### **ORIGINAL ARTICLE**



# Mixture deconvolution by massively parallel sequencing of microhaplotypes

Lindsay Bennett<sup>1</sup> · Fabio Oldoni<sup>2</sup> · Kelly Long<sup>2</sup> · Selena Cisana<sup>2</sup> · Katrina Madella<sup>2</sup> · Sharon Wootton<sup>3</sup> · Joseph Chang<sup>3</sup> · Ryo Hasegawa<sup>3</sup> · Robert Lagacé<sup>3</sup> · Kenneth K. Kidd<sup>4</sup> · Daniele Podini<sup>2</sup>

Received: 29 January 2018 / Accepted: 23 January 2019 © Springer-Verlag GmbH Germany, part of Springer Nature 2019

#### Abstract

Short tandem repeat polymorphisms (STRs) are the standard markers for forensic human identification. STRs are highly polymorphic loci analyzed using a direct PCR-to-CE (capillary electrophoresis) approach. However, STRs have limitations particularly when dealing with complex mixtures. These include slippage of the polymerase during amplification causing stutter fragments that can be indistinguishable from minor contributor alleles, preferential amplification of shorter alleles, and limited number of loci that can be effectively co-amplified with CE. Massively parallel sequencing (MPS), by enabling a higher level of multiplexing and actual sequencing of the DNA, provides forensic practitioners an increased power of discrimination offered by the sequence of STR alleles and access to new sequence-based markers. Microhaplotypes (i.e., microhaps or MHs) are emerging multi-allelic loci of two or more SNPs within < 300 bp that are highly polymorphic, have alleles all of the same length, and do not generate stutter fragments. The growing number of loci described in the literature along with initial mixture investigations supports the potential for microhaps to aid in mixture interpretation and the purpose of this study was to demonstrate that practically. A panel of 36 microhaplotypes, selected from a set of over 130 loci, was tested with the Ion S5<sup>TM</sup> MPS platform (Thermo Fisher Scientific) on single-source samples, synthetic two-to-six person mixtures at different concentrations/contributor ratios, and on crime scene-like samples. The panel was tested both in multiplex with STRs and SNPs and individually. The analysis of single-source samples showed that the allele coverage ratio across all loci was  $0.88 \pm 0.08$  which is in line with the peak height ratio of STR alleles in CE. In mixture studies, results showed that the input DNA can be much higher than with conventional CE, without the risk of oversaturating the detection system, enabling an increased sensitivity for the minor contributor in imbalanced mixtures with abundant amounts of DNA. Furthermore, the absence of stutter fragments simplifies the interpretation. On casework-like samples, MPS of MHs enabled the detection of a higher number of alleles from minor donors than MPS and CE of STRs. These results demonstrated that MPS of microhaplotypes can complement STRs and enhance human identification practices when dealing with complex imbalanced mixtures.

**Keywords** Microhaplotype  $\cdot$  Single-nucleotide polymorphism  $\cdot$  Massively parallel sequencing (MPS)  $\cdot$  Mixture deconvolution  $\cdot$  Forensic DNA samples

**Electronic supplementary material** The online version of this article (https://doi.org/10.1007/s00414-019-02010-7) contains supplementary material, which is available to authorized users.

Daniele Podini podini@gwu.edu

- <sup>1</sup> Metro Nashville Police Department Crime Laboratory, 400 Myatt Drive, Madison, TN 37115, USA
- <sup>2</sup> The Department of Forensic Sciences, The George Washington University, 2100 Foxhall Road NW, Washington, D.C. 20007, USA
- <sup>3</sup> Thermo Fisher Scientific, 180 Oyster Point Boulevard, San Francisco, CA 94080, USA
- <sup>4</sup> Department of Genetics, Yale University, 333 Cedar Street, New Haven, CT 06520, USA

## Introduction

Short tandem repeat polymorphisms (a.k.a. STRs) are the standard DNA markers used in human identification due to their high heterozygosity and their ease of amplification and detection via polymerase chain reaction (PCR) and capillary electrophoresis (CE) [1]. Forensic STRs are likely here to stay; yet they have limitations. The number of loci that can be analyzed in a single assay is limited both by the size ranges of the fragments at the various loci and by the number of dyes that CE platforms allow [2]. During PCR amplification of STRs, polymerase slippage can occur, for some loci over 20% of the time, leading to the generation of fragments of

one or two repeat lengths smaller (or larger) that are artifact peaks (stutter) indistinguishable in size from true alleles of the same size [3]. Increased sensitivity of commercial STR PCR Kits has led to an increased number of mixtures obtained from crime scene evidence. Mixtures currently represent one of the greatest challenges in forensic DNA analysis. Sample degradation, allele sharing between contributors, allele drop-out, preferential amplification of smaller alleles, and stutter all contribute to the difficulties of correctly and consistently interpreting complex (i.e., more than two-person) mixtures using standard STR analysis [3]. Probabilistic genotyping software [4, 5] allows the best use of the data and provides consistency to the results although the inherent issues with CE STR data remain. Massively parallel sequencing (MPS) can enhance mixture deconvolution of STR loci by distinguishing alleles that have the same length but different sequences, both in the repeat sequences and in the flanking regions [6-9]. Yet, experimental data have shown that interpretation of STR sequences from MPS platforms is not as straightforward as originally anticipated [10].

With the increased use of MPS, other markers have been suggested that would provide additional information in human identification. Individual SNPs, both autosomal and on the X and Y chromosomes, have been tested for use in human identification [11, 12]. INDELs can also be used for identification of compromised or degraded samples [13, 14]. All of these markers tend to be di-allelic although tri- and tetra-allelic SNPs have been identified and advocated because they provide more information [15, 16]. Microhaplotypes (microhaps or MHs) provide another approach to more informative multiallelic markers [17]; these are loci composed of two or more SNPs within a short distance from each other (generally < 300 nucleotides, i.e., "micro") with three or more allelic combinations (i.e., "haplotypes") [18-21]. Conventional Sanger sequencing of PCR products does not allow the determination of cis/trans relationships among alleles of the individual SNPs (i.e., the phases of the two haplotypes) [22]. MPS methods, instead, when SNPs are in the same amplified fragment (amplicon), allow direct sequencing of individual DNA strands (clonal/single molecule sequencing), thereby yielding unambiguous phase of parental haplotypes of SNPs.

Previous papers on microhaplotypes have published the allele frequencies for different sets of MHs in up to 96 different population samples [21, 23]. These papers have demonstrated three aspects of the potential uses of microhaps in forensics: (1) theoretical ability to detect mixtures with very high probability [24–26], (2) ability to assign a multi-locus profile to one of several distinct biogeographic regions of ancestry [27–29], and (3) generation of extremely small random match probabilities (RMPs) with selected microhaps [17, 24]. Here we have undertaken to move from the theoretical ability to detect mixtures to practical demonstration of mixture detection by typing a panel of 36 microhaps on the Ion S5<sup>TM</sup>

MPS platform (Thermo Fisher Scientific, Waltham, MA, USA). These markers were selected from a larger panel of over 130 microhaps potentially useful for forensic applications [21, 23]. The 36-locus panel was tested on reference samples, on synthetic mixtures, and on crime scene-like samples. As part of demonstrating mixture deconvolution with microhaps, we evaluated the performance of the panel on casework-like samples including a comparison of the microhap results with results using the standard forensic STRs.

# **Material and methods**

#### Selection of DNA samples

DNA samples from anonymous European American and African American individuals of self-identified ancestry were randomly selected from a large set of samples collected and extracted between 1993 and 2003. The use of these samples was declared exempt from IRB approval by The George Washington University's Office of Human Research (OHR-IRB # 090636) given their anonymous origin. From the same sample collection, eight pristine DNA samples of different biogeographic ancestry were selected and used as singlesource or for synthetic mixtures. Studies mimicking the analysis of crime scene samples (henceforth referred to as "forensic samples") included DNA extracts from eight specimens: cigarette butts (n=3), water bottles (n=3), a piece of chewing gum, and an aged bloodstain. Except for the bloodstain, the number of contributors and their genetic profile(s) were unknown before analysis as these samples were collected from the street and trash cans outside the department. Synthetic mixtures were created by combining single-source samples of known concentrations of DNA at the mixture ratios described in the mixture studies section below.

## Selection and genotyping of loci

The specific panels of STR, SNP, and MH loci analyzed in this study are listed in Supplemental Tables 1, 2 and 3, respectively [21, 30–36]. Two separate marker panels were tested: the first one included primer pairs for 30 STRs plus Amelogenin, 45 SNPs (43 SNPs plus one indel) and 36 MHs, and the second panel included only primer pairs for the 36 MHs. The study was conducted using Thermo Fisher Scientific equipment and kits. Library preparation with the first marker panel was performed using an automated library preparation workflow, which allowed processing eight samples at a time on the Ion Chef<sup>TM</sup> System. DNA samples were processed using the Ion AmpliSeq<sup>TM</sup> Kit for Chef DL8 and barcoded with the IonCode<sup>TM</sup> adaptor following manufacturer's guidelines. Libraries for the 36-MH panel were prepared manually using

Ion Xpress<sup>TM</sup> for barcoding up to 40 samples at a time, as previously described [37]. All the resulting libraries were quantified using the Ion Library TaqMan<sup>TM</sup> Quantitation Kit. Emulsion PCR and chip templating were performed on the Ion Chef<sup>TM</sup> System for both marker panels, and followed by sequencing using the Ion S5TM Precision ID Chef & Sequencing Kit according to the manufacturer's instructions, as previously described [37]. Sequencing was performed on either the Ion 521<sup>TM</sup> (subsequently replaced by 530 <sup>TM</sup>) or the Ion 530<sup>TM</sup> chips and with a maximum of 40 samples loaded per chip, as described below.

## **Experiments**

### Sample genotype study

The eight pristine DNA samples of different biogeographic ancestry, selected for use as single-source and in synthetic mixtures, were processed using the marker panel that also includes both STRs and individual SNPs in addition to the microhaps, and then sequenced on the Ion 521<sup>TM</sup> chip.

## Estimating the number of contributors

An estimate of the minimum number of contributors can be determined by the number of different alleles observed in an individual. On the simple assumption that all distinct alleles in a mixture can be detected, the presence of three alleles at a locus indicates the contribution of at least two donors, the presence of five alleles of at least three donors, and so forth. To assess the potential of the 36-microhap set to estimate the number of contributors in a forensic mixture, allelic diversity across the 36 loci was evaluated by two simulation tests. In particular, the simulations considered all combinations of nperson mixture scenarios (n = 1 to 6) using genotypes from a set of 19 African American (AA) individuals (#1 simulation) and 17 European American (EA) individuals (#2 simulation). For each number of contributors, randomly chosen individuals were mixed in silico and then the total number of distinct alleles was counted.

## **Mixture studies**

Two separate mixture studies (A and B) were designed. In mixture study A, we used the first MPS marker assay including primers for STRs, SNPs, and MHs described above. The input amount of template DNA used for MPS analysis was greater than what can be used in PCR for CE where the amplification of 5–7 ng of DNA would yield overloaded electropherograms and likely not interpretable results. As a result, no comparison was made with conventional CE-STR analysis for these samples. Three separate mixtures of one female to multiple males were prepared using 4 ng of female DNA and a

total of 0.4 ng of two, three, and four male DNAs corresponding to three, four, and five person mixtures, respectively. In addition, a six-person mixture was simulated at a ratio of 1:1:1:1:1:1 and included one female and five males. The three-person mixture was of a Hispanic (HS) female and two males, one HS, and one African American (AA); the fourperson mixture was of a HS female and three males, one HS, one AA, and one European American (EA). The fiveperson mixture was of a HS female and four males, two EAs, one HS, and one AA. Finally, the six-person mixture was composed of one HS female, two HSs, two EAs, and one AA male (Table 1).

The MPS assay used for mixture study B included only the 36 MHs. A total of two two-person mixtures (4:1 and 20:1) and one three-person mixture (4:2:1) were amplified with 2 ng and 0.25 ng of input DNA (Table 2). Sequencing of the mixed samples from both mixture studies was performed on the 530<sup>TM</sup> chip. In mixture study A, 24 samples were individually barcoded and loaded on one chip while in mixture study B, 40 samples were barcoded and loaded on another chip.

## **Forensic samples**

Standard extraction of DNA for the eight crime scene-like samples was performed using previously described methods commonly used in forensic laboratories [38, 39]. DNA extracts were quantified using quantitative real-time PCR with the Quantifiler® Trio DNA Quantification Kit on the Applied Biosystems® 7500 Real-Time PCR Systems following the manufacturer's instructions. Samples were analyzed with the automated library preparation process using the panel that combines 30 STRs plus Amelogenin, 45 SNPs, and 36 MHs, and sequenced all on the same 521<sup>TM</sup> chip. To compare the performance to conventional capillary electrophoresis, Amp*F*/STR<sup>TM</sup> Globalfiler® was used on these samples. PCR products were run on an ABI 3130 Genetic Analyzer following manufacturer's instructions.

## **MPS data analysis**

Sequencing results of STR, SNP, and MH loci were analyzed on the Ion Torrent suite server (version 5.0.2) using the HID Genotyper plugin (2.0r11562) as well as the defined targets Globalfiler\_Mixture\_ID\_NGS\_Panel\_targets\_v1.0.bed and hotspot regions Globalfiler\_Mixture\_ID\_NGS\_Panel\_ hotspot\_v1.0.bed. The plugin is designed to cover all marker types in the assay, and the default settings were applied to the analysis. Results were then uploaded to the Applied Biosystems<sup>™</sup> Converge<sup>™</sup> Software version 2.0 beta (Applied Biosystems, Carlsbad, CA) for visualization. The allele coverage ratio (ACR) of sister alleles, corresponding to the peak height ratio (PHR) with conventional CE methods, was calculated by dividing the total number of reads of the **Table 1** Summary of mixturestudy A. F = female sample, M =male sample, in "Ratio" column10 = 4 ng and 1 = 0.1 ng. HS =Hispanic, AA = AfricanAmerican, EA = EuropeanAmerican

Mixture study A	Numb. Of contributors	Sex	Ratio	Ancestry	Numb. alleles detected
Mix 1	3	F/M/M	10:1:1	HS/HS/AA	117
Mix 2	4	F/M/M/M/M	10:1:1:1	HS/HS/AA/EA	129
Mix 3	5	F/M/M/M/M	10:1:1:1:1	HS/HS/AA/EA/EA	138
Mix 4	6	F/M/M/M/M/M	1:1:1:1:1:1	HS/HS/AA/EA/EA	148

lower allele read count by the higher allele read count to obtain a maximum value of 1.0 for an equal number of reads.

# Results

## Single-source samples

To characterize the sequencing efficiency of the 36 MHs in the panel including also the STR and SNP amplicons, an average and standard deviation of the read depth for the single source samples was calculated and is shown in Fig. 1 (brown line). The read depth or allele coverage varied across loci with average allele coverage of 5717 reads  $\pm$  851. The allele coverage ratio was then determined for sister alleles at each heterozygote MH locus and is also depicted in Fig. 1 (blue line). The overall average of ACR observed was 0.88  $\pm$  0.08. The average read depth and ACRs showed no correlation between the two metrics, understandably because the allele coverage is proportional to the redundancy available on the chip and inversely proportional to the number of samples pooled in the run while the ACR should remain unaltered.

## Estimating the number of contributors to a mixture

A total of 44,436 mixtures were simulated in this study, 19,401 with EA samples, and 25,035 with the AA samples. The total number of unique alleles across all markers in the simulated two- to six-person mixtures was plotted separately

for the European American and African American data (Fig. 2a, b, respectively). The distributions differ in that African American individuals displayed, on average, a greater number of alleles in mixtures with the same number of contributors. For example, focusing on the x-axis one can observe that each distribution is shifted to higher numbers for the mixtures of African American individuals (Fig. 2b) compared to the mixtures of European American individuals (Fig. 2a). This agrees with the general observation that African populations have more genetic variation than non-African populations.

## **Mixture studies**

## Mixture study A

In the three-person mixture, 117 individual alleles were detected across the 36 MHs, in the four-person mixture 129 alleles were detected, in the five-person mixture 138 were detected, and in the six-person mixture 148 alleles were detected (Table 1). The maximum number of alleles detected at a single MH locus was five for the three-person mixture, six for both four- and five-person mixtures, and seven for the sixperson mixture.

In the examination of the STR results for mixtures with a 10:1 ratio, it was difficult to distinguish the true alleles from stutter fragments. If simply counting the number of fragments identified and defining them as possible alleles, the locus with the highest number of possible STR alleles in the three-person mixture had nine (9) fragments detected. In the four- and five-

Table 2	Summary of the results
from miz	xture study B

Mixture study B	Expected total numb. of alleles/Minor contributor alleles	Total alleles detected	Numb. unique minor cont. alleles detected	Numb. drop out alleles of minor cont.
4:1 2 ng	85/30	85	30	0
4:1 250 pg	85/30	75	20	10
20:1 2 ng	85/30	80	25	5
20:1 250 pg	85/30	74	19*	11
4:2:1 2 ng	108//18**	108	18	0
4:2:1 250 pg	108/18**	105	15	3

\*Alleles detected in the noise range, profile likely not suitable for comparison for the minor contributor \*\*Unique alleles of second (lowest) minor contributor



Fig. 1 Average allele coverage (left axis) and average allele coverage ratio ACR (right axis) with standard deviations for 36 multi-allelic loci in eight single-source samples

person mixtures, the numbers were 11 and 12, respectively. In the six-person mixture (1:1:1:1:1), given that there is an equal amount of genomic DNA for each contributor (1 ng) and all true alleles should have similar peak heights for comparable genotypes, most stutter peaks are distinguishable from true alleles and the greatest number of true alleles detected at a locus was nine (9). Figure 3 shows an example of the STR profile obtained at D12S391 where 11 potential alleles are present in the four-person mixture where the maximum number of possible alleles is eight (8).

Figure 4 shows the profiles at three microhap loci coamplified and sequenced with the STRs for the same fourperson mixture described in Fig. 3. The first locus, mh18KK-293 (using suggested nomenclature [40]), displays only three alleles with an ACR consistent with at least two contributors; mh13KK-213 shows five alleles consistent with at least three individuals; mh05KK-170 shows six alleles with an ACR consistent with at least three individuals and probably four assuming the major contributor was homozygous for allele "CAGA". Since alleles at these loci were detected in the same mixture, we can conclude that at least three, and more probably four individuals are included in the mixture.

## Mixture study B

The panel used for mixture study B targeted only the 36 MH loci, and the results are summarized in Table 2. In the twoperson mixtures, at the 4:1 ratio with 2 ng of input DNA all the alleles of the minor contributor were detected while with 250 pg a total of 10 alleles of the minor were undetected. At the 20:1 ratio, when the total amount of DNA was 2 ng, five alleles of the minor dropped out. When the input was lowered to 250 pg, 19 of the 30 unique alleles of the minor were detected but with a level of coverage that was in the noise range, thus likely not suitable for comparison if interpretation thresholds were available. In the three-person mixture (4:2:1), when the input was 2 ng, all the 108 possible alleles were detected while when the input was lowered to 250 pg three alleles of the second (lowest) minor contributor were undetected.

## **Forensic samples**

In this small sample set, overall MH MPS analysis shows a greater sensitivity than that of both CE and MPS STRs. An example of the added value offered by MHs in combination with STRs is illustrated by the results obtained from one of the three cigarette butts tested in this study. The CE STR profile obtained was consistent with a single-source profile with the exception of a single locus displaying three peaks where the lowest of the three peaks was in a stutter position yet above the stutter threshold (Fig. 5). The profile could be classified as a mixture of at least two individuals, but the minor contributor would likely be considered unsuitable for comparison. The MPS STR results of the same sample, amplified with the same amount of extract, indicated the presence of (at least) a second contributor at 10 of the loci (Fig. 6a). A thorough validation study of MPS STR data has not yet been performed, and stutter thresholds have not been defined. Thus, only fragments that were not in n-4 stutter position were considered suitable for comparison relative to the minor contributor. If stutter



**Fig. 2** Frequency of number of alleles counted in simulated *n*-person mixtures (n = 1 to 6) using profiles of (**a**) 17 European American (EA) and (**b**) 19 African American (AA) individuals. The x-axis represents the number of alleles identified in the various mixtures while the y-axis shows the frequency of that number in the totality of the *n*-person mixtures. The total number of mixtures simulated with EA samples was 19,401 (17 N = 1, 136 N = 2, 680 N = 3, 2380 N = 4, 6188 N = 5, and 10,000 N = 6) while the total number of mixtures simulated with AA samples was 25,035 (19 N = 1, 171 N = 2, 969 N = 3, 3876 N = 4, 10,000 N = 5, and 10,000 N = 6)

thresholds were available, fragments in stutter positions could potentially be classified as alleles for at least three loci. When analyzing the results from MPS of MHs, more than two alleles were observed in 27 out of 36 loci. Only 21 out of the 27 loci were considered suitable for comparison relative to the minor contributor. The remaining six loci were considered inconclusive also due to the absence of a comprehensive validation study, and derived ACR, to enable the assignment of two alleles to the same contributor. Figure 6b shows six representative MHs with three or four alleles. As a proof of concept, a random match probability (RMP) was calculated for the 21 MHs using PHASE allele frequencies reported in [21, 23]. Again, no stochastic threshold exists because no thorough validation study has yet been done. Consequently, when a single allele from the minor contributor was present, the conservative "2p" rule approach [42] was used to calculate the RMP. If both alleles of the minor contributor were detected, the 2pq rule was applied. The random match probability calculated from the MPS STR profile resulted in a frequency of the profile in the millionths range  $(10^{-7})$  while that calculated from the MH data was in the quadrillionths range  $(10^{-18})$ . Provided that absence of linkage disequilibrium between MH and STR loci is confirmed, the two values can then be multiplied to increase even further the power of discrimination.

## Discussion

The purpose of this study was to demonstrate practically what was previously only postulated, i.e., that massively parallel sequencing of microhaplotypes can aid human identification practices when dealing with mixtures. The efficacy of MH sequencing in a multiplex panel also containing STRs (and SNPs) was also evaluated. Variation in coverage was observed across loci, but the average ACR between sister alleles within a locus was relatively unaffected by the overall locus coverage. The ACR parameter is equivalent to the peak height ratio of conventional CE STR analysis and is important in mixture deconvolution to determine whether two alleles could have originated from the same contributor. The average ACR across all loci was  $0.88 \pm 0.08$ , with all but one of the heterozygous loci ranging between 0.64 and 0.9874. More data are needed to globally evaluate the performance of MHs within and across populations, and, as for STR analysis, it is likely that SNPs located in primer binding sites will cause allele drop-out or preferential amplification of certain alleles in some individuals.

One of the advantages of MPS is the ability to load in the PCR reaction a higher input amount of DNA than with the conventional direct fluorescent PCR-to-CE detection approach. Increasing the input DNA allows higher MPS coverage of all alleles in a mixture potentially raising the read depth of the minor contributions to values exceeding any stochastic threshold. The loading of 5.6 ng of genomic DNA at a 10:1:1:1:1 mixture ratio in a fluorescent PCR reaction coupled with CE detection would overload the system and generate a very noisy and challenging electropherogram to interpret. Conversely, loading a total of 1 ng of the same mixture would mean adding approximately 70 pg of each minor contributors.

An advantage of MHs, compared to STRs, is the absence of polymerase slippage during the amplification of tandem repeats that causes the generation of stutter peaks. The latter adds an undesirable level of complexity to the interpretation of imbalanced mixtures for both the determination of the Fig. 3 Output table (above) from the software Converge for locus D12S391 of the four-person mixture (10:1:1:1). Fragments in stutter position with coverage below the default stutter threshold are classified as stutter (light green) although a validation study to define the correct threshold has yet to be conducted. The output plot below shows a different way in which Converge can display the data. The plot mimics a conventional STR electropherogram with the coverage (like RFU) on the y-axis and the allele call on the x-axis

Allele	Status	Coverage	Sequence
14	STUTTER	560	[AGAT]7 [AGAC]6 [AGAT]1
23	ABOVE_ST	318	[AGAT]14 [AGAC]8 [AGAT]1
15	ABOVE_ST	7588	[AGAT]8 [AGAC]6 [AGAT]1
17	ABOVE_ST	462	[AGAT]10 [AGAC]6 [AGAT]1
18	STUTTER	1236	[AGAT]10 [AGAC]7 [AGAT]1
14	STUTTER	208	[AGAT]8 [AGAC]5 [AGAT]1
19	ABOVE_ST	5911	[AGAT]11 [AGAC]7 [AGAT]1
21	ABOVE_ST	608	[AGAT]13 [AGAC]8
21	ABOVE_ST	335	[AGAT]13 [AGAC]7 [AGAT]1
21	ABOVE_ST	742	[AGAT]14 [AGAC]6 [AGAT]1
18	STUTTER	276	[AGAT]11 [AGAC]6 [AGAT]1
D128 8000 7000 6000 5000 4000 3000 2000 1000	5391		

number of contributors and deconvolution of a mixture. As microhaps are single base pair sequence variations in singlecopy sequences lacking tandem repeat motifs, the molecular basis of stutter is absent, thus making stutter irrelevant. Moreover, all the alleles of a locus have the same amplicon length, which eliminates the issue of preferential amplification of smaller alleles within a locus, as commonly occurs with STRs. This is particularly critical in imbalanced mixtures where stutter fragments cannot be distinguished from true alleles of the minor contributor. Additionally, these new markers have a mutation rate that is several orders of magnitude lower than that of STRs making them also more effective for relationship testing in civil, criminal, and missing person cases [17].

When evaluating the results from the simulation study aimed at predicting the number of contributors, based on the total number of individual alleles detected in the sample, the differences obtained in the allele distribution from EA vs. AA are not surprising. In fact, African populations (hence, African Americans) have higher levels of polymorphism than non-African populations, as documented in multiple studies [43, 44]. In fact, the total number of individual allele counts in each mixture aligns better with the mixture simulation plots generated from AA allele frequencies, (Fig. 2b), than those generated from EA allele frequencies (Fig. 2a). For example, the four-person mixture composed of two HIS individuals, one AA and one EA, yielded 129 individual allele counts. The distribution of individual allele counts in the simulation plots for four-person mixtures ranged from 110 to 135 in the EA population and from 115 to 144 in the AA population. Although this is just one example, the 129 individual allele counts fit better with the latter distribution. These data suggest that counts of total unique alleles seen in a mixture can be a helpful approach to predict the number of contributors in a mixture. However, the prediction model depends on the biogeographic ancestry of contributors because the diversity of alleles seen in the populations of the contributors will be positively correlated with the number of unique alleles seen in a mixture. Another important parameter to consider is the quality of the sample being tested and whether or not allele dropout is expected due to low template of one or more contributors and/or degradation/inhibition issues.

The data generated in this study for the synthetic mixtures and forensic samples support the assumption that microhaplotype profiling can aid mixture deconvolution. For example, in the same four-person mixture (10:1:1:1) mentioned above, the maximum number of alleles per locus detected in MHs was six with an ACR consistent with at least four contributors. In the forensic sample, CE of STRs resulted in a mixture in which the results for the minor contributor would have been considered not suitable for comparison. For the same mixture, MHs showed three or more alleles at

Fig. 4 Converge output plot for three representative MHs for a four-person mixture (10:1:1:1). mh18KK-293 displays only three alleles with an ACR consistent with at least three contributors given significant imbalance between three alleles detected. mh13KK-213 shows five alleles also consistent with at least three individuals given the three minor alleles (at least two individuals) and the two major alleles (at least one). mh05KK-170 shows six alleles with an ACR consistent with at least four individuals given a major contributor homozygous for allele "CAGA" and five minor alleles consistent with at least three individuals





**Fig. 5** CE-STR profile of extract from cigarette amplified with  $AmpF/STR^{TM}$  Globalfiler® Kit [41]. The profile is classified as a DNA mixture solely due to locus vWA where allele 17 in stutter position of allele 18 displays a RFU value above the stutter threshold

Fig. 6 Forensic sample. Six representative STR (a) and MH (b) markers obtained from the amplification of ~1 ng of DNA obtained from a Qiagen Investigator manual extraction of a cigarette butt. Conventional CE analysis of Amp $F/STR^{TM}$ Globalfiler® Kit amplification with the same amount of template yielded a profile with a single detectable allele (in stutter position) of a possible minor contributor



27 loci. The RMP of the inferred minor contributor was calculated using data from 21 out of the 27 loci, those with one or two alleles unambiguously originating from the minor contributor. Results indicated that the probability of randomly selecting an individual with the same MH profile in the European American population is lower than one in over three quintillions. Such value is in line with a full STR profile obtained with conventional commercial kits and would allow for unambiguous source attribution.

When directly comparing CE of STR and MPS of MH markers, MPS analysis, both with STRs and MHs, yielded a greater number of alleles per locus suggesting an increased sensitivity for detecting minor contributors in mixtures. This increased sensitivity could allow interpretable profiles to be obtained from previously tested evidence samples that had yielded mixed profiles interpretable for the minor contributor.

Since a large number of loci can be co-amplified in a single MPS analysis, manual mixture interpretation becomes impractical. As a result, probabilistic genotyping is a valuable bio-statistical approach for minimizing the subjectivity in the interpretative analysis and will likely play a critical role for casework implementation of this new forensic DNA marker. Results of this study demonstrate that MHs have the potential to be an effective tool to enhance current DNA-based methods for human identification. As more studies provide even more documentation on MHs and their value in forensics, we expect that their implementation in casework will follow.

## **Future studies**

Within this study, we have demonstrated that simple inspection of MPS results of microhaplotypes can provide useful information on the different components of a DNA mixture. This supports empirically the original argument [19] that multiallelic microhaplotypes would be effective in mixture detection and deconvolution. What is now needed is a more rigorous statistical framework for evaluating the likelihood of the multi-locus genotypes disentangled from a mixture. The probabilistic genotyping methods being used for STRs [5] can be adapted for the analysis of microhaplotype data, and the absence of stutter peaks may enable straightforward calculations that yield a meaningfully better interpretation of mixtures. **Acknowledgments** The authors thank Dr. Moses S. Schanfield for providing the DNA samples, collected and extracted between 1993 and 2003, and used in this study.

**Funding information** This study was in part supported by National Institute of Justice through grants No. 2017-DN-BX-0164 awarded to Daniele Podini and No. 2015-DN-BX-K023 awarded to Kenneth K. Kidd, and by the Swiss National Science Foundation through grant No. 2017-P2LAP3\_174742 awarded to Fabio Oldoni.

## **Compliance with ethical standards**

**Conflict of interest** The authors declare no conflict of interest.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.

# References

- Butler JM, Buel E, Crivellente F, McCord BR (2004) Forensic DNA typing by capillary electrophoresis using the ABI Prism 310 and 3100 genetic analyzers for STR analysis. Electrophoresis 25: 1397–1412. https://doi.org/10.1002/elps.200305822
- Butler JM (2015) The future of forensic DNA analysis. Philos Trans R Soc Lond Ser B Biol Sci 370:577–579. https://doi.org/10.1098/ rstb.2014.0252
- Gill P, Haned H, Bleka O, Hansson O, Dørum G, Egeland T (2015) Genotyping and interpretation of STR-DNA: low-template, mixtures and database matches—twenty years of research and development. Forensic Sci Int Genet 18:100–117. https://doi.org/10. 1016/j.fsigen.2015.03.014
- Perlin MW, Belrose JL, Duceman BW (2013) New York State TrueAllele® casework validation study. J Forensic Sci 58:1458– 1466. https://doi.org/10.1111/1556-4029.12223
- Bright J-A, Taylor D, McGovern C, Cooper S, Russell L, Abarno D, Buckleton J (2016) Developmental validation of STRmix<sup>™</sup>, expert software for the interpretation of forensic DNA profiles. Forensic Sci Int Genet 23:226–239. https://doi.org/10.1016/j. fsigen.2016.05.007
- Gettings KB, Aponte RA, Kiesler KM, Vallone PM (2015) The next dimension in STR sequencing: polymorphisms in flanking regions and their allelic associations. Forensic Sci Int Genet Suppl Ser 5:e121–e123. https://doi.org/10.1016/j.fsigss.2015.09.049
- Gettings KB, Aponte RA, Vallone PM (2015) STR allele sequence variation: current knowledge and future issues. Forensic Sci Int Genet 18:118–130. https://doi.org/10.1016/j.fsigen.2015.06.005
- Jäger AC, Alvarez ML, Davis CP, Guzmán E, Han Y, Way L, Walichiewicz P, Silva D, Pham N, Caves G, Bruand J, Schlesinger F, Pond SJK, Varlaro J, Stephens KM, Holt CL (2017) Developmental validation of the MiSeq FGx forensic genomics system for targeted next generation sequencing in forensic DNA casework and database laboratories. Forensic Sci Int Genet 28:52–70. https://doi.org/10.1016/j.fsigen.2017.01.011
- van der Gaag KJ, de Leeuw RH, Hoogenboom J, Patel J, Storts DR, Laros JFJ, de Knijff P (2016) Massively parallel sequencing of short tandem repeats—population data and mixture analysis results for the PowerSeq<sup>™</sup> system. Forensic Sci Int Genet 24:86–96. https://doi.org/10.1016/j.fsigen.2016.05.016
- Aponte RA, Gettings KB, Dueller DL, Coble MD, Vallone PM (2015) Sequence-based analysis of stutter at STR loci: characterization and utility. Forensic Sci Int Genet Suppl Ser 5:e456–e458. https://doi.org/10.1016/J.FSIGSS.2015.09.181

- Børsting C, Fordyce SL, Olofsson J, Mogensen HS, Morling N (2014) Evaluation of the Ion Torrent<sup>TM</sup> HID SNP 169-plex: a SNP typing assay developed for human identification by second generation sequencing. Forensic Sci Int Genet 12:144–154. https://doi.org/10.1016/j.fsigen.2014.06.004
- Ambers AD, Churchill JD, King JL, Stoljarova M, Gill-King H, Assidi M, Abu-Elmagd M, Buhmeida A, Budowle B, Budowle B (2016) More comprehensive forensic genetic marker analyses for accurate human remains identification using massively parallel DNA sequencing. BMC Genomics 17:750. https://doi.org/10. 1186/s12864-016-3087-2
- Wendt FR, Warshauer DH, Zeng X, Churchill JD, Novroski NMM, Song B, King JL, LaRue BL, Budowle B (2016) Massively parallel sequencing of 68 insertion/deletion markers identifies novel microhaplotypes for utility in human identity testing. Forensic Sci Int Genet 25:198–209. https://doi.org/10.1016/j.fsigen.2016.09. 005
- Brown H, Thompson R, Murphy G, Peters D, La Rue B, King J, Montgomery AH, Carroll M, Baus J, Sinha S, Wendt FR, Song B, Chakraborty R, Budowle B, Sinha SK (2017) Development and validation of a novel multiplexed DNA analysis system, InnoTyper® 21. Forensic Sci Int Genet 29:80–99. https://doi.org/ 10.1016/j.fsigen.2017.03.017
- Westen AA, Matai AS, Laros JFJ, Meiland HC, Jasper M, de Leeuw WJF, de Knijff P, Sijen T (2009) Tri-allelic SNP markers enable analysis of mixed and degraded DNA samples. Forensic Sci Int Genet 3:233–241. https://doi.org/10.1016/j.fsigen.2009.02.003
- Phillips C, Amigo J, Carracedo Á, Lareu MV (2015) Tetra-allelic SNPs: informative forensic markers compiled from public wholegenome sequence data. Forensic Sci Int Genet 19:100–106. https:// doi.org/10.1016/j.fsigen.2015.06.011
- Oldoni F, Kidd KK, Podini D (2019) Microhaplotypes in forensic genetics. Forensic Sci Int Genet 38:54–69. https://doi.org/10.1016/ j.fsigen.2018.09.009
- Kidd KK, Pakstis AJ, Speed WC, Lagace R, Chang J, Wootton S, Ihuegbu N (2013) Microhaplotype loci are a powerful new type of forensic marker. Forensic Sci Int Genet Suppl Ser 4:e123–e124. https://doi.org/10.1016/J.FSIGSS.2013.10.063
- Kidd KK, Pakstis AJ, Speed WC, Lagacé R, Chang J, Wootton S, Haigh E, Kidd JR (2014) Current sequencing technology makes microhaplotypes a powerful new type of genetic marker for forensics. Forensic Sci Int Genet 12:215–224. https://doi.org/10.1016/j. fsigen.2014.06.014
- Kidd KK, Speed WC (2015) Criteria for selecting microhaplotypes: mixture detection and deconvolution. Investig Genet 6:1. https:// doi.org/10.1186/s13323-014-0018-3
- Kidd KK, Speed WC, Pakstis AJ, Podini DS, Lagacé R, Chang J, Wootton S, Haigh E, Soundararajan U (2017) Evaluating 130 microhaplotypes across a global set of 83 populations. Forensic Sci Int Genet 29:29–37. https://doi.org/10.1016/j.fsigen.2017.03. 014
- Børsting C, Morling N (2015) Next generation sequencing and its applications in forensic genetics. Forensic Sci Int Genet 18:78–89. https://doi.org/10.1016/j.fsigen.2015.02.002
- Bulbul O, Pakstis AJ, Soundararajan U, Gurkan C, Brissenden JE, Roscoe JM, Evsanaa B, Togtokh A, Paschou P, Grigorenko EL, Gurwitz D, Wootton S, Lagace R, Chang J, Speed WC, Kidd KK (2017) Ancestry inference of 96 population samples using microhaplotypes. Int J Legal Med 132:703–711. https://doi.org/ 10.1007/s00414-017-1748-6
- Kidd KK, Pakstis AJ, Speed WC, Lagace R, Wootton S, Chang J (2018) Selecting microhaplotypes optimized for different purposes. Electrophoresis 39:2815–2823. https://doi.org/10.1002/elps. 201800092
- 25. van der Gaag KJ, de Leeuw RH, Laros JFJ, den Dunnen JT, de Knijff P (2018) Short hypervariable microhaplotypes: a novel set

of very short high discriminating power loci without stutter artefacts. Forensic Sci Int Genet 35:169–175. https://doi.org/10.1016/j. fsigen.2018.05.008

- Chen P, Yin C, Li Z, Pu Y, Yu Y, Zhao P, Chen D, Liang W, Zhang L, Chen F (2018) Evaluation of the microhaplotypes panel for DNA mixture analyses. Forensic Sci Int Genet 35:149–155. https://doi. org/10.1016/j.fsigen.2018.05.003
- Oldoni F, Hart R, Long K, Maddela K, Cisana S, Schanfield M, Wootton S, Chang J, Lagace R, Hasegawa R, Kidd K, Podini D (2017) Microhaplotypes for ancestry prediction. Forensic Sci Int Genet Suppl Ser 6:e513–e515. https://doi.org/10.1016/j.fsigss. 2017.09.209
- Zhu J, Lv M, Zhou N, Chen D, Jiang Y, Wang L, He W, Peng D, Li Z, Qu S, Wang Y, Wang H, Luo H, An G, Liang W, Zhang L (2018) Genotyping polymorphic microhaplotype markers through the Illumina® MiSeq platform for forensics. Forensic Sci Int Genet 39:1–7. https://doi.org/10.1016/j.fsigen.2018.11.005
- Chen P, Zhu W, Tong F, Pu Y, Yu Y, Huang S, Li Z, Zhang L, Liang W, Chen F (2018) Identifying novel microhaplotypes for ancestry inference. Int J Legal Med. https://doi.org/10.1007/s00414-018-1881-x
- Butler JM, Decker AE, Kline MC, Reid TM, Vallone PM (2007) New autosomal and Y-chromosome STR loci: characterization and potential uses. In: 18th Int. Symp. Hum. Identification, Hollywood, CA, Promega
- Karafet TM, Mendez FL, Meilerman MB, Underhill PA, Zegura SL, Hammer MF (2008) New binary polymorphisms reshape and increase resolution of the human Y chromosomal haplogroup tree. Genome Res 18:830–838. https://doi.org/10.1101/gr.7172008
- Pakstis AJ, Speed WC, Kidd JR, Kidd KK (2007) Candidate SNPs for a universal individual identification panel. Hum Genet 121:305– 317. https://doi.org/10.1007/s00439-007-0342-2
- Phillips C, Fang R, Ballard D, Fondevila M, Harrison C, Hyland F, Musgrave-Brown E, Proff C, Ramos-Luis E, Sobrino B, Carracedo A, Furtado MR, Court DS, Schneider PM (2007) Evaluation of the Genplex SNP typing system and a 49plex forensic marker panel. Forensic Sci Int Genet 1:180–185. https://doi.org/10.1016/j.fsigen. 2007.02.007
- Holland MM, Fisher DL, Lee DA, Bryson CK, Weedn VW (1993) Short tandem repeat loci: application to forensic and human remains identification. In: Pena SDJ, Chakraborty R, Epplen JT, Jeffreys AJ (eds) DNA Fingerprinting State Sci. Birkhäuser Basel, Basel, pp 267–274. https://doi.org/10.1007/978-3-0348-8583-6\_24
- Cheung KH, Osier MV, Kidd JR, Pakstis AJ, Miller PL, Kidd KK (2000) ALFRED: an allele frequency database for diverse populations and DNA polymorphisms. Nucleic Acids Res 28:361–363 https://www.ncbi.nlm.nih.gov/pubmed/10592274
- Sherry ST, Ward MH, Kholodov M, Baker J, Phan L, Smigielski EM, Sirotkin K (2001) dbSNP: the NCBI database of genetic variation. Nucleic Acids Res 29:308–311 https://www.ncbi.nlm.nih. gov/pubmed/11125122

- Buchard A, Kampmann M-L, Poulsen L, Børsting C, Morling N (2016) ISO 17025 validation of a next-generation sequencing assay for relationship testing. Electrophoresis 37:2822–2831. https://doi. org/10.1002/elps.201600269
- Ng L-K, Ng A, Cholette F, Davis C (2007) Optimization of recovery of human DNA from envelope flaps using DNA IQ<sup>™</sup> system for STR genotyping. Forensic Sci Int Genet 1:283–286. https://doi. org/10.1016/j.fsigen.2007.05.004
- Phillips K, McCallum N, Welch L (2012) A comparison of methods for forensic DNA extraction: Chelex-100® and the QIAGEN DNA Investigator Kit (manual and automated). Forensic Sci Int Genet 6: 282–285. https://doi.org/10.1016/j.fsigen.2011.04.018
- Kidd KK (2016) Proposed nomenclature for microhaplotypes. Hum Genomics 10:16. https://doi.org/10.1186/s40246-016-0078-y
- Wang DY, Gopinath S, Lagacé RE, Norona W, Hennessy LK, Short ML, Mulero JJ (2015) Developmental validation of the GlobalFiler® express PCR amplification kit: a 6-dye multiplex assay for the direct amplification of reference samples. Forensic Sci Int Genet 19:148–155. https://doi.org/10.1016/j.fsigen.2015.07.013
- 42. Budowle B, Giusti AM, Waye JS, Baechtel FS, Fourney RM, Adams DE, Presley LA, Deadman HA, Monson KL (1991) Fixed-bin analysis for statistical evaluation of continuous distributions of allelic data from VNTR loci, for use in forensic comparisons. Am J Hum Genet 48:841–855 https://www.ncbi.nlm.nih.gov/ pubmed/1673286
- T.I.H. 3 Consortium, Altshuler DM, Gibbs RA, Peltonen L, 43. Altshuler DM, Gibbs RA, Peltonen L, Dermitzakis E, Schaffner SF, Yu F, Peltonen L, Dermitzakis E, Bonnen PE, Altshuler DM, Gibbs RA, de Bakker PIW (Co-leader), Deloukas P (Co-leader), Gabriel SB, Gwilliam R, Hunt S, Inouye M (Co-leader), Jia X, Palotie A, Parkin M (Co-leader), Whittaker P, Yu F (Leader), Chang K, Hawes A, Lewis LR, Ren Y, Wheeler D, Gibbs RA, Marie Muzny D, Barnes C, Darvishi K, Hurles M (Co-leader), Korn JM, Kristiansson K, Lee C, McCarroll SA (Co-leader), Nemesh J, Dermitzakis E, Keinan A (Leader), Montgomery SB, Pollack S, Price AL, Soranzo N, Bonnen PE, Gibbs RA, Gonzaga-Jauregui C, Keinan A, Price AL, Yu F (Leader), Anttila V, Brodeur W, Daly MJ, Leslie S, McVean G, Moutsianas L, Nguyen H, Schaffner SF (Leader), Zhang Q, Ghori MJR, McGinnis R (Coleader), McLaren W, Pollack S, Price AL (Co-leader), Schaffner SF (Co-leader), Takeuchi F, Grossman SR, Shlyakhter I, Hostetter EB, Sabeti PC (Leader), Adebamowo CA, MW Foster, Gordon DR, Licinio J, Manca MC, Marshall PA, Matsuda I, Ngare D, Wang VO, Reddy D, Rotimi CN, Royal CD, Sharp RR, Zeng C, Brooks LD, McEwen JE (2010) Integrating common and rare genetic variation in diverse human populations. Nature 467:52. https://doi.org/10.1038/nature09298
- 44. Jorde LB, Wooding SP (2004) Genetic variation, classification and "race". Nat Genet 36:S28–S33. https://doi.org/10.1038/ng1435