



## Research paper

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



M. de la Puente<sup>a</sup>, C. Santos<sup>a</sup>, M. Fondevila<sup>a</sup>, L. Manzo<sup>a</sup>, The EUROFORGEN-NoE Consortium, Á. Carracedo<sup>a,b</sup>, M.V. Lareu<sup>a</sup>, C. Phillips<sup>a,\*</sup>

<sup>a</sup> Forensic Genetics Unit, Institute of Forensic Sciences, University of Santiago de Compostela, Spain

<sup>b</sup> Center of Excellence in Genomic Medicine Research, King Abdulaziz University, Jeddah, Saudi Arabia

## ARTICLE INFO

## Article history:

Received 7 October 2015

Received in revised form 20 January 2016

Accepted 21 January 2016

Available online 25 January 2016

## Keywords:

SNPs

AIMs

Biogeographical ancestry

SNaPshot

Population-specific Divergence

## 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.

© 2016 Elsevier Ireland Ltd. All rights reserved.

## 1. Introduction

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

With the recent availability of bench-top systems for massively parallel sequencing (MPS) that are applicable to forensic DNA analysis, it is now possible to assemble multiplexes of 400–500 markers [3]. Such enlarged forensic multiplexes can include a portion of carefully chosen ancestry informative markers (AIMs), e.g., the Illumina ForenSeq panel [4], or can be exclusively composed of AIMs [5–7]. Both approaches raise the level of geographic resolution that can be obtained from tests that keep the

necessary forensic sensitivity. However, the forensic community will take time to adopt, optimise and validate MPS technology as a routine analysis system. Therefore, it is important to continue to develop small-scale AIM sets suited to short-amplicon marker genotyping with validated, universally applicable capillary electrophoresis (CE) analysis regimes [8–10]. However, one drawback with use of small-scale AIM sets is the potential for over-estimation of co-ancestry proportions in individuals with admixed ancestry, stemming from the analysis of genotypes strongly differentiated for some populations but not others. The phenomenon of biased estimation of co-ancestry components was detected in a study of Bolivian populations [11] using 46 ancestry-informative Indels [8] compared with a much larger panel of 446 AIM-SNPs [12]. The Indel set consistently over-estimated European co-ancestry and under-estimated Native American co-ancestry using STRUCTURE-based analyses, indicating that the higher European differentiation of the Indel genotypes inflated the estimates of European co-ancestry proportions. Bearing in mind this effect, construction of a dedicated AIM-SNP set for MPS by the EUROFORGEN Consortium [5] sought to carefully balance the cumulative population-specific Divergence values for the five continental population groups of Africa, Europe, East Asia, Native America and Oceania.

\* Corresponding author.

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

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

2. Materials and methods

2.1. Reference population SNP genotype data and DNA samples

SNP variation data from representative populations without high levels of admixture was obtained from 1000 Genomes Phase III [13] and from the Stanford University HGDP-CEPH SNP analysis [14] using the SPSmart frequency browser [15]. SNP genotype data was compiled from 108 YRI Africans (AFR: Yoruba in Ibadan, Nigeria); 99 CEU Europeans (EUR: Utah Residents with North and Western European ancestry); 103 CHB East Asians (EAS: Han Chinese in Beijing, China); 28HGDP-CEPH Oceanians (OCE: 17 Papuan from New Guinea and 11 Melanesian from Bougainville); and 64 HGDP-CEPH Native Americans (AMR: 14 Karitiana, 8 Surui from Brazil; 21 Maya, 14 Pima from Mexico; and 7 Piapoco from Colombia). Phase III 1000 Genomes populations were also analysed, comprising: as a test set, 99 AFR LWK (Luhya in Webuye, Kenya); 113 AFR GWD (Gambian in Western Divisions in the Gambia); 85 AFR MSL (Mende in Sierra Leone); 99 AFR ESN (Esan in Nigeria); 107 EUR TSI (Toscani in Italia); 99 EUR FIN (Finnish in Finland); 91 EUR GBR (British in England and Scotland); 107 EUR

IBS (Iberian Population in Spain); 104 EAS JPT (Japanese in Tokyo, Japan); 105 EAS CHS (Southern Han Chinese); 99 EAS KHV (Kinh in Ho Chi Minh City, Vietnam); 93 EAS CDX (Chinese Dai in Xishuangbanna, China); plus admixed populations 61 ASW (Americans of African Ancestry in SW USA); 96 ACB (African Caribbeans in Barbados); 104 PUR (Puerto Ricans from Puerto Rico), 94 CLM (Colombians from Medellin, Colombia); 64 MXL (individuals with Mexican Ancestry from Los Angeles USA); 85 PEL (Peruvians from Lima, Peru).

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

2.2. AIM-SNP selection and SNaPshot assay design

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

From the selected SNPs, a 31-plex SNaPshot® single base extension assay was designed and optimised following established guidelines [16]. Locus details and summary allele frequencies for component SNPs are summarised in Table 1. PCR and single base extension (SBE) primers are detailed in Supplementary Table S1.

Table 1 Description, reference allele frequencies and population-specific/pairwise Divergence values ( $I_n$ ) of the 31 Nano SNPs. Chr: chromosome; RA: reference allele. All positions from genome build 37.1 (GRCh37). SNPs are ranked according to their Pop 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	EAS	OCE	AMR	AFR	EUR	EAS	OCE	AMR	AFR	EUR	EAS	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	0.000	0.000	0.000	0.000	0.000	0.000	0.001
	rs1871534	8	145539681	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	0.000	0.000	0.000	0.000	0.000	
	rs2789823	9	136769888	G	0.935	0.000	0.000	0.000	0.000	0.555	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	0.000	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	0.000	0.000	0.000	0.000	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	0.000	0.000	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	0.000	0.000	0.000	0.000	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	0.000	0.000	0.000	0.000	0.000	0.000
	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	0.000	0.000	0.000	0.000	0.000	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	0.000	0.000	0.000	0.000	0.000	0.000	
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	0.034	0.434	0.251	0.036	0.000	0.000	0.000	0.000	0.000	0.000
	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	0.000	0.000	0.000	0.000	0.000	0.000
	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	0.000	0.000	0.000	0.000	0.000	0.000
	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	0.000	0.000	0.000	0.000	0.000	0.000
	rs12594144	15	64181351	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	0.000	0.000	0.000	0.000	0.000	0.000
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	0.000	0.000	0.000	0.000	0.000	0.000	
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	0.000	0.000	0.000	0.000	0.000	0.000
	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	0.000	0.000	0.000	0.000	0.000	
	rs2139931	1	84560527	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	0.000	0.000	0.000	0.000	0.000	
	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	0.000	0.000	0.000	0.000	0.000	
	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	0.000	0.000	0.000	0.000	0.000	
rs9809818	3	71480566	C	0.019	0.116	0.869	0.982	0.820	0.295	0.204	0.120	0.172	0.227	0.102	0.021	0.446	0.603	0.099	0.319	0.245	0.276	0.006	0.002	0.042	0.000	0.000	0.000	0.000		
AMR	rs12498138	3	121459598	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	0.000	0.000	0.000	0.000	0.000	
	rs10483251	14	21671277	G	0.921	0.798	0.898	0.712	0.024	0.095	0.006	0.040	0.000	0.429	0.016	0.001	0.038	0.497	0.010	0.055	0.370	0.028	0.489	0.301	0.000	0.000	0.000	0.000	0.000	
	rs2080161	7	13331150	T	0.981	0.758	0.889	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	0.368	0.033	0.314	0.462	0.000	0.000	0.000	0.000	0.000	
	rs8137373	22	41729216	G	0.833	0.707	0.927	0.023	0.000	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	0.563	0.000	0.000	0.000	0.000	0.000	
	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	0.000	0.000	0.000	0.000	0.000	
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	0.000	0.000	0.000	0.000	0.000		
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	0.000	0.000	0.000	0.000	0.000		
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	0.000	0.000	0.000	0.000	0.000	
	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	0.204	0.084	0.355	0.147	0.142	0.300	0.130	0.217	0.027	0.098	0.063	0.101	0.000	0.000	0.000	0.000	0.000	
	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	0.136	0.135	0.083	0.124	0.067	0.001	0.085	0.000	0.000	0.000	0.000	0.000	
<b>CUMULATIVE VALUES</b>										<b>4.739</b>	<b>4.404</b>	<b>3.392</b>	<b>3.986</b>	<b>4.374</b>	<b>5.</b>															

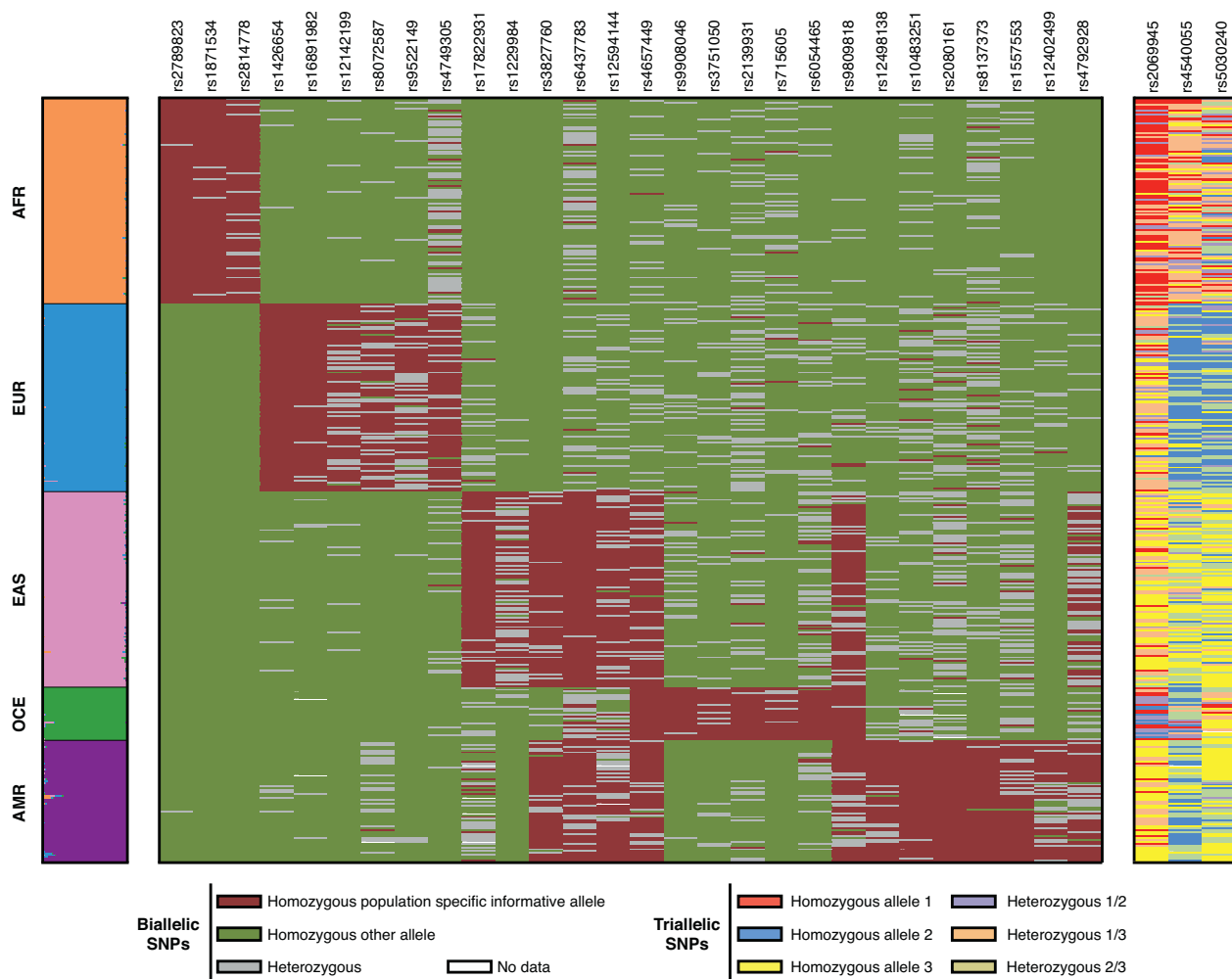
PCR reactions comprised: 1  $\mu$ L Buffer II (100 mM Tris–HCl, pH 8.3, 500 mM KCl); 1.8  $\mu$ L 25 mM MgCl<sub>2</sub>; 0.1  $\mu$ L AmpliTaq Gold<sup>®</sup> DNA Polymerase (at 5 U/ $\mu$ L); 0.4  $\mu$ L of GeneAmp<sup>®</sup> 10 mM dNTP Mix with dTTP (Applied Biosystems, AB); 1  $\mu$ L of 3.2 mg/ml bovine serum albumin; 1.5  $\mu$ L PCR primer mix (Supplementary Table S1); 1 ng of target DNA adjusted to total reaction volume of 10  $\mu$ L. PCR cycling with GeneAmp<sup>®</sup> PCR System 9700 or 2700 (AB) thermocyclers used conditions: 10 min at 95 °C, 32 cycles of 30 s at 95 °C, 40 s at 62 °C and 1 min at 72 °C with a final extension of 20 min at 72 °C. PCR primer clean up combined 2.5  $\mu$ L PCR product with 1  $\mu$ L of 1 in 3 diluted Illustra<sup>™</sup> ExoStar<sup>™</sup> 1-Step (GE Healthcare) then incubation at 37 °C for 45 min and enzyme inactivation at 85 °C for 15 min. SBE reactions comprised: 1.25  $\mu$ L of 1 in 2 diluted SNaPshot<sup>®</sup> Multiplex Ready Reaction Mix; 0.75  $\mu$ L SBE primers (Supplementary Table S1) and 1  $\mu$ L purified PCR product in a total volume of 3  $\mu$ L. SBE cycling used conditions: 33 cycles of 10 s at 96 °C, 5 s at 59 °C and 30 s at 60 °C. SBE primer clean up combined the full SBE volume with 1  $\mu$ L of 1 in 2 diluted Illustra<sup>™</sup> Shrimp Alkaline Phosphatase (GE Healthcare) then incubating at 37 °C for 80 min and enzyme inactivation at 85 °C for 15 min. Purified SBE products were then prepared for CE detection by adding 1  $\mu$ L of product to 9.5  $\mu$ L of Hi-Di<sup>™</sup> Formamide (AB) and 0.25  $\mu$ L of GeneScan<sup>™</sup>-120 LIZ<sup>®</sup> Size Standard (AB). Electrophoresis was performed in an ABI Prism 3130xl Genetic Analyser, with 36 cm capillary arrays and POP-4<sup>™</sup> polymer using standard

conditions. Electropherograms were visualised using AB GeneMapper<sup>®</sup> ID Software v. 3.2.1.

### 2.3. Analysis of population variation in the selected SNPs

Population-specific Divergence and simple pairwise Divergence values were calculated using the Snipper cross-validation option ([http://mathgene.usc.es/snipper/analysispopfile2\\_new.html](http://mathgene.usc.es/snipper/analysispopfile2_new.html)) by marking SNP genotype profiles as AFR and non-AFR, etc., or by comparing each pair of population groups in turn. Output from Snipper lists Shannon's Divergence values for each SNP from the comparisons made by cross-validation [17]. These values were converted to the more widely-used Rosenberg's informativeness-for-assignment metric:  $I_n$  [18], by multiplication with 0.693 (i.e., converting the natural log to log(2)). The Snipper portal was also used to cross-validate the reference population data or calculate classification likelihood ratios (LRs) by uploading an Excel file of reference data (provided ready to use as Supplementary File S1) or by choosing analysis options available in the portal.

Population analyses with STRUCTURE v. 2.3.4 [19] were performed following recommendations outlined previously [20]. Parameters comprised: five iterations (for  $K=1$  to  $K=9$ ) of 100,000 burnin steps and 100,000 MCMC steps, correlated allele frequencies under the Admixture model (no POPFLAG for just reference populations and POPFLAG for analyses of reference



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

populations plus test or admixed populations). The optimum  $K$  value was estimated by computing results with STRUCTURE HARVESTER [21] and following previous guidelines [22]. Ancestry membership plots were constructed with CLUMPAK v. 1.1 [23] or a combination of CLUMPP v. 1.1.2 [24] and Distruct v. 1.1 [25]. PCA analysis was performed using R software v. 3.1.2 [26] and executing a homemade script.  $F_{ST}$  calculations and graphics were computed using Arlequin v. 3.5 [27].

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

### 3. Results

#### 3.1. Characteristics and PSD balance of the Nano SNP set

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

With the reduction in scale from an MPS multiplex of 128 SNPs to the SNaPshot 31-plex, it is important to ensure the cumulative PSD values remain at comparable levels for each group. Reference population comparisons produced the cumulative PSD values listed in Table 1 and summarised in Fig. 2. These indicate the reduced East Asian PSD of 3.39 was noticeably lower than the average of 4.18 and the African PSD of 4.74 was the highest but comparable to three population groups. The reduced differentiation of East Asians from Native American and Oceanian populations is illustrated by their similar allele frequency distributions for many SNPs, exemplified by rs4657449 and rs9809818. The close relationship of East Asian and Native American population variability is highlighted by Fig. 1, with little divergence between the two groups evident in rs3827760, rs6437783 and rs12594144.

This suggests one or two extra East Asian-informative SNPs could address this slight PSD imbalance in future adjustments of the Nano SNP set.

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

#### 3.2. Ancestry inference capabilities of the Nano SNP set

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

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

Analysis of other genotype data from 1000 Genomes which were marked as 'study populations' (POPFLAG=0) gave patterns consistent with analyses made during the development of the 128 Global ancestry set [5] or subsequent analysis of admixed 1000 Genomes ACB, ASW, PEL, MXL, PUR, CLM samples using the same SNPs (Fig. 7 in [2]). The STRUCTURE cluster plots (Fig. 3B) in particular match well with results of both previous analyses using many more SNPs, indicating PEL have the highest Native American co-ancestry proportions and PUR show predominantly European co-ancestry. The PCA plots arranged individually along with STRUCTURE cluster plots for the six admixed 1000 Genomes

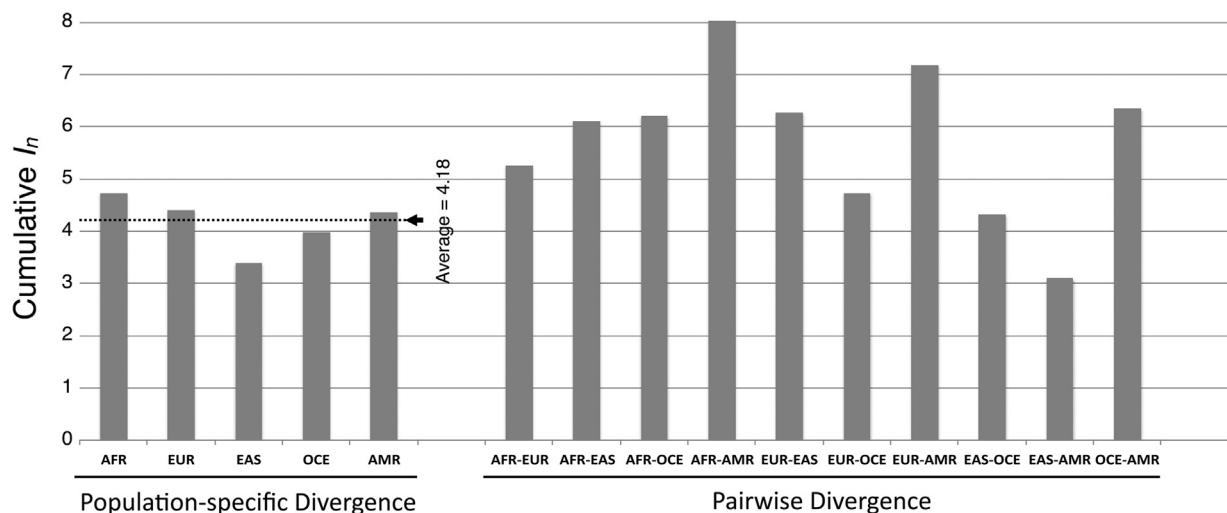
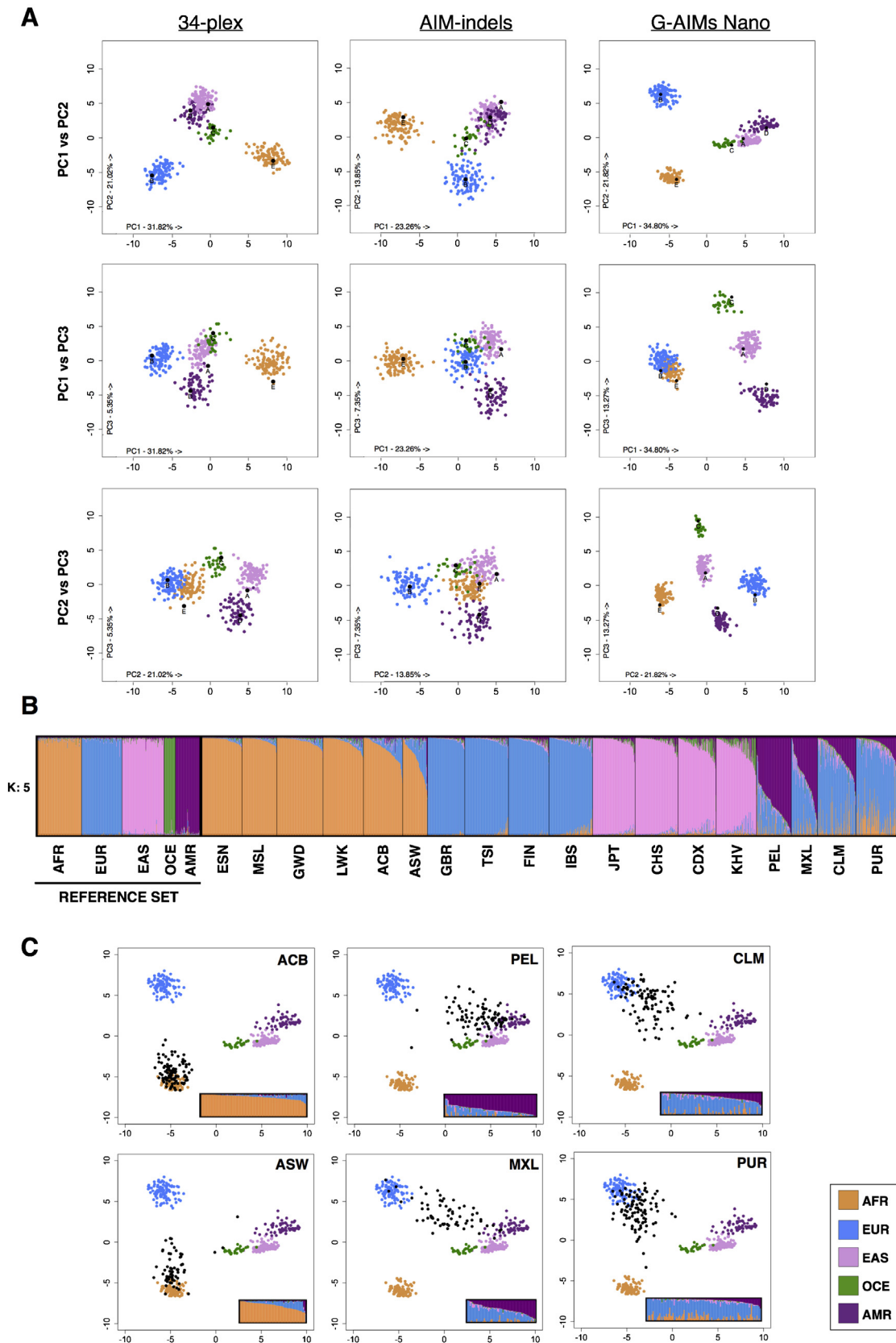


Fig. 2. Bar charts indicating cumulative, population-specific and pairwise  $I_n$  values for Nano SNPs.





**Fig. 3.** (A) Three AIM panels compared with PCA analyses (triallelic SNPs not included in analyses, so Nano = 28 SNPs, 34-plex = 32 SNPs and AIM-indels = 44 markers) for reference samples. PC: principal component; AFR: African; EUR: European; EAS: East Asian; OCE: Oceanian; AMR: Native American. Control samples with known ancestry are shown as black points. (B) STRUCTURE analysis for admixed and non-admixed populations included in the study. (C) Mixed populations are represented below in black in the corresponding PCA analyses (PC1 vs PC2). AFR: African; EUR: European; EAS: East Asian; OCE: Oceanian; AMR: Native American. The 18 individual 1000 Genomes population abbreviations: ACB; ASW; CDX; CHS; CLM; ESN; FIN; GBR; GWD; IBS; JPT; KHV; LWK; MSL; MXL; PEL; PUR; and TSI are described in Section 2.1.

**Table 2**

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

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	5.5E+06	9.7E+06	2.9E+16
		more likely to be AMR <sup>a</sup>	more likely to be EAS	EAS	EAS
B	European (EUR)	4.4E+16	1.7E+11	1.0E+28	1.3E+29
		EUR	EUR	EUR	EUR
C	Oceanian (OCE)	1.0E+07	4.0E+07	1.5E+14	1.5E+16
		OCE	OCE	OCE	OCE
D	Native American (AMR)	1.0E+05	1.2E+09	1.1E+14	4.2E+08
		AMR	AMR	AMR	AMR
E	African (AFR)	6.1E+19	2.8E+21	3.6E+40	2.0E+29
		AFR	AFR	AFR	AFR

<sup>a</sup>Likelihood ratio of AMR/EAS likelihoods.

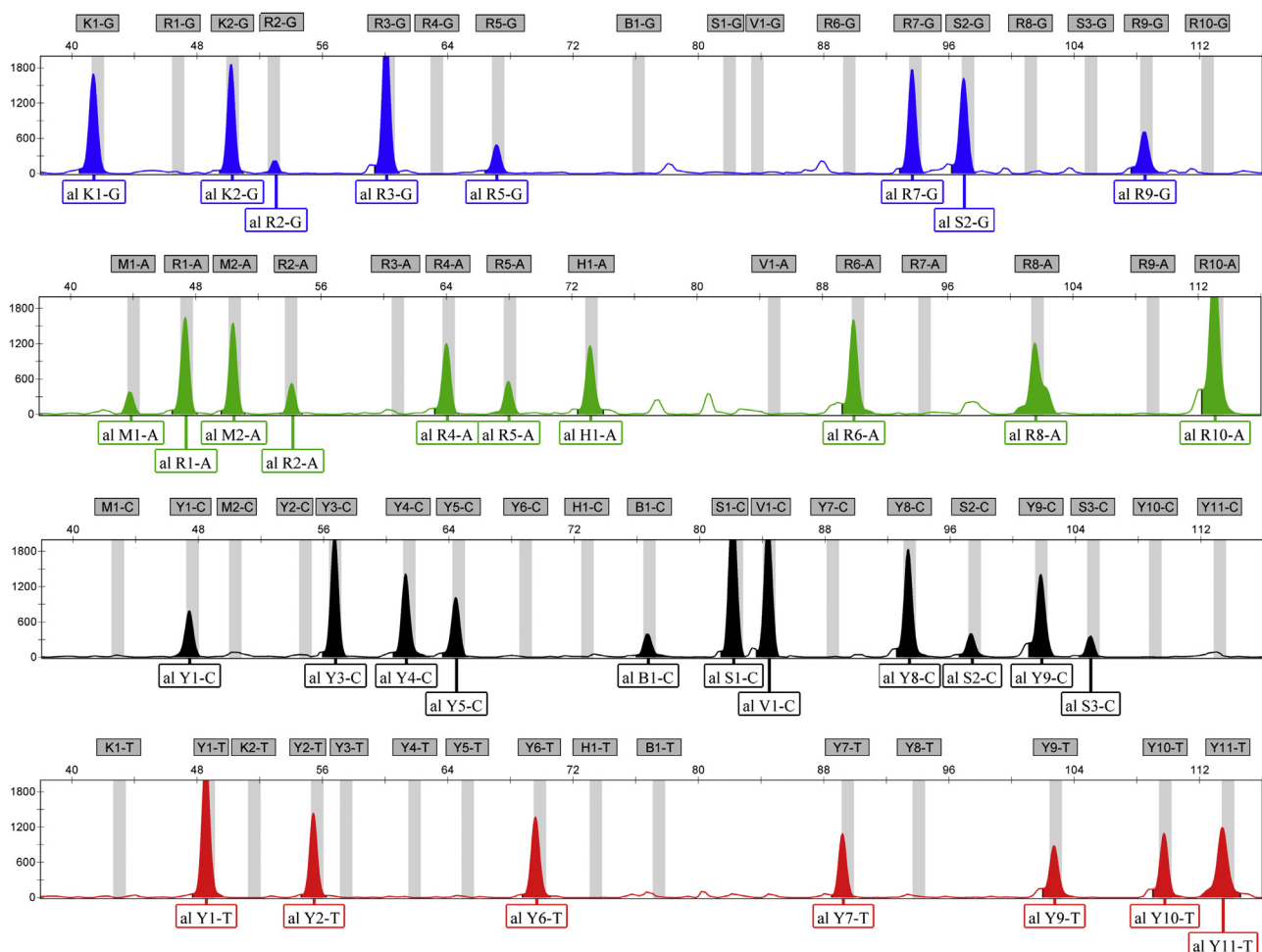
populations in Fig. 3C also give very similar cluster distributions that are positioned between, or close to, the expected population admixture contributor clusters.

As an additional simple gauge of the ancestry informativeness of the 31 Nano SNPs, this assay was used to analyse the control DNAs used in a collaborative EDNAP exercise that assessed the 34 SNPs and 46 Indels described above. The five control DNAs each have confirmed ancestry from one of the five continental regions and they were analysed using Bayes analyses and PCA in Snipper. The five control DNAs are positioned in the middle of their respective population group clusters describing the correct

ancestry in each case. The Bayes analysis likelihoods obtained from Snipper are listed in Table 2. Moreover, as shown in Supplementary Fig. S5, high likelihood ratio values can be obtained from partial profiles even when the fourteen most informative markers are missing.

### 3.3. Forensic performance of the Nano SNPs set

A typical SNaPshot profile from analysis of 1 ng of the 9947A control DNA with the Nano assay is shown in Fig. 4. The dilutions series of 9947A gave full profiles with 0.5, 0.25, 0.125 and 0.064 ng



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

of DNA. Locus and allele drop out occurred with 0.032 and 0.016 ng of input DNA, however >80% profile completeness was obtained for these analyses.

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

#### 4. Discussion

When originally rebuilding a SNP-based forensic ancestry multiplex for MPS analysis, we started to recognise subsets of the most informative markers that would be well suited to development of smaller SNaPshot tests. The present study describes the compilation of 31 SNPs, mainly representing the most informative markers from the full set of 128 Global AIMs. These SNPs maintain the capacity to differentiate five continentally-defined population groups. Therefore, the Nano assay extends the three group comparisons possible with existing SNaPshot forensic ancestry tests [9,28–30] to the two additional population groups of Native Americans and Oceanians. Each of these population groups contribute to admixture patterns commonly seen in large parts of the regions they occupy. For this reason, we prioritised the preservation of balanced PSD, although this was difficult to maintain when accomplishing the 75% reduction in multiplex size. One outcome of this process was a disproportionate lowering of the East Asian cumulative PSD compared to the other groups that we aim to address by careful choice of 1–2 additional AIMs.

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

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

In conclusion, the Nano ancestry assay has brought together a highly informative set of markers with well-balanced population-specific Divergence for the five population groups it is designed to analyse. This characteristic minimises the co-ancestry proportion

estimation bias when analysing admixed samples. Moreover, the single SNaPshot reaction shows high sensitivity (complete profiles obtained down to 64 pg of input DNA), with the potential to analyse degraded samples, making it an ideal forensic ancestry assay for the full range of casework applications where DNA is analysed with capillary electrophoresis.

#### Acknowledgements

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

#### Appendix A. Supplementary data

Supplementary data associated with this article can be found, in the online version, at <http://dx.doi.org/10.1016/j.fsigen.2016.01.015>.

#### References

- [1] M. Kayser, Forensic DNA Phenotyping: predicting human appearance from crime scene material for investigative purposes, *Forensic Sci. Int.: Genet.* 18 (2015) 33–48.
- [2] C. Phillips, Forensic genetic analysis of bio-geographical ancestry, *Forensic Sci. Int.: Genet.* 18 (2015) 49–65.
- [3] B. Budowle, D.H. Warshauer, S.B. Seo, et al., Deep sequencing provides comprehensive multiplex capabilities, *Forensic Sci. Int.: Genet. Suppl. Ser.* 4 (2013) e334–e335.
- [4] Illumina, Forenseq DNA signature prep kit (2014) <http://www.illumina.com/products/forenseq-dna-signature-kit.ilmn>.
- [5] C. Phillips, W. Parson, B. Lundsberg, et al., Building a forensic ancestry panel from the ground up: the EUROFORGEN Global AIM-SNP set, *Forensic Sci. Int. Genet.* 11 (2014) 13–25.
- [6] K.K. Kidd, W.C. Speed, A.J. Pakstis, et al., Progress toward an efficient panel of SNPs for ancestry inference, *Forensic Sci. Int.: Genet.* 10 (2014) 23–32.
- [7] J. Kidd, F. Friedlaender, W. Speed, et al., Analyses of a set of 128 ancestry informative single-nucleotide polymorphisms in a global set of 119 population samples, *Invest. Genet.* 2 (2011) 1–13.
- [8] R. Pereira, C. Phillips, N. Pinto, et al., Straightforward inference of ancestry and admixture proportions through ancestry-informative insertion deletion multiplexing, *PLoS One* 7 (2012) e29684.
- [9] M. Fondevila, C. Phillips, C. Santos, et al., Revision of the SNPforID 34-plex forensic ancestry test: assay enhancements, standard reference sample genotypes and extended population studies, *Forensic Sci. Int. Genet.* 7 (2013) 63–74.
- [10] C. Santos, M. Fondevila, D. Ballard, et al., Forensic ancestry analysis with two capillary electrophoresis ancestry informative marker (AIM) panels: results of a collaborative EDNAP exercise, *Forensic Sci. Int.: Genet.* 19 (2015) 56–67.
- [11] P. Taboada-Echalar, V. Álvarez-Iglesias, T. Heinz, et al., The genetic legacy of the pre-colonial period in contemporary Bolivians, *PLoS One* 8 (2013) e58980.
- [12] J.M. Galanter, J.C. Fernandez-Lopez, C.R. Gignoux, et al., Development of a panel of genome-wide ancestry informative markers to study admixture throughout the Americas, *PLoS Genet.* 8 (2012) e1002554.
- [13] An integrated map of genetic variation from 1,092 human genomes, *Nature* 491, 2012, 56–65.
- [14] J.Z. Li, D.M. Absher, H. Tang, et al., Worldwide human relationships inferred from genome-wide patterns of variation, *Science* 319 (2008) 1100–1104.
- [15] J. Amigo, A. Salas, C. Phillips, et al., SPSmart: adapting population based SNP genotype databases for fast and comprehensive web access, *BMC Bioinform.* 9 (2008) 428.
- [16] J.J. Sanchez, P. Endicott, Developing multiplexed SNP assays with special reference to degraded DNA templates, *Nat. Protocols* 1 (2006) 1370–1378.
- [17] C. Phillips, A. Salas, J.J. Sánchez, et al., Inferring ancestral origin using a single multiplex assay of ancestry-informative marker SNPs, *Forensic Sci. Int.: Genet.* 1 (2007) 273–280.
- [18] N.A. Rosenberg, L.M. Li, R. Ward, et al., Informativeness of genetic markers for inference of ancestry, *Am. J. Human Genet.* 73 (2003) 1402–1422.
- [19] J.K. Pritchard, M. Stephens, P. Donnelly, Inference of population structure using multilocus genotype data, *Genetics* 155 (2000) 945–959.

- [20] L. Porras-Hurtado, Y. Ruiz, C. Santos, et al., An overview of STRUCTURE: applications, parameter settings, and supporting software, *Front. Genet.* 4 (2013) 98.
- [21] D. Earl, B. vonHoldt, STRUCTURE HARVESTER: a website and program for visualizing STRUCTURE output and implementing the Evanno method, *Conserv. Genet. Res.* 4 (2012) 359–361.
- [22] G. Evanno, S. Regnaut, J. Goudet, Detecting the number of clusters of individuals using the software STRUCTURE: a simulation study, *Mol. Ecol.* 14 (2005) 2611–2620.
- [23] N.M. Kopelman, J. Mayzel, M. Jakobsson, et al., Clumpak: a program for identifying clustering modes and packaging population structure inferences across K, *Mol. Ecol. Resour.* 15 (2015) 1179–1191.
- [24] M. Jakobsson, N.A. Rosenberg, CLUMPP: a cluster matching and permutation program for dealing with label switching and multimodality in analysis of population structure, *Bioinformatics* 23 (2007) 1801–1806.
- [25] N.A. Rosenberg, distruct: a program for the graphical display of population structure, *Mol. Ecol. Notes* 4 (2004) 137–138.
- [26] R.C. Team, R: A Language and Environment for Statistical Computing, R Foundation for Statistical Computing, Vienna, Austria, 2014.
- [27] L. Excoffier, H.E. Lischer, Arlequin suite ver 3.5: a new series of programs to perform population genetics analyses under Linux and Windows, *Mol. Ecol. Resour.* 10 (2010) 564–567.
- [28] Y.L. Wei, L. Wei, L. Zhao, et al., A single-tube 27-plex SNP assay for estimating individual ancestry and admixture from three continents, *Int. J. Legal Med.* 130 (2016) 27–37.
- [29] U. Rogalla, E. Rychlicka, M.V. Derenko, et al., Simple and cost-effective 14-loci SNP assay designed for differentiation of European, East Asian and African samples, *Forensic Sci. Int.: Genet.* 14 (2015) 42–49.
- [30] O. Lao, P.M. Vallone, M.D. Coble, et al., Evaluating self-declared ancestry of U.S. Americans with autosomal, Y-chromosomal and mitochondrial DNA, *Human Mutat.* 31 (2010) E1875–E1893.