# UC Davis
## UC Davis Previously Published Works

**Title**

Elucidation of familial relationships using hair shaft proteomics

**Permalink**

**Authors**

Karim, Noreen
Plott, Tempest J
Durbin-Johnson, Blythe P
et al.

**Publication Date**

2021-09-01

**DOI**

Peer reviewed

# Elucidation of Familial Relationships Using Hair Shaft Proteomics

Noreen Karim[a], Tempest J. Plott[a,b], Blythe P. Durbin-Johnson[c], David M. Rocke[c], Michelle Salemi[d], Brett S. Phinney[d], Zachary C. Goecker[a], Marc J. M. Pieterse[a], Glendon J. Parker[a,b], Robert H. Rice[a,b,*]

[a]Department of Environmental Toxicology, University of California, Davis, USA

[b]Forensic Science Program, University of California, Davis, USA

[c]Division of Biostatistics, Department of Public Health Sciences, Clinical and Translational, Science Center Biostatistics Core, University of California, Davis, USA

[d]Proteomics Core Facility, University of California, Davis, USA

*Corresponding author

**Keywords**

Proteomic profiling, genetically variant peptides, human hair, forensic investigation, relationship testing

**Proteomics repository files**

**Acknowledgments**

**Competing Interests**

The authors declare no conflict of interest, with the exception of GJP, who has a patent based on use of genetically variant peptides for human identification (US 8,877,455 B2, Australian Patent 2011229918, Canadian Patent CA 2794248, and European Patent EP11759843.3, GJP inventor). The patent is owned by Parker Proteomics LLC. Protein-Based Identification Technologies LLC (PBIT) has an exclusive license to develop the intellectual property and is co-owned by Utah Valley University and GJP. This ownership of PBIT and associated intellectual property does not alter policies on sharing data and materials. These financial conflicts of interest are administered by the Research Integrity and Compliance Office, Office of Research at the University of California, Davis to ensure compliance with University of California Policy.

1      Elucidation of Familial Relationships Using Hair Shaft Proteomics

2

3      **Abstract**

4      This study examines the potential of hair shaft proteomic analysis to delineate genetic

5      relatedness. Proteomic profiling and amino acid sequence analysis provide information for

6      quantitative and statistically-based analysis of individualization and sample similarity. Protein

7      expression levels are a function of cell-specific transcriptional and translational programs. These

8      programs are greatly influenced by an individual's genetic background, and are therefore

9      influenced by familial relatedness as well as ancestry and genetic disease. Proteomic profiles

10      should therefore be more similar among related individuals than unrelated individuals. Likewise,

11      profiles of genetically variant peptides that contain single amino acid polymorphisms, the result

12      of non-synonymous SNP alleles, should behave similarly. The proteomically-inferred SNP

13      alleles should also provide a basis for calculation of combined paternity and sibship indices. We

14      test these hypotheses using matching proteomic and genetic datasets from a family of two adults

15      and four siblings, one of which has a genetic condition that perturbs hair structure and properties.

16      We demonstrate that related individuals, compared to those who are unrelated, have more similar

17      proteomic profiles, profiles of genetically variant peptides and higher combined paternity indices

18      and combined sibship indices. This study builds on previous analyses of hair shaft protein

19      profiling and genetically variant peptide profiles in different real-world scenarios including

20      different human hair shaft body locations and pigmentation status. It also validates the inclusion

21      of proteomic information with other biomolecular substrates in forensic hair shaft analysis,

22      including mitochondrial and nuclear DNA.

23

24      **Introduction**

25      Hair shafts are a common component of crime scenes and are currently underutilized

26      forensically. Use of morphological patterns in hair shafts is currently considered controversial in

27      forensic science due to the intrinsically subjective nature of pattern matching (Council, 2009).

28      Additionally, nuclear DNA is degraded in hair shafts as part of the natural cornification process

29      (Linch et al, 2001; McNevin et al, 2005). This effectively eliminates the possibility of routinely

30    obtaining identifying STR genotypes. Since the abundant mitochondrial DNA, unlike nuclear

31    DNA, persists in the hair shaft, its matrilineal haplotype analysis is the current best practice for

32    obtaining identifying genetic information from the hair shaft. Recent research has demonstrated

33    that hair shaft protein may also provide forensically relevant identifying information in the form

34    of genetically variant peptides (GVPs) (Goecker et al, 2020; Parker et al, 2016). The forensic

35    utility and scope of proteomic genotyping continues to be extended and demonstrated to be

36    unaffected in forensically relevant, real-world contexts including hair from different body

37    locations (Chu et al, 2019; Milan et al, 2019), different pigmentation states (Franklin et al, 2020),

38    from long term storage (Plott et al, 2020), and even in hair from experimental explosive devices

39    (Chu et al, 2020). This study examines whether the proteomic information in hair shafts is able to

40    delineate familial relationships.

41    Proteomic information in forensic genetics consists of two basic forms, the amino acid sequences

42    themselves and the relative profile of protein expression. The profile, a lineup of the many

43    proteins in the sample and their relative levels of expression, is a function of cell-specific

44    transcriptional and translational programming. In addition to a myriad of physiological,

45    anatomical and biochemical contexts, the genetic background of each individual would also play

46    a significant role. Previous findings with mice (Rice et al, 2012) and humans (Wu et al, 2017)

47    indicate that protein expression levels in the hair shaft are largely genetically determined.

48    However, wide variation is observed among hair samples from individuals in the outbred human

49    population (Laatsch et al, 2014), likely arising from sequence variations in noncoding regions of

50    the genome (Hindorff et al, 2009; Martin-Trujillo et al, 2020), including gene promoters and

51    miRNA binding sites that affect transcription factor binding sites or chromatin accessibility. This

52    background of variation would be predicted to be lower in genetically related individuals, and

53    the proteomic profiles of related individuals would therefore be predicted to be more similar to

54    each other than to those of unrelated individuals (Wu et al, 2017). Since children would be

55    expected to inherit determinants of individual hair protein expression level from each parent,

56    their individual hair protein levels would be expected to mimic those of either parent or to be

57    intermediate between them. Based on this expectation, we test the hypothesis that hair protein

58    profiles in a family are more similar in two-way comparisons between a parent and individual

59    children than between the parents. The family studied in this case has three unaffected offspring

60    and one diagnosed with a rare genetic condition where the hair is brittle and has an unusual

61    protein/lipid ratio (Alsop et al, 2016). This happenstance has permitted the opportunity to

62    determine whether a hair sample appears abnormal within the context of a family.

63    In addition to providing information on protein expression levels, hair shaft proteomic digests

64    also permit analysis of GVPs within those proteins and the development of a proteomically-

65    inferred genotype of non-synonymous single nucleotide polymorphism (SNP) alleles (Parker et

66    al, 2016). This manifestation of allelic differences permits inference of corresponding SNPs in

67    the genomic DNA of hair donors. Although hair protein profiling may have utility in

68    distinguishing individuals, GVPs are more robust and offer a greater power of discrimination.

69    Like any genotype marker system, these profiles would be predicted to be more similar in related

70    individuals, and therefore have the potential also to be exploited to develop measures of genetic

71    relatedness. The present study offers an opportunity to determine kinship indices by analysis of

72    hair shaft digests from a single family compared to nine unrelated individuals.

73    **Materials and Methods**

74    **Sample Collection and Processing**

75    For the current study, six family members of European ancestry were enrolled after obtaining

76    written informed consent either from the individuals or from the parents in the case of minors

77    <18 years of age. The study was conducted in accordance with protocols and procedures

78    approved by the Institutional Review Board of the University of California Davis. The enrolled

79    individuals included mother (M), father (P) and their four children, two sons (S1 and S2) and two

80    daughters (S3 and S4). Hair shafts were collected from each enrolled individual. Abnormalities

81    in hair shaft structure were not visible by light microscopy. For the proteomic analysis three

82    replicates of hair samples from each individual except P and S2 (four and six replicates,

83    respectively) were processed as previously described (Plott et al, 2020), and the randomized

84    protein digests were subjected to LC-MS/MS using a Thermo Scientific Q Exactive Plus

85    Orbitrap mass spectrometric analysis (Wu et al, 2017).

86    **Database Searching and Proteomic Profiling**

87    The data files generated by LC-MS/MS were searched against a Uniprot human database

88    appended with a database containing identical but reversed (decoy) peptides and common human

89    contaminants using X!Tandem (2016.10.15.2). The peptide and protein identifications were

90    validated in Scaffold (version 4.8.2, Proteome Software Inc., Portland). Proteins identified at a

91    minimum of 99% probability and represented by at least two peptides identified at 95%

92    probability were included in the analysis (false discovery rate of 0.7% for proteins and <0.1% for

93    peptides). The weighted spectral count data provided by Scaffold were used for the profiling and

94    statistical analyses after confirming protein presence by exclusive spectral counts. To obtain the

95    number of significant differences between profiles, two-way comparisons were conducted, where

96    the weighted spectral counts were compared separately for each protein from the two subjects

97    (Table S10).  These differential protein expression analyses were conducted using the limma-

98    voom Bioconductor pipeline (Ritchie et al, 2015), which was originally developed for RNA

99    sequencing data (limma version 3.44.3, edgeR version 3.30.3). Normalization factors were

100   calculated using trimmed mean of M values (Robinson and Oshlack, 2010). P-values were

101   adjusted for multiple testing across proteins (Benjamini and Hochberg, 1995). The model used in

102   limma included effects for individual and batch. Analyses were conducted using R version 4.0.0

103   Patched (2020-05-18 r78487). The raw data files and scaffold analysis files are available on the

104   MassIVE repository (https://massive.ucsd.edu) MassIVE # MSV000086665 (reviewer password

105   "Hair Shaft"), ProteomeExchange # = PXD023446.

106   **GVP Analysis**

107   To obtain genetically variant peptides (GVPs) profiles, the raw data files for all the samples were

108   first converted by MSConvertGUI (Proteowizard 2.1 http://proteowizard.sourceforge.net) to

109   MzML format and were subsequently searched using X!Tandem peptide spectra matching

110   algorithm (GPM Fury, X!Tandem Alanine 149 v.3.0 (2017-02-01)). Default search parameters

111   were used except that the search was limited to eukaryotic reference libraries, peptide and

112   protein log(e) scores were set to <-1, fragment mass error of 20 ppm, parent mass error of 100

113   ppm, and point mutations from the refinement specifications were included in the search. The

114   peptides identified by GPM Fury for each sample were subsequently searched for previously

115   identified GVPs using GVP Finder (v 1.2) (https://parkerlab.ucdavis.edu/gvp-finder) where

116   searches for GVPs in the peptide data followed the previously established criteria (Borja et al,

117   2019; Goecker et al, 2020; Plott et al, 2020). Moreover, the .xml files for all the individuals were

118   also explored using the discovery approach by looking for peptides carrying single amino acid

4

119     variations with log(e) scores <-2 with no other chemical or genetic modifications and no peaks

120     representing the alternate amino acid (Borja et al, 2019). The single amino acid variations

121     carrying peptides were evaluated against all human protein sequences in the PROWL web portal

122     (prowl.rockefeller.edu/prowl/) for uniqueness to confirm that they were translated from a single

123     site in the genome. Because of the familial structure of the study, the GVPs were not filtered

124     based on their low allele frequencies contrary to earlier studies (Parker et al, 2016; Borja et al,

125     2020). The obtained GVP profiles of individuals P, M, S1 and S2 were validated from their DNA

126     data. The genetic data of individuals S3 and S4 were not available.

127     **Exome Sequencing Data**

128     Exome DNA sequencing data were provided by the Department of Human Genetics, Radboud

129     University Medical Center, the Netherlands. Data for P, M and S1 were obtained using Illumina

130     HiSeq and those for S2 using SOLiDxl 5500 instrumentation. Genotype information analogous

131     to the detected GVPs were obtained from the exome data for all the four individuals. Data from

132     S2 were not consistent with the M and P at 5 loci encoding GVPs by Mendelian genetics, but

133     proteomic data permitted correction of three of these loci (**Table S1**). The discrepancy reflects

134     the higher error rate in the older SOLiDxl method.

135     **Hierarchical Clustering**

136     Data from a previously published set of unrelated European American individuals (Plott et al,

137     2020) were merged with the GVP list of the currently studied family. A binary format data

138     matrix was generated with 1 representing a GVP detection and 0 a non-detection (**Table S2**).

139     Each row of the matrix represents the GVP information for each individual with the columns

140     representing SNPs. The matrix was used to calculate Euclidean distance between the

141     rows/samples, based on which agglomerative hierarchical clustering was performed, and a

142     dendrogram was plotted for the clustering using the hclust function of R (Version 3.6.3 (2020-

143     02-29)).

144     **Parentage Index and Sibship Index Calculation**

145     The GVPs detected in the samples were used in kinship calculation (parentage indexes and

146     sibship indexes) that can provide a statistical value for the probability of relationship between

147      samples. Likelihood ratios were calculated using the SNP data obtained from exome sequencing

148      corresponding to all the identified GVPs. Moreover, SNPs were inferred from the GVP profiles

149      for all the studied individuals where each locus was treated as homozygous for an allele if a

150      peptide corresponding to only one allele at the locus was detected and heterozygous if both

151      GVPs were detected in the proteomic data. GVPs from the loci where only one GVP was unique

152      were excluded from this analysis. GVPs from different genes were assumed as completely

153      independent whereas complete linkage between loci within a gene was assumed to account for

154      linkage disequilibrium. In cases of more than two GVPs within one gene, the two with the

155      highest allele frequencies of the minor allele were used. Likelihood ratios were calculated as

156      described (Sozer et al, 2010; Wenk et al, 1996) using the formulae in **Table S3**. Relationship

157      indices for each locus were calculated with allele frequencies from the European population

158      (Consortium et al, 2015). Combined paternity and sibship indices were obtained by taking a

159      product of the respective indices for all the loci included in each analysis.

160      **RESULTS**

161      **Proteomic Profiling**

162      The protein levels in hair samples from all six studied individuals were subjected to two way

163      comparisons to evaluate the impact of their genetic relationships. Using the standard significance

164      level of $p<0.05$ (after correction for multiple testing) showed few protein level differences

165      between the parents (**Table 1A**). While unusual, this degree of similarity is occasionally

166      observed among unrelated individuals (Laatsch et al, 2014). Nevertheless, supporting the original

167      hypothesis, the profiles of three offspring (S1, S3, S4) exhibited few proteins whose levels

168      differed from those in the parental hair samples (0-2) or from each other (0) by this criterion. In

169      contrast, however, the profile of one child (S2) with a rare hair phenotype was quite distinct,

170      showing 0-11 proteins differing in level from those in other family members (**Table 1A**).

171      **A**

**A**

| p<0.05 | M | S1 | S2 | S3 | S4 |
|--------|---|----|----|----|----|
| P | 1 | 0 | **11** | 0 | 0 |
| M |  | 1 | **5** | 0 | 2 |
| S1 |  |  | **4** | 0 | 0 |
| S2 |  |  |  | **2** | 0 |
| S3 |  |  |  |  | 0 |

**B**

| p<0.1 | M | S1 | S2 | S3 | S4 |
|-------|---|----|----|----|----|
| P | 13 | 5 (1) | **24 (2)** | 1 (1) | 3 (3) |
| M |  | 1 (1) | **13 (2)** | 0 | 3 (3) |
| S1 |  |  | **6** | 2 | 0 |
| S2 |  |  |  | **2** | **0** |
| S3 |  |  |  |  | 0 |

172

**Table 1.** Proteins with significant differences in hair protein profiles. Values for two-way comparisons between Father (P), Mother (M) and siblings (S1-S4) are tabulated with p<0.05 (A) or p<0.1 (B). In parentheses (B) are numbers of proteins in each case that match those differing between the parents and thus plausibly result through inheritance from the other parent.

To obtain a more expansive view of the proteomic relationships, differences were analyzed at a less stringent significance level of p<0.1. As shown in **Table 1B**, the profiles of mother (M) and father (P) exhibited differences in 13 proteins. Hair from three siblings (S1, S3, S4) exhibited few differences with each parent (0-5), and most of the differences (9 of 13) from one parent were shared with the other parent. Samples from the fourth offspring S2 showed a small number of differences from those of the other siblings (0-6). By contrast, samples from this offspring exhibited numerous differences with the parents (13 and 24), most of which were not evident in comparisons of samples from the parents with each other. The identities of the proteins differing among the parents and offspring are shown in Figure S1.

**Profile of Proteomically-Inferred SNP Genotypes**

Database searching of the samples by GPM Fury identified on average 550 ± 38 proteins with 2390 ± 310 unique peptides per sample (all values given as mean ± std dev), which were then checked for GVPs. A total of 181 GVPs corresponding to 96 loci were identified in datasets of the six studied individuals (**Table S4**). The replicates had on average 52 ± 9 GVPs while the cumulative data of the replicates for each individual showed 75.4 ± 3.6 GVPs. GVPs identified in the individuals P, M, S1 and S2 were validated from the parallel exomic sequencing data, and the GVPs were designated as true positive (TP), true negative (TN), false positive (FP) or false negative (FN) as previously described (Borja et al, 2019; Parker et al, 2016) . The analysis showed a total of 304 (41.7%) TP, 303 (41.6%) TN, 107 (14%) FN, and 14 (1.9%) FP assignments (**Table S4**). The GVPs were also categorized more precisely as undetected when protein regions containing them were not represented due to low yields in the MS run (**Table S5**). Previously  such GVPs were assigned to the false negative category (Borja et al, 2019; Parker et al, 2016). This modification avoids assumptions in cases where no data were provided by the MS scan and increased the negative predictive value (TN/(TN+FN) from 73.9% for data in **Table S4** to 92.8% for data in **Table S5**. The positive predictive value (TP/(TP+FP)) for the data was 95.4%. About 94% of the assumptions made were correct ((TP + TN)/(TP + TN + FP +

203   FN)) when compared to the exomic data. Moreover, because a majority of the homozygous

204   assumptions were made on the major alleles with frequencies >75%, homozygosity was the most

205   conservative assumption.

**Hierarchical Clustering**

207   To evaluate the identifying powers of the GVP profiles, the profiles of the 6 studied family

208   members were compared with those of 9 unrelated individuals. Data from a previously published

209   dataset (the 9 unrelated individuals), processed contemporaneously by the same individual using

210   an identical protocol, were merged with the current GVP dataset (Plott et al, 2020). Each GVP

211   detection was assigned a value of 1 and non-detection a value of 0. The file was then imported

212   into R and Euclidean distances between the samples were calculated. Using agglomerative

213   hierarchical clustering, similar profiles were clustered together based on the Euclidean distances.

214   The clustering showed that the GVP profiles of the 6 related individuals were more closely

215   correlated to each other than to GVP profiles of unrelated subjects (Figure 2), not likely a

216   manifestation of a batch effect of processing (Plott et al, 2020). This was the case even for the

217   sibling with an RPS23 mutation (S2) who manifested a distinct 'wiry' hair shaft phenotype with

218   low lipid levels. The results indicate a high utility of GVP profiling for forensic identification

219   purposes, especially in cases of mass fatalities when samples from the close family members are

220   available for identification, which would likely increase the power of this approach.

221

222



8

**Figure 2.** Hierarchical clustering performed using the GVP data of the currently studied family and nine unrelated individuals. The six family members clustered together (boxed), indicating similarity to each other in contrast to unrelated individuals.

**Relationship Index Using Genotypic Data Corresponding to the Detected GVPs Acquired from Exome Sequencing**

A likelihood ratio (LR) is traditionally used for relationship testing using STRs and/or SNPs where ratios >1 are evidence for individuals to be related, and the higher the LR, the stronger the evidence. However, the value of LR to indicate a relationship conclusively varies among laboratories from 1 to 10 or even 100 (American Association of Blood Banks, 2013; (Ge and Budowle, 2021). The present GVP data were analyzed using several relationship-testing approaches. Initially the corresponding SNP profiles for the GVP profiles of P, M, S1 and S2 were obtained from exome data (**Table S6**) and the profiles of S1 and S2 were tested for the likelihood that they were the offspring of the parents P and M. The likelihood ratios showed combined paternity indexes (CPIs) of 402904 and 5100 and posterior odds of 99.99% and 99.98% calculated with prior probabilities of 0.5 for S1 and S2, respectively. Sibship indexes calculated for the four individuals showed high combined sibship index (CSI) values strongly supporting a relationship for all the genetically related individuals except for M with S2 (7.2) (**Table 2**). This observation reflects a lower number of minor allelic GVPs shared by the siblings with their mother as compared to the father. At 66 of the 96 studied loci, all four analyzed members of the family were homozygous for the same allele. About 90% of this homozygosity was on the major alleles, and these loci added a CSI of 42.1 to the calculations. On the other 30 loci, S1 shared a minor allele with M and P at 6 and 8 loci while S2 shared 4 and 6 such alleles with M and P respectively.

**Table 2:** Combined sibship index values calculated using the genotype data for the GVP loci obtained only from the exome sequencing. P: father, M: mother, S1: sibling 1, S2: sibling 2.

|        | M     | S1     | S2     |
|--------|-------|--------|--------|
| **P**  | <0.01 | 161.30 | 150.97 |
| **M**  |       | 17.39  | 7.22   |
| **S1** |       |        | 47.09  |

## Relationship Index using proteomically-inferred genotypes

The evaluation of the genetic data was proceeded by the same analyses for SNPs inferred from the proteomic data. In this analysis, data from 8 European-American individuals in a previously published cohort (Plott 2020) were included to expand the GVP data from the currently studied family. Loci were assumed heterozygous if peptides encoded by both alleles and homozygous if peptides encoded by only one allele were seen in the proteomic data. Loci where none of the peptides was detected in any of the replicate samples of an individual were called as undetected or uninformative. GVPs for which the frequency of minor allele in European population was 0 were excluded from this analysis. Parentage indexes (ratios based on trio models) using P and M as parents and sibship indexes (ratios based on duo sibling models) for every possible pair of the individuals (from the family and from the additional subjects) were then calculated. The calculations were performed both including **(Table S7 and S8)** and excluding (with the rationale to not include the frequently observed false positive GVPs in real world practices) **(Table 3 and 4)** the false positive GVPs identified in the data. A locus was included in the calculation only if genotype information for that locus could be inferred for both (sibship calculation) or all three (parentage indexes calculation) individuals. Only the four actual children of the couple P and M showed CPIs and posterior probabilities that support the relationship **(Table 3)**. The one locus at which an obligate allele was not found in S2 was rs1455555 in SERPINB5, a false negative **(Table S5)** which was kept out of the CPI calculation for S2 owing to the very low mutation rate per nucleotide ($1\text{-}2 \times 10^{-8}$) (Kong et al, 2012). However, allele dropouts due to technical reasons (e.g. low volatilization of some peptides) that are much more likely in MS based proteomic analyses were taken into account by excluding peptides with a history of false negative detections. The unrelated individuals had at least three loci at which the obligate allele was not present either in the parents or the tested sample except for two (G and H) with two such loci. However, for these individuals the number of loci at which genotype information could be inferred was lower, 21 and 20, respectively.

**Table 3:** Combined paternity indexes and posterior probabilities calculated using the prior odds for the four true offspring and eight random individuals. (For each individual, the chance of being an offspring of the given parents is 4 of a total of 12 individuals or 4/12.) Using P and M as father and mother, the profiles of the 12 individuals were compared. The loci column

278    represents the number of genes used for each analysis with the values in parenthesis indicating

279    numbers of genes with two loci. CPI: combined paternity index.

| CPI when both parents are available | | | | |
|---|---|---|---|---|
| Individual | Loci used in the calculation | Loci with no obligate allele | CPI | Posterior Probability (%) |
| S3 | 24(9) | 0 | 1286.03 | 99.92 |
| S1 | 25(8) | 0 | 1676.36 | 99.94 |
| S4 | 22(8) | 0 | 258.23 | 99.61 |
| S2 | 28(9) | 1 | 1380.90 | 99.92 |
| A | 23(9) | 3 | -- | -- |
| B | 23(9) | 4 | -- | -- |
| C | 21(8) | 5 | -- | -- |
| D | 20(9) | 3 | -- | -- |
| F | 23(8) | 4 | -- | -- |
| G | 21(7) | 2 | -- | -- |
| H | 20(8) | 2 | -- | -- |
| I | 25(7) | 3 | -- | -- |

280

281    Sibship indexes were also calculated for each possible pair of the siblings and eight unrelated

282    individuals belonging to the same population. A total of 91 comparisons were made. The

283    calculated values were >10 for 13 of 14 true sibling pairs. The pair M-S1 was the only one with

284    CSI value <10 (9.75) (**Table 4**). It was observed in the hierarchical clustering that the profiles of

285    S1 and S2 were closer to the father than the mother. The low number of minor alleles shared

286    with the mother could account for the lower relationship index value for this pair. Of the 77

287    unrelated pairs, only two pairs (A-C and D-I) had CSIs >10 falsely supporting a relationship.

288    Consistent with the above calculations, the CPIs including FP-GVPs were more accurate

289    compared to CSIs because of lower genetic similarity in siblings based on Mendelian inheritance

290    patterns (1/4$^{th}$ chance of no allele identical by descent at a given locus) and the inclusion of two

291    profiles (M and P) for comparison in CPI compared to one in CSI (Table S7 and S8). Even

292    though a majority of the calculations were appropriate with a threshold CSI of 10 or greater, a

293    certain threshold for inclusion or exclusion of sibship could not be established in the present

294    study. (However, including more loci to the analysis in the future using optimized techniques

295    (Goecker et al, 2020) should overcome this problem.) The number of loci at which each analysis

11

was made are presented in **Table S9**. Nonetheless, the present findings support the usefulness of GVP profiles in statistically differentiating between related and unrelated individuals.

**Table 4:** Combined sibship index values calculated for the family members and 8 unrelated individuals. The CSI values higher than 10 for unrelated individuals or lower than 10 for true siblings are shown in bold, and the ones that support the relationship in the cases of true relationships are bold italicized.

| | P | S3 | S1 | S4 | S2 | A | B | C | D | F | G | H | I |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **M** | 3.21 | ***198.13*** | 9.75 | ***1383.55*** | ***10.45*** | 0.51 | 5.93 | 1.39 | 6.08 | 0.06 | 8.96 | 0.09 | 2.70 |
| **P** | | ***15.28*** | ***1203.59*** | ***11.04*** | ***287.77*** | 1.76 | 0.15 | 0.02 | 0.05 | 5.54 | 0.08 | 2.20 | 0.18 |
| **S3** | | | ***565.55*** | ***42.90*** | ***24.15*** | 1.41 | 0.77 | 0.34 | 2.67 | 0.25 | 0.45 | 0.07 | 2.83 |
| **S1** | | | | ***15.03*** | ***125.22*** | 0.13 | 0.08 | 0.03 | 0.07 | 0.43 | 0.03 | 0.36 | 1.12 |
| **S4** | | | | | ***35.17*** | 0.66 | 0.86 | 0.06 | 2.32 | 0.07 | 1.18 | 0.15 | 4.19 |
| **S2** | | | | | | 0.17 | 0.02 | 0.00 | 0.79 | 1.96 | 0.73 | 7.39 | 0.09 |
| **A** | | | | | | | 0.83 | **235.11** | 1.07 | 0.14 | 2.19 | 4.99 | 5.78 |
| **B** | | | | | | | | 7.87 | 0.68 | 0.02 | 0.25 | 0.04 | 0.02 |
| **C** | | | | | | | | | 8.77 | 0.02 | 5.27 | 0.95 | 0.23 |
| **D** | | | | | | | | | | 0.00 | 3.27 | 0.13 | **16.74** |
| **F** | | | | | | | | | | | 0.24 | 1.71 | 0.01 |
| **G** | | | | | | | | | | | | 0.19 | 4.95 |
| **H** | | | | | | | | | | | | | 0.03 |
| **I** | | | | | | | | | | | | | |

## DISCUSSION

This study investigates the potential for using proteomic variation, both protein abundance and amino acid sequence information, to compare measures of relatedness within and beyond the family unit. When investigating hair from unidentified remains, reference DNA may be difficult to obtain, while potentially related individuals may be available to investigators. Similarity in protein profiling, or calculations of relationship indices, may be all that can be obtained by investigators. Accordingly, different approaches to measuring relatedness were tested and compared: two way comparison of the proteomic profiles, measurement of correlation distancing using hierarchical clustering of GVP profiles, and indices of relationship using DNA and proteomic genotyping data. Like transcriptional analysis, proteomic profiling is the product of the transcriptional and translation program of each cell. A genetic role in modulating the relative expression and increased similarity in proteomic profiles within the family unit was observed. This was also true for GVP content even though one of the siblings had a genetic condition that affected the hair phenotype and the protein profile. Likewise, paternity and sibling indices using

proteomically-inferred SNP allele genotypes showed elevated scores for related compared to unrelated individuals. This demonstrates the potential for using protein levels and sequences to assist in identification of unidentified remains.

Genetic match is the strongest and most widely accepted evidence for identification of tissues procured from crime scenes, resolving relationship conflicts and/or identification of remains in mass fatalities. To this end, probabilities for marker profiles and relationship indexes from the corresponding population genetics data can be calculated based on laws of Mendelian genetics, hence assigning a statistical value for the degree of match between profiles. As manifestations of allelic differences, permitting inference of corresponding single nucleotide polymorphisms in the genomic DNA, GVP proteomic data permit judging match or mismatch like other genetic marker systems. With random match probabilities as low as one in 640 million (Goecker et al, 2020), they have a greater power of discrimination than protein profiling among related and unrelated samples regardless of the age of individuals/hair samples, anatomic collection sites, chemical treatment and exposure of hairs to extreme conditions (Chu et al, 2020; Franklin et al, 2020; Plott et al, 2020). Thus, proteomics may provide useful information in cases where DNA evidence is insufficient due to age or suboptimal storage.

The combined sibship index values are defined as the likelihood of obtaining the genetic data when the two individuals are related versus unrelated; therefore, a higher LR value supports the relationship and vice versa (Ge and Budowle, 2021). However, the likelihood threshold values for inclusion, exclusion and inconclusive results vary among laboratories. According to American Association of Blood Banks nearly 6% of the laboratories use a LR threshold of 1 and a similar number of laboratories use 10 for inclusion in testing full siblings vs unrelated, while about 20% of the laboratories use 100 (Unit, 2013). In the current data, the likelihood threshold if kept at 10 supported all the relationships, but there were two unrelated pairs with CSI values >10. On the other hand, increasing the minimum LR value supporting a relationship to 100 eliminates the false positives in the data but brings 8 of the true relations into the uncertain range (1 < LR < 100), although not excluding them completely (<1). However, it should be noted that STRs traditionally used for such analyses often exhibit a greater degree of polymorphism (numbers of tandem repeats) than SNP loci, and the number of GVPs used in the current study were lower than the number of SNP loci used earlier in similar testing (Yousefi et al, 2018).

346    Moreover, the detection of GVPs limited the sensitivity of the SNP panel rather than the

347    discriminatory powers of SNPs in a population.

348    Present GVP analysis provided promising results for relationship testing even using only 29-35

349    SNP marker loci for trio parent-child analyses **(Table 3)** and 30-49 loci in 20-28 genes for duo

350    sibship analyses. The former, including data from two individuals (mother and father) for

351    comparison and the obligation for certain alleles to be present, successfully identified the true

352    relationships. Sibship analysis, on the other hand, was less discriminatory as has been seen for

353    DNA analyses of 'duo sibship' cases. The 25% chance that two siblings will have no allele

354    identical by descent at a given diploid locus leads to difficulty solving such identification cases

355    (Lee et al, 2012). Therefore, increasing the amount of data used for the calculation both in terms

356    of GVPs and individual profiles, e.g., comparing the profile of a subject to those of two known

357    true siblings, can better discriminate among the related and unrelated individuals (Lee et al,

358    2012). In the case of STR markers, because of the higher degree of polymorphism at each locus,

359    the number of markers sufficient to discriminate successfully between individuals is relatively

360    low, 13-17, and ~30-40 for resolving second degree relationship status (Fimmers et al, 2008;

361    Presciuttini et al, 2004). This number, due to the low mutation rate and polymorphism, is far

362    higher for SNPs, ~50-150, where including a higher number of markers provides a higher power

363    of discrimination (Chang et al, 2015; Phillips et al, 2008). The same holds true for GVPs in

364    kinship analyses, since GVPs are the expressed manifestation of SNPs in the studied proteomes.

365    Recently published hair sample processing procedures improve the number of GVP

366    identifications by several fold, which will allow for more confident assignment of GVP

367    heterozygosity and will result in higher discrimination and higher indices of relatedness

368    (Goecker et al, 2020).

369    An obvious limitation of the current study is the inference of homozygosity at loci where the

370    alternate allele was not detected. Detection of a peptide encoded by only one allele of a SNP

371    locus provides half the number of markers on which the probability of match/randomness of a

372    profile are calculated provided by DNA sequencing analyses. An intrinsic limitation of GVP

373    detection when using shotgun proteomics is that the presence of an allele can be inferred but no

374    claim can be made concerning alternate alleles. This is currently addressed using genotype

375    frequencies instead of potential homozygotic frequencies for calculating random match

376    probabilities but could lead to inaccuracy if the peptide representing another allele is not detected

377    due, for example, to low volatility. This limitation will be alleviated when GVP quantitation

378    becomes more precise using targeted mass spectrometry. This is a significant issue in analyzing

379    relationships, as the kinship/relationship indexes calculations require data from both alleles at a

380    locus. However, the negative predictive value obtained using the present categorization scheme

381    improved by 20% the value for such data using an earlier approach (Borja et al, 2019, Parker,

382    2016 #2247), thereby increasing confidence in the inferences of SNP alleles. This study makes

383    two assumptions that may change with more study and investigation. This study assumes

384    homozygosity when only one GVP at a locus is detected. Of all the assumptions made for the

385    four individuals whose GVP profiles could be validated, 92.6% were correct when compared

386    with the exome data, whereas 7.2% were incorrect (3.6% were less conservative ($f\mathrm{a}^2 < f2\mathrm{ab}$) and

387    3.6% were more conservative ($f\mathrm{a}^2 > f2\mathrm{ab}$), a balanced outcome). In the future, as genotyping for

388    GVP-inferred loci improves based on proteomic workflows and instruments that are more

389    sensitive and quantitative, this assumption will become moot since the status of the alternate

390    allele based on GVP quantitation could be inferred directly from proteomic detection. The

391    second assumption is that GVP-inferred SNP loci were statistically independent unless they fell

392    in the same gene. Observed SNP locus combinations within a gene were counted in the European

393    population of the 1000 genome project to determine the genotype frequency of the diplotype,

394    consistent with previous studies (Parker et al, 2016). This assumes that linkage disequilibrium

395    dissipates beyond the gene boundary. Although this is worthy of revisiting in the future, it had

396    little impact on final paternity index values in this instance. When calculations were made using

397    both models, treating inferred SNP loci independently or expanding the boundaries of the locus

398    to incorporate the entire reading frame, there was no change in the median sibship indices, and

399    only 2 of 91 changes in concluded relationships (i.e,. sibship index < 10, data not shown). Even

400    so, peptides that are consistently or frequently undetected using a certain protocol should be

401    noted and not included in the analysis. Examples in the current dataset are rs3744786_T in

402    KRT32, rs17843021_A in KRT39, rs2852464_C in KRT83, rs951773_A in KRT 84,

403    rs9636845_T in KRTAP11, and rs13070515_A in LRRC15 **(Table S5)**. Moreover, employing

404    different hair processing protocols, MS instruments or data acquisition strategies will lead to

405    detection of a different set of peptides and proteins, thereby affecting the GVPs detected

406    downstream. Nonetheless, the current study provides a basis and demonstrates feasibility for the

407    use of GVPs in analyzing relationship status.

408 Proteomic profiling, as applied in this study, used label free quantification. Differential protein

409 expression analysis was based on weighted spectral counts obtained from the Scaffold software.

410 This type of label free quantitation is commonly used in proteomics for judging variation in a

411 given protein's level among parallel samples (Dowle et al, 2016; Liu et al, 2004). Consistent

412 with the expected correlation of hair protein profiles within the family, the profiles from three

413 offspring were intermediate between those of the parents in two-way comparisons. Samples from

414 the fourth offspring were distinctly different from both parents, however. The latter finding can

415 possibly be attributed to departure of the hair from an unaffected phenotype due to a *de novo*

416 heterozygous mutation (c.200G>A) in the ribosomal protein RPS23 (Paolini et al, 2017),

417 although a connection to the observed perturbation of hair shaft protein levels in offspring S2 is

418 not obvious. The genetic bases for numerous hair abnormalities are known, and others remain to

419 be discovered (Duverger and Morasso, 2014). We speculate this example could illustrate how a

420 genetic defect could result in an unusual phenotype due to loss of a critical protein or to

421 perturbation of expression levels of a group of proteins in an intracellular signaling pathway.

422 Proteomic analysis could potentially assist in diagnoses or help connect genotype and phenotype

423 if the abnormalities manifested characteristic protein profiles.

## CONCLUSION

425 The major significance of the present work for forensic casework is that GVP analysis of hair

426 evidence offers a viable approach to testing familial relationships. The results obtained

427 complement, and can be combined with, those from mitochondrial DNA analysis. Results from

428 protein profiling, although not readily applicable to calculating random match probabilities,

429 would be expected to support the outcomes of GVP analysis. Discrepancies in protein expression

430 level that do not fit expectation within a family could be indicative of genetic differences not

431 evident by GVP analysis. Such cases may be useful in discovery and characterization of genetic

432 hair abnormalities.

**Supplementary File Legends**

434 **Figure S1.** Two way comparisons of hair protein levels among offspring and parents. Each Venn

435 diagram shows the number of significant differences in samples from the father (P) and mother

436 (M) with each other and with one sibling (S1-S4). Proteins in blue are those significantly

437 different in amount from the mother in samples from sibling and father, while those in red are

438  those different from the father in samples from the mother and sibling. The two way differences

439  between the family members are tabulated in the inset. Note S2 exhibited many more differences

440  than the other siblings with P and M.

441  **Table S1.** Loci at which the genotype obtained from exome data of S2 was not consistent with

442  the parents P and M. Assignments consistent with proteomic data are listed as "corrected".

443  **Table S2.** GVP data matrix used for hierarchical clustering. Each GVP detection was assigned a

444  value of 1 and a non-detection of 0.

445  **Table S3.** Formulae to calculate paternity indices and sibship indices. Capital letters indicate

446  alleles whereas lower case letters indicate the allele frequencies from 1000 Genome Project

447  (Consortium et al, 2015).

448  **Table S4.** Cumulative GVP profiles identified in the six members of the family. The GVPs from

449  P, M, S1 and S2 were validated from the corresponding genomic data. True positive

450  identifications are highlighted in blue, true negative as white, false positive as red and false

451  negative as green.

452  **Table S5.** GVPs identified in the six members of the family. The GVPs from P, M, S1 and S2

453  were validated from the corresponding genomic data. True positive identifications are

454  highlighted in blue, true negative as white, false positive as red and false negative as green.

455  GVPs present in the protein regions that were not sequenced in the MS runs were called as

456  undetected and highlighted as grey.

457  **Table S6.** Genotypes of individual P, M, S1 and S2 for the identified genetically variant

458  peptides. Genotypes at the five dubious loci are highlighted in bold italic.

459  **Table S7.** Combined paternity indexes and posterior probabilities calculated using all the

460  detected GVPs including false positives. The posterior probabilities were obtained using the

461  prior odds of 4/12 for the four true offspring and eight random individuals. Using P and M as

462  father and mother, the profiles of the 12 individuals were compared. CPI: combined paternity

463  index.

464 **Table S8.** Combined sibship index values calculated using all the detected GVPs including false

465 positives for the family members and eight unrelated individuals. The CSI values higher than 10

466 for unrelated individuals or lower than 10 for true siblings are shown in bold, and the ones that

467 support the relationship in the cases of true relationships are bold italicized.

468 **Table S9.** Number of loci at which each comparison was based for CSI calculations. The

469 numbers inside the parentheses represent the genes with two loci included.

470 **Table S10.** Pairwise comparisons of protein levels in samples from the parents and siblings.

471 Shown are the fold difference (FC) for each protein, calculated p values before (P.Value) and

472 after (adj.P.Val) correction for multiple testing, identified proteins, accession numbers in the

473 Uniprot human database and the protein molecular weight for each (MW).

474 **References**

475 1000 Genomes Project Consortium, Auton A, Brooks LD, Durbin RM, Garrison EP, Kang
476 HM, Korbel JO, Marchini JL, McCarthy S, McVean GA, Abecasis GR (2015) A global
477 reference for human genetic variation. Nature 526:68-74. doi: 10.1038/nature15393
478 Alsop RJ, Soomro A, Zhang Y, Pieterse M, Fatona A, Dej K, Rheinstädter MC (2016)
479 Structural abnormalities in the hair of a patient with a novel ribosomopathy. PLoS One
480 11(3):e0149619. doi: 10.1371/journal.pone.0149619
481 American Association of Blood Banks (2013) Annual Report Summary for Testing in 2013.
482 Relationship Testing Program Unit. https://www.aabb.org/docs/default-source/default-document-
483 library/accreditation/2013-relationship-testing-summary-report.pdf?sfvrsn=da4315b2_2
484 Benjamini Y, Hochberg Y (1995) Controlling the false discovery rate: A practical and
485 powerful approach to multiple testing. J Royal Stat Soc, Ser B 57:289-300. doi: 10.1016/0306-
486 9877(95)90228-7
487 Borja T, Karim N, Goecker Z, Salemi M, Phinney BS, Naeem M, Rice RH, Parker GJ (2019)
488 Proteomic genotyping of fingermark donors with genetically variant peptides. Foren Sci Int:
489 Genet 42:21-30. doi: 10.1016/j.fsigen.2019.05.005
490 Chang CC, Chow CC, Tellier LC, Vattikuti S, Purcell SM, Lee JJ (2015) Second-generation
491 PLINK: rising to the challenge of larger and richer datasets. Gigascience 4:7. doi:
492 10.1186/s13742-015-0047-8
493 Chu F, Mason KE, Anex DS, Jones AD, Hart BR (2019) Hair proteome variation at different
494 body locations on genetically variant peptide detection for protein-based human identification.
495 Sci Rep 9(1):7641. doi: 10.1038/s41598-019-44007-7
496 Chu F, Mason KE, Anex DS, Jones AD, Hart BR (2020) Proteomic characterization of
497 damaged single hairs recovered after an explosion for protein-based human identification. J
498 Proteome Res 19:3088-3099. doi: 10.1021/acs.jproteome.0c00102
499 Dowle AA, Wilson J, Thomas JR (2016) Comparing the diagnostic classification accuracy of
500 iTRAQ, peak-area, spectral-counting, and emPAI methods for relative quantification in
501 expression proteomics. J Proteome Res 15:3550-3562. doi: 10.1021/acs.jproteome.6b00308

18

502    Duverger O, Morasso MI (2014) To grow or not to grow: Hair morphogenesis and human
503    genetic hair disorders. Seminars Cell Develop Biol 25-26:22-33. doi:
504    10.1016/j.semcdb.2013.12.006
505    Fimmers R, Baur M, Rabold U, Seifried E, Seidl C (2008) STR-profiling for the
506    differentiation between related and unrelated individuals in cases of citizen rights. For Sci Intl:
507    Genet Suppl Ser 1:510-513. doi: 10.1016/j.fsigss.2008.01.002
508    Franklin RN, Karim N, Goecker ZC, Durbin-Johnson BP, Rice RH P, G J (2020) Proteomic
509    genotyping: Using mass spectrometry to infer SNP genotypes in pigmented and non-pigmented
510    hair. Foren Sci Int 310:110200. doi: 10.1016/j.forsciint.2020.110200
511    Ge J, Budowle B (2021) Forensic investigation approaches of searching relatives in DNA
512    databases. J Forensic Sci 66:430-443. doi: 10.1111/1556-4029.14615
513    Goecker ZC, Salemi MR, Karim N, Phinney BS, Rice RH, Parker GJ (2020) Optimal
514    processing for proteomic genotyping of single human hairs. Foren Sci Int: Genet 47:102314. doi:
515    10.1016/j.fsigen.2020.102314
516    Hindorff LA, Sethupathy P, Junkinsa HA, Ramosa EM, Mehtac JP, Collins FS, Manolio TA
517    (2009) Potential etiologic and functional implications ofgenome-wide association loci for human
518    diseasesand traits. Proc Natl Acad Sci USA 106:9362-9367. doi: 10.1073pnas.0903103106
519    Kong A, Frigge ML, Masson G, Besenbacher S, Sulem P, Magnusson G, Gudjonsson SA,
520    Sigurdsson A, Jonasdottir A, Jonasdottir A, Wong W, Sigurdsson G, Walters GB, Steinberg S,
521    Helgason H, Thorleifsson G, Gudbjartsson DF, Helgason A, Magnusson OT, Thorsteinsdottir U,
522    Stefansson K (2012) Rate of de novo mutations and the importance of father's age to disease risk.
523    Nature 488:471-475. doi: 10.1038/nature11396
524    Laatsch CN, Durbin-Johnson BP, Rocke DM, Mukwana S, Newland AB, Flagler MJ, Davis
525    MG, Eigenheer RA, Phinney BS, Rice RH (2014) Human hair shaft proteomic profiling:
526    individual differences, site specificity and cuticle analysis. PeerJ 2:e506. doi: 10.7717/peerj.506
527    Lee JC, Lin YY, Tsai LC, Lin CY, Huang TY, Chu PC, Yu YJ, Linacre A, Hsieh HM (2012)
528    A novel strategy for sibship determination in trio sibling model. Croat Med J 53:336-342. doi:
529    10.3325/cmj.2012.53.336
530    Linch CA, Whiting DA, Holland MM (2001) Human hair histogenesis for the mitochondrial
531    DNA forensic scientist. J Forensic Sci 46:844-853. doi: 10.1520/JFS15056J
532    Liu H, Sadygov RG, Yates JRI (2004) A model for random sampling and estimation of
533    relative protein abundance in shotgun proteomics. Analyt Chem 76:4193-4201. doi:
534    10.1021/ac0498563
535    Martin-Trujillo A, Patel N, Felix Richter F, Bharati Jadhav B, Garg P, Morton SU, McKean
536    DM, R DS, Goldmuntz E, Gruber D, Kim BR, Jane W. Newburger JW, Porter GAJ, Alessandro
537    Giardini A, Bernstein D, Tristani-Firouzi M, Seidman JG, Seidman CE, Chung WK, Gelb BD,
538    Sharp AJ (2020) Rare genetic variation at transcription factor binding sites modulates local DNA
539    methylation profiles. PLoS Genetics 16(11):e1009189. doi: 10.1371/journal.pgen.1009189
540    McNevin D, Wilson-Wilde L, Robertson J, Kyda J, Lennard C (2005) Short tandem repeat
541    (STR) genotyping of keratinised hair. Part 1. Review of current status and knowledge gaps.
542    Foren Sci Int 153:237-246. doi: 10.1016/j.forsciint.2005.05.005
543    Milan J, Wu P-W, Salemi M, Durbin-Johnson B, Rocke DM, Phinney BS, Rice RH, Parker
544    GJ (2019) Comparison of protein expression levels and proteomically-inferred genotypes using
545    human hair from different body sites. Foren Sci Int: Genet 41:19-23. doi:
546    10.1016/j.fsigen.2019.03.009
547    National Research Council (2009) Strenghtening Forensic Science in the United States: A
548    Path Forward. The National Academies Press, pp 155-161. doi: 10.17226/12589

549      Paolini NA, Attwood M, Sondalle SB, dos Santos Vieira CM, van Adrichem AM, di Summa
550    FM, O'Donohue M-F, Gleizes P-E, Rachuri S, Briggs JW, Fischer R, Ratcliffe PJ, Wlodarski
551    MW, Houtkooper RH, von Lindern M, Kuijpers TW, Dinman JD, Baserga SJ, Cockman ME,
552    MacInnes AW (2017) A ribosomopathy reveals decoding defective ribosomes driving human
553    dysmorphism. Am J Hum Genet 100:506-522. doi: 10.1016/j.ajhg.2017.01.034
554      Parker GJ, Leppert T, Anex DS, Hilmer JK, Matsunami N, Baird L, Stevens J, Parsawar K,
555    Durbin-Johnson BP, Rocke DM, Nelson C, Fairbanks DJ, Wilson AS, Rice RH, Woodward SR,
556    Bothner B, Hart H, Leppert M (2016) Demonstration of protein-based human identification using
557    the hair shaft proteome. PLoS One 11(9):e0160653. doi: 10.1371/journal.pone.0160653
558      Phillips C, Fondevila M, García-Magariños M, Rodriguez A, Salas A, Carracedo A, Lareu
559    MV (2008) Resolving relationship tests that show ambiguous STR results using autosomal SNPs
560    as supplementary markers. Foren Sci Int: Genet 2:198-204. doi: 10.1016/j.fsigen.2008.02.002
561      Plott TJ, Karim N, Durbin-Johnson BP, Swift DP, Youngquist RS, Salemi M, Phinney BS,
562    Rocke DM, Davis MG, Parker GJ, Rice RH (2020) Age-related changes in hair shaft protein
563    profiling and genetically variant peptides Foren Sci Int: Genet 47:102309. doi:
564    10.1016/j.fsigen.2020.102309
565      Presciuttini S, Toni C, Marronia F, Spinetti I, Bailey-Wilson JE, Domenici R (2004) The
566    number of STR markers necessary to resolve relationships in deficiency paternity cases. Int
567    Congr Ser 1261:541-543. doi: doi.org/10.1016/S0531-5131(03)01664-9
568      Rice RH, Bradshaw KM, Durbin-Johnson BP, Rocke DM, Eigenheer RA, Phinney BS,
569    Sundberg JP (2012) Differentiating inbred mouse strains from each other and those with single
570    gene mutations using hair proteomics. PLoS One 7:e51956. doi: 10.1371/journal.pone.0051956
571      Ritchie ME, Phipson B, Wu D, Hu Y, Law CW, Shi W, Smyth GK (2015) limma powers
572    differential expression analyses for RNA-sequencing and microarray studies. Nucl Acids Res
573    43(7):e47. doi: 10.1093/nar/gkv007
574      Robinson MD, Oshlack A (2010) A scaling normalization method for differential expression
575    analysis of RNA-seq data. Genome Biol 11:R25. doi: 10.1186/gb-2010-11-3-r25
576      Sozer A, Baird M, Beckwith M, Harmon B, Lee D, Riley G, Schmitt S (2010) Guidelines for
577    mass fatality DNA identification operations. American Association of Blood Banks.
578    https://www.aabb.org/docs/default-source/default-document-library/about/guidelines-for-mass-
579    fatality-dna-identification-operations.pdf?sfvrsn=af1c96a9_0
580      Wenk RE, Traver M, Chiafari FA (1996) Determination of sibship in any two persons.
581    Transfusion 36:259-262. doi: 10.1046/j.1537-2995.1996.36396182146.x
582      Wu P-W, Mason KE, Durbin-Johnson BP, Salemi M, Phinney BS, Rocke DM, Parker GJ,
583    Rice RH (2017) Proteomic analysis of hair shafts from monozygotic twins: Expression profiles
584    and genetically variant peptides. Proteomics 17:13-14, 1600462. doi: 10.1002/pmic.201600462
585      Yousefi S, Abbassi-Daloii T, Kraaijenbrink T, Vermaat M, Mei H, van 't Hof P, van Iterson
586    M, Zhernakova DV, Claringbould A, Lude Franke L, 't Hart LM, Slieker RC, van der Heijden A,
587    de Knijff P, consortium B, 't Hoen PAC (2018) A SNP panel for identification of DNA and RNA
588    specimens. BMC Genomics 19(1):90. doi: 10.1186/s12864-018-4482-7

# Elucidation of Familial Relationships Using Hair Shaft Proteomics

**Supplementary Material**

**(Figure S1, Tables S1-S9)**



**Figure S1.** Two way comparisons of hair protein levels among offspring and parents. Each Venn diagram shows the number of significant differences in samples from the father (P) and mother (M) with each other and with one sibling (S1-S4). Proteins in blue are those significantly different in amount from the mother in samples from sibling and father, while those in red are those different from the father in samples from the mother and sibling. The two way differences between the family members are tabulated in the inset. Note S2 exhibited many more differences than the other siblings with P and M.

**Table S1** Loci at which the genotype obtained from exome data of S2 was not consistent with the parents P and M. Assignments consistent with proteomic data are listed as "corrected".

| Gene Name | rs# | Reference | P | M | S1 | S2 | S2 corrected |
|---|---|---|---|---|---|---|---|
| KRT32 | rs2604953 | G | TT | TT | TT | GG | TT |
| KRTAP4-1 | rs398825 | C | TT | CT | TT | CC | ?? |
| KRTAP4-9 | rs113059833 | A | AA | AA | AA | AT | AA |
| KRTAP9-2 | rs9902235 | G | GG | GG | GG | CC | ** |
| KRTAP10-6 | rs465279 | G | GG | AA | GA | GG | AG |

?? no peptides
** false positives

Analysis of the DNA profile of S2 revealed 5 loci at which the genotype was not consistent with those of the parents. These included rs2604953 (KRT32), rs398825 (KRTAP4-1), rs113059833 (KRTAP 4-9), rs9902235 (KRTAP9-2) and rs465279 (KRTAP10-6) (Table S5). This problem was attributed to the exome analysis of S2 being performed using an older technique and at a separate time from those of P, M and S1. However, the data obtained from the proteomic analysis at these loci were consistent with the parental genotypes. According to DNA sequencing, both parents were homozygous TT for rs2604953, but S2 was homozygous GG at that position. The proteomic data showed peptides supporting only the T allele in S2, consistent with the parental genotypes. Similarly, for rs113059833 the DNA data of S2 showed a heterozygous AT genotype, but the parents were homozygous for A at that position. The proteomic data for S2 showed translation products only of an A as expected from the genotypes of the parents. For rs465279, the parental genotypes were AA and GG, inconsistent with the sequence of S2 as GG whereas, in the proteomic data, peptides for both the alleles were seen (Table S3). There was no proteomic information at the locus for rs398825 in any of the S2 replicates. The GVP corresponding to rs9902235 in KRTAP9-2 seemed unreliable since the other members of the family had it as a false positive in their GVP profiles.

**Note: Table S2** is at the end of the file after Table S9.

**Table S3:** Formulae to calculate paternity indices and sibship indices. Capital letters indicate alleles whereas lower case letters indicate the allele frequencies from 1000 genome project (1000 Genomes Project Consortium et al, 2015).

| Paternity Indices Calculations | | | |
|---|---|---|---|
| **Parent 1** | **Parent 2** | **Subject** | **Formula** |
| AA | AA | AA | $1/a^2$ |
| AA | BB | AB | $1/2ab$ |
| AA | AB | AA | $1/2a^2$ |
| AA | AB | AB | $1/4ab$ |
| AB | AB | AB | $1/4ab$ |
| AB | AB | AA | $1/4a^2$ |
| AA | BC | AB | $1/4ab$ |
| AB | AC | AA | $1/4a^2$ |
| AB | BC | AB | $1/8ab$ |
| AB | BC | BC | $1/8bc$ |

| Sibship Indices Calculations | | |
|---|---|---|
| **Subject 1** | **Subject 2** | **Formula** |
| AA | AA | $(1+a)^2/(2a)^2$ |
| AA | AB | $(1+a)/4a$ |
| AB | AB | $(1+a+b+2ab)/8ab$ |
| AA | BB | $1/4$ |
| AB | AC | $(1+2a)/8a$ |

# Table S4: Cumulative GVP profiles identified in the six members of the family. The GVPs from P, M, S1 and S2 were validated from the corresponding genomic data. True positive identifications are highlighted in blue, true negative as white, false positive as red and false negative as green.

| Gene Name | rs#_nucleotide | SAP | peptide sequence | P | M | S1 | S2 | S3 | S4 |
|---|---|---|---|---|---|---|---|---|---|
| ALDH2 | rs671_G | E504K | ELGEYGLQAYTEVK | blue | blue | blue | blue | 1 | 1 |
| ALDH2 | rs671_A | E504K | ELGEYGLQAYTk | | | | | | |
| ATG9B | rs7804893_T | N493S | HFNELPHELR | blue | blue | blue | blue | 1 | |
| ATG9B | rs7804893_C | N493S | HFsELPHELR | | | | | | |
| ATP5A1 | rs79011243_C | A32S | VLSIGDGIAR | blue | | blue | blue | 1 | 1 |
| ATP5A1 | rs79011243_A | A32S | VLSIGDGIsR | | | | | | |
| CSRP1 | rs3738283_T | K108I | HEEAPGHRPTTNPNASK | blue | blue | blue | blue | 1 | 1 |
| CSRP1 | rs3738283_A | K108I | HEEAPGHRPTTNPNASiFAQK | | | | | | |
| DSC3 | rs276937_A | S78T | VLNDGSVYTAR | | blue | green | blue | | |
| DSC3 | rs276937_T | S78T | VLNDGtVYTAR | green | | green | green | | |
| DSC3 | rs35296997_T | K180Q | GVDKEPLNLFYIER | green | blue | blue | blue | 1 | |
| DSC3 | rs35296997_G | K180Q | GVDqEPLNLFYIER | | | | | | |
| DSP | rs80325569_G | G939S | NLHSEISGK | blue | blue | blue | blue | 1 | |
| DSP | rs80325569_A | G939S | NLHSEISsK | | | | | | |
| DSP | rs2076299_A | Y1512C | VQYDLQK | green | green | green | blue | | 1 |
| DSP | rs2076299_G | Y1512C | VQcDLQK | green | | | | | |
| DSP | rs28763966_C | N1526K | ANSSATETINK | green | green | green | blue | 1 | |
| DSP | rs28763966_A | N1526K | ANSSATETIk | | | | | | |
| DSP | rs6929069_A | R1738Q | GqSEADSDKNATILELR | | | | | | |
| DSP | rs6929069_G | R1738Q | GRSEADSDKNATILELR/SEADSDKNATILELR | blue | blue | blue | blue | 1 | 1 |
| DSP | rs28763967_C | R1537C | VQEQELTR | blue | blue | blue | blue | 1 | 1 |
| DSP | rs28763967_T | R1537C | VQEQELTcLR | | | | | | |
| FAM83H | rs9969600-C | Q201H | VNLQHVDFLR | | | | | | |
| FAM83H | rs9969600-A/G | Q201H | VNLhHVDFLR | green | green | blue | green | | |
| GSDMA | rs3894194_A | R18Q | QLNPqGDLTPLDSLIDFK | | | | | | |
| GSDMA | rs3894194_G | R18Q | QLNPR/GDLTPLDSLIDFK | blue | blue | blue | blue | 1 | 1 |
| GSDMA | rs7212938_G | V128L | ALETVQER | | | | | | |
| GSDMA | rs7212938_T | V128L | ALETiQER | blue | green | green | green | | 1 |
| GSDMA | rs56030650_A | T314N | GHEVnLEALPK | green | | | | | |
| GSDMA | rs56030650_C | T314N | GHEVTLEALPK | green | green | blue | green | | |
| GSTP1 | rs1138272_C | A114V | YISLIYTNYEAGKDDYVK | blue | green | blue | blue | 1 | 1 |
| GSTP1 | rs1138272_T | A114V | YISLIYTNYEvGKDDYVK | | | | | | |
| GSTP1 | rs1695_A | I105V | YISLIYTNYEAGKDDYVK | blue | blue | blue | blue | | |
| GSTP1 | rs1695_G | I105V | YvSLIYTNYEAGKDDYVK | green | green | | blue | | |
| HEXB | rs10805890_A | I207V | GILIDTSR | blue | blue | blue | blue | 1 | 1 |
| HEXB | rs10805890_G | I207V | GILvDTSR | | | | | | |
| HEXB | rs77499935_A | I420V | LAPGTIVEVWKDSAYPEELSR/LAPGTIVEVWK | blue | blue | blue | blue | 1 | 1 |
| HEXB | rs77499935_G | I420V | LAPGTvVEVWKDSAYPEELSR | | | | | | |
| JUP | rs41283425_C | R142H | SAIVHLINYQDDAELATR | blue | blue | blue | blue | 1 | 1 |
| JUP | rs41283425_T | R142H | SAIVHLINYQDDAELAThALPELTK | | | | | | |
| JUP | rs143043662_C | V648I | NEGTATYAAAVLFR | blue | blue | blue | blue | 1 | 1 |
| JUP | rs143043662_T | V648I | NEGTATYAAAiLFR | | | | | | |
| KRT1 | rs17678945_A | A454S | NKLNDLEDALQQsKEDLAR/LNDLEDALQQsK | | | | | | |
| KRT1 | rs17678945_C | A454S | NKLNDLEDALQQAKEDLAR/LNDLEDALQQAK | blue | green | green | blue | 1 | 1 |
| KRT31 | rs6503627_A | A82V | DNvELENLIR/QLERDNvELENLIR | | blue | | red | 1 | 1 |

| Gene | SNP | Variant | Peptide | | | | | | |
|------|-----|---------|---------|---|---|---|---|---|---|
| KRT32 | rs2071561_T | S222Y | ADLEAQVEyLK | | | | | 1 | |
| KRT32 | rs72830046_C/rs | R280H/C | CQYEAMVEANRR | | | | | 1 | 1 |
| KRT32 | rs72830046_T | R280H | CQYEAMVEANhR | | | | | 1 | |
| KRT32 | rs2604953_G | P427T | SLLENEDCKLPCNPCSTPSCTTCVPSPCVPR/LPCNPCSTPSCT TCVPSPCVPR | | | | | | |
| KRT32 | rs2604953_T | P427T | SLLENEDCKLPCNPCSTPSCTTCVPSPCVtR/LPCNPCSTPSCT TCVPSPCVtR | | | | | 1 | |
| KRT32 | rs3744786_T | Q72R | TYLSSSCQAASGISGSMGPGSWYSEGAFNGNEK | | | | | | |
| KRT32 | rs3744786_C | Q72R | TYLSSSCr | | | | | 1 | |
| KRT32 | rs2071560_A | I171T | MVVNIDNAK | | | | | 1 | 1 |
| KRT32 | rs2071560_G | I171T | MVVNtDNAK | | | | | | |
| KRT32 | rs146792525_C | A255T | LNIEVDAAPPVDLTR | | | | | 1 | 1 |
| KRT32 | rs146792525_T | A255T | LNIEVDtAPPVDLTR | | | | | | |
| KRT33A | rs373657561_C | G33R | PCVPPSCHGCTLPGACNIPANVSNCNWFCEGSFNGSEK | | | | | 1 | 1 |
| KRT33A | rs373657561_T | G33R | PCVPPSCHGCTLPr | | | | | | |
| KRT33A | rs12937519_A | A270V | QVVSSSEQLQSYQvEIIELR | | | | | 1 | 1 |
| KRT34 | rs2071599_T | H348R | DSLENTLTESEAHYSSQLSQVQSLITNVESQLAEIR | | | | | 1 | 1 |
| KRT35 | rs743686_A | S36P | VSAMYSSSSCKLPSLSPVAR | | | | | 1 | 1 |
| KRT35 | rs743686_G | S36P | VSAMYSSSpCKLPSLSPVAR | | | | | 1 | 1 |
| KRT35 | rs200355130-C | E141D | LVVEIDNAK | | | | | 1 | 1 |
| KRT35 | rs200355130-A | E141D | LVVdIDNAK | | | | | 1 | 1 |
| KRT35 | rs138303882_A | R163W | YETEVSLwQLVESDINGLR | | | | | | |
| KRT35 | rs138303882_G | R163W | YETEVSLRQLVESDINGLR/QLVESDINGLR | | | | | 1 | 1 |
| KRT36 | rs75790652_G | A202G | CQLGDRLNVEVDAAPPVDLNK/LNVEVDAAPPVDLNK | | | | | 1 | 1 |
| KRT36 | rs75790652_C | A202G | CQLGDRLNVEVDgAPPVDLNK | | | | | | |
| KRT36 | rs11657323_T | N357T | YSSQLAQMQCLISNVEAQLSEIR | | | | | | |
| KRT36 | rs11657323_G | N357T | YSSQLAQMQCLIStVEAQLSEIR | | | | | | |
| KRT36 | rs9904102_G | R277C | CQYEALVENNR | | | | | | |
| KRT36 | rs9904102_A | R277C | CQYEALVENNcR | | | | | | |
| KRT39 | rs17843021_G | T341M | DSQECILTETEAR | | | | | 1 | 1 |
| KRT39 | rs17843021_A | T341M | DSQECILmETEAR | | | | | | |
| KRT39 | rs7213256_C | R456Q | SGAIESTAPACTSSSPCSLKEHCSACGPLSR/EHCSACGPLSR/E HCSACGPLSRILVK | | | | | 1 | 1 |
| KRT39 | rs7213256_T | R456Q | SGAIESTAPACTSSSPCSLKEHCSACGPLSqLLVK/EHCSACGPL SqLLVK/EHCSACGPLSqILVK | | | | | | |
| KRT39 | rs17843023_G | L383M | QNQEYEILLDVK | | | | | 1 | 1 |
| KRT39 | rs17843023_T | L383M | QNQEYEILmDVK | | | | | | |
| KRT39 | rs112557906_G | S423F | CEPSPWTSCK | | | | | 1 | 1 |
| KRT39 | rs112557906_A | S423F | CEPSPWTfCK | | | | | | |
| KRT39 | rs142154718_C | S86N | FSLDDCSWYGEGINSNEK | | | | | 1 | 1 |
| KRT39 | rs142154718_T | S86N | FSLDDCnWYGEGINSNEK | | | | | | |
| KRT40 | rs2010027_C | R235H | NHEEEVNLLREQLGDR/NHEEEVNLLR | | | | | 1 | 1 |
| KRT40 | rs2010027_T | R235H | NHEEEVNLLhEQLGDR | | | | | | |
| KRT40 | rs140634473_C | R108H | R.SLEETNAELESR | | | | | 1 | 1 |
| KRT40 | s140634473_T | R108H | V h SLEETNAELESR | | | | | | |
| KRT7 | rs6580870_A | H186R | NKYEDEINHRTAAENEFVVLK | | | | | | |

| Gene | SNP | Variant | Peptide | | | | | Count | Count |
|------|-----|---------|---------|---|---|---|---|---|---|
| KRT7 | rs6580870_G | H186R | NKYEDEINrRTAAENEFVVLK | 🟩 | 🟩 | 🟦 | 🟩 | | |
| KRT81 | rs6580873_A | L248R | LYEEEILILQSHISDTSVVVK | 🟦 | 🟦 | 🟦 | 🟥 | 1 | 1 |
| KRT82 | rs2658658_A | T458M | GAFLYEPCGVSmPVLSTGVLR | 🟦 | 🟦 | | 🟦 | | 1 |
| KRT82 | rs2658658_G | T458M | GAFLYEPCGVSTPVLSTGVLR | 🟦 | 🟦 | | | 1 | 1 |
| KRT82 | rs1732263_C | E452D | GAFLYEPCGVSTPVLSTGVLR | 🟦 | 🟩 | 🟩 | 🟦 | | |
| KRT82 | rs1732263_G | E452D | GAFLYdPCGVSTPVLSTGVLR | | | | | | |
| KRT82 | rs1791634_C | E219Q | KYEEELSLRPCVENEFVALK | 🟦 | 🟦 | 🟦 | 🟦 | 1 | 1 |
| KRT82 | rs1791634_G | E219Q | KYEEELSLRPCVqNEFVALK | | | | | | |
| KRT83 | rs61485872_A | C23G | PGNFSCVSACGPR | 🟦 | 🟦 | 🟦 | 🟦 | | 1 |
| KRT83 | rs61485872_C | C23G | PGNFSCVSAgGPR | | | | | | |
| KRT83 | rs2852464_C | I279M | DLNMDCmVAEIK | 🟩 | | 🟩 | 🟩 | | |
| KRT83 | rs2852464_G | I279M | DLNMDCIVAEIK | 🟦 | 🟦 | 🟦 | 🟦 | 1 | 1 |
| KRT83 | rs2857663_G | R149C | LQFYQNR.ECCQSNLEPLFAGYIETLR/LQFYQNR | 🟦 | 🟦 | 🟦 | 🟦 | 1 | 1 |
| KRT83 | rs2857663_A | R149C | LQFYQNcECCQSNLEPLFAGYIETLR | | | | | | |
| KRT84 | RS951773_A | C446R | CEYQELMNAKLGLDIEIATYR | | 🟩 | 🟩 | 🟩 | | |
| KRT84 | RS951773_G | C446R | QLrEYQELMNAKLGLDIEIATYR | 🟦 | 🟦 | 🟦 | 🟦 | 1 | 1 |
| KRT85 | rs2852471-C | W155L | WQFYQNQR | 🟦 | 🟦 | 🟦 | 🟦 | 1 | 1 |
| KRT85 | rs2852471-A | W155L | lQFYQNQR | 🟥 | 🟥 | 🟥 | 🟥 | 1 | 1 |
| KRTAP1-5 | rs148449559_G | T32S | TCCQTSFCGYPSFSISGTCGSSCCQPSCCETSCCQPR | 🟩 | 🟩 | 🟦 | 🟩 | | |
| KRTAP1-5 | rs148449559_C | T32S | TCCQTSFCGYPSFSISGTCGSSCCQPSCCEsSCCQPR | | | | | | |
| KRTAP1-5 | rs62623375_C | C35Y | MTCCQTSFCGYPSFSISGTCGSSCCQPSCCETSCCQPR | 🟩 | 🟦 | 🟦 | 🟦 | 1 | |
| KRTAP1-5 | rs62623375_T | C35Y | MTCCQTSFCGYPSFSISGTCGSSCCQPSCCETSCyQPR | | | | | | |
| KRTAP3-2 | rs9897046_T | S8G | MDCCASRSCSVPTGPATTICSSDKSCR | 🟦 | 🟦 | 🟦 | 🟦 | 1 | 1 |
| KRTAP3-2 | rs3829598_G | R27C | SCSVPTGPATTICSSDKSCR | 🟦 | 🟦 | 🟦 | 🟦 | | |
| KRTAP3-2 | rs3829598_A | R27C | K.SCCCGVCLPSTCPHTVWLLEPTCCDNCPPPCHIPQPCVPTCFLLNSCQPTPGLETLNLTTFTQPCCEPCLPR.G | | | | | | |
| KRTAP3-2 | rs3813050_A | I46T | CGVCLPSTCPHTVWLLEPICCDN | 🟦 | 🟩 | 🟦 | 🟦 | | |
| KRTAP4-1 | rs398825_C | T134A | TTCCRPSCCGSSC- | | 🟦 | | 🟩 | 1 | 1 |
| KRTAP4-1 | rs398825_T | T134A | aTCCRPSCCGSSC- | 🟩 | 🟩 | 🟩 | | 1 | 1 |
| KRTAP4-3 | rs428371_G | P152S | PACCISSCCHPSCCVSSCR | 🟦 | 🟦 | 🟦 | 🟦 | | 1 |
| KRTAP4-3 | rs428371_A | P152S | sACCISSCCHPSCCVSSCR | | | | | | |
| KRTAP4-4 | rs366700_C | R154S | TTCCRPSCCVSRCYR/TTCCRPSCCVSR/TTCCRPSCCVSRCYRPHCGQSLCC- | 🟦 | 🟦 | 🟦 | 🟦 | 1 | 1 |
| KRTAP4-4 | rs366700_G | R154S | TTCCRPSCCVSsCYR/TTCCRPSCCVSsCYRPHCGQSLCC- | | | | | | |
| KRTAP4-4 | rs385055_T | Y25C | VNSCCGSVCSDQGCGLENCCRPSYCQTTCCR | 🟦 | 🟦 | 🟦 | 🟦 | 1 | 1 |
| KRTAP4-4 | rs75030409_T | Q109R | TTCCRPSCCRPQCC | 🟦 | 🟦 | 🟦 | 🟦 | | 1 |
| KRTAP4-4 | rs75030409_C | Q109R | TTCCRPSCCRPr | | | | | | |
| KRTAP4-6 | rs73983172_G | P63S | R.TTCCRPSCCVSSCCRPQCCQSVCCQPTCCRPSCCPSCCQTTCCR.T | 🟦 | 🟩 | 🟩 | 🟦 | | 1 |
| KRTAP4-9 | rs149483591_G | R26H | VSSCCGSVCSDQGCGQDLCQETCCR | 🟩 | 🟩 | 🟩 | 🟦 | | 1 |
| KRTAP4-9 | rs149483591_A | R26H | VSSCCGSVCSDQGCGQDLCQETCChPSCCETTCCR | | | | | | |
| KRTAP4-9 | rs113059833_A | D18V | VSSCCGSVCSDQGCGQDLCQETCCRPSCCETTCCR | 🟩 | 🟩 | 🟩 | 🟦 | | 1 |
| KRTAP4-9 | rs113059833_T | D18V | VSSCCGSVCSDQGCGQvLCQETCCRPSCCETTCCR | | | | 🟩 | | |
| KRTAP5-2 | rs35925287_C | G29R | GCGSGCGGCGSSCGGCGSGCGGCGSGR | 🟦 | 🟦 | 🟦 | 🟦 | | 1 |
| KRTAP5-2 | rs35925287_T | G29R | GCGSGCGGCGSSCGGCGSGCr | | | | | | |
| KRTAP9-2 | rs9902235_C | C56S | CRPTsCQNTCCR | 🟥 | 🟥 | 🟥 | 🟦 | 1 | 1 |

| Gene | rsID | AA | Peptide | | | | | | |
|------|------|-----|---------|---|---|---|---|---|---|
| **KRTAP9-2** | rs9902235_G | C56S | CRPTCCQNTCCR | 🟩 | 🟦 | | | 1 | 1 |
| **KRTAP9-4** | rs2191379_A | S146Y | R.TCYYPTTVCLPGCLNQSCGSNCCQPCCRPACCETTCFQPTCVySCCQPFCC- | 🟩 | 🟩 | 🟦 | | | |
| **KRTAP9-4** | rs2191379_C | S146Y | RTCYYPTTVCLPGCLNQSCGSNCCQPCCRPACCETTCFQPTCVSSCCQPFCC- | | 🟩 | | | | |
| **KRTAP10-6** | rs465279_A | S300P | SSSSVSLLCHPVCK | | 🟦 | 🟦 | 🟥 | 1 | |
| **KRTAP10-12** | rs61745911_G | C236Y | LASCGSLLCR | 🟦 | 🟦 | 🟩 | 🟦 | 1 | |
| **KRTAP10-12** | rs61745911_A | C236Y | LASCGSLLyR | | | | | | |
| **KRTAP10-12** | rs34302939_G | G226S | RVPVPSCCVPTSSCQPSCGR/VPVPSCCVPTSSCQPSCGR | 🟦 | 🟦 | 🟦 | | 1 | |
| **KRTAP10-12** | rs34302939_A | G226S | RVPVPSCCVPTSSCQPSCsR | | | | | | |
| **KRTAP11-1** | rs71321355_C | R72Q | CIVPVAQVTTTSTTDADCLGGICLPSSFQTGSWLLDHCQETCCEPTACQPTCYRR/R.RTSCVSNPCQVTCSR | 🟦 | 🟦 | 🟦 | 🟦 | 1 | 1 |
| **KRTAP11-1** | rs71321355_T | R72Q | CIVPVAQVTTTSTTDADCLGGICLPSSFQTGSWLLDHCQETCCEPTACQPTCYqR | | | | | | |
| **KRTAP11-1** | rs79258920_G | S78F | TSCVSNPCQVTCSR | 🟦 | 🟦 | 🟦 | | 1 | 1 |
| **KRTAP11-1** | rs79258920_A | S78F | TSCVfNPCQVTCSR | | | | | | |
| **KRTAP11-1** | rs9636845_A | C111S | QTTCISNPCSTTYSRPLTFVSSGCQPLGGISSVCQPVGGISTVCQPVGGVSTVCQPACGVSR | 🟦 | 🟦 | 🟦 | | 1 | 1 |
| **KRTAP11-1** | rs9636845_T | C111S | QTTCISNPCSTTYSRPLTFVSSGsQPLGGISSVCQPVGGISTVCQPVGGVSTVCQPACGVSR | | 🟩 | 🟩 | | | |
| **KRTAP16-1** | rs2074285_G | P340R | RCPSVCPEPVSCPSTSCR | 🟦 | 🟦 | 🟦 | | 1 | 1 |
| **KRTAP16-1** | rs2074285_C | P340R | RCrSVCPEPVSCPSTSCR | | 🟦 | | | 1 | |
| **LAMP1** | rs9577230_T | I309T | FFLQGIQLNTILPDAR | 🟦 | 🟦 | 🟦 | | 1 | 1 |
| **LAMP1** | rs9577230_C | I309T | FFLQGIQLNTtLPDAR | | | | | | |
| **LGALS3** | rs10148371_G | R183K | LDNNWGR | 🟦 | 🟦 | 🟦 | | 1 | 1 |
| **LGALS3** | rs10148371_A | R183K | LDNNWGk | | | | | | |
| **LGALS3** | rs11125_A | Q201H | IQVLVEPDHFK | 🟦 | 🟦 | | | 1 | 1 |
| **LGALS3** | rs11125_T | Q201H | IhVLVEPDHFK | | | | | | |
| **LRRC15** | rs13070515_A | P286L | ELSIGIFGPMPNLR | | 🟩 | | | | |
| **LRRC15** | rs13070515_G | P286L | ELSPGIFGPMPNLR | 🟦 | 🟦 | 🟦 | | 1 | |
| **NEU2** | rs2233384_C | S11R | ESVFQSGAHAYR | 🟦 | 🟦 | 🟦 | | 1 | |
| **NEU2** | rs2233384_A | S11R | ASLPVLQKEr | | | | | | |
| **NEU2** | rs2233385_G | R41Q | IPALLYLPGQQSLLAFAEQR | 🟩 | 🟩 | 🟦 | 🟩 | 1 | 1 |
| **NEU2** | rs2233385_A | R41Q | IPALLYLPGQQSLLAFAEQq | | | | | | |
| **NEU2** | rs2233390_G | A145T | DLTDAAIGPAYR | 🟦 | 🟦 | 🟦 | | 1 | 1 |
| **NEU2** | rs2233390_A | A145T | DLTDtAIGPAYR | | | | | | |
| **PKP1** | rs61818256_C | R684W | AAEAARLLLSDMWSSK/LLLSDMWSSK | 🟦 | 🟦 | 🟦 | | 1 | 1 |
| **PKP1** | rs61818256_T | R684W | AAEAAwLLLSDMWSSK | | | | | | |
| **PLCD1** | rs933135_C | R257H | EEAAGPALALSLIER | 🟩 | 🟩 | 🟦 | 🟩 | | |
| **PLCD1** | rs933135_T | R257H | EEAAGPALALSLIEhYEPSETAK | | | | | | |
| **PPL** | rs2037912_C | Q1573E | eNLQLETR | 🟩 | 🟩 | 🟩 | 🟩 | | |
| **PPL** | rs2037912_G | Q1573E | QNLQLETR | 🟩 | 🟩 | 🟦 | | | |
| **PPL** | rs143676756_C | R1457Q | VVLQQDPQQAREHALLR | 🟩 | 🟩 | 🟩 | 🟦 | | |
| **PPL** | rs143676756_T | R1457Q | VVLQQDPQQAqEHALLR | | | | | | |
| **S100A3** | rs36022742_C | R3K | ARPLEQAVAAIVCTFQEYAGR | 🟦 | 🟦 | 🟦 | | 1 | 1 |
| **S100A3** | rs36022742_T | R3K | AkPLEQAVAAIVCTFQEYAGR | | | | | | |

| Gene | rsID | AA change | Peptide | | | | | | |
|------|------|-----------|---------|---|---|---|---|---|---|
| SERPINB5 | rs1455555_A | I319V | GVALSNV**I**HK | blue | | green | blue | | |
| SERPINB5 | rs1455555_G | I319V | GVALSNV**v**HK | | blue | green | green | | |
| SYNGR2 | rs142608913_G | A28S | FLTQPQVV**A**R | blue | blue | blue | green | 1 | |
| SYNGR2 | rs142608913_T | A28S | FLTQPQVV**s**R | | | | | | |
| TCHH | rs2515663_A | L63R | TVDLILELLD**L**DSNGR | | | | | | |
| TCHH | rs2515663_C | L63R | TVDLILELLD**r** | blue | green | green | green | | |
| TGM3 | rs214814_G | S249N | **S**WNGSVEILK | blue | blue | blue | blue | 1 | 1 |
| TGM3 | rs214814_A | S249N | **n**WNGSVEILK | | | | | | |
| TRIM29 | rs11604169_T | Y544C | G**Y**PSLMR | blue | blue | green | blue | 1 | |
| TRIM29 | rs11604169_C | Y544C | G**c**PSLMR | | | | | | |
| VSIG8 | rs62624468_C | V47I | R.LGCPY**V**LDPEDYGPNGLDIEWMQVNSDPAHHR.E | blue | blue | green | blue | | |
| VSIG8 | rs62624468_T | V47I | R.LGCPY**i**LDPEDYGPNGLDIEWMQVNSDPAHHR.E | | | | | | |

| # of observations | |
|---|---|
| Total detected | 235 |
| True positive | 235 |
| False Positive | 0 |
| True Negative | 235 |
| False Negative | 16 |
| Undetected | 138 |

Table S5. GVPs used for CPI and CSI calculation in the six members of the family. The GVPs from P, M, S1 and S2 were validated from the corresponding genomic data. True positive identifications are highlighted in blue, true negative as white, and false negative as green. GVPs present in the protein regions that were not sequenced in the MS runs were called as undetected and highlighted as grey.

| Gene Name | rs#_nucleotide | SAP | Peptide | P | M | S1 | S2 | S3 | S4 |
|---|---|---|---|---|---|---|---|---|---|
| ALDH2 | rs671_G | E504K | ELGEYGLQ | blue | blue | blue | blue | 1 | 1 |
| ALDH2 | rs671_A | E504K | ELGEYGLQ | | | | | | |
| ATG9B | rs7804893_T | N493S | HFNELPHEL | blue | blue | blue | blue | 1 | |
| ATG9B | rs7804893_C | N493S | HFsELPHEL | | | | | | |
| ATP5A1 | rs79011243_C | A32S | VLSIGDGIA | blue | blue | blue | blue | 1 | 1 |
| ATP5A1 | rs79011243_A | A32S | VLSIGDGIs | | | | | | |
| CSRP1 | rs3738283_T | K108I | HEEAPGHR | blue | blue | blue | blue | 1 | 1 |
| CSRP1 | rs3738283_A | K108I | HEEAPGHR | | | | | | |
| DSC3 | rs276937_A | S78T | VLNDGSVY | grey | blue | grey | blue | | |
| DSC3 | rs276937_T | S78T | VLNDGtVYT | grey | white | grey | green | | |
| DSC3 | rs35296997_T | K180Q | GVDKEPLNL | grey | blue | blue | blue | 1 | |
| DSC3 | rs35296997_G | K180Q | GVDqEPLNL | grey | white | | | | |
| DSP | rs80325569_G | G939S | NLHSEISGK | blue | blue | blue | blue | 1 | |
| DSP | rs80325569_A | G939S | NLHSEISsK | | | | | | |
| DSP | rs2076299_A | Y1512C | VQYDLQK | grey | grey | grey | blue | | 1 |
| DSP | rs2076299_G | Y1512C | VQcDLQK | grey | grey | grey | white | | |
| DSP | rs28763966_C | N1526K | ANSSATETI | grey | grey | grey | blue | 1 | |
| DSP | rs28763966_A | N1526K | ANSSATETI | grey | grey | grey | white | | |
| DSP | rs6929069_A | R1738Q | GqSEADSD | grey | grey | grey | | | |
| DSP | rs6929069_G | R1738Q | GRSEADSD | blue | blue | blue | blue | 1 | 1 |
| DSP | rs28763967_C | R1537C | VQEQELTR | blue | blue | blue | blue | 1 | 1 |
| DSP | rs28763967_T | R1537C | VQEQELTcL | | | | | | |
| FAM83H | rs9969600-C | Q201H | VNLQHVDF | grey | grey | white | grey | | |
| FAM83H | rs9969600-A/G | Q201H | VNLhHVDFL | | blue | | | | |
| GSDMA | rs3894194_A | R18Q | QLNPqGDLT | | | | | | |
| GSDMA | rs3894194_G | R18Q | QLNPR/GDL | blue | blue | blue | blue | 1 | 1 |
| GSDMA | rs7212938_G | V128L | ALETVQER | blue | grey | grey | grey | | |
| GSDMA | rs7212938_T | V128L | ALETIQER | blue | | | | | 1 |
| GSDMA | rs56030650_A | T314N | GHEVnLEAL | grey | grey | white | grey | | |
| GSDMA | rs56030650_C | T314N | GHEVTLEAL | | | blue | | | |
| GSTP1 | rs1138272_C | A114V | YISLIYTNYE | blue | grey | grey | grey | 1 | 1 |
| GSTP1 | rs1138272_T | A114V | YISLIYTNYE | grey | grey | grey | grey | | |
| GSTP1 | rs1695_A | I105V | YISLIYTNYE | blue | grey | grey | grey | 1 | |
| GSTP1 | rs1695_G | I105V | YvSLIYTNYE | grey | grey | white | blue | | 1 |
| HEXB | rs10805890_A | I207V | GILIDTSR | blue | blue | blue | blue | 1 | 1 |
| HEXB | rs10805890_G | I207V | GILvDTSR | | | | | | |

| Gene | rs ID | Variant | Peptide | | | |
|------|-------|---------|---------|---|---|---|
| HEXB | rs77499935_A | I420V | LAPGTIVEV | | 1 | 1 |
| HEXB | rs77499935_G | I420V | LAPGTvVEV | | | |
| JUP | rs41283425_C | R142H | SAIVHLINYQ | | 1 | 1 |
| JUP | rs41283425_T | R142H | SAIVHLINYQ | | | |
| JUP | rs143043662_C | V648I | NEGTATYA | | 1 | 1 |
| JUP | rs143043662_T | V648I | NEGTATYA | | | |
| KRT1 | rs17678945_A | A454S | NKLNDLEDA | | | |
| KRT1 | rs17678945_C | A454S | NKLNDLEDA | | 1 | 1 |
| KRT32 | rs72830046_C/rs | R280H/C | CQYEAMVE | | 1 | 1 |
| KRT32 | rs72830046_T | R280H | CQYEAMVE | | 1 | |
| KRT32 | rs2604953_G | P427T | SLLENEDCK | | | |
| KRT32 | rs3744786_C | Q72R | TYLSSSCr | | 1 | |
| KRT32 | rs2071560_A | I171T | MVVNIDNAK | | 1 | 1 |
| KRT32 | rs2071560_G | I171T | MVVNtDNAK | | | |
| KRT32 | rs146792525_C | A255T | LNIEVDAAP | | 1 | 1 |
| KRT32 | rs146792525_T | A255T | LNIEVDtAPP | | | |
| KRT33A | rs373657561_C | G33R | PCVPPSCH | | 1 | 1 |
| KRT33A | rs373657561_T | G33R | PCVPPSCH | | | |
| KRT35 | rs743686_A | S36P | VSAMYSSS | | 1 | 1 |
| KRT35 | rs743686_G | S36P | VSAMYSSS | | 1 | 1 |
| KRT35 | rs200355130-C | E141D | LVVEIDNAK | | 1 | 1 |
| KRT35 | rs138303882_G | R163W | YETEVSLRQ | | 1 | 1 |
| KRT36 | rs75790652_G | A202G | CQLGDRLN | | 1 | 1 |
| KRT36 | rs75790652_C | A202G | CQLGDRLN | | | |
| KRT36 | rs11657323_T | N357T | YSSQLAQM | | | |
| KRT36 | rs11657323_G | N357T | YSSQLAQM | | | |
| KRT36 | rs9904102_G | R277C | CQYEALVE | | | |
| KRT36 | rs9904102_A | R277C | CQYEALVE | | | |
| KRT39 | rs17843021_G | T341M | DSQECILTE | | 1 | 1 |
| KRT39 | rs17843021_A | T341M | DSQECILmE | | | |
| KRT39 | rs7213256_C | R456Q | SGAIESTAP | | 1 | 1 |
| KRT39 | rs7213256_T | R456Q | SGAIESTAP | | | |
| KRT39 | rs17843023_G | L383M | QNQEYEILL | | 1 | 1 |
| KRT39 | rs17843023_T | L383M | QNQEYEILm | | | |
| KRT39 | rs112557906_G | S423F | CEPSPWTS | | 1 | 1 |
| KRT39 | rs112557906_A | S423F | CEPSPWTfC | | | |
| KRT39 | rs142154718_C | S86N | FSLDDCSW | | 1 | 1 |
| KRT39 | rs142154718_T | S86N | FSLDDCnW | | | |
| KRT40 | rs2010027_C | R235H | NHEEEVNLL | | 1 | 1 |
| KRT40 | rs2010027_T | R235H | NHEEEVNLL | | | |
| KRT40 | rs140634473_C | R108H | R.SLEETNA | | 1 | 1 |
| KRT40 | s140634473_T | R108H | V h SLEETN | | | |
| KRT7 | rs6580870_A | H186R | NKYEDEINH | | | |
| KRT7 | rs6580870_G | H186R | NKYEDEINr | | | |

| Gene | RS ID | Variant | Peptide | | |
|---|---|---|---|---|---|
| KRT82 | rs2658658_A | T458M | GAFLYEPC | | 1 |
| KRT82 | rs2658658_G | T458M | GAFLYEPC | 1 | 1 |
| KRT82 | rs1732263_C | E452D | GAFLYEPC | | |
| KRT82 | rs1732263_G | E452D | GAFLYdPC | | |
| KRT82 | rs1791634_C | E219Q | KYEEELSLR | 1 | 1 |
| KRT82 | rs1791634_G | E219Q | KYEEELSLR | | |
| KRT83 | rs61485872_A | C23G | PGNFSCVS | | 1 |
| KRT83 | rs61485872_C | C23G | PGNFSCVS | | |
| KRT83 | rs2852464_C | I279M | DLNMDCmV | | |
| KRT83 | rs2852464_G | I279M | DLNMDCIVA | 1 | 1 |
| KRT83 | rs2857663_G | R149C | LQFYQNR.E | 1 | 1 |
| KRT83 | rs2857663_A | R149C | LQFYQNcE | | |
| KRT84 | RS951773_A | C446R | CEYQELMN | | |
| KRT84 | RS951773_G | C446R | QLrEYQELM | 1 | 1 |
| KRTAP1-5 | rs148449559_G | T32S | TCCQTSFC | | |
| KRTAP1-5 | rs148449559_C | T32S | TCCQTSFC | | |
| KRTAP1-5 | rs62623375_C | C35Y | MTCCQTSF | 1 | |
| KRTAP1-5 | rs62623375_T | C35Y | MTCCQTSF | | |
| KRTAP3-2 | rs3829598_G | R27C | SCSVPTGP | | |
| KRTAP3-2 | rs3829598_A | R27C | K.SCCGV | | |
| KRTAP4-1 | rs398825_C | T134A | TTCCRPSC | 1 | 1 |
| KRTAP4-1 | rs398825_T | T134A | aTCCRPSC | 1 | 1 |
| KRTAP4-3 | rs428371_G | P152S | PACCISSCC | | 1 |
| KRTAP4-3 | rs428371_A | P152S | sACCISSCC | | |
| KRTAP4-4 | rs366700_C | R154S | TTCCRPSC | 1 | 1 |
| KRTAP4-4 | rs366700_G | R154S | TTCCRPSC | | |
| KRTAP4-4 | rs75030409_T | Q109R | TTCCRPSC | | 1 |
| KRTAP4-4 | rs75030409_C | Q109R | TTCCRPSC | | |
| KRTAP4-9 | rs149483591_G | R26H | VSSCCGSV | | 1 |
| KRTAP4-9 | rs149483591_A | R26H | VSSCCGSV | | |
| KRTAP4-9 | rs113059833_A | D18V | VSSCCGSV | | 1 |
| KRTAP4-9 | rs113059833_T | D18V | VSSCCGSV | | |
| KRTAP5-2 | rs35925287_C | G29R | GCGSGCG | | 1 |
| KRTAP5-2 | rs35925287_T | G29R | GCGSGCG | | |
| KRTAP9-4 | rs2191379_A | S146Y | R.TCYYPTT | | |
| KRTAP9-4 | rs2191379_C | S146Y | RTCYYPTTV | | |
| KRTAP10-12 | rs61745911_G | C236Y | LASCGSLLC | 1 | |
| KRTAP10-12 | rs61745911_A | C236Y | LASCGSLLy | | |
| KRTAP10-12 | rs34302939_G | G226S | RVPVPSCC | 1 | |
| KRTAP10-12 | rs34302939_A | G226S | RVPVPSCC | | |
| KRTAP11-1 | rs71321355_C | R72Q | CIVPVAQVT | 1 | 1 |
| KRTAP11-1 | rs71321355_T | R72Q | CIVPVAQVT | | |
| KRTAP11-1 | rs79258920_G | S78F | TSCVSNPC | 1 | 1 |
| KRTAP11-1 | rs79258920_A | S78F | TSCVfNPCQ | | |

| Gene | SNP | Variant | Peptide | | Count | Count |
|------|-----|---------|---------|---|---|---|
| KRTAP11-1 | rs9636845_A | C111S | QTTCISNPC | | 1 | 1 |
| KRTAP11-1 | rs9636845_T | C111S | QTTCISNPC | | | |
| KRTAP16-1 | rs2074285_G | P340R | RCPSVCPE | | 1 | 1 |
| KRTAP16-1 | rs2074285_C | P340R | RCrSVCPEP | | 1 | |
| LAMP1 | rs9577230_T | I309T | FFLQGIQLN | | 1 | 1 |
| LAMP1 | rs9577230_C | I309T | FFLQGIQLN | | | |
| LGALS3 | rs10148371_G | R183K | LDNNWGR | | 1 | 1 |
| LGALS3 | rs10148371_A | R183K | LDNNWGk | | | |
| LGALS3 | rs11125_A | Q201H | IQVLVEPDH | | 1 | 1 |
| LGALS3 | rs11125_T | Q201H | IhVLVEPDH | | | |
| LRRC15 | rs13070515_A | P286L | ELSlGIFGP | | | |
| LRRC15 | rs13070515_G | P286L | ELSPGIFGP | | 1 | |
| NEU2 | rs2233384_C | S11R | ESVFQSGA | | 1 | |
| NEU2 | rs2233384_A | S11R | ASLPVLQKE | | | |
| NEU2 | rs2233385_G | R41Q | IPALLYLPG | | 1 | 1 |
| NEU2 | rs2233385_A | R41Q | IPALLYLPG | | | |
| NEU2 | rs2233390_G | A145T | DLTDAAIGP | | 1 | 1 |
| NEU2 | rs2233390_A | A145T | DLTDtAIGPA | | | |
| PKP1 | rs61818256_C | R684W | AAEAARLLL | | 1 | 1 |
| PKP1 | rs61818256_T | R684W | AAEAAwLLL | | | |
| PLCD1 | rs933135_C | R257H | EEAAGPALA | | | |
| PLCD1 | rs933135_T | R257H | EEAAGPALA | | | |
| PPL | rs2037912_C | Q1573E | eNLQLETR | | | |
| PPL | rs2037912_G | Q1573E | QNLQLETR | | | |
| PPL | rs143676756_C | R1457Q | VVLQQDPQ | | | |
| PPL | rs143676756_T | R1457Q | VVLQQDPQ | | | |
| S100A3 | rs36022742_C | R3K | ARPLEQAV | | 1 | 1 |
| S100A3 | rs36022742_T | R3K | AkPLEQAVA | | | |
| SERPINB5 | rs1455555_A | I319V | GVALSNVIH | | | |
| SERPINB5 | rs1455555_G | I319V | GVALSNVvH | | 1 | |
| SYNGR2 | rs142608913_G | A28S | FLTQPQVV | | 1 | |
| SYNGR2 | rs142608913_T | A28S | FLTQPQVVs | | | |
| TCHH | rs2515663_A | L63R | TVDLILELLD | | | |
| TCHH | rs2515663_C | L63R | TVDLILELLD | | | |
| TGM3 | rs214814_G | S249N | SWNGSVEIL | | 1 | 1 |
| TGM3 | rs214814_A | S249N | nWNGSVEIL | | | |
| TRIM29 | rs11604169_T | Y544C | GYPSLMR | | 1 | |
| TRIM29 | rs11604169_C | Y544C | GcPSLMR | | | |
| VSIG8 | rs62624468_C | V47I | R.LGCPYVL | | | |
| VSIG8 | rs62624468_T | V47I | R.LGCPYiL | | | |

Table S6: Genotypes of individual P, M, S1 and S2 for the identified genetically variant peptides. Genotypes at the five dubious loci are highlighted in bold italic.

| Serial No. | Gene Name | rs# | Reference | P | M | S1 | S2 |
|---|---|---|---|---|---|---|---|
| 1 | ALDH2 | rs671 | G | GG | GG | GG | GG |
| 2 | ATG9B | rs7804893 | T | TT | TT | TT | TT |
| 3 | ATP5A1 | rs79011243 | C | CC | CC | CC | CC |
| 4 | CSRP1 | rs3738283 | T | TT | TT | TT | TT |
| 5 | DSC3 | rs276937 | A | TT | AA | AT | AT |
| 6 | DSC3 | rs35296997 | T | TT | TT | TT | TT |
| 7 | DSP | rs2076299 | A | AG | AA | AA | AA |
| 8 | DSP | rs28763966 | C | CC | CC | CC | CC |
| 9 | DSP | rs28763967 | C | CC | CC | CC | CC |
| 10 | DSP | rs80325569 | G | GG | GG | GG | GG |
| 11 | DSP | rs6929069 | G | GG | GG | GG | GG |
| 12 | FAM83H | rs9969600 | C | GG | GG | GG | GG |
| 13 | GSDMA | rs56030650 | C | CA | CC | CC | CC |
| 14 | GSDMA | rs3894194 | G | GG | GG | GG | GG |
| 15 | GSDMA | rs7212938 | G | TT | TT | TT | TT |
| 16 | GSTP1 | rs1695 | A | AA | AG | AA | AG |
| 17 | GSTP1 | rs1138272 | C | CC | CC | CC | CC |
| 18 | HEXB | rs10805890 | A | AA | AA | AA | AA |
| 19 | HEXB | rs77499935 | A | AA | AA | AA | AA |
| 20 | JUP | rs41283425 | C | CC | CC | CC | CC |
| 21 | JUP | rs143043662 | C | CC | CC | CC | CC |
| 22 | KRT1 | rs17678945 | C | CC | CC | CC | CC |
| 23 | KRT31 | rs6503627 | G | GG | AG | GG | GG |
| 24 | KRT32 | rs2071560 | A | AA | AA | AA | AA |
| 25 | KRT32 | rs72830046 | C | CT | CC | CT | CC |
| 26 | KRT32 | rs146792525 | C | CC | CC | CC | CC |
| 27 | KRT32 | rs2071561 | G | GT | GG | GT | GT |
| 28 | KRT32 | rs3744786 | T | TC | TT | TC | TC |
| 29 | KRT32 | rs2604953 | G | *TT* | *TT* | *TT* | *GG* |
| 30 | KRT33A | rs373657561 | C | CC | CC | CC | CC |
| 31 | KRT33A | rs12937519 | G | GG | GA | GG | GG |
| 32 | KRT34 | rs2071599 | T | TT | TC | TT | TT |
| 33 | KRT35 | rs743686 | A | GG | AG | GG | GG |
| 34 | KRT35 | rs200355130 | C | CC | CC | CC | CC |
| 35 | KRT35 | rs138303882 | G | GG | GG | GG | GG |
| 36 | KRT36 | rs75790652 | G | GG | GG | GG | GG |
| 37 | KRT36 | rs9904102 | G | GG | GG | GG | GG |
| 38 | KRT36 | rs11657323 | T | TG | TG | TG | TG |
| 39 | KRT39 | rs7213256 | C | CC | CC | CC | CC |
| 40 | KRT39 | rs142154718 | C | CC | CC | CC | CC |
| 41 | KRT39 | rs17843021 | G | GA | GG | GG | GG |
| 42 | KRT39 | rs17843023 | G | GG | GG | GG | GG |

| 43 | KRT39 | rs112557906 | G | | GG | GG | GG | GG |
| 44 | KRT40 | rs2010027 | C | | CT | CC | CC | CC |
| 45 | KRT40 | rs140634473 | C | | CC | CC | CC | CC |
| 46 | KRT7 | rs6580870 | A | | AG | GG | GG | AG |
| 47 | KRT81 | rs6580873 | A | | AC | AC | AA | CC |
| 48 | KRT82 | rs1732263 | C | | CC | CC | CC | CC |
| 49 | KRT82 | rs1791634 | C | | CC | CC | CC | CC |
| 50 | KRT82 | rs2658658 | G | | GA | GA | GG | AA |
| 51 | KRT83 | rs61485872 | A | | AA | AA | AA | AA |
| 52 | KRT83 | rs2852464 | G | | GC | GG | GG | GC |
| 53 | KRT83 | rs2857663 | G | | GG | GG | GG | GG |
| 54 | KRT84 | RS951773 | A | | GG | AG | AG | GG |
| 55 | KRT85 | rs2852471 | C | | CC | CC | CC | CC |
| 56 | KRTAP10-12 | rs61745911 | G | | GG | GG | GG | GG |
| 57 | KRTAP10-12 | rs34302939 | G | | GG | GG | GG | GG |
| 58 | KRTAP10-6 | rs465279 | G | | *GG* | *AA* | *GA* | *GG* |
| 59 | KRTAP11-1 | rs9636845 | A | | AA | AT | AT | AA |
| 60 | KRTAP11-1 | rs71321355 | C | | CC | CC | CC | CC |
| 61 | KRTAP11-1 | rs79258920 | G | | GG | GG | GG | GG |
| 62 | KRTAP1-5 | rs62623375 | C | | CC | CC | CC | CC |
| 63 | KRTAP1-5 | rs148449559 | G | | GG | GG | GG | GG |
| 64 | KRTAP16-1 | rs2074285 | G | | GG | GC | GG | GG |
| 65 | KRTAP3-2 | rs3813050 | A | | AA | AA | AA | AA |
| 66 | KRTAP3-2 | rs3829598 | G | | GG | GG | GG | GG |
| 67 | KRTAP3-2 | rs9897046 | T | | TT | TT | TT | TT |
| 68 | KRTAP4-1 | rs398825 | C | | *TT* | *CT* | *TT* | *CC* |
| 69 | KRTAP4-3 | rs428371 | G | | GG | GG | GG | GG |
| 70 | KRTAP4-4 | rs366700 | C | | CC | CC | CC | CC |
| 71 | KRTAP4-4 | rs385055 | T | | TT | TT | TT | TT |
| 72 | KRTAP4-4 | rs75030409 | T | | TT | TT | TT | TT |
| 73 | KRTAP4-6 | rs73983172 | G | | GG | GG | GG | GG |
| 74 | KRTAP4-9 | rs149483591 | G | | GG | GG | GG | GG |
| 75 | KRTAP4-9 | rs113059833 | A | | *AA* | *AA* | *AA* | *AT* |
| 76 | KRTAP5-2 | rs35925287 | C | | CC | CC | CC | CC |
| 77 | KRTAP9-2 | rs9902235 | **G** | | *GG* | *GG* | *GG* | *CC* |
| 78 | KRTAP9-4 | rs2191379 | C | | AA | CA | AA | AA |
| 79 | LAMP1 | rs9577230 | T | | TT | TT | TT | TT |
| 80 | LGALS3 | rs11125 | A | | AA | AA | AA | AA |
| 81 | LGALS3 | rs10148371 | G | | GG | GG | GG | GG |
| 82 | LRRC15 | rs13070515 | G | | GG | GA | GG | GG |
| 83 | NEU2 | rs2233384 | C | | CC | CC | CC | CC |
| 84 | NEU2 | rs2233385 | G | | GG | GG | GG | GG |
| 85 | NEU2 | rs2233390 | G | | GG | GG | GG | GG |

| 86 | **PKP1** | rs61818256 | C | | CC | CC | CC | CC |
| 87 | **PLCD1** | rs933135 | C | | CC | CC | CC | CC |
| 88 | **PPL** | rs143676756 | C | | CC | CC | CC | CC |
| 89 | **PPL** | rs2037912 | G | | GC | GC | GC | CC |
| 90 | **S100A3** | rs36022742 | C | | CC | CC | CC | CC |
| 91 | **SERPINB5** | rs1455555 | A | | AA | GG | AG | AG |
| 92 | **SYNGR2** | rs142608913 | G | | GG | GG | GG | GG |
| 93 | **TCHH** | rs2515663 | A | | CC | CC | CC | CC |
| 94 | **TGM3** | rs214814 | G | | GG | GG | GG | GG |
| 95 | **TRIM29** | rs11604169 | T | | TT | TT | TT | TT |
| 96 | **VSIG8** | rs62624468 | C | | CC | CC | CC | CC |

**Table S7:** Combined paternity indexes and posterior probabilities calculated using all the detected GVPs including false positives. The posterior probabilities were obtained using the prior odds of 4/12 for the four true offspring of the couple P and M and eight random individuals. CPI: combined paternity index.

| CPI when both parents are available | | | | |
|---|---|---|---|---|
| Individual | Loci used in the calculation | Loci with no obligate allele | CPI | Posterior Probability (%) |
| S3 | 35 | 0 | 1696.75 | 99.94 |
| S1 | 36 | 0 | 1506.12 | 99.93 |
| S4 | 33 | 0 | 340.70 | 99.70 |
| S2 | 39 | 1 | 1239.76 | 99.91 |
| A | 34 | 3 | -- | -- |
| B | 34 | 4 | -- | -- |
| C | 31 | 6 | -- | -- |
| D | 31 | 4 | -- | -- |
| F | 33 | 4 | -- | -- |
| G | 30 | 2 | -- | -- |
| H | 30 | 2 | -- | -- |
| I | 35 | 3 | -- | -- |

**Table S8:** Combined sibship index values calculated using all the detected GVPs including false positives for the family members and eight unrelated individuals. The CSI values higher than 10 for unrelated individuals or lower than 10 for true siblings are shown in bold, and the ones that support the relationship in the cases of true relationships are bold italicized.

| CSI | P | S3 | S1 | S4 | S2 | A | B | C | D | F | G | H | IYO |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| M | 3.21 | *198.13* | 5.71 | *1383.55* | 6.12 | 0.51 | 7.32 | 2.18 | 9.54 | 0.04 | 14.07 | 0.05 | 1.58 |
| P | | *6.16* | *1648.28* | *10.83* | *394.09* | 1.76 | 0.06 | 0.02 | 0.05 | 7.58 | 0.05 | 3.01 | 0.24 |
| S3 | | | *330.91* | *67.33* | *14.13* | 1.45 | 1.00 | 0.56 | 4.40 | 0.15 | 0.71 | 0.04 | 1.70 |
| S1 | | | | *8.80* | *171.49* | 0.07 | 0.02 | 0.03 | 0.04 | 0.59 | 0.01 | 0.50 | 1.53 |
| S4 | | | | | *60.96* | 0.66 | 1.07 | 0.10 | 3.65 | 0.04 | 1.86 | 0.09 | 2.45 |
| S2 | | | | | | 0.11 | 0.01 | 0.00 | 0.51 | 2.83 | 0.45 | **10.65** | 0.13 |
| A | | | | | | | 1.36 | **388.00** | 1.77 | 0.08 | 1.35 | 2.99 | 3.46 |
| B | | | | | | | | **10.74** | 0.93 | 0.00 | 0.34 | 0.01 | 0.01 |
| C | | | | | | | | | **15.22** | 0.01 | 9.15 | 0.59 | 0.14 |
| D | | | | | | | | | | 0.00 | 5.67 | 0.08 | **10.30** |
| F | | | | | | | | | | | 0.15 | 2.40 | 0.01 |
| GH | | | | | | | | | | | | 0.12 | 3.05 |
| H | | | | | | | | | | | | | 0.04 |
| IYO | | | | | | | | | | | | | |

**Table S9:** Number of loci at which each comparison was based for CSI calculations. The numbers inside the parentheses represent the genes with two loci included.

| CSI | P | S3 | S1 | S4 | S2 | A | B | C | D | F | G | H | I |
|-----|---|----|----|----|----|---|---|---|---|---|---|---|---|
| M | 30(9) | 29(9) | 28(8) | 24(8) | 31(9) | 25(9) | 24(9) | 22(8) | 21(9) | 24(8) | 22(7) | 21(8) | 25(8) |
| P | | 27(9) | 27(9) | 24(8) | 31(11) | 26(11) | 25(11) | 22(10) | 21(10) | 24(10) | 22(9) | 21(10) | 26(10) |
| S3 | | | 27(10) | 24(8) | 28(9) | 24(10) | 23(10) | 21(8) | 21(9) | 23(10) | 22(9) | 21(10) | 24(11) |
| S1 | | | | 24(8) | 30(9) | 25(9) | 25(10) | 22(8) | 22(9) | 25(9) | 25(9) | 23(9) | 25(9) |
| S4 | | | | | 25(8) | 24(9) | 22(8) | 19(7) | 20(8) | 22(7) | 22(7) | 19(7) | 23(8) |
| S2 | | | | | | 27(12) | 26(11) | 27(10) | 22(11) | 27(12) | 26(11) | 24(12) | 27(12) |
| A | | | | | | | 29(12) | 31(12) | 23(9) | 25(13) | 28(12) | 26(13) | 26(12) |
| B | | | | | | | | 27(13) | 26(10) | 26(13) | 24(12) | 25(13) | 25(12) |
| C | | | | | | | | | 25(9) | 27(11) | 30(10) | 30(11) | 27(11) |
| D | | | | | | | | | | 22(11) | 23(9) | 23(11) | 24(11) |
| F | | | | | | | | | | | 26(13) | 27(15) | 25(13) |
| G | | | | | | | | | | | | 29(14)) | 25(12) |
| H | | | | | | | | | | | | | 25(13) |
| I | | | | | | | | | | | | | |

Table S2. GVP data matrix used for hierarchical clustering. Each GVP detection was assigned a value of 1 and a non-detection of 0.

| Individuals | rs2 229528_T | rs2 229528_C | rs6 71_ G | rs6 71_ A | rs1 784522 | rs7 804893 6_A | rs7 804893 _T | rs7 901124 _C | rs7 901124 3_C | rs3 738283 3_A | rs3 738283 _T | rs2 76937_ _A | rs2 76937_ A | rs3 529699 T | rs3 529699 7_T | rs8 032556 7_G | rs8 032556 9_G | rs2 076299 9_A | rs2 076299 _A | rs2 876396 _G | rs2 876396 6_C | rs2 876396 6_A | rs2 876396 7_C | rs6 929069 7_T | rs6 929069 _A | rs9 969600- _G | rs9 969600- C | rs1 155069 | rs1 155069 A/G | rs1 785602 9_A | rs1 785602 9_G | rs3 536328 4_C | rs3 536328 4_A | rs3 894194 7_C | rs3 894194 7_T | rs7 212938 _A | rs7 212938 _G | rs5 603065 0_A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| S3 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| S4 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 |
| S1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| P | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 |
| M | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| S2 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| A | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| B | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 |
| C | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| D | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| E | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| F | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| G | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| H | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| I | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |

| rs56030650_C | rs1138272_C | rs1138272_T | rs1695_A0_A | rs1695_G0_G | rs10805895_A | rs10805895_G | rs7749993_C | rs7749993_T | rs67612765_C | rs67612765_T | rs412834262_C | rs412834262_T | rs14304365_A | rs14304365_C | rs1767894_A | rs1767894_T | rs65036271_G | rs20715616_T | rs7398345_G | rs7283004_T | rs2604953_T | rs2604953_C | rs3744786_A | rs3744786_G | rs207156025_C | rs207156025_T | rs146792561_C | rs146792561_T | rs37365759_A | rs3736575_A | rs1293751_T | rs22397108_G | rs20715998_A | rs6174066A | rs6174066G | rs743686_A | rs743686_G | rs2003551 30-C |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | 1 |
| 0 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 |
| 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 |
| 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 |
| 0 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 |
| 0 | 1 | 0 | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 0 | 0 |
| 0 | 1 | 0 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 0 |
| 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 |
| 0 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 |
| 0 | 1 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 0 |

| rs2003551 30-A | rs1383038 82_A | rs1383038 82_G | rs1245165 2_C | rs1245165 2_T | rs2071601 _C | rs2071601 _G | rs7579065 2_G | rs7579065 2_C | rs1165732 3_T | rs1165732 3_G | rs9904102 _G | rs9904102 _A | rs9916484 _T | rs9916484 _C | rs9916475 _T | rs9916475 _A | rs1696681 1_A | rs8974 16_ A | rs8974 16_ G | rs1784302 1_G | rs1784302 1_A | rs7213256 _C | rs7213256 _T | rs1784302 3_G | rs1784302 3_T | rs1125579 06_ G | rs1125579 06_ A | rs1421547 18_ C | rs1421547 18_ T | rs2010027 _C | rs2010027 _T | rs1406344 73_ | s14063447 3_T C | rs7219 57_ T | rs7219 57_ | rs6580870 _A | rs6580870 _G | rs2232393 _A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 |
| 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 1 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |
| 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 |

| rs2 | rs6 | rs2 | rs2 | rs1 | rs1 | rs1 | rs1 | rs6 | rs6 | rs2 | rs2 | rs2 | rs2 | rs2 | rs2 | RS9 | RS9 | rs6 | rs6 | rs2 | rs2 | rs5 | rs1 | rs1 | rs1 | rs1 | rs1 | rs6 | rs6 | rs1 | rs1 | rs9 | rs3 | rs3 | rs3 | rs3 | rs3 | rs6 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 232 | 580 | 658 | 658 | 732 | 732 | 791 | 791 | 148 | 148 | 852 | 852 | 857 | 857 | 857 | 857 | 517 | 517 | 163 | 163 | 852 | 852 | 871 | 502 | 502 | 382 | 484 | 484 | 262 | 262 | 387 | 387 | 897 | 829 | 829 | 813 | 988 | 988 | 206 |
| 393 | 873 | 658 | 658 | 263 | 263 | 634 | 634 | 587 | 587 | 464 | 464 | 671 | 671 | 663 | 663 | 73_ | 73_ | 000 | 000 | 471- | 471- | 726 | 184 | 184 | 008 | 495 | 495 | 337 | 337 | 587 | 587 | 046 | 598 | 598 | 050 | 25_ | 25_ | 729 |
| _G | _A | _A | _G | _C | _G | _C | _G | 2_A | 2_C | _C | _G | _G | _A | _G | _A | A | G | 4_C | 4_T | C | A | 6_C | 95_G | 95_C | 23_C | 59_G | 59_C | 5_C | 5_T | 76_T | 76_C | _T | _G | _A | _A | C | T | 2_G |
| 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 |
| 0 | 1 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 |
| 0 | 1 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 |
| 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 |
| 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 |
| 0 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 |
| 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 |

| rs62067292_C | rs389784_T | rs389784_C | rs428371_G | rs428371_A | rs366700_C | rs366700_G | rs385055_T | rs444509_A | rs444509_T | rs75030409_T | rs75030409_C | rs1497383_A | rs73983172_G | rs201814486_G | rs201814486_C | rs138296121_C | rs138296121_T | rs149483591_G | rs149483591_A | rs113059833_A | rs113059833_T | rs9897031_C | rs9897031_T | rs113376601_C | rs35925287_C | rs35925287_T | rs9902235_C | rs9902235_G | rs2191379_A | rs2191379_C | rs12938692_A | rs12938692_G | rs576405629_A | rs576405629_G | rs537301040_T | rs537301040_C | rs12938374_G | rs12938374_A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 |
| 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 |
| 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |

| rs2332 52_C | rs2332 52_T | rs4643 91_C | rs4643 91_G | rs4652 79_A | rs1116 686 37_G | rs1116 686 37_A | rs4112 54_A | rs4112 54_G | rs9980 129_T | rs9980 129_C | rs4818 950_G | rs4818 950_A | rs6174 591 1_G | rs6174 591 1_A | rs3430 293 9_G | rs3430 293 9_A | rs7132 135 5_C | rs7132 135 5_T | rs7925 892 0_G | rs7925 892 0_A | rs9636 845_A | rs9636 845_T | rs3804 010_G | rs3804 010_C | rs2074 285_G | rs2074 285_C | rs9577 230_T | rs9577 230_C | rs1014 837 1_G | rs1014 837 1_A | rs1112 5_A | rs1112 5_T | rs1306 062 7_C | rs1306 062 7_T | rs1307 051 5_A | rs1307 051 5_G | rs2233 384_C | rs2233 384_A |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 |
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 |
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 |
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 |
| 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 |
| 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 1 | 1 | 0 |
| 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 1 | 0 |
| 1 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 1 | 1 | 0 | 1 | 0 | 1 | 1 | 0 | 0 | 0 | 1 | 1 | 0 |

| rs2233385_G | rs2233385_A | rs2233390_G | rs2233390_A | rs2233391_A | rs2233391_C | rs412750_G | rs412750_A | rs61818256_C | rs61818256_T | rs10920171_C | rs10920171_T | rs6753929_G | rs6753929_C | rs933135_C | rs933135_T | rs2037912_C | rs2037912_G | rs143676756_C | rs143676756_T | rs116208483_G | rs36022742_C | rs36022742_T | rs41265164_G | rs41265164_T/C | rs1455555_A | rs1455555_G | rs142608913_G | rs142608913_T | rs2515663_A | rs2515663_C | rs214803_A | rs214803_C | rs214814_G | rs214814_A | rs11604169_T | rs11604169_C | rs62624468_C | rs62624468_T |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 |
| 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 |
| 0 | 0 | 1 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 1 | 0 | 0 | 1 | 0 | 1 | 0 | 1 | 0 |
| 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 |
| 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 0 | 0 | 0 |
| 1 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 1 | 1 | 0 | 1 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 | 0 |