

Serveur Académique Lausannois SERVAL [serval.unil.ch](http://serval.unil.ch)

## Author Manuscript

Faculty of Biology and Medicine Publication

This paper has been peer-reviewed but does not include the final publisher proof-corrections or journal pagination.

Published in final edited form as:

**Title:** Inferring biogeographic ancestry with compound markers of slow and fast evolving polymorphisms.

**Authors:** Moriot A, Santos C, Freire-Aradas A, Phillips C, Hall D

**Journal:** European journal of human genetics : EJHG

**Year:** 2018 Nov

**Issue:** 26

**Volume:** 11

**Pages:** 1697-1707

**DOI:** 10.1038/s41431-018-0215-2

In the absence of a copyright statement, users should assume that standard copyright protection applies, unless the article contains an explicit statement to the contrary. In case of doubt, contact the journal publisher to verify the copyright status of an article.

1 **Inferring biogeographic ancestry with compound markers of slow and fast**  
2 **evolving polymorphisms**

3

4

5 Amandine Moriot<sup>1¶</sup>, Carla Santos<sup>2¶</sup>, Ana Freire-Aradas<sup>2</sup>, Christopher Phillips<sup>2</sup>, Diana  
6 Hall<sup>1\*</sup>

7

8

9 <sup>1</sup> Unité de Génétique Forensique, Centre Universitaire Romand de Médecine Légale,  
10 Centre Hospitalier Universitaire Vaudois et Université de Lausanne, Lausanne,  
11 Switzerland

12

13 <sup>2</sup> Forensic Genetics Unit, Institute of Forensic Science, University of Santiago de  
14 Compostela, Santiago de Compostela, Spain

15

16 <sup>¶</sup>These authors contributed equally to this work.

17

18

19 <sup>\*</sup>Corresponding author

20 Diana Hall

21 Centre Universitaire Romand de Médecine Légale, Centre Hospitalier Universitaire  
22 Vaudois et Université de Lausanne

23 Ch. de la Vulliette 4

24 1000, Lausanne

25 Switzerland

26

27 Telephone: +41 021 314 7081, FAX: +41 021 314 7090

28 Email: [Diana.Hall@chuv.ch](mailto:Diana.Hall@chuv.ch)

29

30 **Conflict of interest** The authors declare that they have no conflict of interest.

31 **Abstract**

32 Bio-geographic ancestry is an area of considerable interest in the medical genetics,  
33 anthropology and forensics. Although genome-wide panels are ideal as they provide  
34 dense genotyping data, small sets of ancestry informative marker provide a cost-  
35 effective way to investigate genetic ancestry and population structure. Here, we  
36 investigate the performance of a reduced marker set that combine different types of  
37 autosomal markers through haplotype analysis. In particular, recently described DIP-  
38 STR markers should offer the advantage of comprising both, low mutation rate Indels  
39 (DIPs), to study human history over longer time scale; and high mutation rate STRs, to  
40 trace relatively recent demographic events.

41 In this study, we assessed the ability of an initial set of 23 DIP-STRs to distinguish  
42 major population groups using the HGDP-CEPH reference samples. The results  
43 obtained applying the STRUCTURE algorithm show that the discrimination capacity of  
44 the DIP-STRs is comparable to currently used small-scale ancestry informative markers  
45 by approaching seven major demographic groups. Yet, the DIP-STRs show an  
46 improved success rate in assigning individuals to populations of Europe and Middle  
47 East.

48 These data show a remarkable ability of a preliminary set of 23 DIP-STR markers  
49 to infer major biogeographic origins. A novel set of DIP-STRs preselected to contain  
50 ancestry information should lead to further improvements.

51

52 **Key Words**

53 DIP-STR, ancestry inference, HGDP-CEPH, population structure, Indels

## 54 **Introduction**

55 Bio-geographic ancestry inference has largely contributed to controlling for  
56 population structure in disease or trait association studies and to infer human  
57 evolutionary history<sup>1-4</sup>. Sets of small-scale ancestry informative markers (AIM) are  
58 valuable to provide most of the information in a cost-effective manner when sample sets  
59 have not been typed with genome-wide arrays of SNPs or when DNA amount is limited  
60 as in forensic science. Targeted association studies – such as candidate gene studies or  
61 replication studies following up genome-wide scans typically analyse a much smaller  
62 number of markers than genome-wide scans, making it difficult to infer ancestry in  
63 order to correct for stratification. AIMs are also used to select samples for follow-up  
64 studies before genome-wide scans are performed. In criminal investigations or missing  
65 persons identifications, in the absence of any other investigative leads (no database or  
66 suspect match), AIM genotypes obtained from evidential material could indicate the  
67 likely ancestry of the donor, and therefore help direct the course of investigations<sup>5,6</sup>.

68 Numerous AIM panels composed of dozens of autosomal SNPs, Indels or STRs  
69 have been published<sup>7-15</sup>. Briefly, these and other studies have shown that populations  
70 clustering patterns are robust, provided that at least about 60-150 markers are used<sup>16-18</sup>,  
71 or about 30 or fewer if markers are preselected to have a high information content about  
72 ancestry<sup>8,12,15,19</sup>. Most of small AIM sets are able to structure global populations into  
73 five major geographic regions: America, Sub-Saharan Africa, East Asia, Oceania and a  
74 cluster composed of Europe, Central-South Asia and Middle East populations within  
75 Eurasia<sup>7</sup>.

76 Nevertheless, different types of markers offer different population structure  
77 resolution. In the field of molecular evolution, it has long been recognized that markers

78 with relatively low mutation rates (SNP, Alu, Indel) serve as best loci for the analysis of  
79 human history over longer time scales (therefore to provide a biogeographical resolution  
80 at the level of continents), whereas rapidly evolving markers (STR and mtDNA)  
81 provide the greatest resolution over shorter time scales (regional resolution) <sup>20</sup>.  
82 Haplotypes composed of both slow- and fast-evolving loci should combine the benefits  
83 of two types of markers. For example, in the case of linked SNP-STR, the instability of  
84 the STR should generate many alleles, in proportion to the populations' divergence  
85 time, while the more stable flanking SNP should allow greater certainty in tracing the  
86 lineage of each haplotype. In addition, the number of possible haplotypes formed by the  
87 combination of these markers is much greater than the number of alleles at each  
88 individual locus. Such high variability increases the likelihood of rare haplotypes, that  
89 are easily lost during population founding and bottleneck events.

90 The first example of compound genetic markers are the SNP-STRs residing in the  
91 nonrecombining region of the human Y chromosome <sup>20</sup>. However, the uniparental  
92 transmission of the Y chromosome, like the mitochondrial genome, has the limit of  
93 reflecting only a fraction of individual's ancestry with the risk of misinterpreting the  
94 overall ancestry in admixed individuals <sup>6,21</sup>. To overcome such limitation Mountain et  
95 al. <sup>22</sup> proposed autosomal SNP-STR markers. In a pilot study they were able to show  
96 that two such compound markers were sufficient to provide support for the "out of  
97 Africa" human evolutionary hypothesis. Unfortunately, the application of SNP-STR for  
98 ancestry inference found limited success because of the difficulties initially encountered  
99 in genotyping SNPs linked STRs.

100 The current high throughput DNA sequencing technology finally enables genetic  
101 practitioners to consider multiple polymorphisms grouped into haplotypes. Many

102 benchtop DNA sequencing platforms provide today continuous runs of a hundred base  
103 pairs or more on a single DNA molecule that allows to directly inferring the phase of  
104 the multiple markers within a small DNA segment. Minihaplotypes and  
105 microhaplotypes encompassing two to four SNPs spanning less than 10 kb and 200 bp,  
106 respectively, were recently described<sup>23,24</sup>. Haplotype systems based on multiple SNPs  
107 have proven a forensically useful DNA marker for family or lineage inference<sup>24,25</sup> and  
108 in anthropology for population relationships<sup>19,24,26,27</sup>. A recent study showed that  
109 assignment of individuals to candidate populations significantly improves through  
110 combining linked SNPs. For example, incorrectly assigned individuals from empirical  
111 data of French and German populations can be decreased by 73% using haplotypes<sup>28</sup>.

112 In this study, we evaluated the ability of an additional compound genetic marker  
113 recently described named DIP-STR, to serve as useful marker for ancestry inference.  
114 Originally, DIP-STRs were proposed as typing method for the characterization of DNA  
115 mixtures from two individuals present in very different proportions<sup>29-33</sup>. These mixtures  
116 are commonly encountered in forensic investigations during mixed trace analyses; in  
117 cases of peripheral blood DNA microchimerism during pregnancy or induced by solid  
118 organ transplant. Because of PCR amplification bias, the genetic identification of a  
119 DNA that contributes trace amounts to a mixed sample (minor DNA) represents a  
120 tremendous challenge. The DIP-STRs have the innovative feature of combining the  
121 analysis of a DIP (Deletion/Insertion Polymorphism) and a closely linked STR  
122 polymorphism. To target the analysis of the minor DNA contributor of the mixture,  
123 PCR primers overlap the deleted/inserted sequence (S or L) to produce allele-specific  
124 amplifications of one or two minor DIP-STR haplotypes comprising the DIP allele that  
125 is not shared with the major DNA. Allele-specific amplifications of DIP-STR

126 haplotypes enable the characterization a minor DNA in the presence of more than 1000-  
127 fold excess of a major DNA contributor.

128 Here we present the results of 23 DIP-STR markers typed in the global HGDP-  
129 CEPH reference samples. A comparative HGDP-CEPH analysis is shown using existing  
130 AIM SNP and AIM Indel marker sets.

131

## 132 **Materials and Methods**

### 133 **DIP-STR marker selection**

134 DIP-STRs were selected from three groups: 14 from the first two sets of DIP-STRs  
135 developed for the analysis of unbalanced DNA mixtures (Table 1). Three markers from  
136 the first set were not included because of the large amplicon size (over 600 bp) and one  
137 appeared not interesting based on balanced allele frequencies of the DIP across the  
138 major geographic regions as reported by data from the 1000 Genomes project. Five  
139 additional DIP-STRs were rescued from the preliminary marker selection of the study  
140 <sup>33</sup>. In this previous study, several markers were initially eliminated because of the  
141 negative results of the first amplification assay; here they produced positive results with  
142 newly designed PCR primers. Among these candidates, we selected those markers with  
143 larger DIP allele frequencies differences between continents using the 1000 Genomes  
144 dataset. Finally, four DIP-STRs were selected considering previously described sets of  
145 AIM Indels <sup>12-14,34-39</sup>. DIPs of interest were located near a repeated sequence, showed  
146 similar genotyping conditions and skewed global allele frequencies among continents.

### 147 **Population samples**

148 The CEPH Human Genome Diversity panel (CEPH-HGDP) contains 1,064  
149 individuals from African, European, North African/Middle Eastern, Central-South

150 Asian, East Asian, Native American and Oceanian populations<sup>40</sup>. For all data analyses  
151 purposes we considered only 952 individuals (H952 subset) after exclusion of  
152 duplicates, first- and second-degree relatives<sup>41</sup>. Populations were combined into  
153 continental-based groups which have been previously established<sup>16</sup> with the following  
154 composite populations, sample sizes and labels: 6 African (105 AFR), 8 European (158  
155 EUR), 4 North African/Middle Eastern (162 ME), 9 Central-South Asian (202 CSA), 17  
156 East Asian (230 EAS), 2 Oceanian (28 OCE) and 5 Native American (64 NAM). For all  
157 subjects, blood cell samples were obtained according to protocols and informed-consent  
158 procedures approved by institutional review boards, and were labelled with an  
159 anonymous code number linked only to demographic information and sex.

#### 160 **Marker typing**

161 DNA samples were genotyped for the 23 DIP-STRs and also for the 23 DIPs alone  
162 in order to further control for allele-specificity of the S- and L-DIP-STR amplifications.  
163 PCR reactions were performed in 20  $\mu$ L final volume. This contained 1 $\times$  PCR Buffer  
164 containing 1.5 mM MgCl<sub>2</sub> (Thermo Fisher Écublens, Switzerland), 250  $\mu$ M dNTP  
165 (Thermo Fisher Écublens, Switzerland), 1.2 U AmpliTaq Gold DNA Polymerase  
166 (Thermo Fisher Écublens, Switzerland) and 0.5 ng DNA. Primers' sequences, quantities  
167 and multiplexes are indicated in Supplementary Table 1 and Supplementary Table 2.  
168 PCR thermal cycling conditions were: 5 min at 95°C, 1 min at 94°C, 1 min at annealing  
169 temperature specific to the markers set to be genotyped, 1 min at 72°C for a number of  
170 PCR cycles that also varied across multiplex and a final extension of 30 min at 72°C.  
171 Annealing temperatures and number of cycles are indicated in Supplementary Table 2.

172 PCR fragments were separated by capillary electrophoresis after adding 1  $\mu$ L PCR  
173 amplicon to 8.5  $\mu$ L deionized formamide HI-DI (Thermo Fisher Écublens, Switzerland)



174 and to 0.5  $\mu$ L 600 LIZ size standard (Thermo Fisher Écublens, Switzerland). Capillary  
175 electrophoresis was performed using an ABI PRISM 3130xl Genetic Analyzer (Thermo  
176 Fisher Écublens, Switzerland) according to the manufacturer's instruction and analyzed  
177 using the GeneMapper<sup>®</sup> ID v3.2.1 software (Thermo Fisher Écublens, Switzerland),  
178 with a minimum peak height threshold of 50 RFU. The commercial DNA CEPH 1347-  
179 02 (Thermo Fisher Écublens, Switzerland) was added to two empty positions in each  
180 PCR plate as positive control of amplification and internal standard for allele calls, at  
181 least one empty well per plate was used a negative control of amplification. Markers  
182 information and genotypes are available at the HGDP-CEPH database  
183 ([http://www.cephb.fr/en/hgdp\\_panel.php#basedonnees](http://www.cephb.fr/en/hgdp_panel.php#basedonnees)).

#### 184 **Data analysis**

185 In order to assess the predictive value of the 23 DIP-STR marker set, cluster  
186 analysis was performed using the STRUCTURE program version 2.3.4<sup>42,43</sup>. Runs  
187 consisted of 50 000 Markov Chain steps after a burn-in of length 50 000 with ten  
188 replicates for K=5 and K=7, using the *admixture* ancestry model and *correlated allele*  
189 *frequencies*. The DIP-STR allele names indicated the DIP variant, either S (deletion) or  
190 L (insertion) and the STR allele expressed in DNA fragment size, this last corresponds  
191 to the STR allele size that one would obtain using primers located around the repeated  
192 sequence (allele names can be found in Figure 1 and Supplementary Figure 1). For the  
193 STRUCTURE analysis DIP-STR alleles were recoded with serial numbers from 1 to *n*  
194 starting from the shortest L-DIP-STR to the longest S-DIP-STR. The bar plot was  
195 prepared using CLUMPAK (<http://clumpak.tau.ac.il/>).

196 To estimate the success rate in assigning the population of origin of each individual  
197 considering the genetic information of 23 DIP-STRs the likelihood-based approach

198 implemented by Phillips et al.<sup>44</sup> in the SNIPPER App Suite  
199 (<http://mathgene.usc.es/snipper/>) was used. With this tool, a single, unknown genetic  
200 profile can be compared to a set of reference populations, the “training set”. The  
201 software calculates individual maximum likelihoods estimates for the inclusion of the  
202 unknown sample into each reference population. A cross-validation has been performed  
203 using the frequency-based classifier in *Snipper* app suite. Each sample was tested in  
204 turn as unknown sample against the training set containing all the remaining samples.  
205 No marker deviated from the Hardy–Weinberg equilibrium and linkage disequilibrium.  
206 For comparative analysis, both STRUCTURE and *Snipper* based analyses were  
207 repeated for two previously published AIM markers sets, these include a 34 AIM SNP  
208 set<sup>44,45</sup> and a 46 AIM Indel set<sup>37</sup>.

209

## 210 **Results**

### 211 **Patterns of DIP-STR variability**

212 The number of DIP-STR markers studied here is 23 (Table 1). Of the 368 alleles  
213 present more than once in the dataset, 27.5% appeared in all major regions represented,  
214 these are Africa, Europe, Middle East, Central-South Asia, East Asia, Oceania and  
215 Native America. Those exclusive to one region were 8.7%; region-specific alleles  
216 showed a median relative frequency of 1.4% in their region of occurrence.

217 As previously observed, Africa was the most variable region<sup>46,47</sup>. Number of  
218 alleles, and mean number of private alleles followed a common trend: they were highest  
219 in the African samples, were somewhat lower in Europeans and East Asians, and were  
220 lowest in Amerindians and Oceanians. This was especially true for 13 markers  
221 (Supplementary Figure 1). Note that Oceanians population size is about half the Native

222 Americans and between one fourth and one eighth of the other major geographic regions.  
223 In Figure 1a is showed the allele frequency distribution of marker MID1013-D5S490  
224 representative of the pattern described above. Conversely, seven markers showed a  
225 higher number of alleles outside Africa. For these markers the DIP variant is not  
226 polymorphic in Africa therefore all the DIP-STR haplotypes containing the DIP allele  
227 that probably appeared outside Africa, are missing in this group. Here as well, the  
228 number of observed alleles decreases with increasing distance from Eurasia. An  
229 example is reported in Figure 1b, marker rs112604544-STR.

### 230 **Population clustering analysis**

231 We assessed the ability of the DIP-STR genotypes to cluster the global HGDP-  
232 CEPH reference population samples applying the widely used STRUCTURE Bayesian  
233 grouping algorithm. Previous studies including several types of ancestry informative  
234 markers suggest that when limited numbers of markers are analyzed it is appropriate to  
235 aim to assign individuals to five major population groups in the first instance, these are  
236 Africans, Europeans, East Asians, Oceanians and Native Americans. To obtain the  
237 clearest pattern of group membership from STRUCTURE, the study population  
238 complexity was reduced by excluding geographically close populations such as Middle  
239 East and Central-South Asia. These are known to show higher misclassification rates  
240 since they occupy regions in the middle of a continuum of variability. The  
241 STRUCTURE results with the highest likelihood at K= 5 shows a clear pattern of five  
242 clusters corresponding to major geographical regions (Figure 2a). These results are  
243 comparable to currently used AIM markers sets comprising 34 AIM SNPs (Figure 2b)  
244 or 46 AIM Indels (Figure 2c). The results at K=7 including the complete HGDP-CEPH  
245 dataset show that although a “private” ancestry component is present in Middle East and

246 Central South Asia, clustering patterns are less distinct and Europe loses clear definition  
247 incorporating the additional inferred cluster as partial degrees of ancestry (Figure 3a).  
248 These data are comparable to the results obtained using the 34 AIM SNP set (Figure  
249 3b), yet they show an improved clustering capacity of the European and Middle East  
250 groups when compared to the results obtained with the 46 AIM Indels (Figure 3c).

### 251 **Individual ancestry assignment analysis**

252 *Snipper* cross-validation classification success values for five and seven HGDP-  
253 CEPH reference samples obtained using the 23 DIP-STRs considered in this study are  
254 in agreement with the clustering obtained using the STRUCTURE classification  
255 algorithm. For the five group analysis the classification success rates are higher than  
256 99% for all populations. These results are similar to those obtained with 34 AIM SNPs  
257 and 46 AIM Indels, with a somewhat better classification success rate for the East Asian  
258 group (Table 2). For the seven group analysis, samples with Eurasians origin were more  
259 difficult to classify, yet about 96% of Europeans 86% of Middle Easterns are correctly  
260 classified (Table 3). These values are lower when using 34 AIM SNPs (70% and 62%)  
261 and 46 AIM Indels (40% and 46%).

### 262 **Separate STRUCTURE analysis of the 23 DIP set and the 23 STR set**

263 The independent contribution to the clustering of the HGDP-CEPH populations of  
264 the 23 DIPs and the 23 STRs composing the haplotypes analyzed before, was  
265 investigated (Figure 4). This is possible because haplotypes are named based on the  
266 respective DIP and STR comprising alleles. The data at K=5 after excluding  
267 geographically close populations of Middle East and Central South Asia, show similar  
268 results for the DIPs and the STRs. The populations of Africa are distinguished from the  
269 other main geographical regions while part of Native Americans are clustered with the

270 Oceanians. The populations of Eurasia form a unique cluster with a high degree of noise  
271 and some indication of admixture with the Native Americans, especially when using the  
272 DIPs.

273

## 274 **Discussion**

275 The aim of this study was to explore the contribution to small-scale marker sets  
276 based biogeographic inference of DIP-STR haplotype markers. We hypothesized that  
277 the compound nature of DIP-STRs may represent an attractive feature not only for DNA  
278 mixture resolutions as we originally proposed, but also for studies of population  
279 structure. Briefly, each DIP-STR haplotype is provided with a slow and fast mutating  
280 variant that should confer higher biogeographic information compared to the use of  
281 each single type of polymorphism. Here, the survey of the HGDP-CEPH global  
282 reference samples with an initial set of 23 DIP-STRs allowed us to determine the  
283 relative value of our markers with respect to other validated AIM marker sets.

284 Overall, the clustering patterns observed with the STRUCTURE algorithm and the  
285 population assignments obtained with the Snipper program were in good correlation  
286 with each other as well as with worldwide population structure<sup>8,12,15,19,37</sup>. The prediction  
287 of bio-geographic ancestry was achieved for five major populations using a smaller set  
288 of markers (23 DIP-STRs *versus* 34 to 46 markers of validated AIM SNP and AIM  
289 Indel marker sets). Moreover, the seven group analysis shows that with these markers  
290 the global pattern is approaching seven clusters and some additional selected DIP-STRs  
291 may produce stable and reproducible distinction of the populations of Middle East and  
292 Central-South Asia within Eurasia. The discrimination of these groups still represents a  
293 challenge for most available AIM panels<sup>7</sup>, as also indicated by our comparative

294 analysis using 34 AIM SNPs and 46 AIM Indels. However, according to Snipper cross-  
295 validation success values the DIP-STRs show an improved success rate in assigning  
296 individuals to populations of Europe and Middle East. It should be noted that the results  
297 showed here, were obtained with DIP-STR markers not selected for distinguishing  
298 continental or intra-continental structures. Yet, these results are comparable to those  
299 obtained with larger sets of SNPs and Indels produced after refined selections using  
300 allele frequency data from worldwide population surveys. Finally, the fact that the  
301 obtained resolution is not due to single composing markers (either DIPs or STRs) with  
302 no advantage of considering the combined information, is also showed by the  
303 corresponding cluster analysis using the DIP marker set and the STR marker set,  
304 separately.

305 It is generally accepted that identifying SNPs that can distinguish among four or  
306 five continental groups of populations is not difficult <sup>48</sup>. Many different non-  
307 overlapping sets of SNPs capable of inferring continental ancestry were described (see  
308 references in the Introduction). However, previous studies showed the difficulty of  
309 developing a single panel comprising a limited number of SNPs or Indels that is capable  
310 of differentiating populations both globally and within regions. Note that, the HGDP  
311 collection has the drawback of not sampling densely within each geographic region;  
312 however, the use of a large number of markers showed the possibility of a sub-  
313 continental estimation of bio-geographic ancestry.

314 To reach such a result with DIP-STR markers; first of all, they need to be selected  
315 for ancestry purposes. Large whole genome sequencing data will provide the basis for  
316 annotating novel DIP-STRs and estimating allele frequencies across populations from  
317 long read data; alternatively, additional marker selecting criteria can be formulated

318 based on annotated lists of unphased DIP and STR markers. For example: a tentative  
319 approach could be a marker selection based on DIP skewed allele frequencies among  
320 large geographical regions (this would include a mixture of old and young DIP  
321 mutations) linked to highly polymorphic STRs.

322 Finally, the value of a small ancestry panel is also measured based on the capacity  
323 of analyzing trace amount of DNA with a cost-effective method easy to implement in  
324 routine laboratories. Most of the DIP-STRs genotyped here were analyzed in short  
325 fragments by standard PCR and capillary electrophoresis analysis and were validated  
326 for forensic contact traces <sup>32</sup>. Yet, several multiplex were used to produce the data. This  
327 is because: on one side, markers were progressively added to the study; on the other  
328 side, the typing of 23 DIP-STRs conserving S- and the L-specific amplifications require  
329 some developmental efforts. We plan to work on a compact multiplex method once  
330 final DIP-STRs and minimum number of marker are identified.

331 As stated before, significant progress has been made in the research of ancestry  
332 marker sets; therefore the effort that the development of AIM-DIP-STRs requires is  
333 justified only by an advantage over existing methods. Besides the hope of providing  
334 better (more refined and robust) resolution of bio-geographic ancestry that can be  
335 assessed with simple and cost-effective methods, these markers have the attractive  
336 feature of combining ancestry estimates and unbalanced DNA mixture resolution,  
337 especially useful in forensics. Although, ancestry and identity DIP-STRs are may not to  
338 be the same marker sets, one can apply DIP-STR markers to first, detect the minor DNA  
339 of the trace sample; for example, the DNA of a man who sexually assaulted a woman,  
340 retrieved from a gynecological sample. Second, if the DIP-STR profile does not find a  
341 match in the list of suspects, additional ancestry DIP-STRs can be used as intelligence

342 tool to further guide the investigation. A caveat to this is that the ancestry inference in  
343 the context of mixture resolution would produce results for a reduced number of  
344 markers (those that contain DIP alleles unique to the minor DNA) and therefore the  
345 validation of redundant DIP-STRs may be required to have sufficient markers for  
346 ancestry resolution.

347 Our study provides another example of forensic markers capable of providing  
348 substantial biogeographic information. The fact that individual identifiability and  
349 population identifiability may be correlated has been observed before for forensic  
350 microsatellite markers <sup>49</sup> and it has been extensively studied by Algee-Hewitt and  
351 colleagues <sup>50</sup>. These authors identified in the high level of polymorphism the potential  
352 for high population identifiability. In the case of DIP-STRs the number of possible  
353 haplotypes formed by the combination of two markers is indeed much greater than the  
354 number of alleles at each individual locus, but the combination of slow- and fast-  
355 evolving loci is also expected to play a role.

356 In conclusion, our study demonstrates that DIP-STRs can provide a useful adjunct  
357 to AIM panels aiming at providing a better resolution and, in forensic science, at  
358 combining the function of mixture deconvolution to bio-geographic ancestry estimates.  
359 We showed here that continental regions can be readily distinguished, while more  
360 markers are necessary to improve the classification of closely related populations such  
361 as Eurasians.

362 This represents an exceptional result for the first evaluation of this set of marker  
363 that was not selected to contain information about ancestry. The research described here  
364 should be considered a pilot study to encourage more efforts on DIP-STRs marker



365 discovery and multiplex development. Finally, these data provide references of global  
366 DIP-STR allele frequencies including the set proposed for forensic casework analysis.

367

368 **Acknowledgments** This work was financially supported by funding from the  
369 University Center of Legal Medicine of the University Hospital of Lausanne and the  
370 Faculty of Biology and Medicine of the University of Lausanne (Pro-Femmes  
371 fellowship to Diana Hall).

372

373 **Conflict of interest** The authors declare that they have no conflict of interest.

374

#### 375 **References**

- 376 1. Hellenthal G, Busby GB, Band G *et al*: A genetic atlas of human admixture  
377 history. *Science* 2014; **343**: 747-751.  
378
- 379 2. Li JZ, Absher DM, Tang H *et al*: Worldwide human relationships inferred from  
380 genome-wide patterns of variation. *Science* 2008; **319**: 1100-1104.  
381
- 382 3. Tishkoff SA, Kidd KK: Implications of biogeography of human populations for  
383 'race' and medicine. *Nat Genet* 2004; **36**: S21-27.  
384
- 385 4. Elhaik E, Tatarinova T, Chebotarev D *et al*: Geographic population structure  
386 analysis of worldwide human populations infers their biogeographical origins.  
387 *Nat Commun* 2014; **5**: 3513.  
388
- 389 5. Kayser M, de Knijff P: Improving human forensics through advances in  
390 genetics, genomics and molecular biology. *Nature reviews Genetics* 2011; **12**:  
391 179-192.  
392
- 393 6. Phillips C, Prieto L, Fondevila M *et al*: Ancestry analysis in the 11-M Madrid  
394 bomb attack investigation. *PLoS One* 2009; **4**: e6583.  
395
- 396 7. Phillips C: Forensic genetic analysis of bio-geographical ancestry. *Forensic  
397 science international Genetics* 2015; **18**: 49-65.  
398
- 399 8. de la Puente M, Santos C, Fondevila M *et al*: The Global AIMs Nano set: A 31-  
400 plex SNaPshot assay of ancestry-informative SNPs. *Forensic science  
401 international Genetics* 2016; **22**: 81-88.

- 402  
403 9. Pakstis AJ, Kang L, Liu L *et al*: Increasing the reference populations for the 55  
404 AISNP panel: the need and benefits. *Int J Legal Med* 2017.  
405
- 406 10. Phillips C, Parson W, Lundsberg B *et al*: Building a forensic ancestry panel  
407 from the ground up: The EUROFORGEN Global AIM-SNP set. *Forensic science international Genetics* 2014; **11**: 13-25.  
408  
409
- 410 11. Lao O, van Duijn K, Kersbergen P, de Knijff P, Kayser M: Proportioning  
411 whole-genome single-nucleotide-polymorphism diversity for the identification  
412 of geographic population structure and genetic ancestry. *Am J Hum Genet* 2006;  
413 **78**: 680-690.  
414
- 415 12. Santos C, Phillips C, Oldoni F *et al*: Completion of a worldwide reference panel  
416 of samples for an ancestry informative Indel assay. *Forensic science international Genetics* 2015; **17**: 75-80.  
417  
418
- 419 13. Zaumsegel D, Rothschild MA, Schneider PM: A 21 marker insertion deletion  
420 polymorphism panel to study biogeographic ancestry. *Forensic science international Genetics* 2013; **7**: 305-312.  
421  
422
- 423 14. Yang N, Li HZ, Criswell LA *et al*: Examination of ancestry and ethnic  
424 affiliation using highly informative diallelic DNA markers: application to  
425 diverse and admixed populations and implications for clinical epidemiology and  
426 forensic medicine. *Human Genetics* 2005; **118**: 382-392.  
427
- 428 15. Londin ER, Keller MA, Maista C *et al*: CoAIMs: a cost-effective panel of  
429 ancestry informative markers for determining continental origins. *PLoS One*  
430 2010; **5**: e13443.  
431
- 432 16. Rosenberg NA, Pritchard JK, Weber JL *et al*: Genetic structure of human  
433 populations. *Science* 2002; **298**: 2381-2385.  
434
- 435 17. Halder I, Shriver M, Thomas M, Fernandez JR, Frudakis T: A panel of ancestry  
436 informative markers for estimating individual biogeographical ancestry and  
437 admixture from four continents: Utility and applications. *Human Mutation* 2008;  
438 **29**: 648-658.  
439
- 440 18. Kosoy R, Nassir R, Tian C *et al*: Ancestry informative marker sets for  
441 determining continental origin and admixture proportions in common  
442 populations in America. *Hum Mutat* 2009; **30**: 69-78.  
443
- 444 19. Kidd KK, Pakstis AJ, Speed WC *et al*: Current sequencing technology makes  
445 microhaplotypes a powerful new type of genetic marker for forensics. *Forensic science international Genetics* 2014; **12**: 215-224.  
446  
447

- 448 20. de Knijff P: Messages through bottlenecks: On the combined use of slow and  
449 fast evolving polymorphic markers on the human Y chromosome. *American*  
450 *Journal of Human Genetics* 2000; **67**: 1055-1061.
- 451 21. King TE, Parkin EJ, Swinfield G *et al*: Africans in Yorkshire? The deepest-  
452 rooting clade of the Y phylogeny within an English genealogy. *Eur J Hum*  
453 *Genet* 2007; **15**: 288-293.
- 454 22. Mountain JL, Knight A, Jobin M *et al*: SNPSTRs: Empirically derived, rapidly  
455 typed, autosomal Haplotypes for inference of population history and mutational  
456 processes. *Genome Research* 2002; **12**: 1766-1772.
- 457 23. Kidd KK, Speed WC, Pakstis AJ *et al*: Evaluating 130 microhaplotypes across a  
458 global set of 83 populations. *Forensic science international Genetics* 2017; **29**:  
459 29-37.
- 460 24. Pakstis AJ, Fang R, Furtado MR, Kidd JR, Kidd KK: Mini-haplotypes as lineage  
461 informative SNPs and ancestry inference SNPs. *Eur J Hum Genet* 2012; **20**:  
462 1148-1154.
- 463 25. Kidd KK, Pakstis AJ, Speed WC *et al*: Microhaplotype loci are a powerful new  
464 type of forensic marker. *Forensic Science International: Genetics Supplement*  
465 *Series* 2013; **4**: e123-e124.
- 466 26. Kidd JR, Friedlaender F, Pakstis AJ *et al*: Single nucleotide polymorphisms and  
467 haplotypes in Native American populations. *Am J Phys Anthropol* 2011; **146**:  
468 495-502.
- 469 27. Schlebusch CM, Soodyall H: Extensive population structure in San, Khoe, and  
470 mixed ancestry populations from southern Africa revealed by 44 short 5-SNP  
471 haplotypes. *Hum Biol* 2012; **84**: 695-724.
- 472 28. Gattepaille LM, Jakobsson M: Combining markers into haplotypes can improve  
473 population structure inference. *Genetics* 2012; **190**: 159-174.
- 474 29. Castella V, Gervaix J, Hall D: DIP-STR: highly sensitive markers for the  
475 analysis of unbalanced genomic mixtures. *Hum Mutat* 2013; **34**: 644-654.
- 476 30. Cereda G, Biedermann A, Hall D, Taroni F: An investigation of the potential of  
477 DIP-STR markers for DNA mixture analyses. *Forensic science international*  
478 *Genetics* 2014; **11**: 229-240.
- 479 31. Cereda G, Biedermann A, Hall D, Taroni F: Object-oriented Bayesian networks  
480 for evaluating DIP-STR profiling results from unbalanced DNA mixtures.  
481 *Forensic science international Genetics* 2014; **8**: 159-169.
- 482  
483  
484  
485  
486  
487  
488  
489  
490  
491  
492  
493

- 494 32. Oldoni F, Castella V, Grosjean F, Hall D: Sensitive DIP-STR markers for the  
495 analysis of unbalanced mixtures from "touch" DNA samples. *Forensic science*  
496 *international Genetics* 2017; **28**: 111-117.  
497
- 498 33. Oldoni F, Castella V, Hall D: A novel set of DIP-STR markers for improved  
499 analysis of challenging DNA mixtures. *Forensic science international Genetics*  
500 2015; **19**: 156-164.  
501
- 502 34. Bastos-Rodrigues L, Pimenta JR, Pena SDJ: The genetic structure of human  
503 populations studied through short insertion-deletion polymorphisms. *Annals of*  
504 *Human Genetics* 2006; **70**: 658-665.  
505
- 506 35. Francez PA, Ribeiro-Rodrigues EM, dos Santos SE: Allelic frequencies and  
507 statistical data obtained from 48 AIM INDEL loci in an admixed population  
508 from the Brazilian Amazon. *Forensic science international Genetics* 2012; **6**:  
509 132-135.  
510
- 511 36. Pereira R, Phillips C, Alves C, Amorim A, Carracedo A, Gusmao L: A new  
512 multiplex for human identification using insertion/deletion polymorphisms.  
513 *Electrophoresis* 2009; **30**: 3682-3690.  
514
- 515 37. Pereira R, Phillips C, Pinto N *et al*: Straightforward inference of ancestry and  
516 admixture proportions through ancestry-informative insertion deletion  
517 multiplexing. *PLoS One* 2012; **7**: e29684.  
518
- 519 38. Rosenberg NA, Mahajan S, Ramachandran S, Zhao CF, Pritchard JK, Feldman  
520 MW: Clines, clusters, and the effect of study design on the inference of human  
521 population structure. *Plos Genetics* 2005; **1**: 660-671.  
522
- 523 39. Santos NPC, Ribeiro-Rodrigues EM, Ribeiro-dos-Santos AKC *et al*: Assessing  
524 Individual Interethnic Admixture and Population Substructure Using a 48-  
525 Insertion-Deletion (INSEL) Ancestry-Informative Marker (AIM) Panel. *Human*  
526 *Mutation* 2010; **31**: 184-190.  
527
- 528 40. Cann HM, de Toma C, Cazes L *et al*: A human genome diversity cell line panel.  
529 *Science* 2002; **296**: 261-262.  
530
- 531 41. Rosenberg NA: Standardized subsets of the HGDP-CEPH Human Genome  
532 Diversity Cell Line Panel, accounting for atypical and duplicated samples and  
533 pairs of close relatives. *Ann Hum Genet* 2006; **70**: 841-847.  
534
- 535 42. Falush D, Stephens M, Pritchard JK: Inference of population structure using  
536 multilocus genotype data: Linked loci and correlated allele frequencies. *Genetics*  
537 2003; **164**: 1567-1587.  
538
- 539 43. Pritchard JK, Stephens M, Donnelly P: Inference of population structure using  
540 multilocus genotype data. *Genetics* 2000; **155**: 945-959.  
541

- 542 44. Phillips C, Salas A, Sanchez JJ *et al*: Inferring ancestral origin using a single  
543 multiplex assay of ancestry-informative marker SNPs. *Forensic Science*  
544 *International-Genetics* 2007; **1**: 273-280.  
545
- 546 45. Fondevila M, Phillips C, Santos C *et al*: Revision of the SNPforID 34-plex  
547 forensic ancestry test: Assay enhancements, standard reference sample  
548 genotypes and extended population studies. *Forensic science international*  
549 *Genetics* 2013; **7**: 63-74.  
550
- 551 46. Bowcock AM, Ruiz-Linares A, Tomfohrde J, Minch E, Kidd JR, Cavalli-Sforza  
552 LL: High resolution of human evolutionary trees with polymorphic  
553 microsatellites. *Nature* 1994; **368**: 455-457.  
554
- 555 47. Calafell F, Shuster A, Speed WC, Kidd JR, Kidd KK: Short tandem repeat  
556 polymorphism evolution in humans. *Eur J Hum Genet* 1998; **6**: 38-49.  
557
- 558 48. Soundararajan U, Yun L, Shi M, Kidd KK: Minimal SNP overlap among  
559 multiple panels of ancestry informative markers argues for more international  
560 collaboration. *Forensic science international Genetics* 2016; **23**: 25-32.  
561
- 562 49. Phillips C, Fernandez-Formoso L, Garcia-Magarinos M *et al*: Analysis of global  
563 variability in 15 established and 5 new European Standard Set (ESS) STRs using  
564 the CEPH human genome diversity panel. *Forensic science international*  
565 *Genetics* 2011; **5**: 155-169.  
566
- 567 50. Algee-Hewitt BF, Edge MD, Kim J, Li JZ, Rosenberg NA: Individual  
568 Identifiability Predicts Population Identifiability in Forensic Microsatellite  
569 Markers. *Curr Biol* 2016; **26**: 935-942.  
570  
571
- 572

573 **Titles and legends to the figures**

574 **Figure 1** Allele frequency distributions of two representative DIP-STR markers  
575 estimated for the seven groups of Africa (AFR), Europe (EUR), Middle East (ME),  
576 Central-South Asia (CSA), East Asians (EAS), Oceania (OCE) and Native America  
577 (NAM). (a) MID1013-D4S490. (b) rs112604544-STR

578

579 **Figure 2** STRUCTURE analyses of the HGDP-CEPH reference populations at  $K = 5$ ,  
580 after exclusion of Middle Eastern and Central-South Asian groups. Analyses were  
581 computed using the *admixture* ancestry model and *correlated allele frequencies*. (a)  
582 Results obtained using 23 DIP-STRs. (b) 34 AIM SNPs. (c) 46 AIM Indels.

583

584 **Figure 3** STRUCTURE analyses of the HGDP-CEPH reference populations at  $K = 7$ .  
585 Analyses were computed using the *admixture* ancestry model and *correlated allele*  
586 *frequencies*. (a) Results obtained using 23 DIP-STRs. (b) 34 AIM SNPs. (c) 46 AIM  
587 Indels.

588

589 **Figure 4** STRUCTURE analyses of the HGDP-CEPH reference populations at  $K = 5$ ,  
590 after exclusion of Middle Eastern and Central-South Asian groups using 23 markers. (a)  
591 Results obtained using only the 23 DIP genotypes of the DIP-STRs. (b) Results  
592 obtained using only the 23 STRs genotypes of the DIP-STRs.

593

594 **Supporting information**

595 **Supplementary Table 1** DIP-STR primers

596 **Supplementary Table 2** Marker information, PCR primer concentration, number of

597 cycles and multiplex groups

598 **Supplementary Table 3** STRUCTURE individual ancestry proportions of HGDP-

599 CEPH individuals analyzed for three marker sets: 23 DIP-STRs, 34 AIM SNPs and 46

600 AIM Indels at K=5 and K=7

601 **Supplementary Figure 1** DIP-STRs haplotype frequencies for seven HGDP-CEPH

602 major population groups

**Table 1 DIP-STR marker list**

DIP-STR	Chr. Loc.	DIP S/L sequence	STR repeat	DIP-STR size (bases)	Reference
MID1013 <sup>a</sup> -D5S490	5q23.2	-/CCAG	GT	299-345	Castella et al. 2013
rs11277790-D10S530	10q25.1	-/TCCAAC	GT	340-371	Castella et al. 2013
rs60194384-D15S1514	15q26.2	-/TCTTAA	TATC	283-325	Castella et al. 2013
rs66679498-D2S342	2q32.3	-/CCAAC	CA	331-359	Castella et al. 2013
rs35032587-STR	15q26.1	-/TATT	AGAT	239-271	Oldoni et al. 2015
rs142543564-STR	2q34	-/TACT	ATAA	210-238	Oldoni et al. 2015
rs34212659-STR	7p14.1	-/AGG	TGAA	182-199	Oldoni et al. 2015
rs112604544-STR	1q25.3	-/TTTAA	TTCC	134-204	Oldoni et al. 2015
rs111478323-STR	2p25.3	-/GAGA	TTTA	229-265	Oldoni et al. 2015
rs146332920-STR	9q31.3	-/AGG	TAAA	179-207	Oldoni et al. 2015
rs71070706-STR	1p12	-/TGT	AAAG	212-264	Oldoni et al. 2015
rs72406828-STR	4q21.3	-/ATTG	AATTT	178-250	Oldoni et al. 2015
rs145423446-STR	16p13.2	-/AGTC	GATA	230-256	Oldoni et al. 2015
rs2308142-STR	20p13	-/ATT	TTA	205-238	Oldoni et al. 2015
rs71725104-STR	13q31.3	-/ATAG	AAAT	211-235	
rs72534187-STR	5p13.1	-/ACAGGCC	ATAG	208-236	
rs139592446-STR	2q24.2	-/ACTTAGTC	CATC	154-174	
rs36194161-STR	2q32.1	-/CTC	TTTA	138-178	
rs138331044-STR	1p12	-/CATATGC	AGAT	266-302	
MID473 <sup>a</sup> -STR	6q16.1	-/TTACATTT	AGGA	179-227	Francez et al. 2012 <sup>b</sup> Rosemberg et al. 2005 <sup>b</sup> Santos et al. 2010 <sup>b</sup> Santos et al. 2015 <sup>b</sup>
MID2538 <sup>a</sup> -STR	15q25.3	-/TGTT	AC	299-311	Santos et al. 2015 <sup>b</sup>
MID2824 <sup>a</sup> -STR	11p13	-/AGGACT	AAAC	197-222	Zaumsegel et al. 2013 <sup>b</sup>
MID73 <sup>a</sup> -STR	22q12.3	-/GAA	CCACT	362-442	Rosemberg et al. 2005 <sup>b</sup>

<sup>a</sup> Marker name is from the Marshfield database and corresponds to rs1611095, rs140762, rs3054057, rs11278940, rs16365, respectively

<sup>b</sup> References for the DIP marker only



**Table 2 *Snipper* cross-validation classification success values for five HGDP-CEPH major population groups using 23 DIP-STRs, 34 AIM SNPs and 46 AIM Indels.**

Population of origin		Population assignment				
		AFR	EUR	EAS	OCE	NAM
23 DIP-STRs						
	AFR	<b>100.00%</b>	-	-	-	-
	EUR	-	<b>99.37%</b>	0.63%	-	-
	EAS	-	0.43%	<b>99.57%</b>	-	-
	OCE	-	-	-	<b>100.00%</b>	-
	NAM	-	-	-	-	<b>100.00%</b>
34 AIM SNPs						
	AFR	<b>100.00%</b>	-	-	-	-
	EUR	-	<b>99.37%</b>	-	-	0.63%
	EAS	-	0.43%	<b>94.71%</b>	0.44%	4.85%
	OCE	-	-	-	<b>100.00%</b>	-
	NAM	-	-	-	-	<b>100.00%</b>
46 AIM Indels						
	AFR	<b>100.00%</b>	-	-	-	-
	EUR	-	<b>100%</b>	-	-	-
	EAS	-	0.44%	<b>96.07%</b>	2.62%	0.87%
	OCE	-	-	-	<b>100.00%</b>	-
	NAM	-	-	-	-	<b>100.00%</b>

The most likely major component of ancestry was considered independently of the admixture proportions

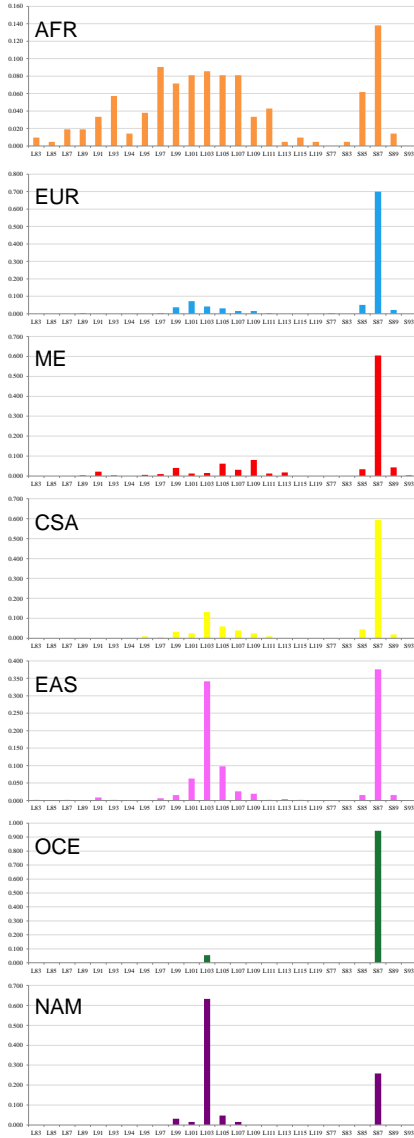
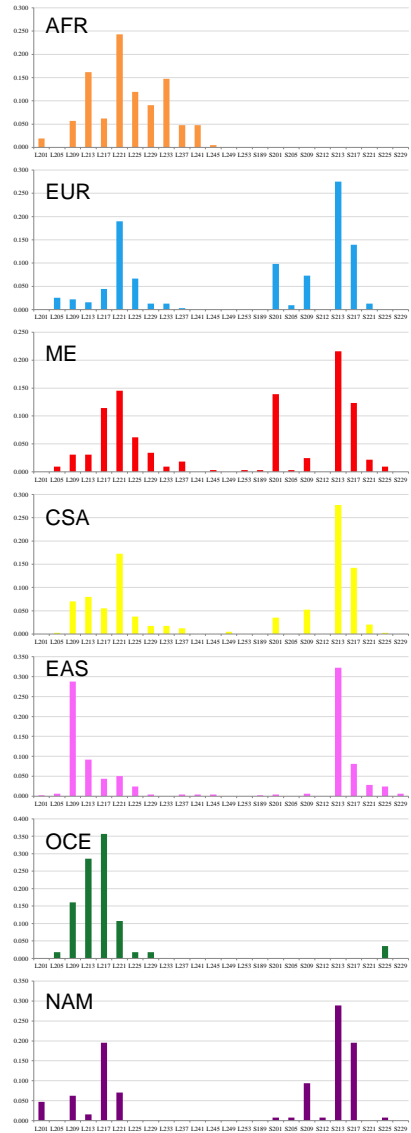
Populations are AFR-Africa, EUR-Europe, EAS-East Asia, OCE-Oceania and NAM- Native America

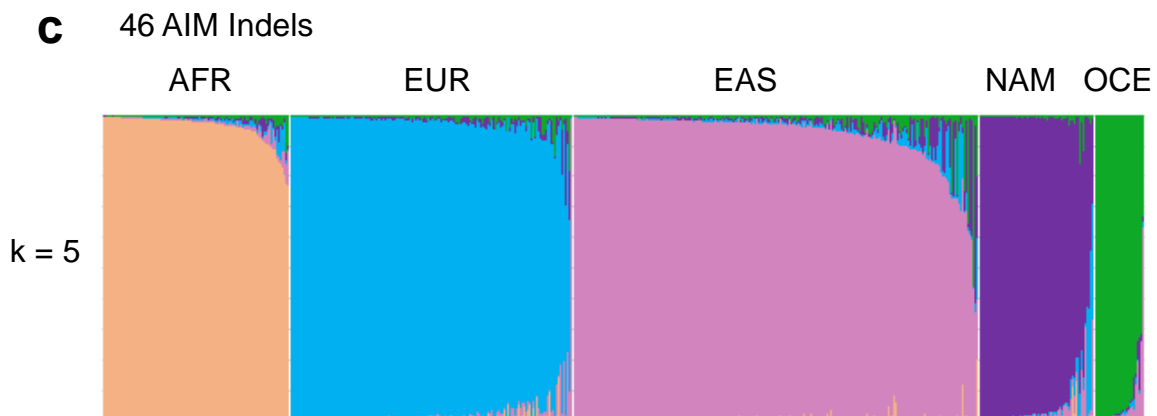
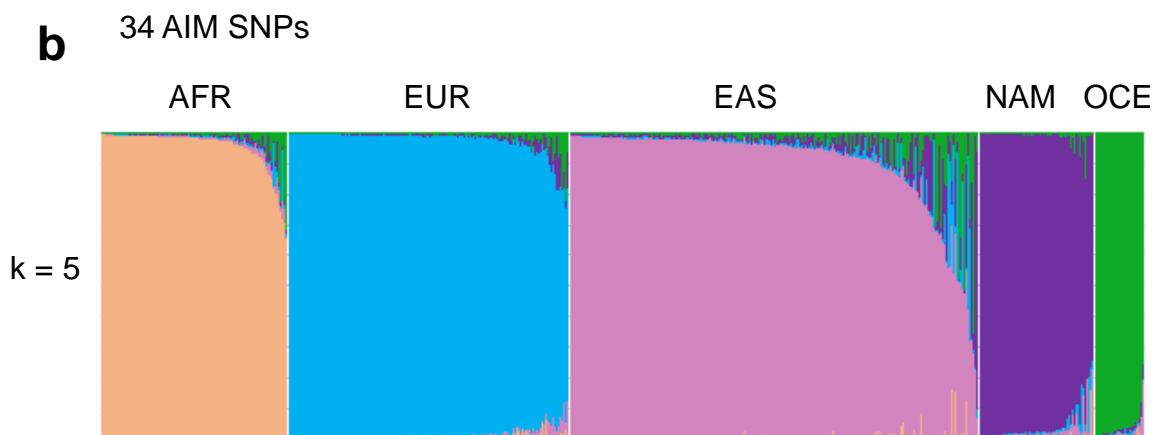
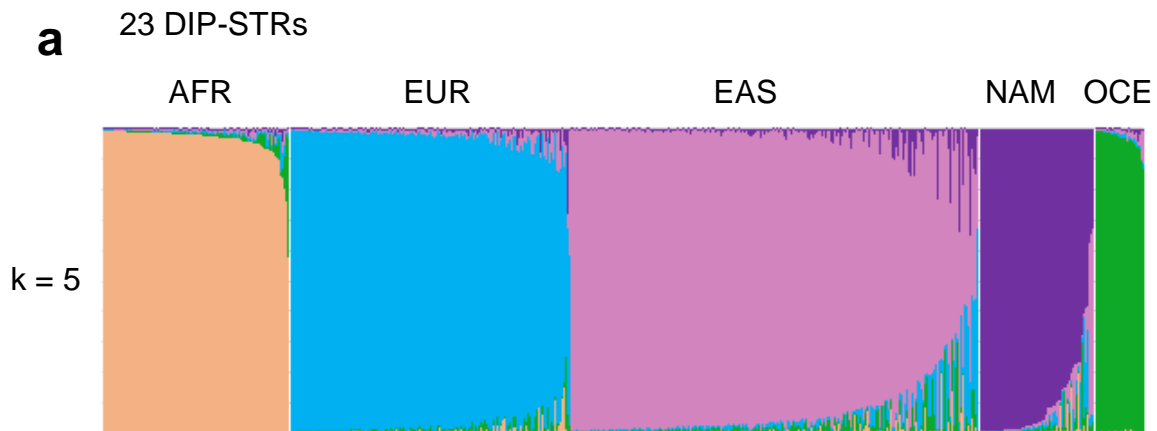
**Table 3 *Snipper* cross-validation classification success values for seven HGDP-CEPH major population groups using 23 DIP-STRs, 34 AIM SNPs and 46 AIM Indels.**

Population of origin		Population assignment							
		AFR	EUR	ME	CSA	EAS	OCE	NAM	
23 DIP-STRs		AFR	<b>100.00%</b>	-	-	-	-	-	-
	EUR	-	<b>95.57%</b>	1.90%	1.90%	0.63%	-	-	-
	ME	1.23%	5.56%	<b>86.42%</b>	6.17%	0.62%	-	-	-
	CSA	-	7.92%	7.43%	<b>82.18%</b>	2.48%	-	-	-
	EAS	-	0.43%	-	0.87%	<b>98.70%</b>	-	-	-
	OCE	-	-	-	-	-	<b>100.00%</b>	-	-
	NAM	-	-	-	-	-	-	<b>100.00%</b>	-
34 AIM SNPs		AFR	<b>100.00%</b>	-	-	-	-	-	-
	EUR	-	<b>69.62%</b>	25.95%	4.43%	-	-	-	-
	ME	0.62%	1.85%	<b>61.73%</b>	35.19%	0.62%	-	-	-
	CSA	-	4.95%	3.47%	<b>84.16%</b>	3.96%	1.49%	1.98%	-
	EAS	-	-	0.43%	3.04%	<b>86.96%</b>	1.74%	7.83%	-
	OCE	-	-	-	-	-	<b>96.43%</b>	3.57%	-
	NAM	-	-	-	-	-	-	<b>100.00%</b>	-
46 AIM Indels		AFR	<b>100.00%</b>	-	-	-	-	-	-
	EUR	-	<b>40.51%</b>	39.87%	19.62%	-	-	-	-
	ME	1.23%	0.61%	<b>46.01%</b>	52.15%	-	-	-	-
	CSA	-	-	0.50%	<b>90.10%</b>	7.43%	1.49%	0.50%	-
	EAS	-	-	-	-	<b>96.94%</b>	0.87%	2.18%	-
	OCE	-	-	-	-	-	<b>100.00%</b>	-	-
	NAM	-	-	-	1.56%	-	-	<b>98.44%</b>	-

The most likely major component of ancestry was considered independently of the admixture proportions

Populations are AFR-Africa, EUR-Europe, ME-Middle East, CSA-Central-South Asia, EAS-East Asia, OCE-Oceania and NAM- Native America

**a****b**

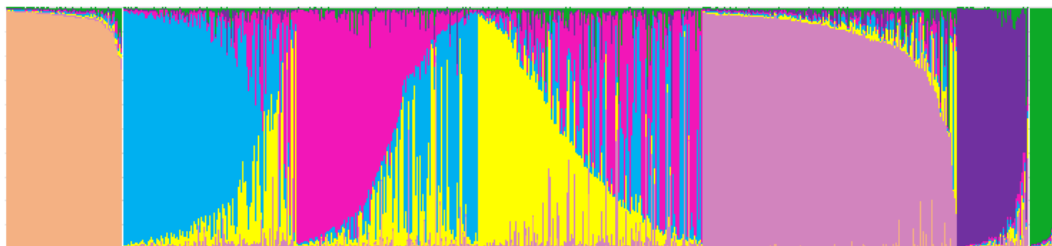


**a**

23 DIP-STRs

AFR EUR ME CSA EAS NAM OCE

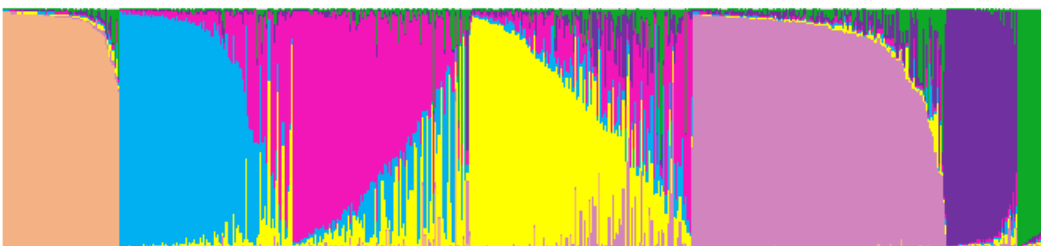
k = 7

**b**

34 AIM SNPs

AFR EUR ME CSA EAS NAM OCE

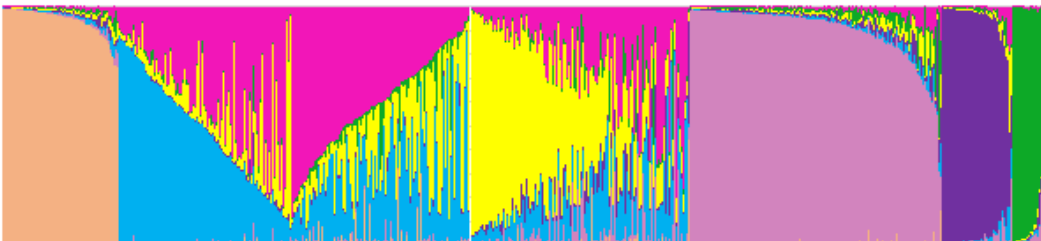
k = 7

**c**

46 AIM Indels

AFR EUR ME CSA EAS NAM OCE

k = 7



**a**

23 DIPs

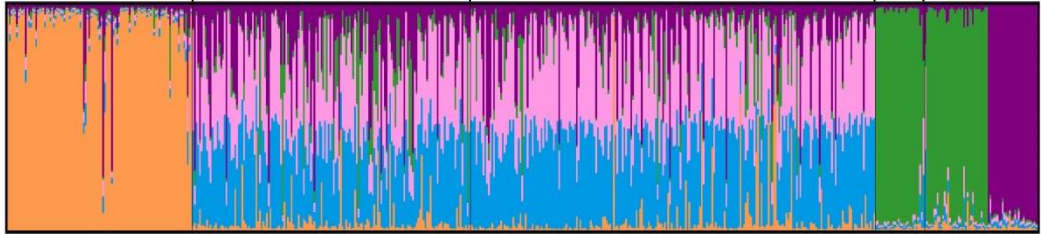
AFR

EUR

EAS

OCE NAM

k = 5



**b**

23 STRs

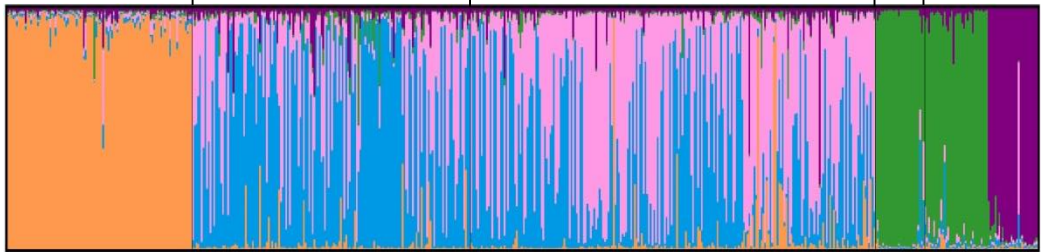
AFR

EUR

EAS

OCE NAM

k = 5



**Supplementary Table 1 DIP-STR primers**

Marker	DIP primers	STR primers	S-DIP allele-specific primer L-DIP allele-specific primer DIP-STR reverse primer
rs71725104-STR	*TTTTTGCCACAAAATAAATT GCAGCTCTGCAAAAATTC	*GTGTGGTGGTAGCTGGGACT TTTTGTGGCAAAAACTTGG	CCTTCCTTCTATTCTTGCTTTAT.TTT CCTTCCTTCTATTCTTGCTTTAT <u>CTA</u> *GTGTGGTGGTAGCTGGGACT
rs72534187-STR	*GCTCATGCAATTGATCAAACC GCTTTGGGCTTGATACAGAAA	*GGTATGCAATCTATCCTGATGTGA TGCATGAGCCAATTTATCTGA	TCTCTGGTTCTCTAGCTTGTAGAT.TAC CTCTAGCTTGTAGAT <u>GGCCTGTTA</u> *AGGCCAAAATTGACATTATAGTTTA
rs139592446-STR	*ACTGGGAAAACATGAGAAACAAA CCTTCCTTTTATCTTCTATCACACA	*GATTTAGGAGGGGATGTGGT TTCCCAAGTGTCTGCTCAA	CCATTTTGCCCCACTAGT.TC TGCCCCACTAGT <u>GACTAAGTTC</u> *TAGCCTTCTGCCCAAACATC
rs36194161-STR	*CCAAGATTGTGCCACTGC TCACCAGGTCTTGGGTATGTT	*AAAACATACCCAAGACCTGGTG TTTAAACCTCTTCTGCTTGC	AAATATTACTAGTTGTATTAGTCTGTT.ACG TATTACTAGTTGTATTAGTCTGTT <u>CTCAC</u> *TGCAGTGAGCAGGGTGAC
rs13831044-STR	*ACAATCGCTGCTCACTGAAG GCCGAAGCAGGTGTATTCTT	*AGCACATAGCAGGCACTAGC GCGATTGTGCCACTACACAG	ATTAGCTGGGCTTAGTG.CCTGT TAGCTGGGCTTAGTGG <u>CATATG</u> *GCACTAGCTGTTAGTTCTTTTCTG
MID473 <sup>a</sup> -STR	*AAATGTTAAGCCTCCCTGTG CCACTGACAGCAACAACCAA	*CCTTGCTTGGTTGTTGCTG TGCAGGCAGATTTTAAAGGAA	TGGGCTTTCTA.TTACATTTTTAGC TGGGCTTTCTA <u>TACATTTT</u> TACAT *TGCAGGCAGATTTTAAAGGAA
MID2538 <sup>a</sup> -STR	*ACAATCTTGGCACCCATTT GCTCGCAAAGTAGGCAAGTT	*TCATTACCTTCTCTGCATTGGA GTGCCAAGATTGTTGGTGTG	GTTCAAAATCACAATCACTCA.TTT TCAAAATCACAATCACTCA <u>AAACA</u> *TGGAATCACTCATTACCTTCTCTG
MID2824 <sup>a</sup> -STR	*TGTTCACCTTCTGCCATGTG TCTAGTGGGGTTTGCAGAG	*ACTTGGGAGGCTGAGACAGA CACATGGCAGAAGTGAACA	TCCAAGATGAGCACTG.GGC CCTCCAAGATGAGCACTG <u>AGTC</u> *CCAGCCTGGCAACAGAGTA
MID73 <sup>a</sup> -STR	*TGTGTTTCTAAGGAGCGCTGT CACAGTGAGGAGAAGGAAGGA	*CCATTCTCTCCTTCTCTCC CCTGGTGCCAGAGCAT	CATACTCAGAAGTGCCTT.GAAAAG CATACTCAGAAGTGCCTT <u>GAAAGAA</u> *GCACATGGCTCTTTAATACACTG

<sup>a</sup> Marker name is from the Marshfield database and corresponds to rs140762, rs3054057, rs11278940, rs16365, respectively

\* Fluorescent labeled primers. In the last column dots indicate the insertion/deletion point and underlined is the inserted sequence. Primers for markers previously published were not changed

**Supplementary Table 2 Marker information, PCR primer concentration, number of cycles and multiplex groups**

Marker	dbSNP ID of the DIP marker	Genome hg38 position of the DIP marker	Genbank ID of the STR marker	Genome hg38 position of the STR marker
MID1013-D5S490	rs1611095:insCCAG	chr5:g.127507347_127507348insCCAG	Z23637 [GT]	chr5:127507078-127507107
rs2308142-STR	rs2308142:insATT	chr20:g.3484768_3484769insATT	G08041 [TTA]	chr20:3484910-3484948
rs11277790-D10S530	rs11277790:delTCCAAC	chr10:g.105760904_105760910delTCCAAC	Z23432 [GT]	chr10:105761139-105761184
rs112604544-STR	rs112604544:insTTTAA	chr1:g.185716220_185716221insTTTAA	Simple Tandem Repeat [TTCC]	chr1:185716115-185716167
rs34212659-STR	rs34212659:delAGG	chr7:g.38563258_38563260delAGG	Simple Tandem Repeat [TGAA]	chr7:38563129-38563173
rs145423446-STR	rs145423446:delAGTC	chr16:g.8094810_8094813delAGTC	Simple Tandem Repeat [GATA]	chr16:8094861-8094992
rs111478323-STR	rs111478323:insGAGA	chr2:g.2619471_2619472insGAGA	Simple Tandem Repeat [TTTA]	chr2:2619638-2619680
rs146332920-STR	rs146332920:delAGG	chr9:g.111781629_111781631delAGG	Simple Tandem Repeat [TAAA]	chr9:111781725-111781753
rs35032587-STR	rs35032587:insTATT	chr15:g.93295596_93295597insTATT	Simple Tandem Repeat [AGAT]	chr15:93295413-93295457
rs142543564-STR	rs142543564:delTACT	chr2:g.211170402_211170405delTACT	Simple Tandem Repeat [ATAA]	chr2:211170264-211170317
rs72406828-STR	rs72406828:delATTG	chr4:g.86985935_86985938delATTG	Simple Tandem Repeat [AATTT]	chr4:86986026-86986089
rs71070706-STR	rs71070706:delTGT	chr1:g.117804313_117804315delTGT	Simple Tandem Repeat [AAAG]	chr1:117804149-117804266
MID473-STR	rs140762:delTTACATTT	chr6:g.94529411_94529418delTTACATTT	Simple Tandem Repeat [AGGA]	chr6:94529500-94529560
MID2538-STR	rs3054057:delTGTT	chr15:g.85467307_85467310delTGTT	Simple Tandem Repeat [AC]	chr15:85467076-85467101
MID2824-STR	rs11278940:delAGGACT	chr11:g.36029406_36029411delAGGACT	Simple Tandem Repeat [AAAC]	chr11:36029245-36029271
rs60194384-D15S1514	rs60194384:delTCTTAA	chr15:g.95010636_95010641delTCTTAA	G10630 [TATC]	chr15:95010830-95010878
rs66679498-D2S342	rs66679498:delCCAAC	chr2:g.195019736_195019750delCCAAC	Z23993 [CA]	chr2:195019569-195019613
rs139592446-STR	rs139592446:delACTTAGTC	chr2:g.162477942_162477949delACTTAGTC	Simple Tandem Repeat [CATC]	chr2:162477813-162477859
rs71725104-STR	rs71725104:delATAG	chr13:g.93633710_93633713delATAG	Simple Tandem Repeat [AAAT]	chr13:93633636-93633677
rs36194161-STR	rs36194161:delCTC	chr2:g.187252184_187252186delCTC	Simple Tandem Repeat [TTTA]	chr2:187252243-187252287
rs72534187-STR	rs72534187:delACAGGCC	chr5:g.42202012_42202018delACAGGCC	Simple Tandem Repeat [ATAG]	chr5:42201874-42201943
rs138331044-STR	rs138331044:delCATATGC	chr1:g.119135915_119135921delCATATGC	G07791 [AGAT]	chr1:119135695-119135762
MID73-STR	rs16365:delGAA	chr22:g.34308579_34308581delGAA	Simple Tandem Repeat [CCATC]	chr22:34308746-34308839

<sup>a</sup>Annealing temperatures and number of cycles are the same for the S- and L-specific DIP-STR amplifications



Supplementary Table 2 (continued)

Marker	Primers DIP			Annealing		Primers S-DIP-STR		Primers L-DIP-STR		Annealing <sup>a</sup>	
	multiplex group	concentration (nM)	temperature (°C)	N cycles	multiplex group	concentration (nM)	multiplex group	concentration (nM)	temperature (°C)	N cycles <sup>a</sup>	
MID1013-D5S490	1	100	52	28	6	100	5	100	52	34	
rs2308142-STR	1	200	52	28	5	100	6	100	52	34	
rs11277790-D10S530	1	100	52	28	5	100	6	100	52	34	
rs112604544-STR	2	100	55	34	7	100	8	150	55	34	
rs34212659-STR	2	100	55	34	7	100	8	100	55	34	
rs145423446-STR	2	100	55	34	7	100	8	100	55	34	
rs111478323-STR	2	50	55	34	7	100	8	100	55	34	
rs146332920-STR	2	50	55	34	7	100	8	100	55	34	
rs35032587-STR	2	200	55	34	9	100	10	100	55	34	
rs142543564-STR	2	100	55	34	9	100	10	100	55	34	
rs72406828-STR	2	150	55	34	9	100	10	100	55	34	
rs71070706-STR	2	50	55	34	9	200	10	200	55	34	
MID473-STR	3	100	55	30	11	50	12	50	55	34	
MID2538-STR	3	100	55	30	11	200	12	200	55	34	
MID2824-STR	3	400	55	30	11	50	12	50	55	34	
rs60194384-D15S1514	1	200	52	28	11	200	12	200	55	34	
rs66679498-D2S342	1	200	52	28	11	100	12	500	55	34	
rs139592446-STR	4	150	55	34	13	50	14	37	58	34	
rs71725104-STR	4	150	55	34	13	75	14	50	58	34	
rs36194161-STR	4	150	55	34	13	100	14	150	58	34	
rs72534187-STR	4	75	55	34	13	100	14	100	58	34	
rs138331044-STR	4	120	55	34	13	100	14	500	58	34	
MID73-STR	3	400	55	30	15	200	16	100	59	34	

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

.....

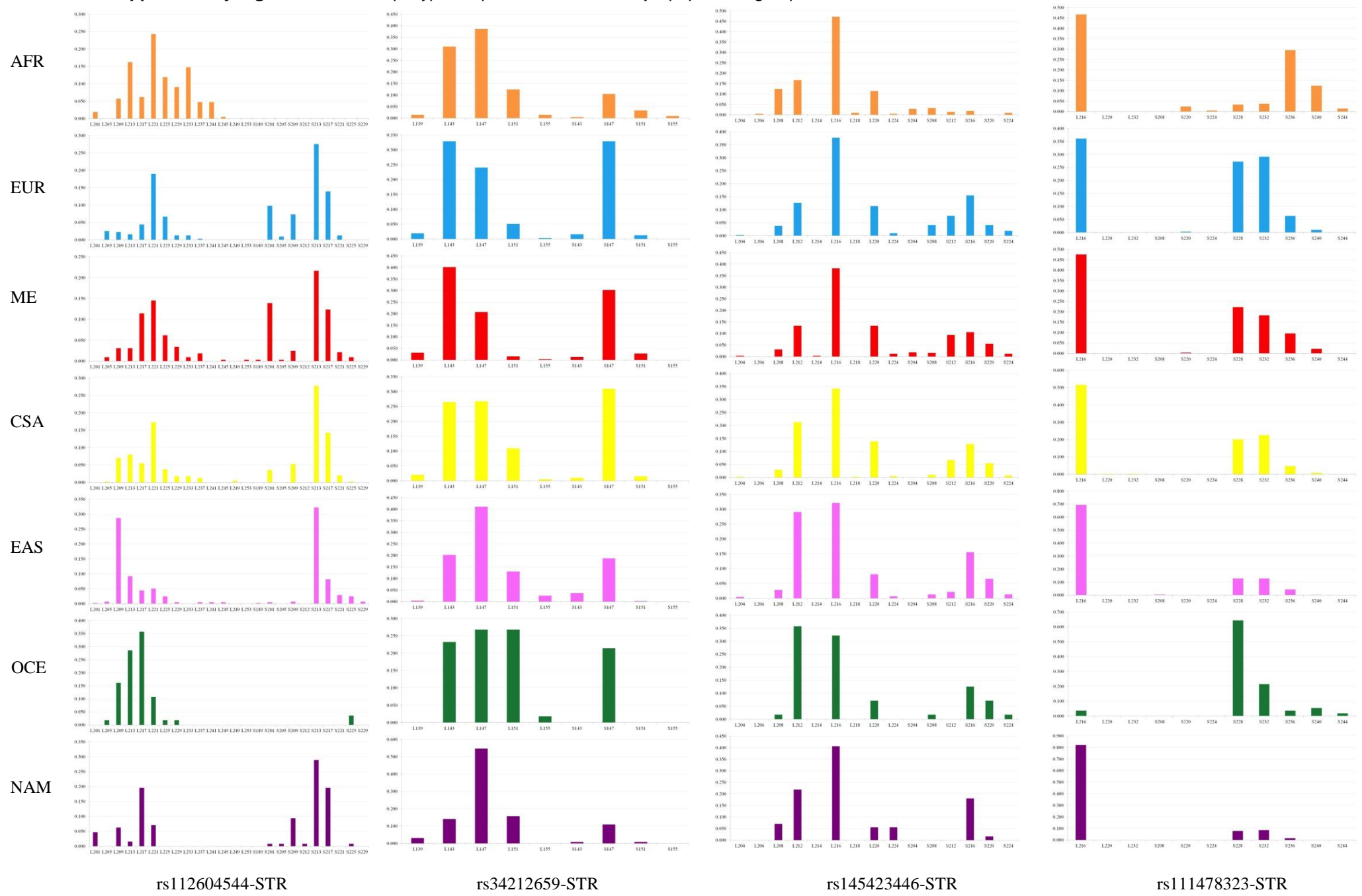
.....

.....

1. Introduction  
2. Literature Review  
3. Methodology  
4. Results  
5. Discussion  
6. Conclusion

The following text is a placeholder for the main body of the document, which would contain the detailed analysis and findings of the study.

# Supplementary Fig. 1 DIP-STRs haplotype frequencies in seven major population groups



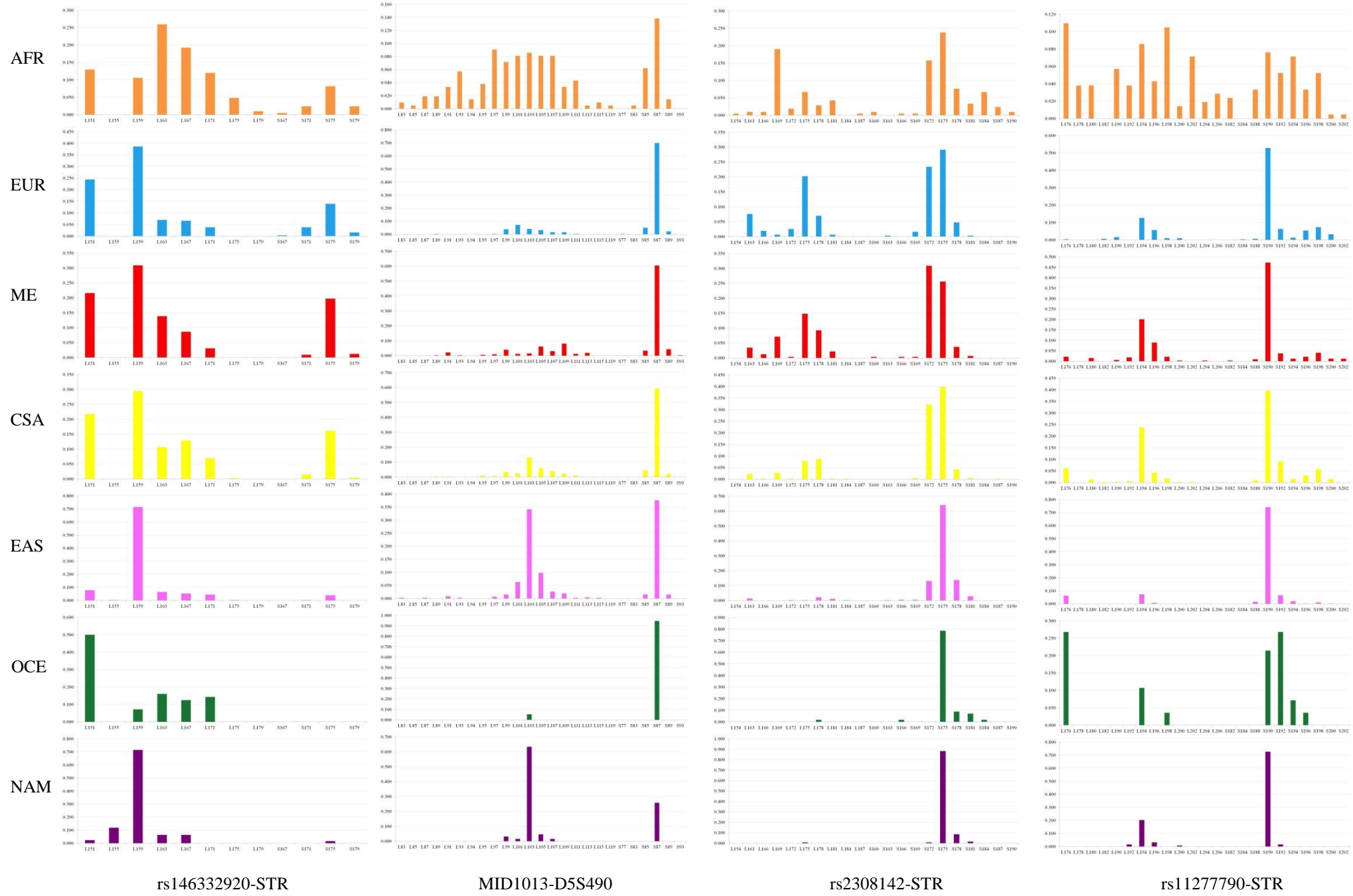
rs112604544-STR

rs34212659-STR

rs145423446-STR

rs111478323-STR

# Supplementary Fig. 1 (continued)



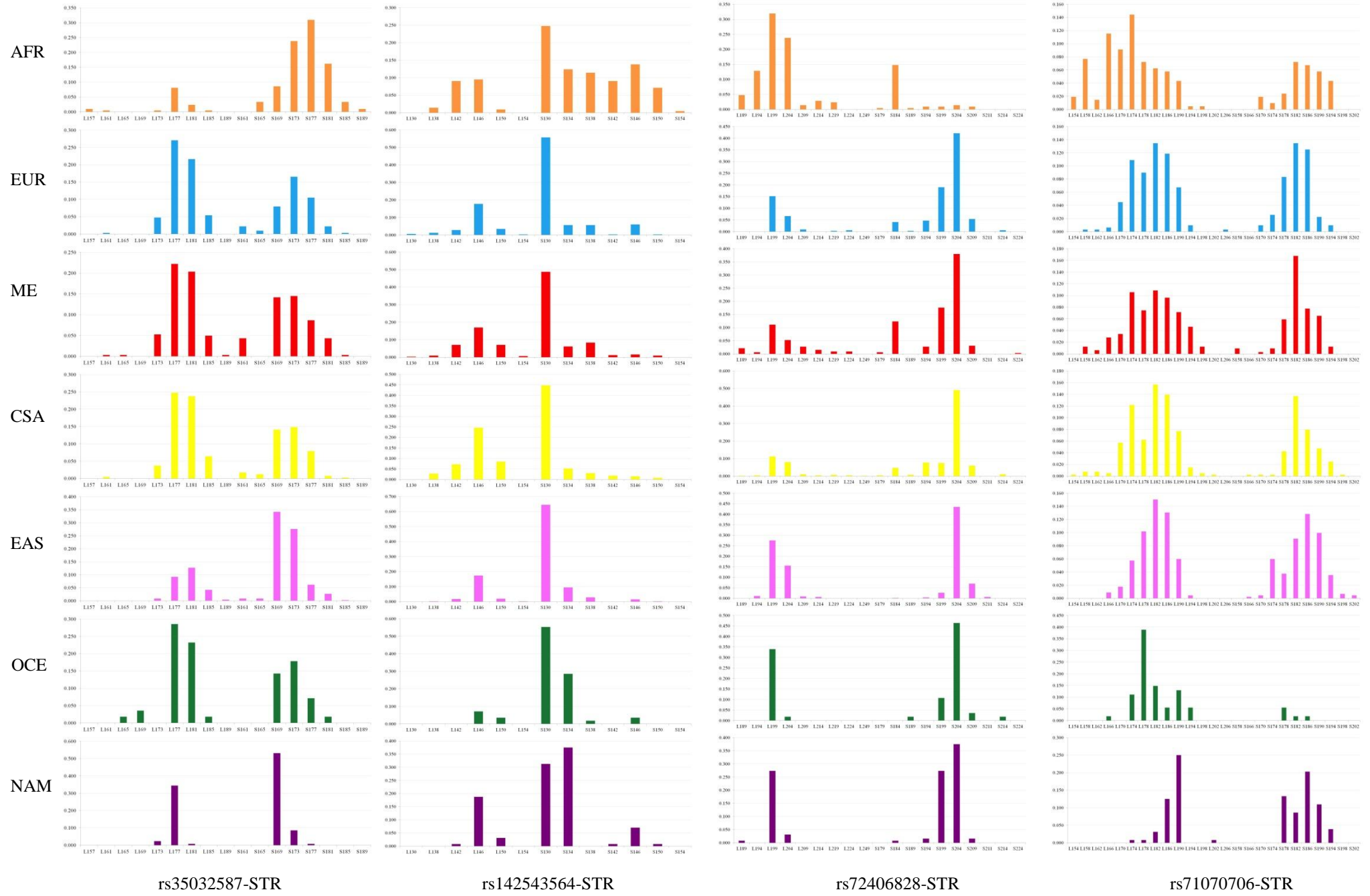
rs146332920-STR

MID1013-D5S490

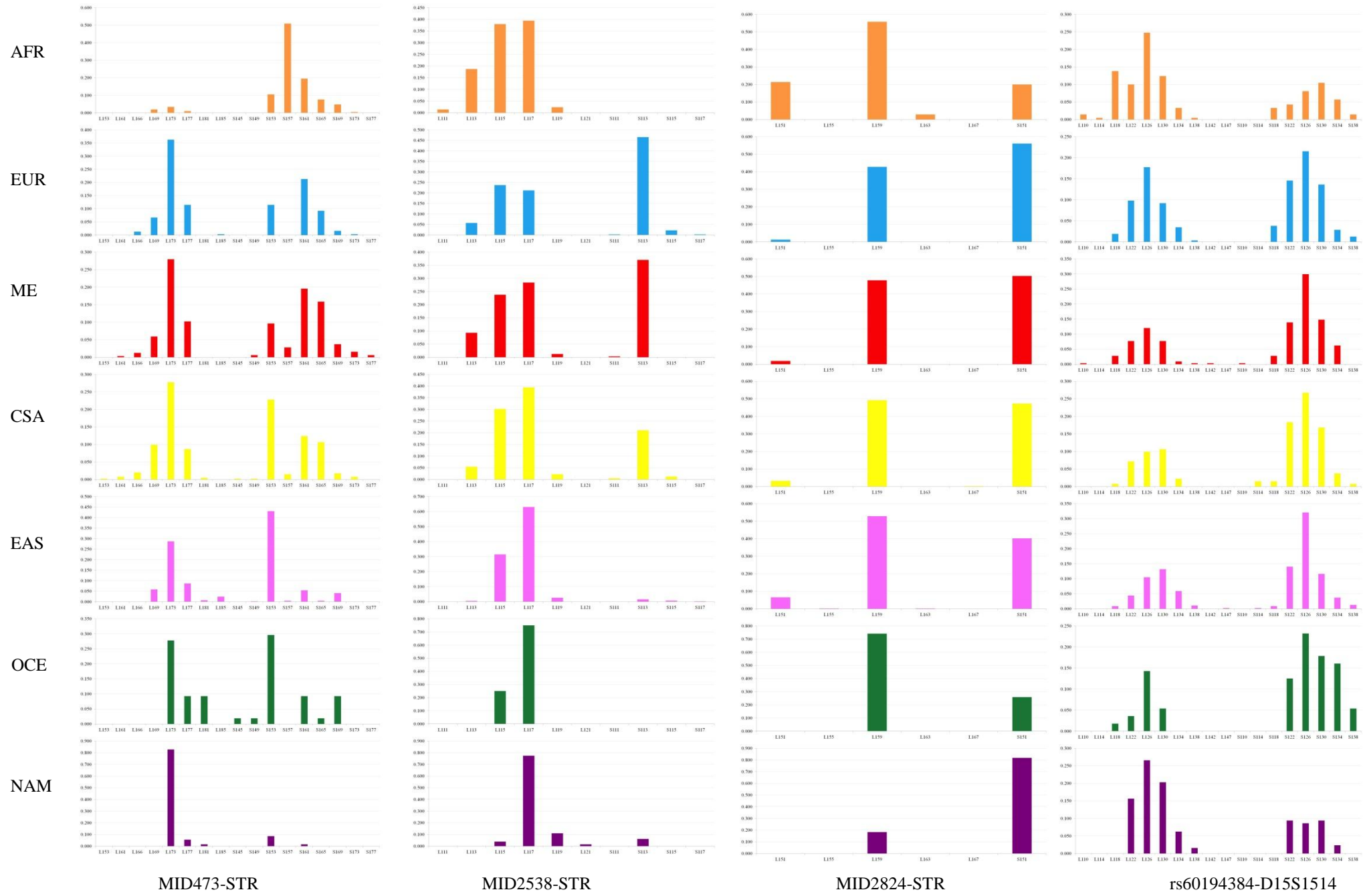
rs2308142-STR

rs1127790-STR

# Supplementary Fig. 1 (continued)



# Supplementary Fig. 1 (continued)



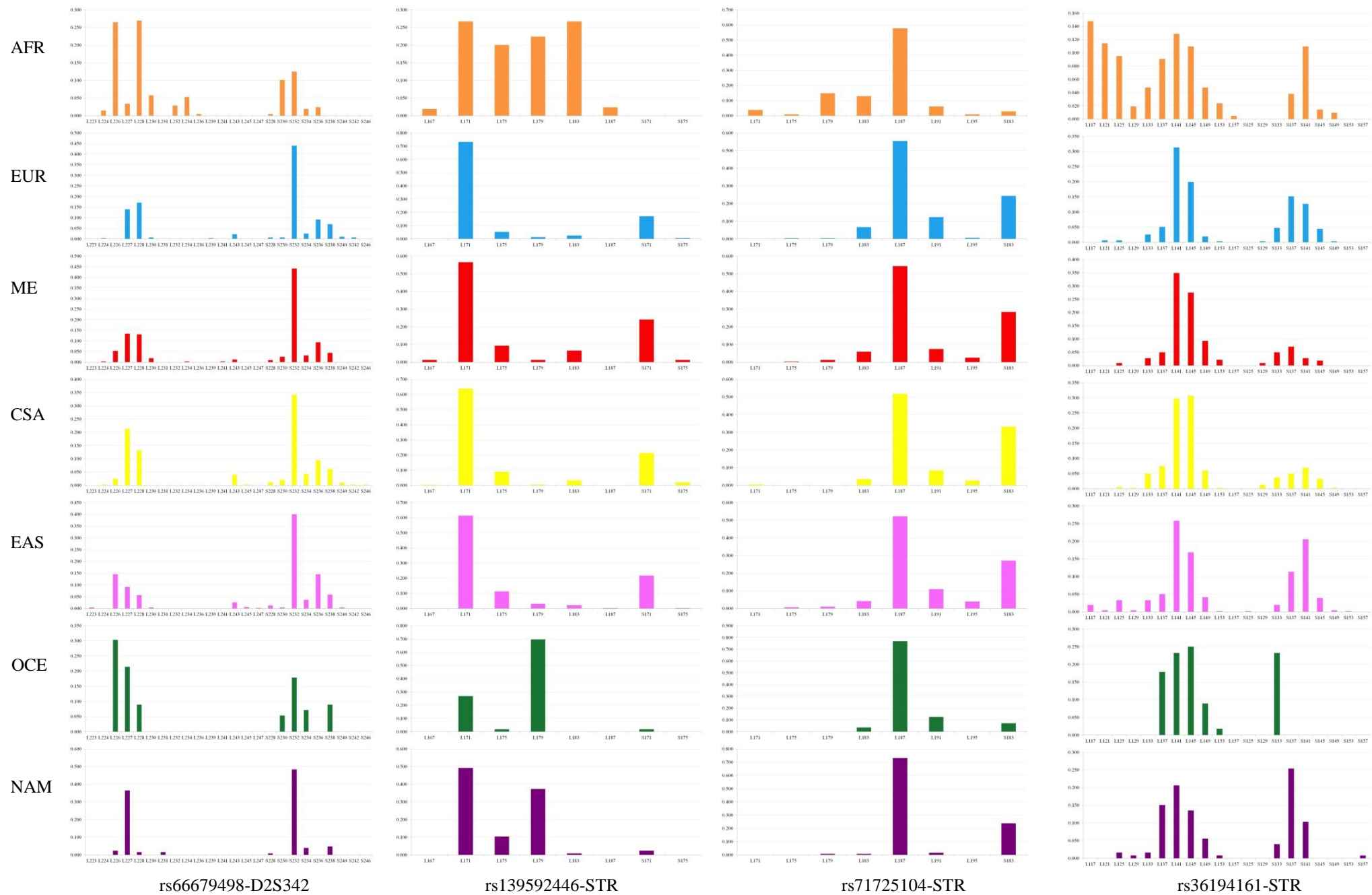
MID473-STR

MID2538-STR

MID2824-STR

rs60194384-D15S1514

# Supplementary Fig. 1 (continued)



rs66679498-D2S342

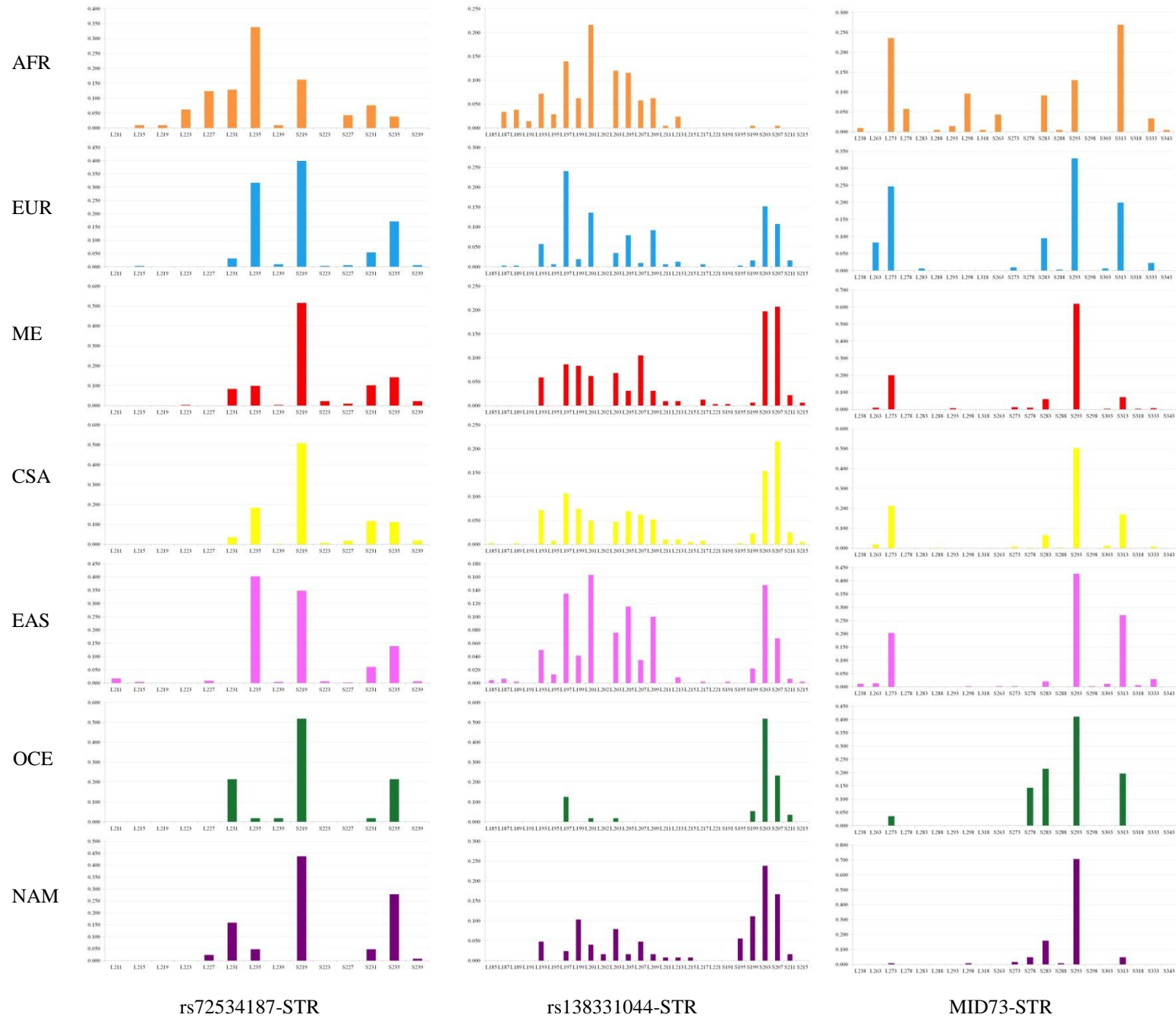
rs139592446-STR

rs71725104-STR

rs36194161-STR



Supplementary Fig. 1 (continued)



rs72534187-STR

rs138331044-STR

MID73-STR