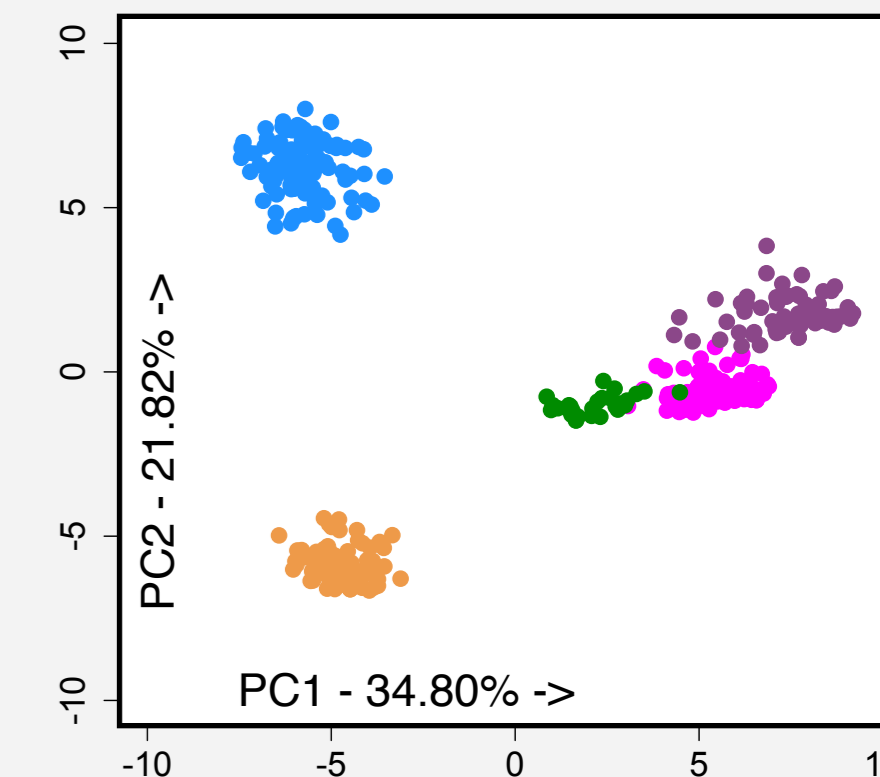
G-AIMs Nano

- ✓ Sensitivity 0.064 ng of DNA
- ✓ Degraded samples

5 GROUPS



Average = 4.18



The Global AIMs Nano set: a 31-plex SNaPshot assay of ancestry-informative SNPs

- A 31-plex AIM-SNP set was assembled from the most differentiated markers in the Global forensic ancestry panel.
- The SNPs extend population differentiations to Native Americans and Oceanians.
- SNPs preserve balanced population-specific Divergence, reducing estimation bias of individual co-ancestry.
- Assessments of the assay's ability to type degraded or low level DNA indicate good sensitivity for routine forensic use.
- The SNaPshot assay gives a compact but powerful multiplex to infer ancestry for labs yet to adopt massively parallel sequencing regimes.

28 1. Introduction

29

30 Although STR profiling has been successfully applied to the majority of
31 forensic DNA analyses for many years, there are still situations when STR
32 typing is unable to inform criminal investigations, for example, with no matching
33 profile found in DNA database searches or when no suspect is apprehended.
34 For this reason, there is interest in developing DNA tests that can provide
35 investigative leads, focused on panels of single nucleotide polymorphisms
36 (SNPs) to predict external visible characteristics (ECVs), including common
37 variation in pigmentation [1], or to infer an individual's biogeographical ancestry
38 [2].

39

40 With the recent availability of bench-top systems for massively parallel
41 sequencing (MPS) that are applicable to forensic DNA analysis, it is now
42 possible to assemble multiplexes of 400-500 markers [3]. Such enlarged
43 forensic multiplexes can include a portion of carefully chosen ancestry
44 informative markers (AIMs), e.g. the Illumina ForenSeq panel [4], or can be
45 exclusively composed of AIMs [5-7]. Both approaches raise the level of
46 geographic resolution that can be obtained from tests that keep the necessary
47 forensic sensitivity. However, the forensic community will take time to adopt,
48 optimise and validate MPS technology as a routine analysis system. Therefore,
49 it is important to continue to develop small-scale AIM sets suited to short-
50 amplicon marker genotyping with validated, universally applicable capillary
51 electrophoresis (CE) analysis regimes [8-10]. However, one drawback with use
52 of small-scale AIM sets is the potential for over-estimation of co-ancestry
53 proportions in individuals with admixed ancestry, stemming from the analysis of
54 genotypes strongly differentiated for some populations but not others. The
55 phenomenon of biased estimation of co-ancestry components was detected in a
56 study of Bolivian populations [11] using 46 ancestry-informative Indels [8]
57 compared with a much larger panel of 446 AIM-SNPs [12]. The Indel set
58 consistently over-estimated European co-ancestry and under-estimated Native
59 American co-ancestry using STRUCTURE-based analyses, indicating that the
60 higher European differentiation of the Indel genotypes inflated the estimates of
61 European co-ancestry proportions. Bearing in mind this effect, construction of a

62 dedicated AIM-SNP set for MPS by the EUROFORGEN Consortium [5] sought
63 to carefully balance the cumulative population-specific Divergence values for
64 the five continental population groups of Africa, Europe, East Asia, Native
65 America and Oceania.

66

67 The emphasis on keeping balanced population group Divergence values
68 provided the main focus for the new ancestry informative SNP panel reported
69 here. We took the most differentiated AIM-SNPs from the EUROFORGEN
70 Global AIMS panel [5] and assembled a compact 31-plex assay genotyped with
71 SNaPshot[®] single base extension technology. The SNP set, named 'Global
72 AIMS Nano' (herein Nano) was designed to be applicable to forensic analyses
73 where several different admixture combinations may be commonly
74 encountered, e.g. in Australia; where comparisons of European, Oceanian and
75 East Asian co-ancestry components will be routinely necessary. As well as
76 preserving a comparable level of differentiation amongst the five population
77 groups, the Nano assay aimed to provide a single CE-based test that is
78 sufficiently informative for all five groups.

79

80 **2. Materials and methods**

81

82 *2.1. Reference population SNP genotype data and DNA samples*

83

84 SNP variation data from representative populations without high levels of
85 admixture was obtained from 1000 Genomes Phase III [13] and from the
86 Stanford University HGDP-CEPH SNP analysis [14] using the SPSmart
87 frequency browser [15]. SNP genotype data was compiled from 108 YRI
88 Africans (AFR: Yoruba in Ibadan, Nigeria); 99 CEU Europeans (EUR: Utah
89 Residents with North and Western European ancestry); 103 CHB East Asians
90 (EAS: Han Chinese in Beijing, China); 28 HGDP-CEPH Oceanians (OCE: 17
91 Papuan from New Guinea and 11 Melanesian from Bougainville); and 64
92 HGDP-CEPH Native Americans (AMR: 14 Karitiana, 8 Surui from Brazil; 21
93 Maya, 14 Pima from Mexico; and 7 Piapoco from Colombia). Phase III 1000
94 Genomes populations were also analysed, comprising: as a test set, 99 AFR
95 LWK (Luhya in Webuye, Kenya); 113 AFR GWD (Gambian in Western

96 Divisions in the Gambia); 85 AFR MSL (Mende in Sierra Leone); 99 AFR ESN
97 (Esan in Nigeria); 107 EUR TSI (Toscani in Italia); 99 EUR FIN (Finnish in
98 Finland); 91 EUR GBR (British in England and Scotland); 107 EUR IBS (Iberian
99 Population in Spain); 104 EAS JPT (Japanese in Tokyo, Japan); 105 EAS CHS
100 (Southern Han Chinese); 99 EAS KHV (Kinh in Ho Chi Minh City, Vietnam); 93
101 EAS CDX (Chinese Dai in Xishuangbanna, China); plus admixed populations
102 61 ASW (Americans of African Ancestry in SW USA); 96 ACB (African
103 Caribbeans in Barbados); 104 PUR (Puerto Ricans from Puerto Rico), 94 CLM
104 (Colombians from Medellin, Colombia); 64 MXL (individuals with Mexican
105 Ancestry from Los Angeles USA); 85 PEL (Peruvians from Lima, Peru).

106

107 To evaluate the forensic sensitivity of the Nano assay, challenging
108 casework samples plus control DNAs were analysed, comprising: (i) five DNA
109 samples each from separate population groups, previously used in an ancestry
110 analysis collaborative exercise [10]; (ii) highly degraded skeletal DNA extracts;
111 (iii) a doubling dilution series of 1 ng/ μ L; 0.5 ng/ μ L; 0.25 ng/ μ L; 0.125 ng/ μ L;
112 0.064 ng/ μ L; 0.032 ng/ μ L; and 0.016 ng/ μ L of the 9947A forensic kit DNA
113 standard.

114

115 *2.2. AIM-SNP selection and SNaPshot assay design*

116

117 Ancestry-informative SNPs were selected directly from the
118 EUROFORGEN Global AIM-SNP set according to the following criteria: (i)
119 differentiation of five population groups to comparable levels to produce
120 population-specific Divergence (PSD) values as balanced as possible (use of
121 capitalised Divergence distinguishes the metric from the phenomenon of
122 population divergence); (ii) inclusion of certain informative tri-allelic SNPs to
123 allow a level of mixed DNA detection; (iii) genomic separation of component
124 SNPs by a minimum inter-marker distance of 1 Mb to minimise the effects of
125 linkage on likelihood calculations that assume independence for the loci
126 analysed.

127

128 From the selected SNPs, a 31-plex SNaPshot[®] single base extension
129 assay was designed and optimised following established guidelines [16]. Locus

130 details and summary allele frequencies for component SNPs are summarised in
131 Table 1. PCR and single base extension (SBE) primers are detailed in
132 Supplementary Table S1.

133

134 PCR reactions comprised: 1 μ L Buffer II (100 mM Tris-HCl, pH 8.3, 500
135 mM KCl); 1.8 μ L 25 mM MgCl₂; 0.1 μ L AmpliTaq Gold[®] DNA Polymerase (at 5
136 U/ μ L); 0.4 μ L of GeneAmp[®] 10 mM dNTP Mix with dTTP (Applied Biosystems,
137 AB); 1 μ L of 3.2 mg/ml bovine serum albumin; 1.5 μ L PCR primer mix
138 (Supplementary Table S1); 1 ng of target DNA adjusted to total reaction volume
139 of 10 μ L. PCR cycling with GeneAmp[®] PCR System 9700 or 2700 (AB)
140 thermocyclers used conditions: 10 mins at 95°C, 32 cycles of 30 secs at 95°C,
141 40 secs at 62°C and 1 min at 72°C with a final extension of 20 mins at 72°C.
142 PCR primer clean up combined 2.5 μ L PCR product with 1 μ L of 1 in 3 diluted
143 Illustra[™] ExoStar[™] 1-Step (GE Healthcare) then incubation at 37°C for 45
144 mins and enzyme inactivation at 85°C for 15 mins. SBE reactions comprised:
145 1.25 μ L of 1 in 2 diluted SNaPshot[®] Multiplex Ready Reaction Mix; 0.75 μ L SBE
146 primers (Supplementary Table S1) and 1 μ L purified PCR product in a total
147 volume of 3 μ L. SBE cycling used conditions: 33 cycles of 10 secs at 96°C, 5
148 secs at 59°C and 30 secs at 60°C. SBE primer clean up combined the full SBE
149 volume with 1 μ L of 1 in 2 diluted Illustra[™] Shrimp Alkaline Phosphatase (GE
150 Healthcare) then incubating at 37°C for 80 mins and enzyme inactivation at
151 85°C for 15 mins. Purified SBE products were then prepared for CE detection
152 by adding 1 μ L of product to 9.5 μ L of Hi-Di[™] Formamide (AB) and 0.25 μ L of
153 GeneScan[™]-120 LIZ[®] Size Standard (AB). Electrophoresis was performed in
154 an ABI Prism 3130xl Genetic Analyser, with 36 cm capillary arrays and POP-
155 4[™] polymer using standard conditions. Electropherograms were visualised
156 using AB GeneMapper[®] ID Software v. 3.2.1.

157

158 *2.3. Analysis of population variation in the selected SNPs*

159

160 Population-specific Divergence and simple pairwise Divergence values
161 were calculated using the Snipper cross-validation option
162 (http://mathgene.usc.es/snipper/analysispopfile2_new.html) by marking SNP
163 genotype profiles as AFR and non-AFR, etc., or by comparing each pair of

164 population groups in turn. Output from Snipper lists Shannon's Divergence
165 values for each SNP from the comparisons made by cross-validation [17].
166 These values were converted to the more widely-used Rosenberg's
167 informativeness-for-assignment metric: I_n [18], by multiplication with 0.693 (i.e.
168 converting the natural log to $\ln(2)$). The Snipper portal was also used to cross-
169 validate the reference population data or calculate classification likelihood ratios
170 (LRs) by uploading an Excel file of reference data (provided ready to use as
171 Supplementary File S1) or by choosing analysis options available in the portal.

172

173 Population analyses with STRUCTURE v. 2.3.4 [19] were performed
174 following recommendations outlined previously [20]. Parameters comprised: five
175 iterations (for $K=1$ to $K=9$) of 100,000 burnin steps and 100,000 MCMC steps,
176 correlated allele frequencies under the Admixture model (no POPFLAG for just
177 reference populations and POPFLAG for analyses of reference populations plus
178 test or admixed populations). The optimum K value was estimated by
179 computing results with STRUCTURE HARVESTER [21] and following previous
180 guidelines [22]. Ancestry membership plots were constructed with CLUMPAK v.
181 1.1 [23] or a combination of CLUMPP v. 1.1.2 [24] and distruct v. 1.1 [25]. PCA
182 analysis was performed using R software v. 3.1.2 [26] and executing a
183 homemade script. F_{ST} calculations and graphics were computed using Arlequin
184 v. 3.5 [27].

185

186 To assess the ancestry inference performance of the Nano SNP set,
187 comparisons were made with two previously developed biallelic AIM sets
188 comprising 46 Indels [8] and 34 SNPs [9], by applying STRUCTURE and PCA
189 analyses of the same 1000 Genomes African, European and East Asian
190 genotypes plus HGDP-CEPH Native American and Oceanian genotypes
191 compiled from each set (data used from 44 of 46 Indels currently listed by 1000
192 Genomes Phase III).

193

194 **3. Results**

195

196 *3.1. Characteristics and PSD balance of the Nano SNP set*

197

198 The 31 SNPs selected showed highly contrasting allele frequency
199 distributions. In each of the 28 biallelic SNPs one allele was close to fixation
200 (allele frequencies between 0.9 and 1) in at least one population group, as
201 shown by the raster plot of Fig. 1. Population group summary allele frequency
202 pie charts are also shown in Supplementary Fig. S1. All SNPs were distributed
203 in the genome with sufficient distance between syntenic marker sets to be free
204 from the effects of close physical linkage (Supplementary Fig. S2).

205

206 With the reduction in scale from an MPS multiplex of 128 SNPs to the
207 SNaPshot 31-plex, it is important to ensure the cumulative PSD values remain
208 at comparable levels for each group. Reference population comparisons
209 produced the cumulative PSD values listed in Table 1 and summarised in Fig.
210 2. These indicate the reduced East Asian PSD of 3.39 was noticeably lower
211 than the average of 4.18 and the African PSD of 4.74 was the highest but
212 comparable to three population groups. The reduced differentiation of East
213 Asians from Native American and Oceanian populations is illustrated by their
214 similar allele frequency distributions for many SNPs, exemplified by rs4657449
215 and rs9809818. The close relationship of East Asian and Native American
216 population variability is highlighted by Fig. 1, with little divergence between the
217 two groups evident in rs3827760, rs6437783 and rs12594144. This suggests
218 one or two extra East Asian-informative SNPs could address this slight PSD
219 imbalance in future adjustments of the Nano SNP set.

220

221 The F_{ST} and pairwise genotype difference data from Arlequin analyses
222 are summarised in Supplementary Fig. S3. Results indicate a high average
223 number of pairwise genotype differences between the population groups and
224 consequently F_{ST} values are low within-groups and high between-groups.
225 Admixed populations from 1000 Genomes give high within-population average
226 number of pairwise genotype differences and F_{ST} values, as would be expected
227 from the complex patterns of variation that are characteristic of population
228 admixture.

229

230 *3.2. Ancestry inference capabilities of the Nano SNP set*

231

232 Cross-validation of reference populations gave 100% ancestry
233 assignment success (Supplementary Table S2). Additionally, assignment
234 success remained at 100% for all populations when excluding the 14 most
235 informative SNPs (those with highest overall Divergence values), indicating this
236 marker set maintains a high level of informativeness even when many
237 components fail to be reliably genotyped, e.g. analysing low-level DNA.

238

239 STRUCTURE analysis of reference population data (no POPFLAG)
240 produced a pattern of five distinct clusters matching the known origin of the
241 individuals. The 34 SNP and 46 Indel forensic ancestry assays with which Nano
242 was compared also distinguish the five genetic clusters, but in contrast, both
243 give estimated optimum numbers of clusters below five (Supplementary Fig.
244 S4). The PCA plots shown in Fig. 3A reveal improved separation and
245 clusteredness of populations from the Nano set compared to the other two
246 assays. This is especially evident in the PC2 vs. PC3 plots for Nano genotypes,
247 where no set of points overlap and the five population groups have almost
248 equidistant cluster positions.

249

250 Analysis of other genotype data from 1000 Genomes which were marked
251 as 'study populations' (POPFLAG=0) gave patterns consistent with analyses
252 made during the development of the 128 Global ancestry set [5] or subsequent
253 analysis of admixed 1000 Genomes ACB, ASW, PEL, MXL, PUR, CLM
254 samples using the same SNPs (Fig. 7 in [2]). The STRUCTURE cluster plots
255 (Fig. 3B) in particular match well with results of both previous analyses using
256 many more SNPs, indicating PEL have the highest Native American co-
257 ancestry proportions and PUR show predominantly European co-ancestry. The
258 PCA plots arranged individually along with STRUCTURE cluster plots for the six
259 admixed 1000 Genomes populations in Fig. 3C also give very similar cluster
260 distributions that are positioned between, or close to, the expected population
261 admixture contributor clusters.

262

263 As an additional simple gauge of the ancestry informativeness of the 31
264 Nano SNPs, this assay was used to analyse the control DNAs used in a
265 collaborative EDNAP exercise that assessed the 34 SNPs and 46 Indels

266 described above. The five control DNAs each have confirmed ancestry from
267 one of the five continental regions and they were analysed using Bayes
268 analyses and PCA in Snipper. The five control DNAs are positioned in the
269 middle of their respective population group clusters describing the correct
270 ancestry in each case. The Bayes analysis likelihoods obtained from Snipper
271 are listed in Table 2. Moreover, as shown in Supplementary Fig. S5, high
272 likelihood ratio values can be obtained from partial profiles even when the
273 fourteen most informative markers are missing.

274

275 *3.3. Forensic performance of the Nano SNPs set*

276

277 A typical SNaPshot profile from analysis of 1 ng of the 9947A control
278 DNA with the Nano assay is shown in Fig. 4. The dilutions series of 9947A
279 gave full profiles with 0.5, 0.25, 0.125 and 0.064 ng of DNA. Locus and allele
280 drop out occurred with 0.032 and 0.016 ng of input DNA, however >80% profile
281 completeness was obtained for these analyses.

282

283 The sensitivity of the Nano assay was assessed with paternity test
284 samples, comprising biopsy, bones (cranium, femur and tibia) and teeth
285 identified as degraded or PCR-inhibited (~35% to 95% STR profile
286 completeness). Nano profile completeness ranged from ~20% to 70% and
287 produced ancestry assignment likelihoods above 60,000.

288

289 **4. Discussion**

290

291 When originally rebuilding a SNP-based forensic ancestry multiplex for
292 MPS analysis, we started to recognise subsets of the most informative markers
293 that would be well suited to development of smaller SNaPshot tests. The
294 present study describes the compilation of 31 SNPs, mainly representing the
295 most informative markers from the full set of 128 Global AIMs. These SNPs
296 maintain the capacity to differentiate five continentally-defined population
297 groups. Therefore, the Nano assay extends the three group comparisons
298 possible with existing SNaPshot forensic ancestry tests [9,28-30] to the two
299 additional population groups of Native Americans and Oceanians. Each of

300 these population groups contribute to admixture patterns commonly seen in
301 large parts of the regions they occupy. For this reason, we prioritised the
302 preservation of balanced PSD, although this was difficult to maintain when
303 accomplishing the 75% reduction in multiplex size. One outcome of this process
304 was a disproportionate lowering of the East Asian cumulative PSD compared to
305 the other groups that we aim to address by careful choice of 1-2 additional
306 AIMs.

307

308 The selection of SNPs informative for Native American and Oceanian
309 populations requires use of much smaller population sample sizes from the
310 HGDP-CEPH panel compared to those of 1000 Genomes. Therefore, SNP
311 ascertainment bias could reduce the power of the 31 SNPs to differentiate novel
312 populations not yet characterised from each of these regions. However, such
313 bias is unlikely to lead to the discovery of new SNPs as divergent as the 22
314 Native American informative and 28 Oceanian informative SNPs assembled in
315 the original 128-plex set. The inclusion in the 31-plex Nano set of the five most
316 informative SNPs for both Native Americans and Oceanians, comprising SNPs
317 with alleles near to fixation, ensures the Nano set is almost equally informative
318 for all five groups and the assignment likelihoods obtained for populations from
319 the two additional population groups exceed the values possible with previous
320 AIM-SNP sets developed for SNaPshot genotyping.

321

322 When the same control samples with known ancestries are tested with
323 established multiplexes of 34 SNPs, 46 Indels and the Nano 31-plex, higher
324 assignment likelihoods are obtained for the 31 SNPs than 80 markers combined
325 in three of five population groups, and all likelihoods exceed those obtained
326 from 34-plex SNP data by considerable margins (between 3 and 16 orders of
327 magnitude). Therefore, for the bulk of forensic samples that require an ancestry
328 analysis in laboratories without MPS systems in place, the Nano SNP set
329 represents the best option as a stand-alone CE test. The use of the Nano SNPs
330 alongside Indels, with their enhanced capacity to detect mixed DNA profiles [8],
331 will provide a particularly powerful approach to forensic ancestry analysis from
332 the use of conventional capillary electrophoresis techniques, already optimised
333 for routine DNA analyses in every forensic laboratory.

334

335 In conclusion, the Nano ancestry assay has brought together a highly
336 informative set of markers with well-balanced population-specific Divergence for
337 the five population groups it is designed to analyse. This characteristic
338 minimises the co-ancestry proportion estimation bias when analysing admixed
339 samples. Moreover, the single SNaPshot reaction shows high sensitivity
340 (complete profiles obtained down to 64 pg of input DNA), with the potential to
341 analyse degraded samples, making it an ideal forensic ancestry assay for the
342 full range of casework applications where DNA is analysed with capillary
343 electrophoresis.

344

345 **Acknowledgements**

346

347 This work was funded by the EUROFORGEN Node of Excellence (Grant
348 Agreement No. 285487). MdIP is supported by funding awarded by the
349 Consellería de Cultura, Educación e Ordenación Universitaria of the Xunta de
350 Galicia as part of the Plan Galego de Investigación, Innovación e Crecemento
351 2011-2015 (Plan I2C). CS is supported by a PhD grant (SFRH/BD/75627/2010)
352 awarded by the Portuguese Foundation for Science and Technology (FCT) and
353 co-financed by the European Social Fund (Human Potential Thematic
354 Operational Program). MVL was supported by funding from Xunta de Galicia,
355 Incite 09208163 PR.

356

357 **References**

- 358 [1] M. Kayser, Forensic DNA Phenotyping: Predicting human appearance from
359 crime scene material for investigative purposes, *Forensic Science International: Genetics* 18 (2015) 33-48.
360
- 361 [2] C. Phillips, Forensic genetic analysis of bio-geographical ancestry, *Forensic Science International: Genetics* 18 (2015) 49-65.
362
- 363 [3] B. Budowle, D.H. Warshauer, S.B. Seo, et al., Deep sequencing provides
364 comprehensive multiplex capabilities, *Forensic Science International: Genetics Supplement Series* 4 (2013) e334-e335.
365
- 366 [4] Illumina, Foreseq DNA signature prep kit (2014)
367 <http://www.illumina.com/products/foreseq-dna-signature-kit.ilmn>.
- 368 [5] C. Phillips, W. Parson, B. Lundsberg, et al., Building a forensic ancestry
369 panel from the ground up: The EUROFORGEN Global AIM-SNP set, *Forensic Sci Int Genet* 11 (2014) 13-25.
370
- 371 [6] K.K. Kidd, W.C. Speed, A.J. Pakstis, et al., Progress toward an efficient
372 panel of SNPs for ancestry inference, *Forensic Science International: Genetics* 10 (2014) 23-32.
373
- 374 [7] J. Kidd, F. Friedlaender, W. Speed, et al., Analyses of a set of 128 ancestry
375 informative single-nucleotide polymorphisms in a global set of 119 population
376 samples, *Investigative Genetics* 2 (2011) 1-13.
- 377 [8] R. Pereira, C. Phillips, N. Pinto, et al., Straightforward inference of ancestry
378 and admixture proportions through ancestry-informative insertion deletion
379 multiplexing, *PLoS One* 7 (2012) e29684.
- 380 [9] M. Fondevila, C. Phillips, C. Santos, et al., Revision of the SNPforID 34-plex
381 forensic ancestry test: Assay enhancements, standard reference sample
382 genotypes and extended population studies, *Forensic Sci Int Genet* 7 (2013)
383 63-74.
- 384 [10] C. Santos, M. Fondevila, D. Ballard, et al., Forensic ancestry analysis with
385 two capillary electrophoresis ancestry informative marker (AIM) panels: Results
386 of a collaborative EDNAP exercise, *Forensic Science International: Genetics* 19
387 (2015) 56-67.
- 388 [11] P. Taboada-Echalar, V. Álvarez-Iglesias, T. Heinz, et al., The genetic
389 legacy of the pre-colonial period in contemporary Bolivians, *PLoS ONE* 8 (2013)
390 e58980.
- 391 [12] J.M. Galanter, J.C. Fernandez-Lopez, C.R. Gignoux, et al., Development of
392 a panel of genome-wide ancestry informative markers to study admixture
393 throughout the Americas, *PLoS Genet* 8 (2012) e1002554.
- 394 [13] An integrated map of genetic variation from 1,092 human genomes, *Nature*
395 491 (2012) 56-65.
- 396 [14] J.Z. Li, D.M. Absher, H. Tang, et al., Worldwide human relationships
397 inferred from genome-wide patterns of variation, *Science* 319 (2008) 1100-
398 1104.
- 399 [15] J. Amigo, A. Salas, C. Phillips, et al., SPSmart: adapting population based
400 SNP genotype databases for fast and comprehensive web access, *BMC*
401 *Bioinformatics* 9 (2008) 428.
- 402 [16] J.J. Sanchez, P. Endicott, Developing multiplexed SNP assays with special
403 reference to degraded DNA templates, *Nat Protocols* 1 (2006) 1370-1378.
- 404 [17] C. Phillips, A. Salas, J.J. Sánchez, et al., Inferring ancestral origin using a
405 single multiplex assay of ancestry-informative marker SNPs, *Forensic Science International: Genetics* 1 (2007) 273-280.
406

407 [18] N.A. Rosenberg, L.M. Li, R. Ward, et al., Informativeness of Genetic
408 Markers for Inference of Ancestry, *American Journal of Human Genetics* 73
409 (2003) 1402-1422.

410 [19] J.K. Pritchard, M. Stephens, P. Donnelly, Inference of population structure
411 using multilocus genotype data, *Genetics* 155 (2000) 945-959.

412 [20] L. Porras-Hurtado, Y. Ruiz, C. Santos, et al., An overview of STRUCTURE:
413 applications, parameter settings, and supporting software, *Front Genet* 4 (2013)
414 98.

415 [21] D. Earl, B. vonHoldt, STRUCTURE HARVESTER: a website and program
416 for visualizing STRUCTURE output and implementing the Evanno method,
417 *Conservation Genetics Resources* 4 (2012) 359-361.

418 [22] G. Evanno, S. Regnaut, J. Goudet, Detecting the number of clusters of
419 individuals using the software STRUCTURE: a simulation study, *Mol Ecol* 14
420 (2005) 2611-2620.

421 [23] N.M. Kopelman, J. Mayzel, M. Jakobsson, et al., Clumpak: a program for
422 identifying clustering modes and packaging population structure inferences
423 across K, *Mol Ecol Resour* (2015) doi: 10.1111/1755-0998.12387. [Epub ahead
424 of print].

425 [24] M. Jakobsson, N.A. Rosenberg, CLUMPP: a cluster matching and
426 permutation program for dealing with label switching and multimodality in
427 analysis of population structure, *Bioinformatics* 23 (2007) 1801-1806.

428 [25] N.A. Rosenberg, distruct: a program for the graphical display of population
429 structure, *Molecular Ecology Notes* 4 (2004) 137-138.

430 [26] R.C. Team. R: A language and environment for statistical computing.
431 Vienna, Austria: R Foundation for Statistical Computing; 2014.

432 [27] L. Excoffier, H.E. Lischer, Arlequin suite ver 3.5: a new series of programs
433 to perform population genetics analyses under Linux and Windows, *Mol Ecol*
434 *Resour* 10 (2010) 564-567.

435 [28] Y.-L. Wei, L. Wei, L. Zhao, et al., A single-tube 27-plex SNP assay for
436 estimating individual ancestry and admixture from three continents, *International*
437 *Journal of Legal Medicine* (2015) 1-11.

438 [29] U. Rogalla, E. Rychlicka, M.V. Derenko, et al., Simple and cost-effective
439 14-loci SNP assay designed for differentiation of European, East Asian and
440 African samples, *Forensic Science International: Genetics* 14 (2015) 42-49.

441 [30] O. Lao, P.M. Vallone, M.D. Coble, et al., Evaluating Self-declared Ancestry
442 of U.S. Americans with Autosomal, Y-chromosomal and Mitochondrial DNA,
443 *Human Mutation* 31 (2010) E1875-E1893.

444

445

446 **Figure legends**

447

448 **Fig. 1.** Raster plot summarising the allele frequency distributions of the Nano
449 SNPs in five population groups. AFR: African; EUR: European; EAS: East
450 Asian; OCE: Oceanian; AMR: Native American.

451

452 **Fig. 2.** Bar charts indicating cumulative, population-specific and pairwise I_n
453 values for Nano SNPs.

454

455 **Fig. 3.** (A) Three AIM panels compared with PCA analyses (triallelic SNPs not
456 included in analyses, so Nano=28 SNPs, 34-plex=32 SNPs and AIM-indels=44
457 markers) for reference samples. PC: principal component; AFR: African; EUR:
458 European; EAS: East Asian; OCE: Oceanian; AMR: Native American. Control
459 samples with known ancestry are shown as black points. (B) STRUCTURE
460 analysis for admixed and non-admixed populations included in the study. (C)
461 Mixed populations are represented below in black in the corresponding PCA
462 analyses (PC1 vs PC2). AFR: African; EUR: European; EAS: East Asian; OCE:
463 Oceanian; AMR: Native American; ESN: Esan in Nigeria; MSL: Mende in Sierra
464 Leone; GWD: Gambian in Western Divisions in the Gambia; LWK: Luhya in
465 Webuye, Kenya; ACB: African Caribbeans in Barbados; ASW: Americans of
466 African Ancestry in SW USA; GBR: British in England and Scotland; TSI:
467 Toscani in Italia; FIN: Finnish in Finland; IBS: Iberian Population in Spain; JPT:
468 Japanese in Tokyo, Japan; CHS: Southern Han Chinese; CDX: Chinese Dai in
469 Xishuangbanna, China; KHV: Kinh in Ho Chi Minh City, Vietnam; PEL:
470 Peruvians from Lima, Peru); MXL: Mexican Ancestry from Los Angeles, USA;
471 CLM: Colombians from Medellin, Colombia and PUR: Puerto Ricans from
472 Puerto Rico.

473

474 **Fig. 4.** Electropherogram of 31 Nano SNP genotypes obtained from the control
475 DNA 9947A.

476 **Table 1.** Description, reference allele frequencies and population-specific/pairwise Divergence values (I_n) of the 31 Nano SNPs.
 477 Chr: chromosome; RA: reference allele. All positions from genome build 37.1 (GRCh37). SNPs are ranked according to their Pop
 478 vs. Other Pop I_n (highlighted in grey) inside each population informative (Informat.) group.

SNP details					Reference allele frequency					Population-specific Divergence					Pairwise Divergence									
Informat.	SNP ID	Chr	Position	RA	AFR	EUR	EAS	OCE	AMR	AFR	EUR	EAS	OCE	AMR	AFR - EUR	AFR - EAS	AFR - OCE	AFR - AMR	EUR - EAS	EUR - OCE	EUR - AMR	EAS - OCE	EAS - AMR	OCE - AMR
AFR	rs2814778	1	159174683	A	0.005	1.000	1.000	1.000	0.992	0.672	0.131	0.134	0.083	0.108	0.663	0.663	0.634	0.656	0.000	0.672	0.000	0.000	0.000	0.001
	rs1871534	8	145639681	C	0.981	0.000	0.000	0.000	0.000	0.641	0.128	0.131	0.081	0.107	0.631	0.632	0.603	0.624	0.000	0.641	0.000	0.000	0.000	0.000
	rs2789823	9	136769888	G	0.935	0.000	0.000	0.000	0.000	0.565	0.121	0.123	0.076	0.100	0.556	0.557	0.527	0.548	0.000	0.565	0.000	0.000	0.000	0.000
EUR	rs1426654	15	48426484	A	0.014	1.000	0.029	0.000	0.039	0.117	0.622	0.093	0.081	0.069	0.641	0.001	0.000	0.003	0.611	0.117	0.595	0.001	0.000	0.002
	rs16891982	5	33951693	C	1.000	0.020	0.985	1.000	0.984	0.126	0.620	0.104	0.073	0.087	0.629	0.001	0.000	0.002	0.606	0.126	0.603	0.000	0.000	0.000
	rs12142199	1	1249187	G	0.977	0.177	0.971	1.000	1.000	0.076	0.402	0.068	0.061	0.083	0.393	0.000	0.000	0.002	0.383	0.076	0.423	0.001	0.003	0.000
	rs8072587	17	19211073	C	0.986	0.182	1.000	1.000	0.817	0.099	0.358	0.113	0.069	0.003	0.405	0.001	0.000	0.047	0.425	0.099	0.218	0.000	0.058	0.043
	rs9522149	13	111827167	T	0.972	0.237	0.995	1.000	0.977	0.063	0.353	0.092	0.055	0.056	0.334	0.004	0.001	0.000	0.377	0.063	0.341	0.002	0.003	0.000
	rs4749305	10	28391596	A	0.389	0.909	0.078	0.036	0.008	0.001	0.309	0.095	0.100	0.155	0.162	0.072	0.105	0.141	0.404	0.001	0.514	0.004	0.017	0.005
EAS	rs17822931	16	48258198	C	1.000	0.869	0.029	0.875	0.650	0.190	0.053	0.432	0.039	0.000	0.039	0.613	0.037	0.129	0.428	0.190	0.004	0.434	0.251	0.036
	rs1229984	4	100239319	A	0.000	0.015	0.709	0.071	0.000	0.089	0.070	0.320	0.018	0.071	0.002	0.336	0.018	0.000	0.313	0.089	0.001	0.238	0.328	0.015
	rs3827760	2	109513601	T	1.000	1.000	0.063	0.946	0.109	0.219	0.209	0.319	0.098	0.203	0.000	0.559	0.012	0.500	0.558	0.219	0.499	0.471	0.003	0.415
	rs6437783	3	108172817	C	0.259	0.146	0.995	0.589	0.891	0.078	0.153	0.268	0.001	0.104	0.010	0.359	0.057	0.223	0.459	0.078	0.311	0.157	0.031	0.062
	rs12594144	15	64161351	C	1.000	0.889	0.121	0.607	0.177	0.237	0.097	0.222	0.000	0.130	0.032	0.487	0.149	0.430	0.334	0.237	0.283	0.136	0.003	0.101
	rs4657449	1	165465281	G	0.912	0.909	0.102	0.000	0.117	0.176	0.164	0.177	0.210	0.130	0.000	0.380	0.497	0.364	0.376	0.176	0.360	0.017	0.000	0.022
OCE	rs9908046	17	53563782	C	0.958	0.929	0.883	0.018	0.992	0.020	0.008	0.000	0.528	0.042	0.002	0.010	0.562	0.006	0.003	0.020	0.015	0.463	0.030	0.626
	rs3751050	11	9091244	A	0.972	0.924	0.966	0.089	0.961	0.020	0.002	0.016	0.451	0.011	0.006	0.000	0.477	0.000	0.004	0.020	0.003	0.467	0.000	0.459
	rs2139931	1	84590527	A	0.898	0.753	0.879	0.018	0.898	0.017	0.002	0.011	0.433	0.014	0.019	0.000	0.480	0.000	0.013	0.017	0.019	0.458	0.000	0.480
	rs715605	22	30640308	T	0.866	0.914	0.985	0.089	1.000	0.000	0.003	0.039	0.422	0.041	0.003	0.030	0.345	0.036	0.015	0.000	0.020	0.502	0.001	0.516
	rs6054465	20	6673018	T	0.972	0.859	0.743	0.036	0.859	0.062	0.005	0.005	0.408	0.004	0.022	0.061	0.553	0.022	0.011	0.062	0.000	0.306	0.011	0.408
	rs9809818	3	71480566	C	0.019	0.116	0.869	0.982	0.820	0.245	0.120	0.172	0.227	0.102	0.021	0.446	0.603	0.399	0.319	0.245	0.276	0.026	0.002	0.042
AMR	rs12498138	3	121459589	G	1.000	0.949	0.922	0.911	0.094	0.085	0.034	0.020	0.011	0.443	0.011	0.020	0.024	0.519	0.002	0.085	0.437	0.000	0.401	0.387
	rs10483251	14	21671277	G	0.921	0.798	0.898	0.712	0.024	0.055	0.006	0.040	0.000	0.429	0.016	0.001	0.038	0.497	0.010	0.055	0.370	0.028	0.469	0.301
	rs2080161	7	13331150	T	0.981	0.758	0.689	0.820	0.000	0.144	0.007	0.001	0.036	0.424	0.064	0.091	0.044	0.624	0.003	0.144	0.366	0.033	0.314	0.462
	rs8137373	22	41729216	G	0.833	0.707	0.927	0.982	0.023	0.020	0.000	0.068	0.093	0.406	0.011	0.011	0.037	0.402	0.043	0.020	0.298	0.009	0.506	0.593
	rs1557553	22	44760984	C	0.949	0.904	0.714	0.786	0.094	0.076	0.042	0.000	0.002	0.325	0.004	0.053	0.031	0.436	0.030	0.076	0.380	0.003	0.219	0.270
	rs12402499	1	101528954	G	1.000	0.919	1.000	1.000	0.258	0.061	0.007	0.059	0.032	0.324	0.021	0.000	0.000	0.361	0.021	0.061	0.252	0.000	0.361	0.335
	rs4792928	17	42105174	T	1.000	0.960	0.345	0.804	0.195	0.171	0.111	0.108	0.011	0.178	0.008	0.297	0.064	0.413	0.239	0.171	0.350	0.113	0.014	0.199
Triallelic	rs2069945	20	33761837	CG	0.153 / 0.796	0.480 / 0.420	0.680 / 0.277	0.036 / 0.536	0.766 / 0.234	0.114	0.002	0.045	0.201	0.087	0.078	0.156	0.118	0.210	0.022	0.114	0.065	0.289	0.017	0.384
	rs4540055	4	38803255	AC	0.069 / 0.514	0.793 / 0.010	0.301 / 0.068	0.536 / 0.250	0.605 / 0.000	0.217	0.148	0.055	0.024	0.084	0.355	0.147	0.142	0.300	0.130	0.217	0.027	0.098	0.063	0.101
	rs5030240	11	32424389	CA	0.278 / 0.389	0.712 / 0.055	0.228 / 0.039	0.093 / 0.278	0.258 / 0.023	0.083	0.117	0.062	0.062	0.054	0.125	0.123	0.052	0.135	0.135	0.083	0.124	0.067	0.001	0.085
CUMULATIVE VALUES										4.739	4.404	3.392	3.986	4.374	5.263	6.111	6.210	8.029	6.274	4.739	7.184	4.323	3.106	6.350

480 **Table 2.** Ancestry assignment likelihood ratios from a five-population group
 481 comparison in Snipper for control samples A-E of known ancestry. All
 482 assignments were correct except for Sample A misclassified as AMR with 34-
 483 plex data. Data for 34-plex, AIM-indel markers and the 80 markers combined
 484 was obtained directly using the Snipper portal options and for Nano AIMs, by
 485 applying the data in Supplementary File S1 as a custom training set.
 486

Sample	Known Ancestry	Assignment likelihood ratios (from two highest likelihoods)			
		34-plex	AIM indels	80 Markers	G-AIMs Nano
A	East Asian (EAS)	6.8E+00 more likely to be AMR*	5.5E+06 more likely to be EAS	9.7E+06 EAS	2.9E+16 EAS
B	European (EUR)	4.4E+16 EUR	1.7E+11 EUR	1.0E+28 EUR	1.3E+29 EUR
C	Oceanian (OCE)	1.0E+07 OCE	4.0E+07 OCE	1.5E+14 OCE	1.5E+16 OCE
D	Native American (AMR)	1.0E+05 AMR	1.2E+09 AMR	1.1E+14 AMR	4.2E+08 AMR
E	African (AFR)	6.1E+19 AFR	2.8E+21 AFR	3.6E+40 AFR	2.0E+29 AFR

487 *Likelihood ratio of AMR / EAS likelihoods

488

489 **Supplementary material legends**

490

491 **Supplementary Fig. S1.** Pie charts indicating the frequency of each Nano SNP
492 allele in five population groups (AFR: Africans; EUR: Europeans; EAS: East
493 Asians; OCE: Oceanians; AMR: Native Americans).

494

495 **Supplementary Fig. S2.** Autosomal chromosome ideograms representing the
496 genomic positions of the Nano SNPs.

497

498 **Supplementary Fig. S3.** Pairwise F_{ST} and number of pairwise genotype
499 differences between and within populations. AFR: African; EUR: European;
500 EAS: East Asian; OCE: Oceanian; AMR: American; ESN: Esan in Nigeria; MSL:
501 Mende in Sierra Leone; GWD: Gambian in Western Divisions in the Gambia;
502 LWK: Luhya in Webuye, Kenya; ACB: African Caribbeans in Barbados; ASW:
503 Americans of African Ancestry in SW USA; GBR: British in England and
504 Scotland; TSI: Toscani in Italia; FIN: Finnish in Finland; IBS: Iberian Population
505 in Spain; JPT: Japanese in Tokyo, Japan; CHS: Southern Han Chinese; CDX:
506 Chinese Dai in Xishuangbanna, China; KHV: Kinh in Ho Chi Minh City, Vietnam;
507 PEL: Peruvians from Lima, Peru); MXL: Mexican Ancestry from Los Angeles,
508 USA; CLM: Colombians from Medellin, Colombia and PUR: Puerto Ricans from
509 Puerto Rico. Reference populations are highlighted in red, and admixed
510 populations in blue.

511

512 **Supplementary Fig. S4.** Three panel comparison of STRUCTURE analysis
513 (Nano - 31 SNPs, 34-plex - 34 SNPs and AIM-indels – 44 of 46 markers) for the
514 reference population samples. Graphics on the right indicate the mean of
515 estimated Ln probability-of-data and Delta K according to Evanno's method
516 across different numbers of clusters (K). Red arrows indicate the most probable
517 number of clusters taking into account both graphs and black arrows indicate
518 other probable numbers of clusters represented on the left. K: number of
519 clusters; AFR: African; EUR: European; EAS: East Asian; OCE: Oceanian;
520 AMR: Native American.

521

522 **Supplementary Fig. S5.** Likelihoods ratios of ancestry estimation for control
523 samples of known ancestry (A-E) and 9947A, removing sequentially the most
524 informative markers. The grey box represents a threshold likelihood value of
525 1,000 and the red line indicates the most informative fourteen SNPs.

526

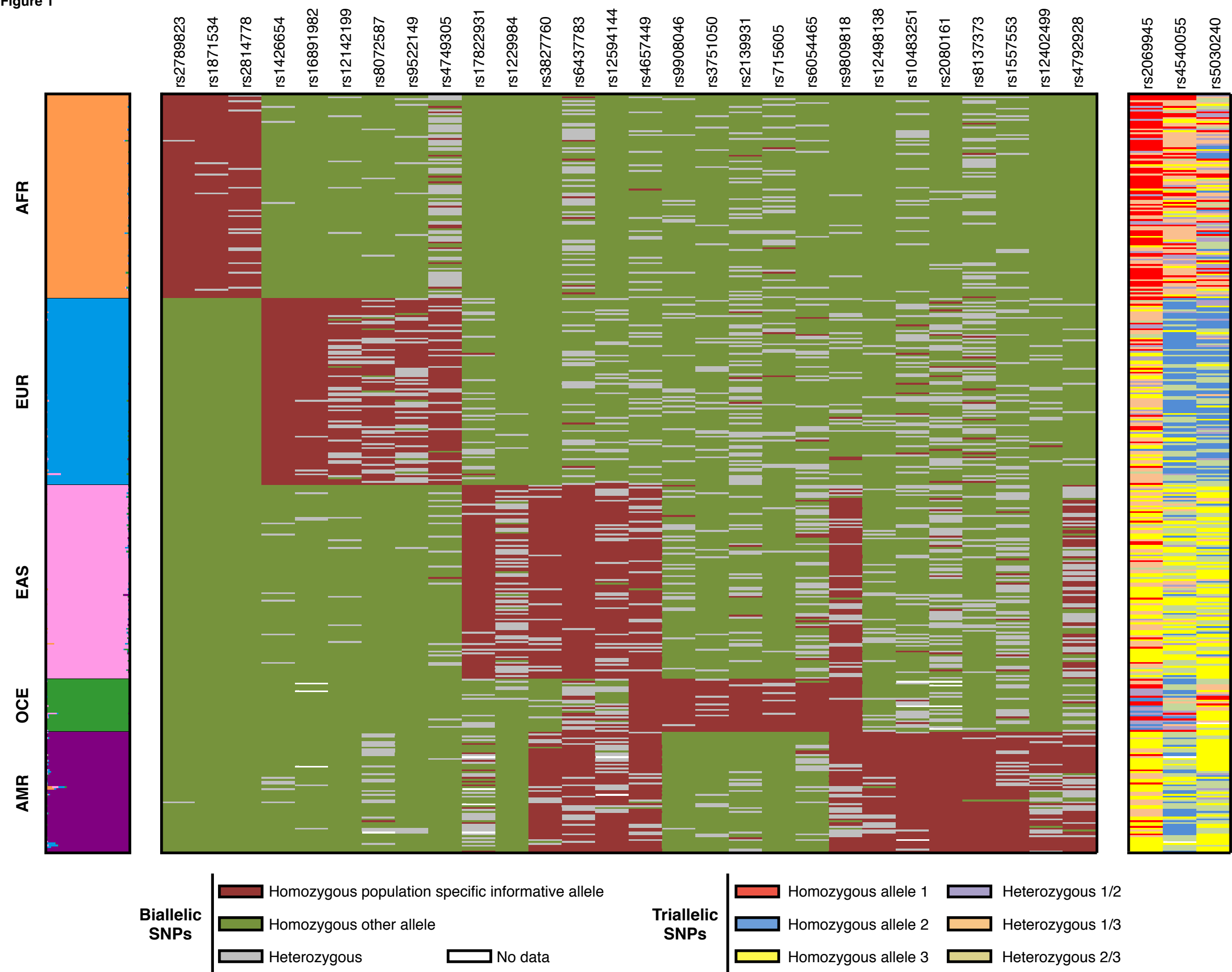
527 **Supplementary Table S1.** PCR and SBE primer sequences and their mixture
528 concentrations for the 31-SNP Nano multiplex reactions. Internal locus
529 identifiers use the IUPAC code for the extended base and a number indicating
530 the position in the multiplex. The SBE primer sequence and non-specific
531 mobility modifying tails are distinguished by upper and lower case nucleotides
532 respectively.

533

534 **Supplementary Table S2.** Cross-validation analyses from Snipper for the 31
535 Nano SNPs. Grey-highlighted values correspond to the percentage of correctly
536 classified individuals from the reference training set of each population.

537

538 **Supplementary File S1.** Custom training set file that can be uploaded
539 unmodified for five-group analyses in Snipper. SNP rs-numbers and internal
540 codes are provided as identifiers of the markers, genotypes are provided in the
541 strand direction that allows a direct classification when using the primers listed
542 in Supplementary Table S1. To classify a profile: choose the '*Classification with
543 a custom Excel file of populations*' option; upload this Excel file in '*Data input
544 (population)*'; choose '*Naive Bayes (Hardy-Weinberg principle applies)*' as the
545 classifier and input the genotype profile of the individual to classify in the same
546 SNP order as this file.



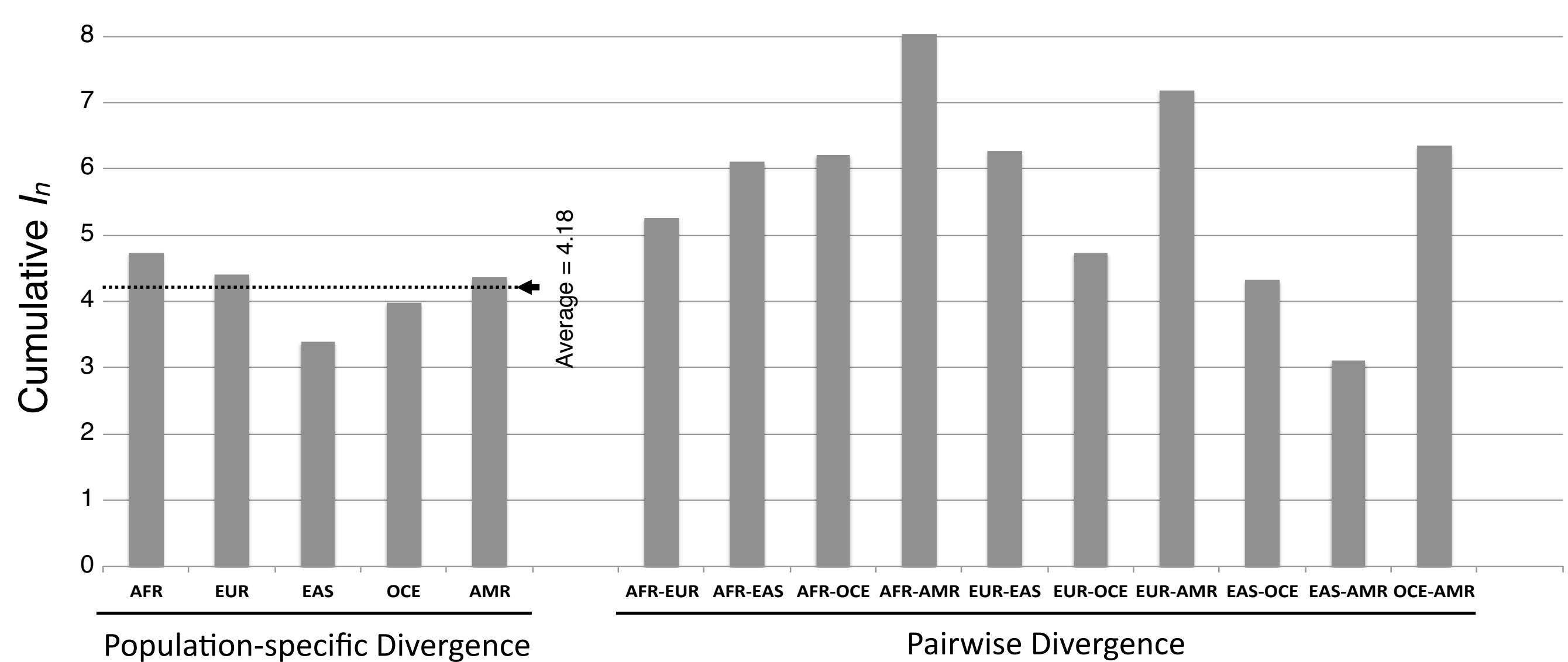
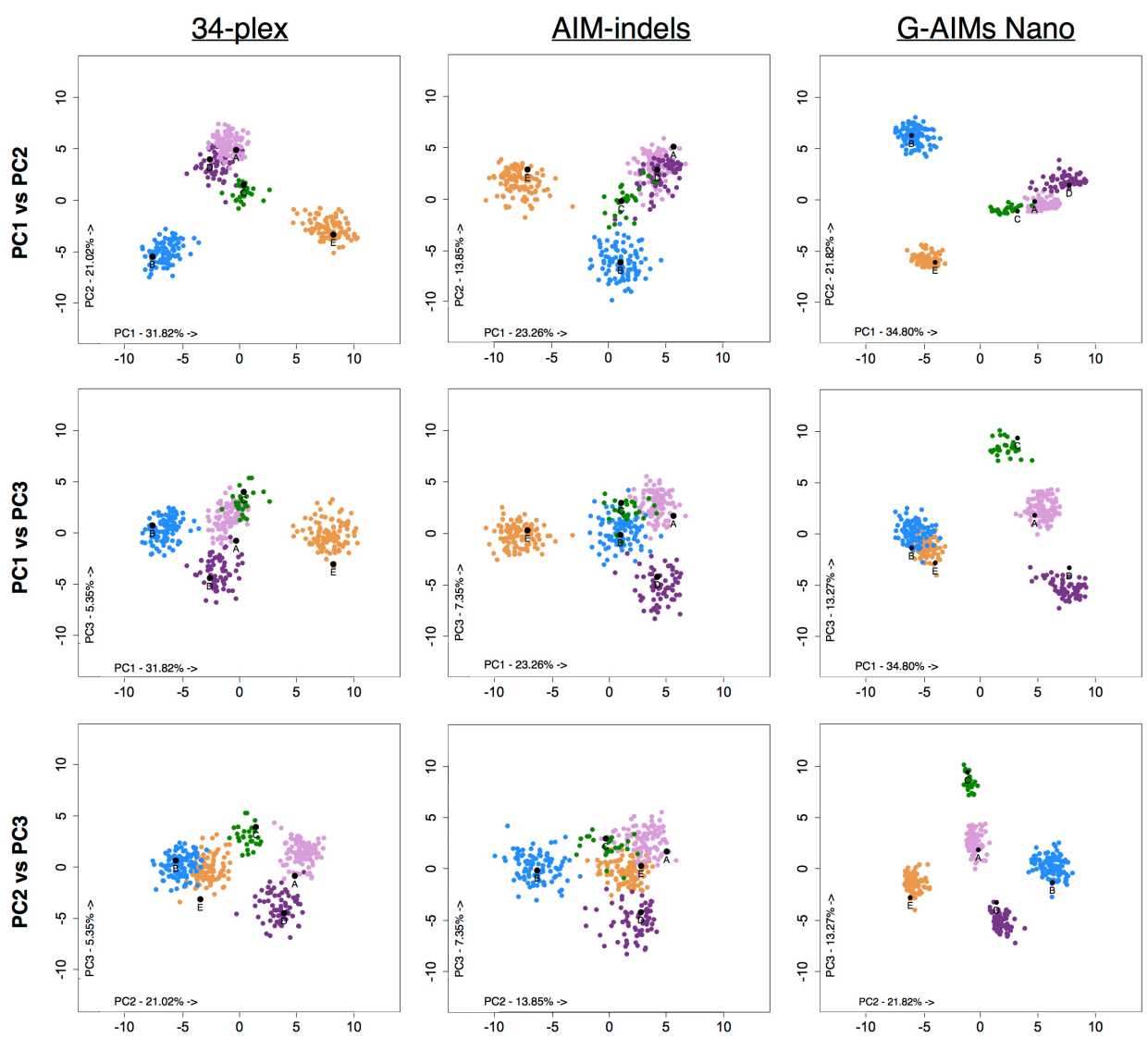
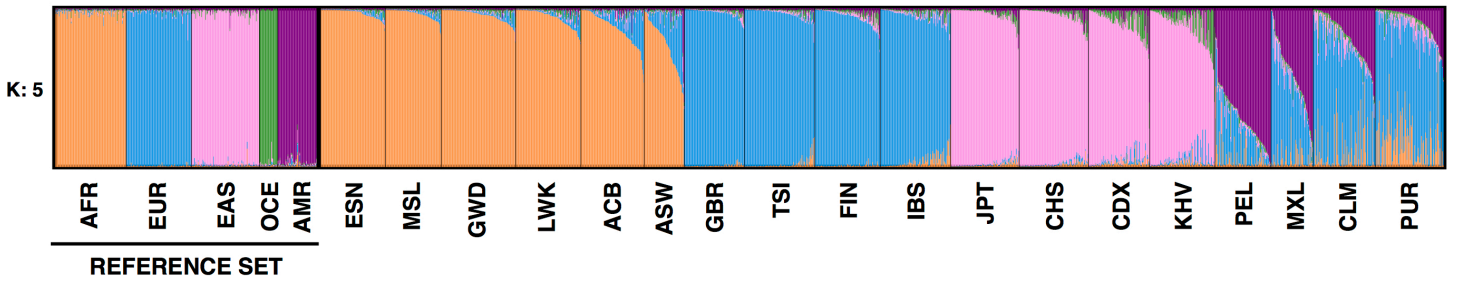


Figure 3



B



C

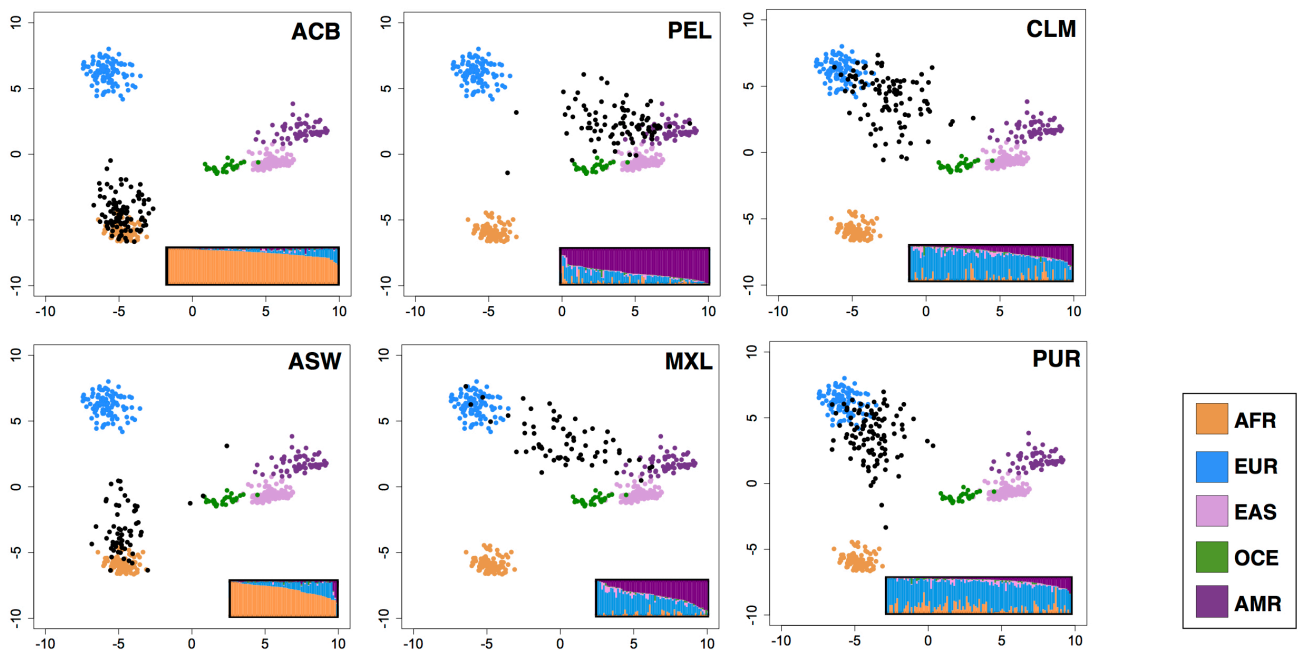
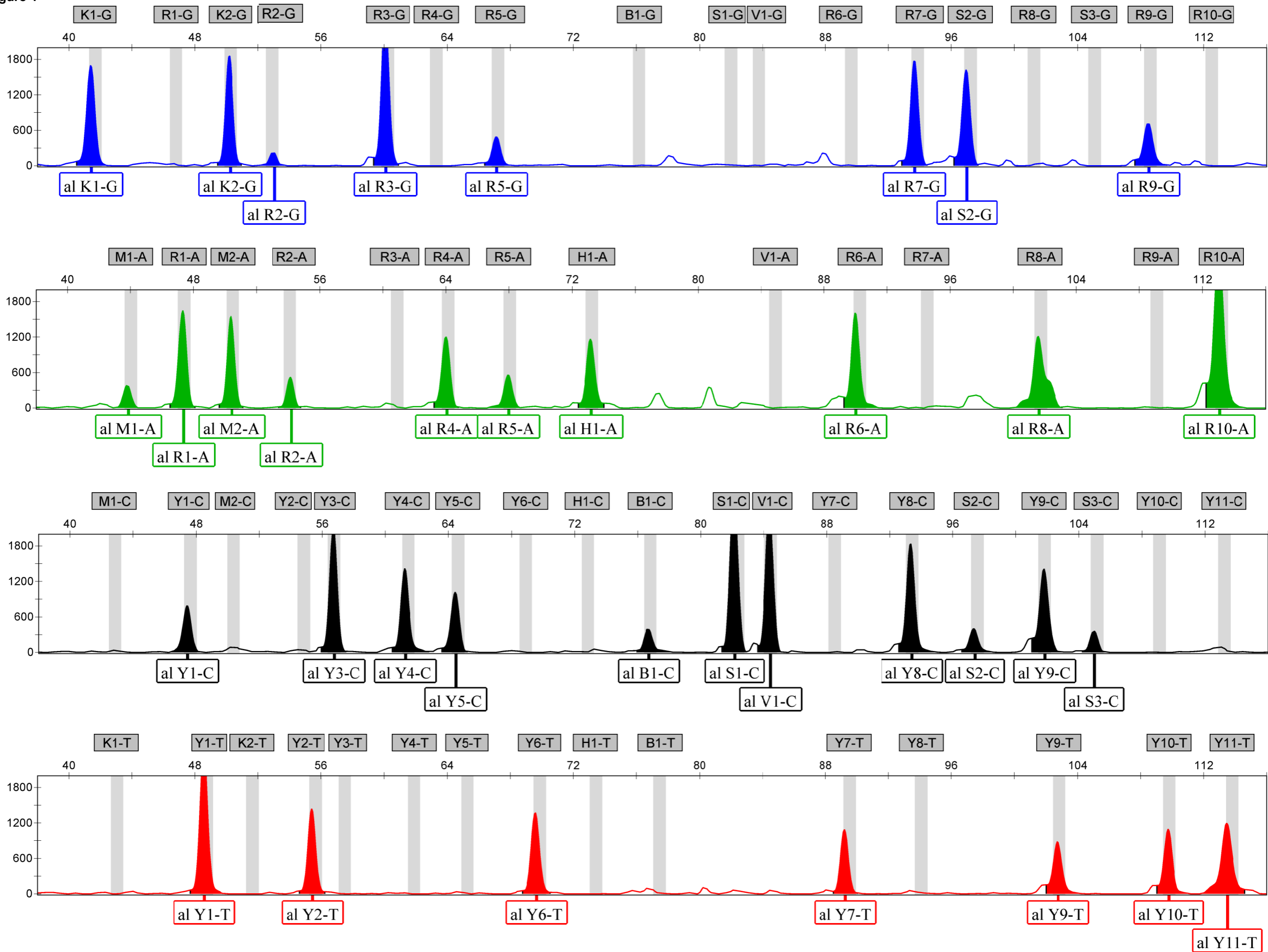


Figure 4



Supplementary Table 2. Cross-validation analyses from Snipper for the 31 Nano SNPs. Grey-highlighted values correspond to the percentage of correctly classified individuals from the reference training set of each population.

	AFR	EUR	EAS	OCE	AMR
Population of AFR origin	100.00 %	0.00 %	0.00 %	0.00 %	0.00 %
Population of EUR origin	0.00 %	100.00 %	0.00 %	0.00 %	0.00 %
Population of EAS origin	0.00 %	0.00 %	100.00 %	0.00 %	0.00 %
Population of OCE origin	0.00 %	0.00 %	0.00 %	100.00 %	0.00 %
Population of AMR origin	0.00 %	0.00 %	0.00 %	0.00 %	100.00 %

Supplementary File S1[Click here to download e-component: Supplementary File S1.xlsx](#)

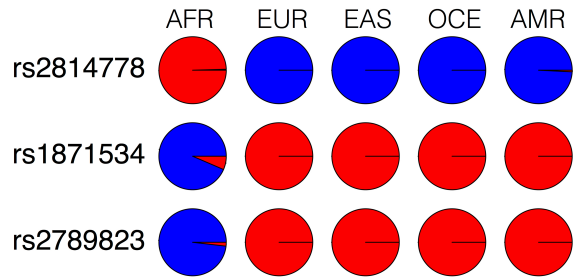
402	31	5	rs3751050_R1	rs12402499_R2
(sample no.)	(SNPs)	(populations)	1	2

Supplementary File S1. Custom training set file that can be uploaded unmodified for

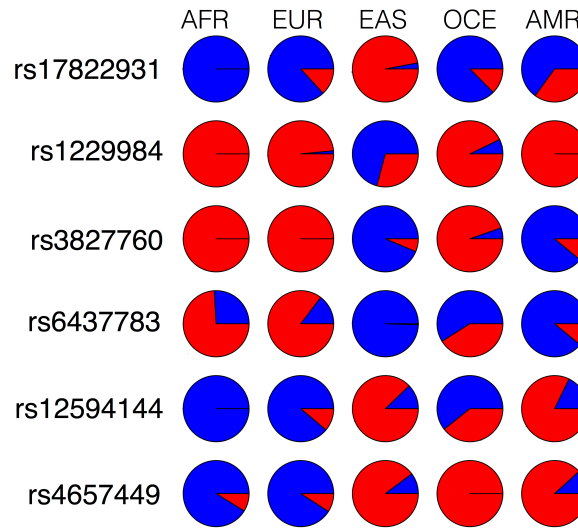
1	AFR	NA18486	AA	GG
2	AFR	NA18488	AA	GG
3	AFR	NA18489	AA	GG
4	AFR	NA18498	AA	GG
5	AFR	NA18499	AA	GG
6	AFR	NA18501	AA	GG
7	AFR	NA18502	AA	GG
8	AFR	NA18504	AA	GG
9	AFR	NA18505	AA	GG
10	AFR	NA18507	AA	GG
11	AFR	NA18508	AA	GG
12	AFR	NA18510	AG	GG
13	AFR	NA18511	AA	GG
14	AFR	NA18516	AA	GG
15	AFR	NA18517	AA	GG
16	AFR	NA18519	AA	GG
17	AFR	NA18520	AA	GG
18	AFR	NA18522	AG	GG
19	AFR	NA18523	AA	GG
20	AFR	NA18853	AA	GG
21	AFR	NA18856	AA	GG
22	AFR	NA18858	AA	GG
23	AFR	NA18861	AA	GG
24	AFR	NA18864	AA	GG
25	AFR	NA18865	AA	GG
26	AFR	NA18867	AA	GG
27	AFR	NA18868	AA	GG
28	AFR	NA18870	AA	GG
29	AFR	NA18871	AA	GG
30	AFR	NA18873	AG	GG
31	AFR	NA18874	AA	GG
32	AFR	NA18876	AA	GG
33	AFR	NA18877	AA	GG
34	AFR	NA18878	AA	GG
35	AFR	NA18879	AA	GG
36	AFR	NA18881	AA	GG

Supplementary Fig. S1 Pie charts indicating the frequency of each Nano SNP allele in five population groups (AFR: Africans; EUR: Europeans; EAS: East Asians; OCE: Oceanians; AMR: Native Americans).

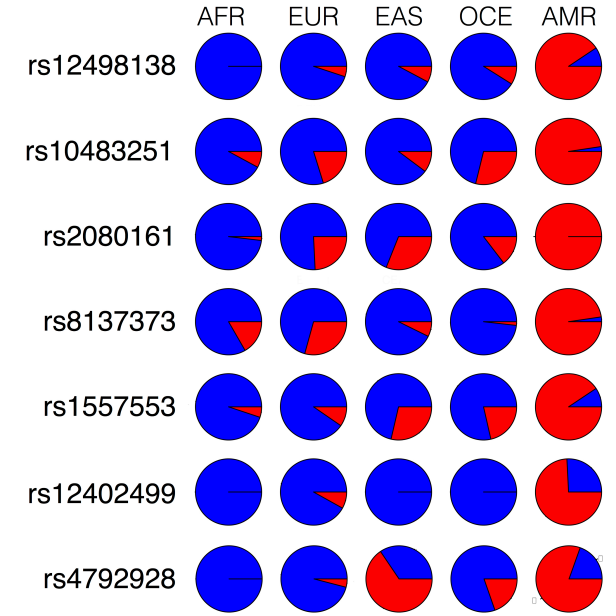
African informative



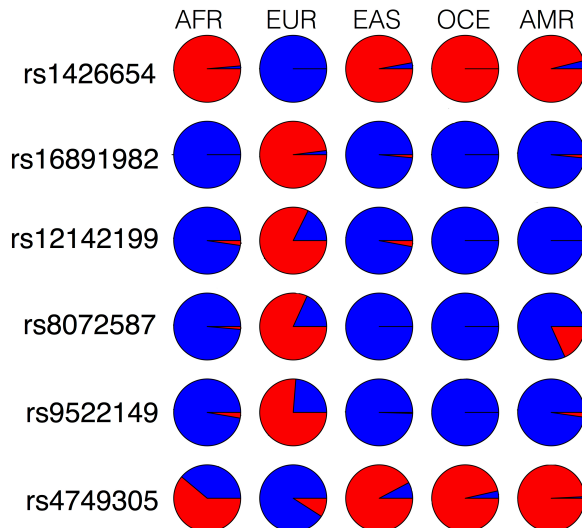
East Asian informative



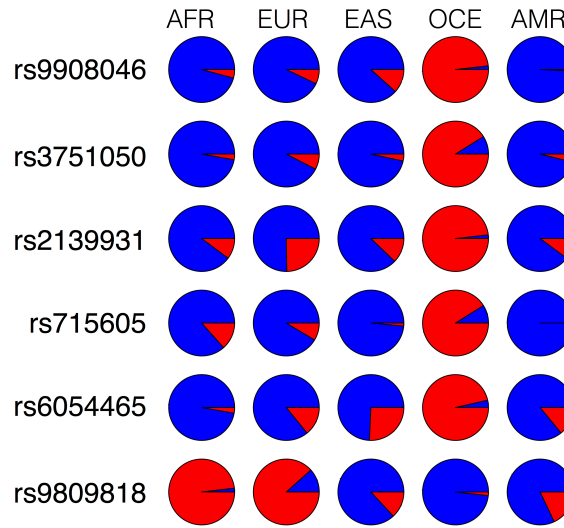
Native American informative



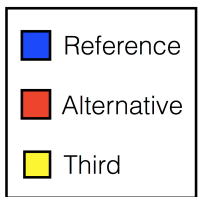
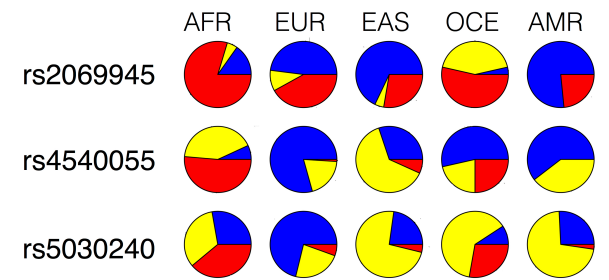
European informative



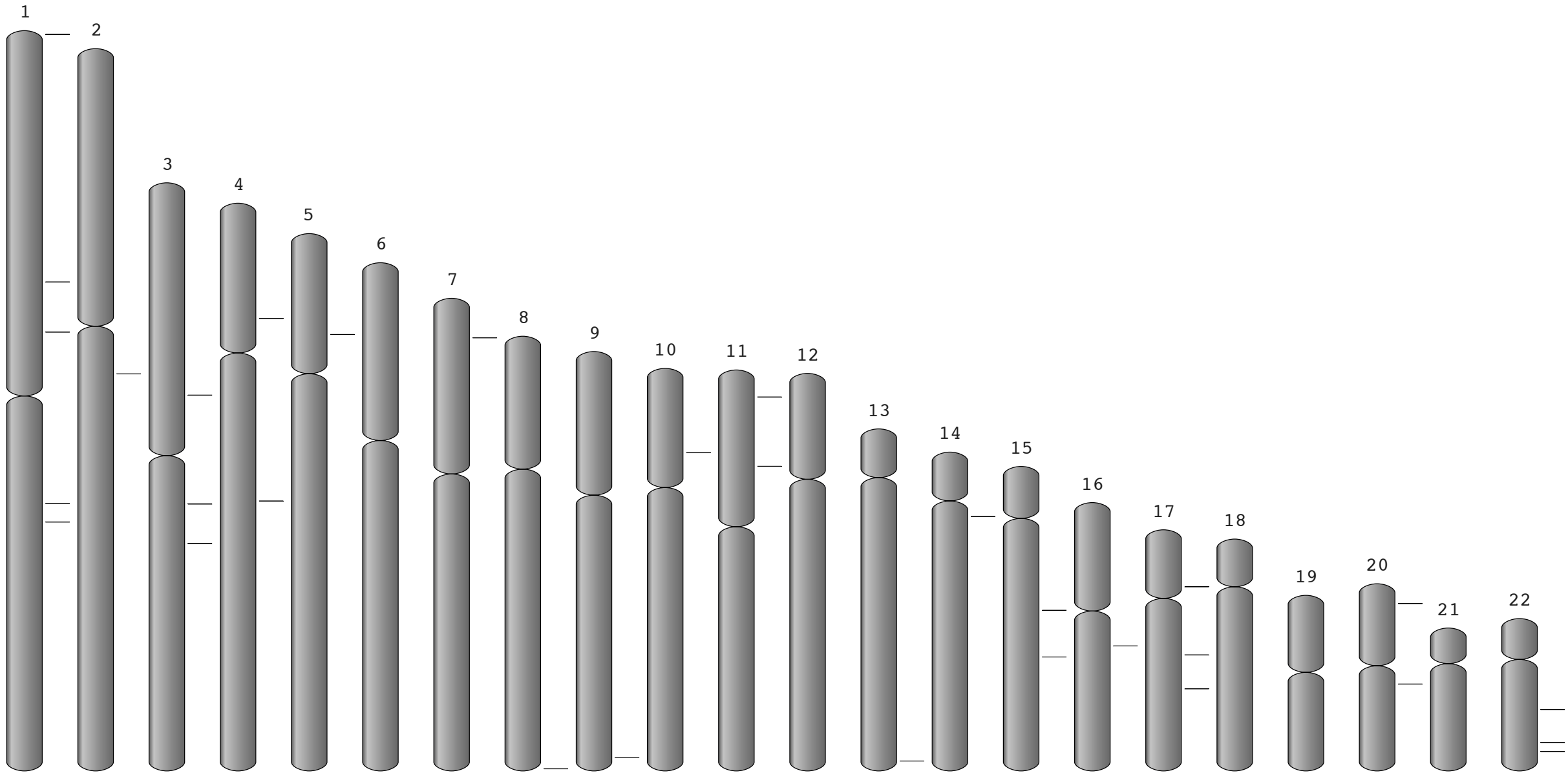
Oceanian informative



Triallelic

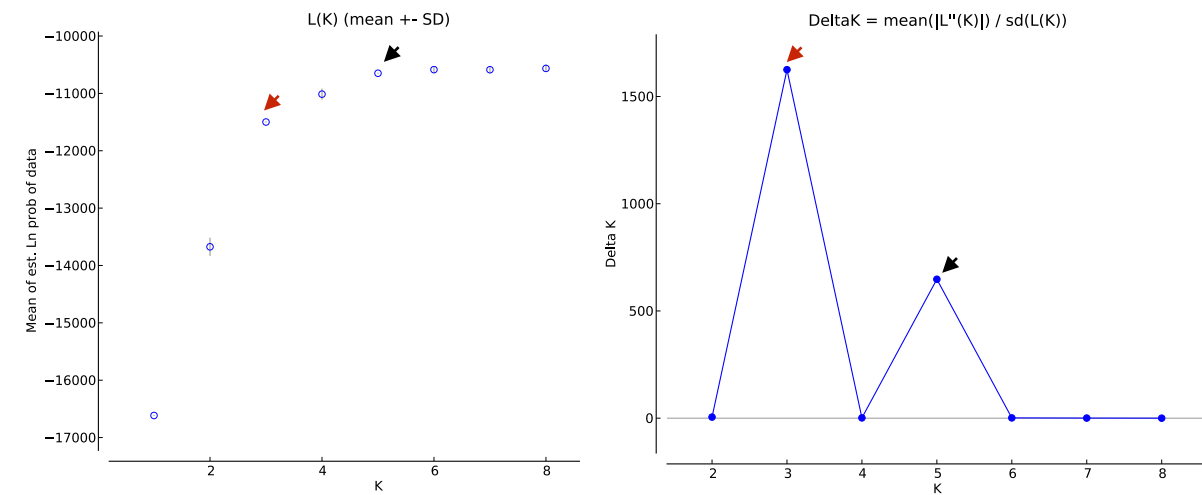
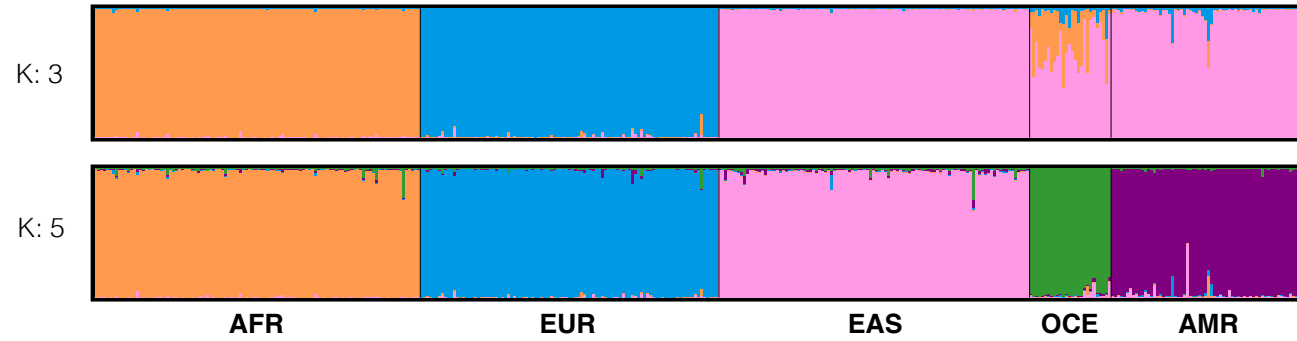


Supplementary Fig. S2. Autosomal chromosome ideograms representing the positions of the Nano SNPs.

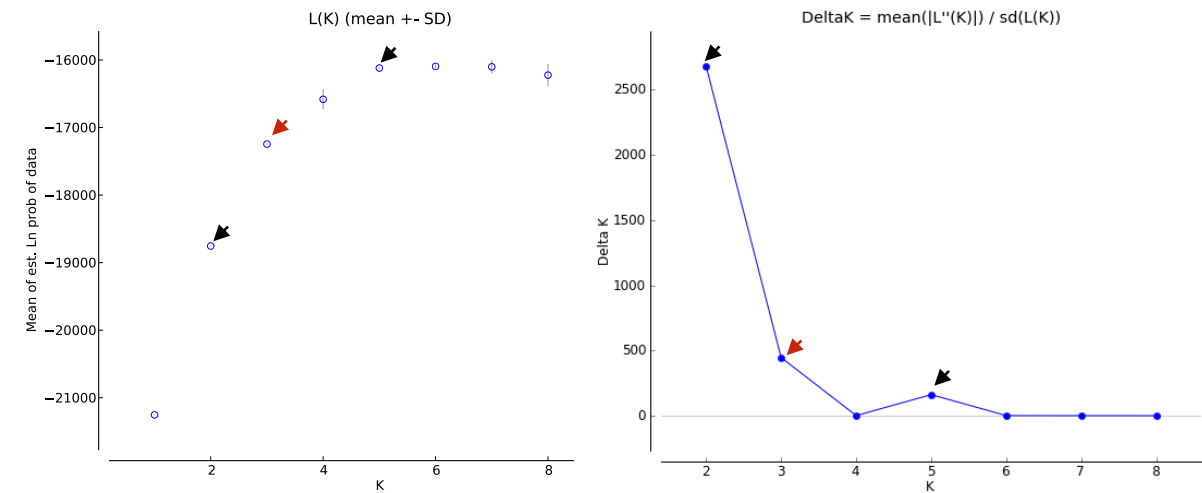
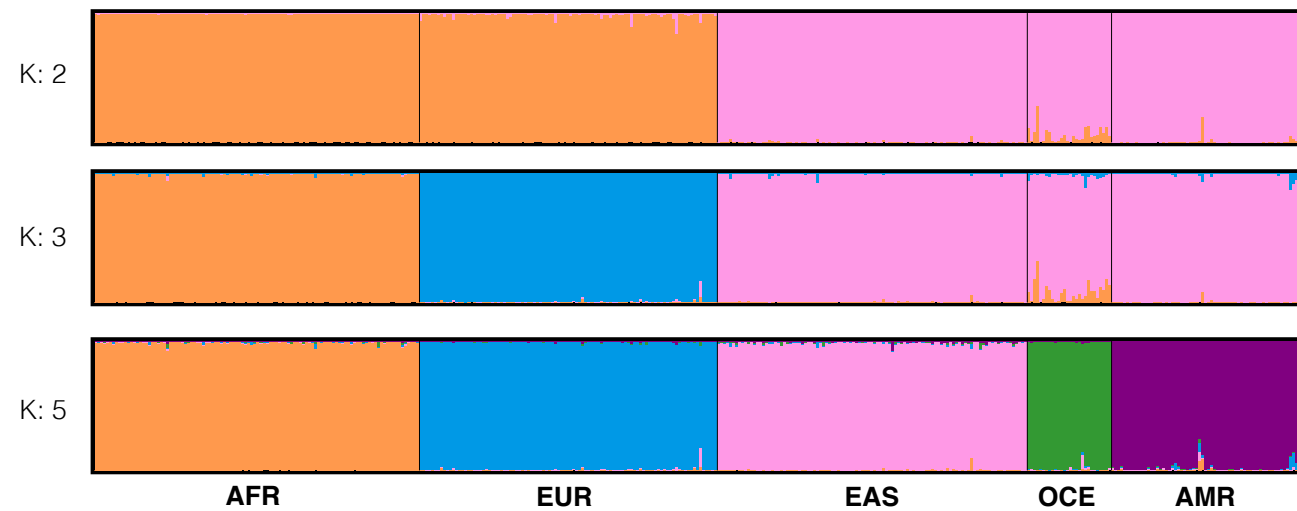


Supplementary Fig. S4. Three panel comparison of STRUCTURE analysis (Nano - 31 SNPs, 34-plex - 34 SNPs and AIM-indels - 44 markers) for the reference population samples. Graphics on the right indicate mean of estimated Ln probability of data and Delta K according to Evanno's method across different numbers of clusters (K). Red arrows indicate the most probable number of clusters taking into account both graphs and black arrows indicate other probable numbers of clusters represented on the left. K: number of clusters; AFR: African; EUR: European; EAS: East Asian; OCE: Oceanian; AMR: Native American.

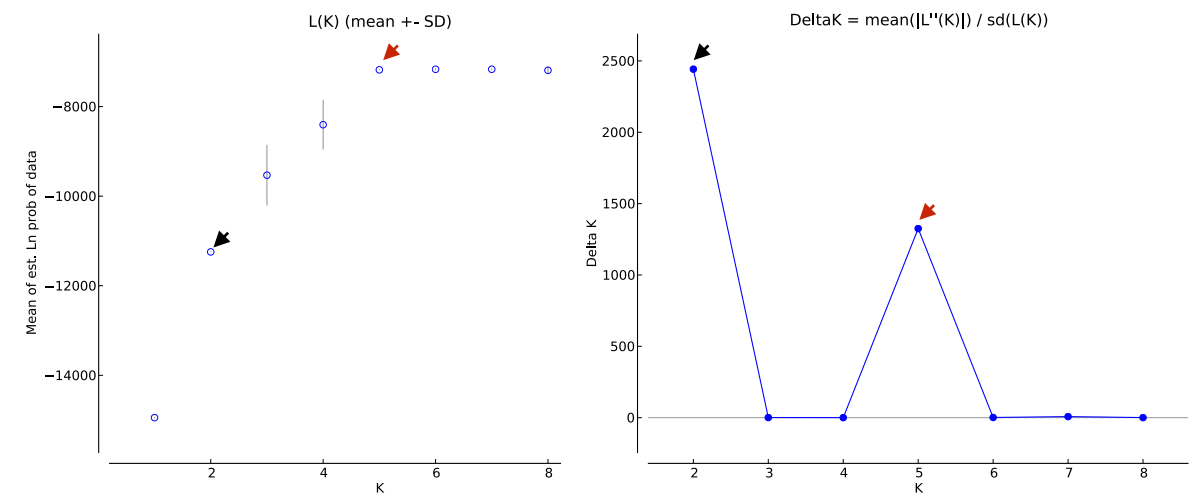
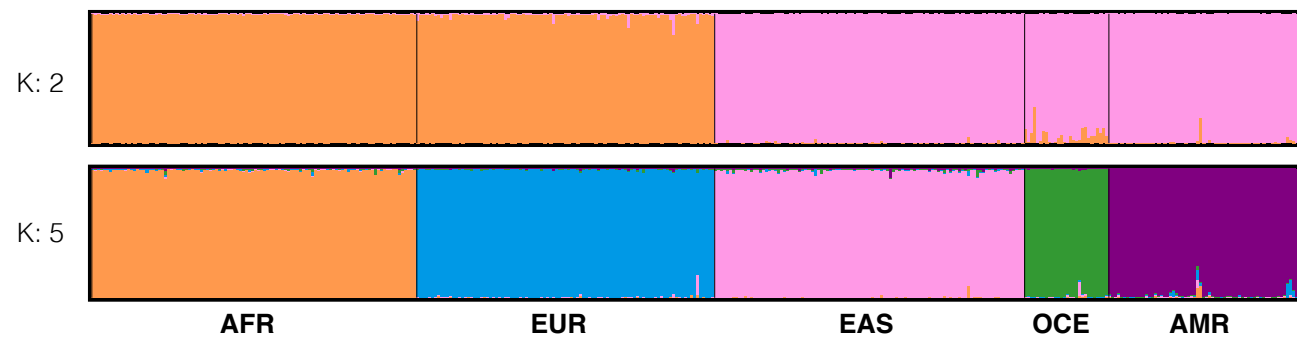
34-plex



AIM-indels



G-AIMs Nano



Supplementary Fig. S5. Likelihood ratios of ancestry estimation for control samples of known ancestry (A-E) and 9947A, removing sequentially the most informative markers. Grey box represents a threshold likelihood value of 1,000 and the red line indicates the most informative fourteen SNPs.

