1    Comprehensive analysis of microsatellite polymorphisms in human populations

2

3    Leo Gochi[1], Yosuke Kawai[2], and Akihiro Fujimoto[1]

4

5    1; Department of Human Genetics, The University of Tokyo, Graduate School of Medicine,

6    Japan

7    2; Genome Medical Science Project, National Center for Global Health and Medicine

8

9    Corresponding author:

10

11    Akihiro Fujimoto

12    Department of Human Genetics, The University of Tokyo, Graduate School of Medicine, 113-

13    0003, Japan;

14    E-mail: afujimoto@m.u-tokyo.ac.jp

15

16 **Abstract**

17 Microsatellites (MS) are tandem repeats of short units and have been used for population

18 genetics, individual identification, and medical genetics. However, studies of MS on a whole

19 genome level are limited, and genotyping methods for MS have yet to be established. Here, we

20 analyzed approximately 8.5 million MS regions using a previously developed MS caller

21 (MIVcall method) for three large publicly available human genome sequencing data sets: the

22 Korean Personal Genome Project (KPGP), Simons Genome Diversity Project (SGDP), and

23 Human Genome Diversity Project (HGDP). Our analysis identified 253,114 polymorphic MS. A

24 comparison among different populations suggests that MS in the coding region evolved by

25 random genetic drift and natural selection. In an analysis of genetic structures, MS clearly

26 revealed population structures as SNPs and detected clusters that were not found by SNPs in

27 African and Oceanian populations. Based on the MS polymorphisms, we selected an effective

28 MS set for individual identification. We also showed that our MS analysis method can be applied

29 to ancient DNA samples. This study provides a comprehensive picture of MS polymorphisms

30 and application to human population studies.

31

32 **Introduction**

33 Repetitive sequences account for more than two-thirds of the human genome (1). Among them,

34 sequences consisting of tandem repeats of short units are classified as microsatellites (MS). The

35 mutation rate of MS is higher than these of other genomic regions, and MS have high diversity

36 among individuals (1).

37 Due to their high heterozygosity and multiallelic nature, MS have been widely used as

38 genetic markers in population studies (2–5). Previous studies analyzed dozens to several

39 hundreds of MS and revealed the genetic structure of modern human populations, the

40 relationship of genetic and linguistic variations, and the trace of natural selection (3–6). These

41 studies have made a great contribution to understanding human evolutionary history. However,

42 conventional PCR-based MS genotyping methods are not suitable for high-throughput

43 genotyping; therefore, MS has largely been replaced with single nucleotide polymorphism (SNP)

44 in related studies. Indeed, in the past decade, genome-wide SNP analysis had become a standard

45 method for human population genetics (7,8).

46 In addition to population studies, MS has been widely used for personal identification in

47 forensic science and paternity testing (9,10) since MS are multiallelic and have more

48 discriminative power than SNP (2). Established common MS sets, such as the Globalfiler kit,

49 have been used for many years (11). Although these MS marker sets have sufficient power for

50 most cases, they are not selected from entire genomes and there may be other more efficient MS

51 in the human genome.

52 Next generation sequencing technologies (NGS) enable whole genome sequencing

53 (WGS). In the past decade, applications of NGS and the development of algorithms for the

54 analysis have successfully identified genetic variations, including single nucleotide variations,

55 insertions and deletions, copy number variations, and structural variations (12). However, due to

56 the short read length and the high sequencing error rate in repeat regions, the identification of

57 mutations and polymorphisms in MS regions have been difficult. Previously, several groups

58 including ours have developed MS genotyping tools from WGS data (13–15). These methods

59 identified somatic mutations or germline polymorphisms in MS in the human population and

60 revealed genome-wide patterns of MS polymorphisms, factors that determine the mutation rate

61 of MS, and functional roles of MS on gene expressions (13–17). Although these studies provide

62 important information on MS, the MS genotyping method is not perfect, and only few studies

63 have been conducted for MS polymorphisms. Indeed, the patterns of genome-wide MS

64    polymorphisms have not been well analyzed in various human populations. Nor has the amount

65    of genetic variation in MS among different human populations been compared in detail.

66    Furthermore, clustering based on principal component analysis (PCA) with MS polymorphisms

67    has presented unclear results compared to that with SNPs (14,18), and the efficiency of genome-

68    wide MS for analyzing population structures requires more study.

69          Here, we analyzed approximately nine million MS regions using a previously developed

70    MS caller for three large publicly available human genome sequencing data sets: Korean

71    Personal Genome Project (KPGP), Simons Genome Diversity Project (SGDP), and Human

72    Genome Diversity Project (HGDP) (13,18,19). We revealed the pattern of MS polymorphisms,

73    analyzed the genetic structure of the populations with several dimensionality reduction methods,

74    and identified useful candidate MS for individual identification. Additionally, we analyzed MS

75    of an ancient DNA sample (20). Our analysis provides a comprehensive picture of MS

76    polymorphisms and their application to human population studies.

77

78

79 **Results**

80 *Establishment of MS calling*

81 We identified genotypes of MS with the MIVcall method (13). MIVcall outputs the

82 $\log_{10}$(likelihood) and number of reads for each MS locus, which can be used to evaluate the

83 reliability of MS genotypes. We examined these parameters using monozygotic twins in the

84 KPGP (KPGP-00088 and KPGP-00089). Because monozygotic twins have completely same

85 genotypes, we considered all disconcordant genotypes between the twins as errors. We compared

86 the genotypes of 8,343,174 MS in the twins and classified them into concordant homozygote,

87 concordant heterozygote, disconcordance of two alleles, and disconcordance of one allele (S1

88 Table). Based on the result, the cutoff $-\log_{10}$(likelihood) and minimum number of reads for allele

89 detection were set to -4 and 3, respectively, in this study (S1 Table).

90

91 *Selection of samples and MS*

92 We selected samples for the analysis. In this study, MS covered by less than 10 reads were

93 considered insufficient depth, and we excluded samples if more than 4% had insufficient depth

94 of MS. As a result, 277 samples from the SGDP, 692 samples from the HGDP, and 81 samples

95 from the KPGP (excluding one of the monozygotic twins) were selected.

96    We next selected MS loci for this study using the 81 Korean samples from the KPGP. We

97 selected MS that were genotyped (number of reads ≥ 10) in more than 85% of the 81 samples,

98 leaving 8,468,218 MS. We also tested deviation from the Hardy-Weinberg equilibrium (HWE)

99 in the KPGP, but no MS showed significant deviation.

100    Of the selected 8,468,218 MS, 7,740,569 MS were monomorphic. For each MS, minor allele

101 frequency (MAF) was calculated by 1 – (major allele frequency). Of these, 727,649 MS had

102 variations in at least one sample, and 253,114 had a MAF ≥ 1% (S1 Fig). Of all MS, 71,040 MS

103 were in coding sequences (CDS), and 8,397,178 MS were in non-CDS (S2 Table). Among the

104 CDS MS, 893 MS had variations in at least one sample, and 71 MS had a MAF ≥ 1% (S2

105 Table).

106

107 *Genome-wide pattern of MS polymorphisms*

108 The length of MS was negatively correlated to the number of samples with insufficient depth (p-

109 value < $10^{-200}$ Kruskal-Wallis test) but positively correlated to the number of alleles (p-value <

110 $10^{-200}$ Kruskal-Wallis test) (Fig 1A and B). The increase in the number of alleles was more

111    gradual in longer MS (Fig 1B). The longer MS had fewer reads fully covering them compared to

112    shorter MS and thus a lower detection sensitivity. The number and heterozygosity of different

113    repeat units showed that MS with higher AT content had significantly higher heterozygosity

114    (Pearson's uncorrelated test p-value = $1.47 \times 10^{-61}$), suggesting that the mutability of MS is

115    affected by the base composition (Fig 1 C and D).

116       A comparison between CDS and non-CDS showed that MS in CDS had a higher GC content

117    (p-value = $6.49 \times 10^{-144}$ Wilcoxon rank sum test), shorter length (p-value = $1.25 \times 10^{-29}$ Wilcoxon

118    rank sum test), smaller number of alleles (p-value = $1.80 \times 10^{-47}$ Wilcoxon rank sum test), and

119    lower heterozygosity (p-value = $1.42 \times 10^{-89}$ Wilcoxon rank sum test) (Fig 2 A-D). A comparison

120    of distributions of the number MS in each unit length showed that polymorphic $3n$ MS (3 and 6

121    bp) were frequent in the CDS region (p-value = $5.22 \times 10^{-232}$ Fisher's exact test) (Fig 2 E-H, S3

122    Table). In the CDS, $3n$ MS had higher heterozygosity than non-$3n$ MS (1, 2, 4, 5 and 7 bp) (p-

123    value = $5.23 \times 10^{-13}$ Wilcoxon rank sum test), suggesting that non-$3n$ MS was strongly influenced

124    by negative selection (Fig 2I).

125       To compare the genetic variation among human populations, we calculated the distribution of

126    the heterozygosity of all MS, MS in CDS, MS in non-CDS, $3n$ MS in CDS, and non-$3n$ MS in

127    CDS (Fig 3 A-E and S4 Table). We also calculated the ratio of the mean heterozygosity of non-

128    $3n$ MS to mean heterozygosity of $3n$ MS in CDS among each region (Fig 3F). Africa had the

129    highest heterozygosity and lowest ratio of mean heterozygosity (Fig 3 and S4 Table).

130       We observed that 689 genes had 724 polymorphic non-$3n$ MS (S4 Table, S5 Table). We

131    performed a pathway analysis for these genes but did not find any significantly over-represented

132    pathway (data not shown).

133

134    *Analysis of population structure*

135    To analyze the population structure, we conducted five dimensionality reduction methods for MS

136    polymorphisms: PCA, t-Distributed Stochastic Neighbor Embedding (t-SNE), Uniform Manifold

137    Approximation and Projection (UMAP), PCA-t-SNE, and PCA-UMAP using MS and SNP. For

138    this analysis, we used MS and SNP with MAF $\geq$ 1%. In the MS, 253,114 MS had MAF $\geq$ 1% in

139    all samples, 340,114 in Africa, 176,214 in America, 229,636 in Central Asia and Siberia,

140    197,081 in East Asia, 231,706 in Oceania, 220,737 in South Asia, and 209,204 in West Eurasia

141    (Supplementary data1). In the SNPs, 486,579 SNPs had MAF $\geq$ 1% in all samples, 515,456 in

142    Africa, 358,538 in America, 416,191 in Central Asia and Siberia, 403,015 in East Asia, 389,358

143    in Oceania, 437,410 in South Asia, and 432,163 in West Eurasia (S2 Table).

144    To apply dimensionality reduction methods, we converted MS genotypes to numerical values

145    using two methods, multiallelic method and average method (see Methods), and performed PCA

146    with both methods for all samples. Although the results were not significantly different between

147    the two, the average method had a higher contribution rate (PC1 = 5.80%) than the multiallelic

148    method (PC1 = 4.97%) (S2 Fig). Therefore, we selected the average method in this study.

149    PCA was conducted for all regions and each individual region  (Fig 4, S3 Fig). We found

150    similar patterns between MS and SNPs for all samples (Fig 4 AB). However, patterns were

151    slightly different in African and Oceanian populations between MS and SNPs (Fig 4 C-F). These

152    results suggest that genome-wide MS has compatible resolution to SNPs for the genetic structure

153    of human populations and that MS can be used to find new genetic structures. t-SNE, UMAP,

154    PCA-t-SNE, and PCA-UMAP were also conducted for all samples (S4 Fig). In most of these

155    analyses MS did not detect novel clusters; however, MS discriminated African populations from

156    the others in the t-SNE analysis (S4 Fig A).

157

158    *MS marker set for individual identification*

159    We selected a set of 22 MS and calculated the discriminative power. From MS with 4-bp units,

160    22 MS with the highest heterozygosity were selected from each chromosome (Table 1).

161    For the selected 22 loci, we calculated the discriminative power using allele frequencies of 255

162    Japanese ICGC datasets (Table 1). The discriminative power for this MS set was estimated to be

163    $1.0 \times 10^{-17}$.

164

165    *MS analysis for an ancient DNA sample*

166    Although the genome-wide MS analysis of ancient DNA samples has yet to be conducted, an

167    analysis of MS variation in ancient DNA samples may contribute to clarifying the genetic

168    structure of ancient populations. Since low-quality sequencing data of ancient DNA samples can

169    result in incorrect results, we selected an ancient DNA sample with high sequence depth (sample

170    ID; F23) (20) and analyzed the distribution of the variant allele frequency (VAF) for MS (Fig 5

171    A-D). VAF was calculated by (number of reads with second most frequent pattern)/(total number

172    of reads) for heterozygous MS. The distributions of VAF were quite different between F23 and

173    modern human samples in MS with lengths < 3 bp, suggesting that genotypes contained many

174     errors. However, in MS with lengths ≥ 3 bp, the distributions of VAF were not different (Fig 5

175     A-D). Therefore, we used MS with lengths ≥ 3 bp and performed PCA with HGDP and SGDP

176     samples (All, East Asia, America, Oceania, Central Asia and Siberia, and South Asia) (Fig 5

177     E,F). F23 was clustered close to Central Asia and Siberia and East Asia populations in the PCA

178     plot.

179

180     **Discussion**

181     Population genetic studies strongly depend on variant calling. Thus, like other types of variants,

182     MS analysis is affected by genotyping methods. Since most MS calling methods analyze only

183     predefined MS regions, the numbers of target MS are different among studies (700,000 loci in

184     Gymrek et al., 2017 and Willems et al., 2014, and 1.6 million loci in Jakubosky, Smith, et al.,

185     2020) (15,18). In this study, we analyzed a larger number of MS regions (8,468,218 MS) than

186     previous studies for a comprehensive analysis of human MS polymorphisms. Compared to

187     conventional PCR-based MS studies(17), this study has another advantage: MS were not

188     influenced by ascertainment bias. Most conventional studies have analyzed pre-screened MS

189     marker sets (21), which are influenced by the MS marker selection. On the other hand, we did

190     not select MS based on the allele frequency in certain populations and could analyze features of

191     MS and compare the variation among different populations without the influence of

192     ascertainment bias.

193     We first selected parameters for the analysis based on the evaluation of monozygotic twins in

194     the KPGP. The concordant rate of this parameter set was estimated to be 99.87%, which is

195     sufficiently accurate for population studies (S1 Table). We then selected 8,468,218 MS based on

196     the call rate and HWE in the KPGP for 81 Korean individuals. Our MS calling identified

197     253,114 MS with MAF ≥ 1% in the SGDP and HGDP (S1 Fig). Since the SGDP and HGDP

198     datasets represent human genome diversity, these MS polymorphisms can be used for future

199     population studies (S1 data). Previous genome-wide studies reported that MS polymorphisms

200     could influence gene expression patterns and the risk of human diseases (17,22). Therefore

201     applications of our MS set and our MS calling method may contribute to discovering novel

202     disease susceptibility genes. In particular, 696 genes with non-$3n$ MS are good targets for disease

203     studies (S5 Table).

204     The amount of genetic variation reflects the effective population size ($N$). A comparison of

205     heterozygosities across regions showed Africa with the highest (Fig 3, S4 Table), and America,

206  which was estimated to have a very small effective population size, showed the lowest (Fig 3, S4

207  Table) (23). Such patterns were observed in other genetic variations and MS in a previous study

208  (14,24), suggesting that the heterozygosity of MS reflects the size of each population.

209  Theoretical population genetics also predicts that the effectiveness of natural selection depends

210  on the selection coefficient ($s$) of the genetic variation and population size (25). Therefore, a

211  comparison of genetic variations of MS in the CDS may provide additional information about the

212  evolution of MS. Most of the 3$n$ MS should not cause severe damage to protein functions and

213  have neutral or nearly neutral effects, whereas non-3$n$ MS cause a frameshift and should have

214  deleterious effects. We attempted to evaluate the strength of negative selection among

215  populations. For this purpose, we compared the average heterozygosity of non-3$n$ MS in CDS

216  with the average heterozygosity of 3$n$ MS in CDS among populations (Fig 3F, S4 Table). The

217  African population showed the lowest ratio, and populations with lower heterozygosity tended to

218  have a higher ratio (Fig 3F, S4 Table). This pattern indicates that the African population has the

219  largest effective population size and that stronger natural selection has acted to remove

220  deleterious non-3$n$ MS. In Central Asia and Siberia, the heterozygosity was not the lowest, but

221  the heterozygosity ratio of non-3$n$ MS to 3$n$ MS in CDS was the highest (Fig 3F, S4 Table). A

222  previous study showed that subdivided populations show a higher effective population size and

223  lower selection coefficient (26). The Central Asia and Siberia population may be composed of

224  subpopulations, which may affect the selection pressure against MS. These results indicate that

225  MS is evolved by the combination of population history and natural selection.

226      MS has been used to infer genetic structures because of high genetic diversity (27–30). In a

227  previous study, PCA was conducted using 53,002 MS, but the genetic structures by the MS PCA

228  were less clear than that by SNP PCA (18). In the present study, we used a larger number of MS

229  (253,114 MS with MAF $\geq$ 1 %) and obtained highly concordant results with SNPs (Fig 4 A,B).

230  Although the overall patterns were similar between MS and SNP (Fig 4, S3 Fig), small

231  differences were observed in African and Oceanian populations (Fig 4 C-F). In Oceanian

232  populations, NAN-Melanesians (NAN; Non-Austronesian) and Bougainville, who belong to

233  Melanesians, were clustered in the SNP PCA but not in the MS PCA (Fig 4 EF). In the African

234  populations, Biaka and Mbuti populations showed different patterns between the SNP and MS

235  PCAs (Fig 4 CD), which may be caused by hidden population structures. Although the efficiency

236  of using MS for genetic structures should be evaluated by a larger number of samples, these

237 results indicate that MS can be an additional marker set and may detect hidden population

238 structures in the human population.

239 In addition to the modern human samples, we analyzed a deep sequenced ancient human

240 sample (F23) (20). In ancient genome sequencing, the DNA fragmentation and library

241 construction process should affect the quality of the sequence reads. To evaluate the quality of

242 the MS call, we compared the distribution of the VAF of this sample with that of modern human

243 samples (Fig 5 A-D). The clear skew of the VAF was observed in MS with unit lengths ≤ 2 bp,

244 suggesting that MS with a short unit are strongly affected by the quality of DNA samples.

245 However, the distributions of unit lengths ≥ 3 bp were not different, and therefore we used these

246 MS for the analysis. In the PCA, F23 was close to East Asians, which is consistent with the SNP

247 PCA in a previous study (Kanzawa-Kiriyama et al., 2019) (Fig 5 E,F). This result suggests the

248 applicability of MS to ancient human samples.

249 We found 22 novel highly polymorphic MS for the personal identification. Using the allele

250 frequencies in a Japanese population, the discriminative power was estimated to be $1\times10^{-17}$,

251 which is sufficient for personal identification. Although the discriminative power of our MS set

252 is slightly lower than that of the Globalfiler kit, which is a standard MS set, for a Japanese

253 population ($5.6\times10^{18}$) (31), the length of our MS was shorter and can be genotyped by short-read

254 sequencers. Additionally, the PCR success rate of MS is known to be affected by the length of

255 the MS (32), and our shorter MS may be robust to DNA degradation.

256 This study provides a comprehensive catalog of MS in human populations and shows the

257 applicability of MS to modern and ancient human population studies. Nevertheless, our study has

258 several limitations. First, the genotyping of MS needs reads that cover MS regions. Therefore,

259 the amount of data and read length strongly affect the results. For example, we removed 824,459

260 MS and 395 samples from the SGDP and HGDP due to insufficient depth. Deeper sequence data

261 would improve the quality of the MS calling. Second, long MS cannot be analyzed using short-

262 read data. A recent study using a long-read sequencer reported high genetic variation in long

263 repeat regions (33). In the future, the application of our algorithm to long-read data should detect

264 a larger number of polymorphic MS.

265 To conclude, here we analyzed MS polymorphisms using large publicly available human

266 genome sequencing datasets. This study revealed a pattern of MS polymorphisms and identified

267 polymorphic MS in the human population. The comparison of the heterozygosity among

268 populations suggests that MS have evolved by random genetic drift and negative selection. PCA

269    suggests that MS detect the genetic structures of human populations. Currently, large-scale

270    sequencing projects are ongoing worldwide, in which the analysis of MS, in addition to SNPs,

271    should provide deeper understanding of human genetic variations and benefit genome medicine.

272

273    **Materials and Methods**

274    *Data*

275    We downloaded the following publicly available sequencing datasets: the Korean Personal

276    Genome Project (KPGP), Simons Genome Diversity Project (SGDP)(18), and Human Genome

277    Diversity Project (HGDP)(19), Japanese samples from the International Cancer Genome

278    Consortium (ICGC)(34), and an ancient DNA sample (20) (S6 Table, S7 Table). The SGDP (n =

279    300) and HGDP (n = 1064) samples were collected from various populations throughout the

280    world. The KPGP sequenced 107 Koreans. The ICGC performed WGS of cancer and matched

281    normal samples; in this study, we used the WGS data of normal Japanese samples (n = 255). A

282    deep sequenced ancient genome dataset (F23) was also analyzed (20).

283       KPGP data were used for the MS selection and parameter optimization of the MS calling.

284    The KPGP has data from monozygotic twins (KPGP-00088 and KPGP-00089), which were used

285    for the parameter optimization of the MS calling (S6 Table). One of the monozygotic twins and

286    other Korean samples (in total n = 81) were used for the MS selection. The HGDP and SGDP

287    samples were merged and used as a single dataset (S6 Table). When population names were

288    inconsistent between the SGDP and HGDP, we adopted the population names of the SGDP (S6

289    Table).

290       As a result of the sample selection (see below), we selected 81 Korean samples from the

291    KPGP and 969 samples from the HGDP and SGDP (138 samples from Africa, 70 from America,

292    48 from Central Asia and Siberia, 193 from East Asia, 38 from Oceania, 195 from South Asia,

293    and 287 from West Eurasia). The quality of the ICGC sequencing data was not constant among

294    samples; therefore, Japanese samples from the ICGC were used to estimate the allele frequencies

295    of our MS marker set for personal identification (S7 Table).

296       The downloaded bam files were results of the mapping to the GRCh37; therefore, our analysis

297    was based on the GRCh37.

298

299    *MS genotyping using MIVcall method*

300    Target MS regions were selected in our previous study (13) using three software packages：

301    MSDetector, Tandem Repeat Finder, and MISA software (35–37). Regions were filtered based

302    on the uniqueness of the flanking sequences and the distance to other MS. Insertions and

303    deletions in a target MS were detected using the MIVcall method (13); MIVcall counts the

304    length of each MS in each read. When multiple lengths are observed in a MS locus in a sample,

305    the most frequent pattern is assumed to be present, and the second most frequent pattern is

306    examined. The likelihood value was calculated based on the number of reads. Genotypes are

307    determined based on the likelihood value, the number of reads that support the pattern, and the

308    VAF.

309

310    *Establishing the MS detection method*

311    In our previous study (13), the optimal criteria of likelihood and number of reads (the minimum -

312    $\log_{10}$(likelihood) value and minimum number of reads for allele detection) were chosen for

313    analyzing somatic mutations. To obtain the optimal criteria for a polymorphism, we used

314    monozygotic twins in the KPGP (KPGP-00088 and KPGP-00089). Since all genotypes of

315    monozygotic twins are identical, we tested various parameter sets and compared the concordance

316    rates of genotypes between twins.

317

318    *Sample selection and MS filtering*

319    Since MS are susceptible to sequencing errors, selecting high-quality samples is necessary. Thus,

320    MS covered by less than 10 reads were considered MS with insufficient depth. We excluded

321    samples if more than 4% of MS loci had insufficient depth.

322        Next, we selected MS loci from the 9,292,677 MS selected in our previous study (13). Using

323    the 81 Korean samples in the KPGP, we counted the number of samples with insufficient depth

324    for each MS and removed samples if the percentage of MS with insufficient depth was ≥ 15%.

325    Additionally, we tested deviations from the HWE with Fisher's exact test for 2 x $n$ contingency

326    table ($n$; number of genotypes) in the KPGP ($\alpha = 0.0001$).

327

328    *Genome-wide pattern of MS*

329    To reveal the landscape of MS polymorphisms, we analyzed the features of MS. This analysis

330    was performed using HGDP and SGDP samples. We analyzed the association of the length of a

331    MS region in the reference genome with the number of samples with insufficient depth and

332    heterozygosity. We also examined the number and heterozygosity of MS for repeat units with

333    different sequences (for example, A, G, and AC). In this analysis, we merged MS with different

334    unit sequences if reverse-complement (e.g., GT to CA) or reverse (e.g., TA to AT) generated the

335    same sequences.

336        We then focused on MS in CDS and in non-CDS. We compared GC content, lengths of MS,

337    number of alleles, and heterozygosity between CDS and non-CDS MS. In CDS, MS were

338    classified into $3n$ MS (3 and 6 bp) and non-$3n$ MS (1, 2, 4, 5 and 7 bp). We compared

339    heterozygosity among all MS, $3n$ MS, non-$3n$ MS, MS in CDS, and MS in non-CDS. The ratios

340    of the mean heterozygosity of non-$3n$ MS to the mean heterozygosity of $3n$ MS were compared

341    among the different populations. We also performed a pathway analysis for genes with multi-

342    allelic non-$3n$ MS using the Reactome database.

343

344    *Analysis of population structure with MS*

345    We conducted five dimensionality reduction methods: Principal Component Analysis (PCA), t-

346    Distributed Stochastic Neighbor Embedding (t-SNE), Uniform Manifold Approximation and

347    Projection (UMAP), PCA-t-SNE, and PCA-UMAP. For these analyses, the genotypes of MS and

348    SNPs were converted to numerical values. SNP genotypes were converted to numerical values

349    by counting the number of minor alleles. For MS, we used two methods, the multiallelic method

350    and average method. In the multiallelic method, each allele at a MS locus is treated as a different

351    marker (if three samples have the following genotypes, [8/8, 13/14, 8/14], we converted them

352    into 3 independent pseudo-loci: [2, 0, 1] (8 or else), [0, 1, 0] (13 or else), and [0, 1, 1] (14 or

353    else)) (38). In the average method, we calculated the average length of two alleles in each

354    individual (in the previous example, the average method converts the genotypes to [8, 13.5, 11])

355    (S5 Fig). We performed a PCA with both methods and compared the results.

356        We used the MS and SNP with a MAF $\geq$ 1% in the SGDP and HGDP samples. MAF was

357    calculated by 1 – (major allele frequency). A PCA was conducted for all samples and samples in

358    each region (Africa, America, Central Asia and Siberia, East Asia, Oceania, South Asia, and

359    West Eurasia). Other dimensionality reduction methods (t-SNE, UMAP, PCA-t-SNE, and PCA-

360    UMAP) were applied to all samples only.

361

362    *MS set for personal identification*

363     Individual identification with MS is an important topic in human genetics. We selected a set of

364     22 MS and estimated the discriminative power. We calculated the heterozygosity of MS with

365     repeat unit lengths = 4 and selected 22 MS with the highest heterozygosity in each chromosome.

366     The allele frequencies of the selected 22 loci were estimated using 255 Japanese samples from

367     the ICGC data. The discriminative power was calculated as the product of frequencies of the

368     most frequent genotype in each locus (31).

369

370     *Analysis of ancient DNA samples*

371     To analyze MS variation in ancient DNA, we analyzed one ancient DNA sample with a higher

372     depth of coverage from a previous study (Sample ID; F23) (20). To examine the quality of the

373     variant calling, we calculated the VAF of each MS of F23 and compared it with the average VAF

374     of 10 randomly selected HGDP samples. We then conducted a PCA for F23 sample with the

375     SGDP and HGDP samples.

376

377     *Programming languages*

378     We used Python (https://www.python.org) for this study. PCA, TSNE, UMAP, and Decision

379     Tree were conducted with the sklearn package.

380

387

388     **Author Contributions**

389     Study design: A.F. Data analysis: L.G. and A. F. Manuscript writing: L.G. and A. F.

390     Interpretation of data: L.G., Y. K. and A. F.

391

392     **URLs.**

393     KPGP sequencing data: http://kpgp.kr.

394     The Reactome database: https://reactome.org/PathwayBrowser/

395

396 **Disclosure declarations**

397 The authors declare that they have no competing interests.

398

399 **FIGURE LEGENDS**

400

401 **Fig. 1.**

402 Features of microsatellites (MS) in the human genome. (A) MS length and number of samples

403 with insufficient depth. The MS length was negatively correlated to the number of samples with

404 insufficient depth (p-value $< 10^{-200}$ Kruskal-Wallis test). (B) MS length and heterozygosity. MS

405 length was positively correlated to the number of alleles (p-value $< 10^{-200}$ Kruskal-Wallis test).

406 (C) Proportion of MS unit types. (D) Number of MS and heterozygosity of different units. MS

407 with higher AT content had significantly higher heterozygosity (Pearson's uncorrelated test p-

408 value $= 1.47 \times 10^{-61}$).

409

410 **Fig. 2.**

411 Features of MS in CDS and non-CDS regions. (A) MS in CDS had higher GC content (p-value $=$

412 $6.49 \times 10^{-144}$ Wilcoxon rank sum test). (B) MS in CDS had shorter length (p-value $= 1.25 \times 10^{-29}$

413 Wilcoxon rank sum test). (C) MS in CDS had fewer alleles (p-value $= 1.80 \times 10^{-47}$ Wilcoxon rank

414 sum test). (D) MS in CDS had lower heterozygosity (p-value $= 1.42 \times 10^{-89}$ Wilcoxon rank sum

415 test). (E) Total number of MS loci per unit length in non-CDS. (F) Total number of MS loci per

416 unit length in CDS. (G) Number of multiallelic MS of each unit length in CDS. (H) Number of

417 multiallelic MS of each unit length in non-CDS. (I) Heterozygosity of MS in non-CDS and CDS.

418 Non-3$n$ MS had lower heterozygosity than 3$n$ MS (p-value $= 5.23 \times 10^{-13}$ Wilcoxon rank sum test).

419

420 **Fig. 3.**

421 Heterozygosity of MS among demographic regions. (A) Heterozygosity of all MS. (B)

422 Heterozygosity of MS in CDS. (C) Heterozygosity of MS in non-CDS. (D) Heterozygosity of 3$n$

423 MS in CDS. (E) Heterozygosity of non-3$n$ MS in CDS. In (A)-(E), regions were sorted by their

424 mean heterozygosity. (F) Ratios of mean heterozygosity of non-3$n$ MS to 3$n$ MS in CDS.

425

426 **Fig. 4.**

427    PCA using MS and SNPs. (A) PCA for all samples using MS. (B) PCA for all samples using

428    SNPs. (C) PCA for African populations using MS. (D) PCA for African populations using SNPs.

429    PC2 values of the Mbuti and Biaka populations were different between MS and SNP. (E) PCA

430    for Oceanian populations using MS. (F) PCA for Oceanian populations using SNPs.

431    NAN_Melanesuans and Bougainville were clustered together in the SNP but separated in the MS.

432

433    **Fig. 5.**

434    Analysis of an ancient sample (F23). Distribution of VAF in F23 and HGDP samples for unit

435    lengths ≥ 1 bp (A), ≥ 2 bp (B), ≥ 3 bp (C), and ≥ 4 bp (D). The distributions of VAF were quite

436    different between F23 and modern human samples in MS with lengths ≤ 2 bp. (E) PCA using

437    MS for F23 and all modern human samples. (F) PCA using MS for F23 and modern human

438    samples in South Asia, Oceania, America, East Asia, and Central Asia and Siberia.

439

440    **References**

441    1.    de Koning APJ, Gu W, Castoe TA, Batzer MA, Pollock DD. Repetitive elements may

442          comprise over Two-Thirds of the human genome. PLoS Genetics. 2011;7(12).

443    2.    Ellegren H. Microsatellites: Simple sequences with complex evolution. Nature Reviews

444          Genetics. 2004;5(6):435–45.

445    3.    Shriver MD, Smith MW, Jin L, Marcini A, Akey JM, Deka R, et al. Ethnic-affiliation

446          estimation by use of population-specific DNA markers. American Journal of Human

447          Genetics. 1997;60(4):957–64.

448    4.    Rosenberg NA, Pritchard JK, Weber JL, Cann HM, Kidd KK, Zhivotovsky LA, et al.

449          Genetic structure of human populations. Science (1979). 2002 Dec;298(5602):2381–5.

450    5.    Pemberton TJ, DeGiorgio M, Rosenberg NA. Population structure in a comprehensive

451          genomic data set on human microsatellite variation. G3: Genes, Genomes, Genetics.

452          2013;3(5):891–907.

453    6.    de Filippo C, Bostoen K, Stoneking M, Pakendorf B. Bringing together linguistic and

454          genetic evidence to test the Bantu expansion. Proceedings of the Royal Society B:

455          Biological Sciences. 2012;279(1741):3256–63.

456    7.    Hofer T, Ray N, Wegmann D, Excoffier L. Large allele frequency differences between

457          human continental groups are more likely to have occurred by drift during range

458          expansions than by selection. Annals of Human Genetics. 2009;73(1):95–108.

459    8.    Novembre J, Johnson T, Bryc K, Kutalik Z, Boyko AR, Auton A, et al. Genes mirror
460           geography within Europe. Nature. 2008;456(7218):98–101.

461    9.    Butler JM. Genetics and genomics of core short tandem repeat loci used in human identity
462           testing. Vol. 51, Journal of Forensic Sciences. 2006. p. 253–65.

463    10.   Ruitberg CM, Reeder DJ, Butler JM. STRBase: A short tandem repeat DNA database for
464           the human identity testing community. Nucleic Acids Research. 2001;29(1):320–2.

465    11.   Ludeman MJ, Zhong C, Mulero JJ, Lagacé RE, Hennessy LK, Short ML, et al.
466           Developmental validation of GlobalFiler$^{TM}$ PCR amplification kit: a 6-dye multiplex assay
467           designed for amplification of casework samples. International Journal of Legal Medicine.
468           2018;132(6):1555–73.

469    12.   Nagasaki M, Yasuda J, Katsuoka F, Nariai N, Kojima K, Kawai Y, et al. Rare variant
470           discovery by deep whole-genome sequencing of 1,070 Japanese individuals. Nature
471           Communications. 2015;6:1–13.

472    13.   Fujimoto A, Fujita M, Hasegawa T, Wong JH, Maejima K, Oku-Sasaki A, et al.
473           Comprehensive Analysis of Indels in Whole-genome Microsatellite Regions and
474           Microsatellite Instability across 21 Cancer Types. Genome Reserch [Internet].
475           2020;(30):334–46. Available from:
476           https://genome.cshlp.org/content/30/3/334.full.pdf+html

477    14.   Willems T, Gymrek M, Highnam G, Mittelman D, Erlich Y. The landscape of human STR
478           variation. Genome Research. 2014;24(11):1894–904.

479    15.   Jakubosky D, Smith EN, D'Antonio M, Jan Bonder M, Young Greenwald WW,
480           D'Antonio-Chronowska A, et al. Discovery and quality analysis of a comprehensive set of
481           structural variants and short tandem repeats. Nature Communications. 2020;11(1).

482    16.   Gymrek M, Willems T, Reich D, Erlich Y. Interpreting short tandem repeat variations in
483           humans using mutational constraint. Nature Genetics. 2017;49(10):1495–501.

484    17.   Gymrek M, Willems T, Guilmatre A, Zeng H, Markus B, Georgiev S, et al. Abundant
485           contribution of short tandem repeats to gene expression variation in humans. Nature
486           Genetics [Internet]. 2015;48(1):22–9. Available from: http://dx.doi.org/10.1038/ng.3461

487    18.   Mallick S, Li H, Lipson M, Mathieson I, Gymrek M, Racimo F, et al. The Simons
488           Genome Diversity Project: 300 genomes from 142 diverse populations. Nature.
489           2016;538(7624):201–6.

490    19.    Cavalli-sforza LL. The Human Genome Diversity Project: past , present and future.
491            2005;6(April):3–10.

492    20.    Kanzawa-Kiriyama H, Jinam TA, Kawai Y, Sato T, Hosomichi K, Tajima A, et al. Late
493            jomon male and female genome sequences from the funadomari site in Hokkaido, Japan.
494            Anthropological Science. 2019;127(2):83–108.

495    21.    Cacciò S, Homan W, Camilli R, Traldi G, Kortbeek T, Pozio E. A microsatellite marker
496            reveals population heterogeneity within human and animal genotypes of Cryptosporidium
497            parvum. Parasitology. 2000;120(3):237–44.

498    22.    Jakubosky D, D'Antonio M, Bonder MJ, Smail C, Donovan MKR, Young Greenwald
499            WW, et al. Properties of structural variants and short tandem repeats associated with gene
500            expression and complex traits. Nature Communications. 2020;11(1):1–15.

501    23.    Hey J. On the number of new world founders: A population genetic portrait of the
502            peopling of the Americas. PLoS Biology. 2005;3(6):0965–75.

503    24.    Auton A, Abecasis GR, Altshuler DM, Durbin RM, Bentley DR, Chakravarti A, et al. A
504            global reference for human genetic variation. Vol. 526, Nature. Nature Publishing Group;
505            2015. p. 68–74.

506    25.    Charlesworth B. Fundamental concepts in genetics: Effective population size and patterns
507            of molecular evolution and variation. Vol. 10, Nature Reviews Genetics. 2009. p. 195–205.

508    26.    Cherry JL, Wakeley J. A Diffusion Approximation for Selection and Drift in a Subdivided
509            Population [Internet]. 2003. Available from:
510            https://academic.oup.com/genetics/article/163/1/421/6052796

511    27.    Vieira MLC, Santini L, Diniz AL, Munhoz C de F. Microsatellite markers: What they
512            mean and why they are so useful. Genetics and Molecular Biology. 2016;39(3):312–28.

513    28.    Aimé C, Verdu P, Ségurel L, Martinez-Cruz B, Hegay T, Heyer E, et al. Microsatellite
514            data show recent demographic expansions in sedentary but not in nomadic human
515            populations in Africa and Eurasia. European Journal of Human Genetics.
516            2014;22(10):1201–7.

517    29.    Shriver MD, Jin L, Ferrell RE, Deka R. Microsatellite data support an early population
518            expansion in Africa. Genome Research. 1997;7(6):586–91.

519    30.    L. J. McIver. Population-scale analysis of human microsatellites reveals novel sources of
520            exonic variation. NIH [Internet]. 2013;23(1):1–7. Available from:
521            https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3624763/pdf/nihms412728.pdf

522   31.   Fujii K, Watahiki H, Mita Y, Iwashima Y, Kitayama T, Nakahara H, et al. Allele
523         frequencies for 21 autosomal short tandem repeat loci obtained using GlobalFiler in a
524         sample of 1501 individuals from the Japanese population. Legal Medicine [Internet].
525         2015;17(5):306–8. Available from: http://dx.doi.org/10.1016/j.legalmed.2015.08.007

526   32.   Schneider PM, Bender K, Mayr WR, Parson W, Hoste B, Decorte R, et al. STR analysis
527         of artificially degraded DNA - Results of a collaborative European exercise. Forensic
528         Science International. 2004;139(2–3):123–34.

529   33.   Audano PA, Sulovari A, Graves-Lindsay TA, Cantsilieris S, Sorensen M, Welch AME, et
530         al. Characterizing the Major Structural Variant Alleles of the Human Genome. Cell. 2019
531         Jan 24;176(3):663-675.e19.

532   34.   Campbell PJ, Getz G, Korbel JO, Stuart JM, Jennings JL, Stein LD, et al. Pan-cancer
533         analysis of whole genomes. Nature. 2020;578(7793):82–93.

534   35.   Benson G. Tandem repeats finder: A program to analyze DNA sequences. Nucleic Acids
535         Research. 1999;27(2):573–80.

536   36.   Girgis HZ, Sheetlin SL. MsDetector: Toward a standard computational tool for DNA
537         microsatellites detection. Nucleic Acids Research. 2013;41(1).

538   37.   Hause RJ, Pritchard CC, Shendure J, Salipante SJ. Classification and characterization of
539         microsatellite instability across 18 cancer types. Nature Medicine. 2016;22(11):1342–50.

540   38.   Putman AI, Carbone I. Challenges in analysis and interpretation of microsatellite data for
541         population genetic studies. Ecology and Evolution. 2014;4(22):4399–428.

542

543

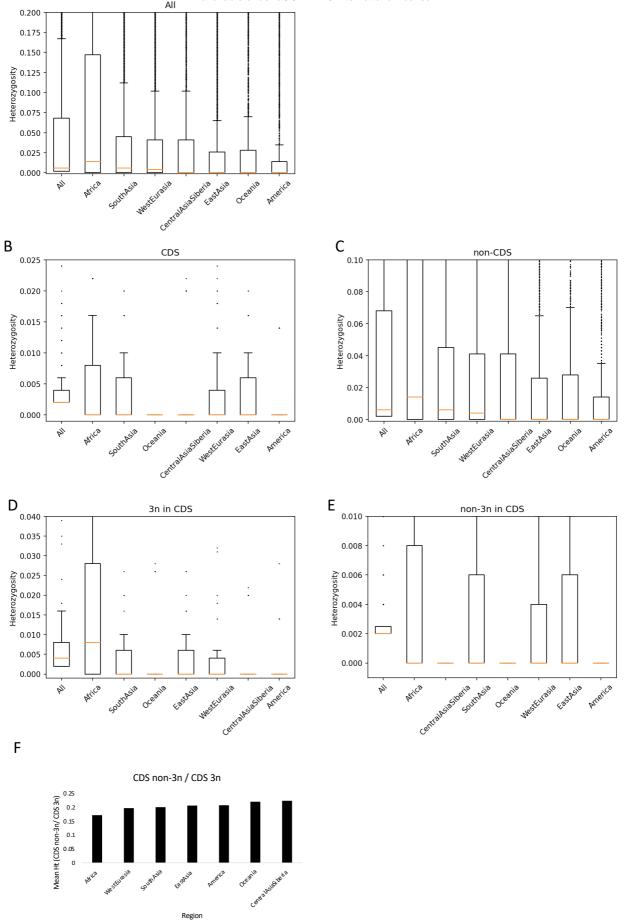Table1 MS set for individual identification

| Chromosome | Start position | End position | Type | Max length of MS in ICGC Japanese samples (bp) | Max length of MS in HGDP+SGDP (bp) | Number of alleles in ICGC Japanese samples | Number of alleles in HGDP+SGDP | Match probability in ICGC Japanese samples | Match probability in HGDP+SGDP |
|---|---|---|---|---|---|---|---|---|---|
| 1 | 105083241 | 105083272 | (TATC)n | 48 | 55 | 12 | 20 | 0.074 | 0.076 |
| 2 | 11032690 | 11032717 | (CTTC)n | 44 | 56 | 5 | 10 | 0.270 | 0.134 |
| 3 | 76233694 | 76233746 | (AGAT)n | 61 | 62 | 11 | 19 | 0.075 | 0.107 |
| 4 | 117885601 | 117885630 | (TATC)n | 49 | 53 | 11 | 16 | 0.085 | 0.054 |
| 5 | 82900377 | 82900424 | (TCTA)n | 55 | 55 | 10 | 14 | 0.154 | 0.095 |
| 6 | 49948157 | 49948210 | (ATCT)n | 66 | 70 | 10 | 11 | 0.142 | 0.091 |
| 7 | 67348144 | 67348171 | (TTTC)n | 44 | 60 | 7 | 17 | 0.318 | 0.150 |
| 8 | 108769208 | 108769243 | (TAAG)n | 56 | 60 | 7 | 13 | 0.247 | 0.165 |
| 9 | 101627740 | 101627784 | (ATCT)n | 57 | 61 | 7 | 8 | 0.231 | 0.178 |
| 10 | 2918258 | 2918301 | (TAGA)n | 52 | 48 | 7 | 9 | 0.175 | 0.095 |
| 11 | 114544813 | 114544850 | (AAAC)n | 47 | 47 | 10 | 14 | 0.209 | 0.149 |
| 12 | 121657632 | 121657667 | (TAGA)n | 52 | 52 | 9 | 11 | 0.189 | 0.116 |
| 13 | 107091029 | 107091073 | (ATGG)n | 50 | 54 | 5 | 9 | 0.187 | 0.148 |
| 14 | 38581137 | 38581184 | (ATAG)n | 58 | 58 | 12 | 15 | 0.171 | 0.122 |
| 15 | 47067031 | 47067072 | (GAAT)n | 50 | 46 | 7 | 12 | 0.202 | 0.143 |
| 16 | 86386308 | 86386351 | (GATA)n | 56 | 60 | 7 | 9 | 0.156 | 0.155 |
| 17 | 72561544 | 72561595 | (CTAT)n | 56 | 56 | 5 | 10 | 0.204 | 0.155 |
| 18 | 5249011 | 5249068 | (AGAT)n | 70 | 74 | 12 | 21 | 0.180 | 0.103 |
| 19 | 15754216 | 15754267 | (ATAG)n | 60 | 60 | 8 | 15 | 0.164 | 0.116 |
| 20 | 5596482 | 5596528 | (TATG)n | 47 | 51 | 7 | 9 | 0.178 | 0.126 |
| 21 | 42031280 | 42031319 | (TATC)n | 52 | 52 | 9 | 13 | 0.136 | 0.157 |
| 22 | 36513967 | 36514013 | (TATC)n | 57 | 75 | 13 | 31 | 0.194 | 0.114 |
| | | | | | | | Product of the match probability of all MS | $1.0 \times 10^{-17}$ | $6.3 \times 10^{-21}$ |

Fig 1

Fig 2

Fig 3

Fig 4

## Fig 5

A



B



C



D



E



F