

The author(s) shown below used Federal funds provided by the U.S. Department of Justice and prepared the following final report:

Document Title: A SNP-Based Microarray Technology for Use in Forensic Applications” Final Technical Report

Author: Giulia C. Kennedy

Document No.: 223977

Date Received: September 2008

Award Number: 2005-DA-BX-K101

This report has not been published by the U.S. Department of Justice. To provide better customer service, NCJRS has made this Federally-funded grant final report available electronically in addition to traditional paper copies.

Opinions or points of view expressed are those of the author(s) and do not necessarily reflect the official position or policies of the U.S. Department of Justice.

NIJ Award No. 2005-DA-BX-K101 (PI: Giulia C. Kennedy)
“A SNP-Based Microarray Technology for Use in Forensic Applications”
Final Technical Report

Abstract

The purpose of this project is to enable the detection and analysis of forensically-relevant single-nucleotide polymorphisms (SNPs), which will enhance and complement the genomic information currently being used for forensic identification. The project is focused on developing an accurate, affordable, microarray-based forensic DNA analysis assay, capable of rapid, simultaneous, SNP genotyping for human identification testing for three specific forensic sample types: 1) low template single donor samples; 2) degraded samples and 3) mixtures containing two or more DNA sources. During the two year project period, we conducted assay development for forensic applications. Specifically, we have focused our efforts in 2 areas with the following results:

- 1) **Assay development:** Using Affymetrix 500K nuclear SNP genotyping arrays, we demonstrated that:
 - a. Pre-amplification using whole genome amplification (WGA) on low template samples yields successful genotyping for 80-90% of the markers on the 500K arrays. These results represent approximately 5-10% lower call rates than the non-pre-amplified standard assay.
 - b. Degraded samples prepared by DNase I digestion provide genotype data on many thousands of SNPs in the 500K assay
 - c. Using the standard assay, we detect mixed samples composed of two, three or four components. In mixtures with two components, we can detect as low as a 5% contribution from the minor component.
- 2) **Blind studies on forensic samples provided by Bode Technologies.** We ran the 500K SNP assay on all samples and analyzed the data using the standard algorithm. Our results show that:
 - a. All samples except for hair shaft (no root) yielded useable genotype data
 - b. Genders were correctly assigned in all cases
 - c. Mixtures were detected and the two components of the mixture were identified
 - d. Semen and buccal samples yielded high call rates comparable to control samples.
 - e. Vaginal samples had 5-10% lower call rates than semen and buccal
 - f. WGA samples yielded 5-10% lower call rates than unamplified genomic DNA
 - g. All ethnicities were correctly identified, including a sample with ancestry from a Taiwanese aborigine.

Table of Contents	Page
Abstract	1
Executive Summary	3
Main Body	
Introduction	4
Methods	7
Results	10
Conclusions	36
Implications for Policy and Practice	39
Implications for Future Research	39
References Cited	40
Dissemination of Research Findings	42

Executive Summary

A thirteen-locus STR panel is sufficient for human identification provided that the DNA sample is intact or minimally degraded and at sufficient concentration. As conditions become less optimal, STR markers fail to genotype, compromising the power of the marker set to identify or exclude individuals with sufficient match probabilities. We hypothesize that the functional power of a high-density SNP marker set (e.g. 500,000 SNPs) to identify individuals is retained even when large numbers of markers drop out due to degradation, low template amounts or other assay failures. We set out to explore the feasibility of using 500K SNPs to analyze three types of forensically relevant samples: low copy, degraded and mixtures of two or more contributors. First we tested the ability to amplify nanogram and picogram quantities of starting DNA using isothermal whole genome amplification (WGA) with phi29 polymerase. Second, we generated degraded samples by DNase I digestion and assessed the effects on genotyping quality. Finally, we constructed a series of DNA mixtures from two, three and four contributors and genotyped them with 500K SNP arrays. We determined the ethnicity of each contributor using STRUCTURE.

As the overall goal of this work is to use high-density SNP assays to extract accurate genotypic information from typical forensics samples, we began by laying a solid foundation of assay and algorithm development work which then allowed us to successfully test blinded forensics samples supplied by Bode Technologies. We concluded the following:

- The assay development successfully introduced a whole genome amplification (WGA) step to adapt smaller sample quantities to the standard 500K protocol. While call rates are generally lower with pre-amplified samples, nonetheless sufficient data are collected on hundreds of thousands of SNPs to provide critical genotype information such as relatedness and ancestry.
- We generated DNase I-degraded samples and obtained useful genotype information on a subset of SNPs, using the standard algorithm.
- We generated panels of defined mixtures of two, three and four-component DNAs at varying proportions. We processed them on the standard 500K assay and determined that a series of simple metrics are statistically sufficient to detect mixtures down to a 5% contribution of the minor component. Ethnicity of contributors can also be determined at a mixture level down to 25:75.

The implications of this work are three-fold. As low template samples often pose challenges in casework, the ability to extract genotype information from picogram quantities (~20 genome equivalents) would significantly improve the ability to solve cases with minimal evidence. Degraded samples also pose a difficult challenge for criminalists; as the number of additional STR loci that can be developed to work with current panels is finite, methods that seek out SNP information in these samples would be highly useful. Finally, methods that can be used not only to identify mixtures, but to estimate the number of contributors and also to identify the contributors from known reference samples would be very powerful. The results described in this report indicate that high-density SNP panels can be used with success on low template, degraded and mixed samples. The next step will be for a forensics lab to test actual samples and enhance the information provided by STRs to help exclude or include individuals in the identification process.

I. Introduction

Statement of the problem The demand for forensic DNA analysis of evidentiary samples far exceeds current capabilities. In addition to a staggering backlog of samples in the criminal justice system yet to be analyzed, there are additional limitations in analyzing DNA from highly degraded samples. While mtDNA analysis is useful for samples containing compromised biological evidence, only data from the maternal line is generated, limiting its statistical value in human identification. Furthermore, current forensic genotyping technologies use simple tandem repeat (STR) markers, and while they have been successful in identity matching, they provide limited additional information about DNA contributors, e.g. geographic ancestry or kinship. Thus, there is a need to identify novel genetic variants to provide this and other information. A rapid expansion of genomic sequence information has led to the identification of >5 million single-nucleotide polymorphisms (SNPs) in the human genome. SNPs are effective markers for human genetic variation and can be measured using a variety of methods. Despite this treasure of genetic markers, the ability to genotype them in large numbers in a cost-effective manner has been challenging. Recognizing the potential for synthetic microarrays (i.e. “DNA chips”) to overcome significant bottlenecks in DNA analysis, we have developed methods for generating large amounts of genetic information from a single array as small as 5 x 5 mm. This array contains enough DNA probes to genotype many thousands of SNPs, all at the same time. Because of the vast “overkill” of content available on the microarrays, we may be able to extract partial, but significant, information from a degraded, mixed or compromised sample where previous methods have failed. In addition to providing large amounts of genetic information, these tiny chips can be processed in a variety of configurations depending on need: by large forensics laboratories with an automated walk-away system, or by an individual scientist in a small local crime laboratory.

Literature citations and review

Current methods used in forensic DNA analysis

Before describing Affymetrix genotyping technologies and how they may impact forensic science, we summarize currently used technologies in the forensics community.

STRs Simple tandem repeats, or STRs, are useful markers for scoring human genetic variation and are the mainstay in forensic identity testing (Gill, 2002). STR analysis requires PCR amplification by sequence-specific primers followed by size discrimination on a gel-based platform (<http://www.appliedbiosystems.com>). A panel composed of 13 or 16 STRs is used extensively in forensic science for identity matching between test and reference sample DNAs (<http://www.promega.com/applications/hmnid/>). The amplification of all loci is performed in one reaction as a multiplex PCR. One limitation of such a system is that the addition of more markers (for example to add ancestry-informative markers, to include mitochondrial or plant DNA sequences) would require re-optimization of the multiplex PCR reaction. While high levels of STR multiplexing by PCR (100-1000 fold) are achievable with extensive optimization, this approach is difficult to scale to large numbers of loci, due to limited space on the gel for resolving additional fragments. A further bottleneck is that individual profiles must be examined and checked by highly-trained personnel and often reviewed by a second individual because of stutter peaks and sizing reproducibility that confound interpretation of results.

Mitochondrial DNA sequencing As there are hundreds to thousands of copies of the mitochondrial genome in a cell, compared to two copies of the nuclear genome, mitochondrial DNA (mtDNA) analysis can be used on samples with little or no intact nuclear DNA, such as teeth, hair, skeletal remains, etc. (Budowle et al., 2003; Ginther et al., 1992; Higuchi et al., 1988; Holland et al., 1993). Currently, mtDNA analysis begins with locus-specific amplification of two hypervariable regions from sample DNA followed by standard dideoxysequencing sequencing (Holland and Parsons, 1999; Stoneking et al., 1991). Because mtDNA is maternally inherited and

is not subject to recombination, it is particularly useful in tracing ancestry of samples; however its high rate of mutation and mode of inheritance precludes its use as a unique identifier because all individuals within a maternal lineage will share the same mtDNA sequence.

SNPs Single-nucleotide polymorphisms, or SNPs, have become desirable as genetic markers because of their abundance and genetic stability (Xiong and Jin, 1999). Current SNP genotyping technologies use locus-specific PCR to amplify specific SNPs, followed by a variety of allele discrimination methodologies (Kwok, 2001; Syvanen, 2001). Several SNP-based tests have been commercialized for eye color (Frudakis et al., 2003), paternity (<http://www.dnprint.com>) and ancestry (Frudakis et al., 2003). These latter SNPs quantify mixtures of geographic ancestry from four groups, African, Indo-European, Native American and East Asian. As it may be desirable to specify geographical ancestry at higher resolution, i.e. determination of sub-population group, etc., more SNPs will be needed. The technology we developed is a high-throughput cost-effective method for genotyping large numbers of SNPs in large numbers of forensics samples.

Recent advances in Affymetrix DNA analysis

WGSA Recognizing the potential for synthetic microarrays to overcome significant bottlenecks in genotyping technology, we developed a simple but powerful approach, termed whole genome sampling analysis (WGSA), to genotype simultaneously thousands of SNPs in complex DNA without locus-specific primers or the need for automation (Kennedy et al., 2003). We devised a generic sample preparation method that uses a single oligonucleotide primer for amplification, coupled to allele discrimination on synthetic DNA microarrays (Figure 2). Our method amplifies highly reproducible fractions of the genome across multiple DNA samples and calls genotypes at >99.8% accuracy. Each interrogated SNP is comprised of 40 features, 20 on each strand. Both the A and B alleles are represented in 25mer oligonucleotides staggered at various positions relative to the SNP base, along with a single basepair mismatch sequence to assess specificity. This technology led to the commercialization of an Affymetrix 10,000 SNP (10K) chip and because of its inherent potential for scalability, has resulted in commercialization of 100,000 SNP arrays followed by 500,000 SNP arrays and finally, nearly 1 million:

<http://www.affymetrix.com/products/arrays/>

In this report, we use the two-array 500K SNP set and the standard genotyping software GTYPE which implements the Bayesian Robust Linear Model with Mahalanobis distance (BRLMM) genotyping algorithm [Affymetrix (2006). Product Update: BRLMM Analysis Tool (Affymetrix, Inc.).

URL: http://www.affymetrix.com/support/technical/product_updates/brlmm_algorithm.affx]

Detecting DNA Mixtures STRs have been used extensively to identify DNA mixtures. Evett *et al.* (1991) first investigated DNA mixtures at the single locus level. Weir *et al.* (1997) proposed formulae for likelihood calculations for mixed samples. Curran *et al.* (1998) and Fung & Hu (2000, 2002a, 2002b) incorporated population substructure into DNA mixture interpretation. Hu & Fung (2003) further considered the case when contributors in the mixture are relatives. In all these studies, the number of contributors in the DNA mixture is assumed to be known, i.e., set by prior information (Weir 1997) or bounded by some threshold (Lauritzen and Mortera 2002). Stockmarr (2000) first estimated the number of contributors by maximizing the likelihood for a single STR locus. However, detection of DNA mixtures and estimation of number of contributors with SNP markers has not yet been thoroughly investigated. Chakraborty *et al.* (2008) have used simulated data and two maximum-likelihood methods to detect mixtures, and also to estimate the exact number of contributors. We include some of his unpublished data simulations on 2 and 3 component mixtures, computed for a range of inbreeding coefficient values (Tables 4-7 below). In this report, we also provide preliminary 500K SNP data studying empirically derived mixtures

at different ratios of contributors markers using simple genotype concordance and heterozygosity measures (Figures 7-12 and Table 3).

Statement of hypothesis or rationale for the research

The purpose of this project is to enable the detection and analysis of forensically-relevant single-nucleotide polymorphisms (SNPs), which will enhance and complement the genomic information currently being used for forensic identification. The project is focused on developing an accurate, affordable, microarray-based forensic DNA analysis assay, capable of rapid, simultaneous, SNP genotyping for human identification testing for three specific forensic sample types: 1) low template single donor samples; 2) degraded samples and 3) mixtures containing two or more DNA sources.

II. Methods

Sample Processing

Bode Technologies provided blinded forensic samples to Affymetrix (Table 8). The sample manifest provided by Bode included: 1) DNA concentration (ng/ul), 2) Source (blood, hair, semen, buccal, vaginal), 3) Gender. All samples were amplified by Whole Genome Amplification using commercially available methods (Qiagen). The 2 array set (Human Mapping 250K Nsp and Human Mapping 250K Sty array) was used for this preliminary work. For simplicity, only data from the Sty chip is presented here, as it was highly similar to the Nsp data. Standard mapping protocols were followed, as discussed below.

Nuclear Genome Information from Mapping Experiments

Whole Genome Amplification: The standard mapping assay requires 500ng of starting material (250 ng per array). Because this amount of starting material is not forensically useful, we utilized a preamplification step prior to array target preparation. Up to 10 ng DNA was utilized for each WGA reaction, according to the recommended Qiagen Repli-g Midi Kit standard protocol (Qiagen, 2007). For several of the Bode samples, the amount of genomic DNA provided was less than 10 ng. We were still able to derive partial useful data on these samples despite the lower template amount.

Whole Genome Sampling Assay (WGSA): Following the WGA reaction, DNA target was prepared using the Whole Genome Sampling Assay (WGSA). Affymetrix investigators have demonstrated the utility of high-density oligonucleotide microarrays for simultaneous genotyping of thousands of SNPs in complex DNA without locus-specific primers or the need for automation (Kennedy et al., 2003). Whole genome sampling analysis (WGSA) is a generic sample preparation method that uses a single oligonucleotide primer for amplification, coupled to allele discrimination on high-density oligonucleotide microarrays. This method, illustrated in Figure 2, amplifies highly reproducible fractions of the genome across multiple DNA samples and calls genotypes at >99.5% accuracy.

The GeneChip® Human Mapping 500K product has undergone extensive product validation. As discussed above, the commercial 500K assay has been designed for 250ng single-source DNA input per array. Analysis software and product validation were therefore optimized for these conditions. A large focus of our NIH-funded work has been to develop the assay and software for specific use by the forensics community. Therefore parameters and data analysis performed by the standard commercial GTYPE algorithms, for example, would need to be modified in the future to analyze low template and mixed samples.

Data Analysis Workflow and Software Tools

Affymetrix GeneChip® Genotyping Analysis Software (GTYPE) was used to analyze the data. GTYPE is commercially available software that is currently used to analyze data from all Affymetrix Mapping arrays. GTYPE uses the Dynamic Model (DM) (Di et al., 2005) and Bayesian Robust Linear with Mahalanobis (BRLMM) algorithm (Affymetrix, 2006; Rabbee and Speed, 2005; Rabbee and Speed, 2006) to make genotype calls for the 500K Array data. GTYPE and the standard genotyping algorithms were used to call genotypes. With this information, summary statistics, ethnicity, sex and relatedness were determined.

Summary Statistics Using Standard Analysis

Failed samples: Using standard assay conditions for 500K, some of the samples were designated as “failed” due to the non-standard appearance of PCR products on a gel, which did not pass standard metrics for size distribution. Despite this failure at the PCR gel visualization stage, the

samples still generated reasonable call rates, which still provide genotypes on hundreds of thousands of SNPs.

Call rate: The call rate represents the number of called genotypes divided by the number of possible genotypes multiplied by 100. Put another way, of the possible SNPs, the call rate reflects how many of the possible SNPs the algorithm actually made a decision on. Standard Affymetrix products have specifications associated with call rate. For example, using the Dynamic Model algorithm for 500K arrays, the specification is 93% call rate. The call rates we observed on the Bode samples are above the 93% specifications for most of the samples.

MDR-MCR: A high MDR-MCR metric is suggestive of cross-sample contamination or mixtures.

Percent Heterozygote (%AB): The AB metric reflects the number of SNPs called as heterozygotes over the total number of SNPs on the array. If this heterozygosity metric is significantly lower than the average, this could indicate possible allele dropout. Average heterozygosity for SNPs on the 500K Mapping arrays is 30%, as measured on a panel of ethnically diverse individuals. Heterozygosity for a given ethnic group may differ from this value. In this study, most of the Bode samples show heterozygosity between 22-26%.

Sex: Sex is determined by %AB call rate on the X chromosome. If % AB is less than 7%, the GTYPE software calls the sex as a male.

Ethnicity Determination

Nuclear SNPs have been used previously to classify individuals according to geographic ethnicity. Several publicly-available computer programs (e.g. STRUCTURE (Chicago, 2007; Falush et al., 2003; Pritchard et al., 2000)) are available to assign ethnicity, and while it is appreciated that larger the numbers of SNPs afford greater resolution in distinguishing closely related ethnic groups, the software has not kept pace with advances in SNP genotyping technology. At present, it is extremely computationally intensive to use 500,000 SNPs to classify individuals according to geographic ethnicity using publicly-available programs. Keeping the computational constraints in mind, we identified a set of 3700 SNPs to conduct ethnicity determinations.

An overview of the experimental workflow is shown in Figure 1. We used up to 10 ng of genomic DNA as input into whole genome amplification (WGA) reactions (Qiagen). The resulting amplified product (250 ng) was then used as input into the 500K Mapping assay (Figure 2). In sensitivity titrations, lower amounts were used in the WGA step (see Figure 3).

The 500K mapping assay uses whole-genome sampling analysis (Kennedy et al., 2003) which simultaneously genotypes 500,000 human SNPs using a single generic primer. Validation studies have shown the 500K SNP genotyping assay to be >99.5% accurate when used in standard assay and software mode.

Standard summary statistics are reported for each chip as follows:

- **% Call rate**: indicates the number of called genotypes divided by the number of possible genotypes x 100. For the standard assay, a call rate of >93% meets QC specifications
- **% AB**: Indicates the level of heterozygosity computed across all called SNPs on the array. The average heterozygosity for the 500K SNPs is 30% and was computed on a large multiethnic panel of DNA samples. A lower than average %AB indicates possible allele drop-out. A higher than average %AB indicates a possible mixture.

- **% MDR-MCR:** this value is the difference between the detection rate and the call rate on a small set of QC SNPs present on the array. A high value (>15%) is suggestive of sample mixtures.
- **Sex:** The %AB is computed for SNPs on the X chromosome. A value of <7% AB is called a male.

III. Results

Statement of Results

Pre-amplification using Whole Genome Amplification (WGA): The standard mapping assay requires 500ng of starting material (250 ng per array). Because this amount of starting material is not forensically useful, we utilized a preamplification step prior to array target preparation. We compared the results of WGA vs non-WGA templates and consistently found about 5-10% lower call rates when pre-amplifying DNA prior to the 500K assay. Table 1 shows a direct comparison of WGA vs non-WGA for a set of samples. The results shown indicate that 10 ng of starting genomic DNA followed by WGA provides sufficient quality and number of genotypes to support most forensics applications when coupled with the 500K assay. We also tested lower amounts of starting material for the preamplification step to determine whether ultra-low copy templates could be detected on the arrays. Figure 3 shows the raw intensity signal for a titration series of the standard amount of template (250 ng or 250,000 pg) titrated down to 1 pg. A negative control (0 ug template) shows the signal obtained; as expected it is extremely low ie concentrated in the lower left quadrant. Even at 1 pg of starting material, the arrays begin to detect signal, primarily in the AA and BB homozygote clusters. As the template is increased, AB calls increase and the middle cluster becomes more populated, until it reaches the expected pattern of the standard input DNA (Figure 3, lower right). By removing the signal in the lower left quadrant (negative control) for all low template samples, we remove some of the genotyping errors associated with the BRLMM algorithm, which was not designed to call genotypes in low copy samples. Table 2 shows the call rate and concordancy (surrogate for accuracy in this case) for a titration series. Clearly the genotype accuracy decreases as template is decreased; a future challenge will be to determine which of the thousands of genotypes remaining are accurate.

Degraded DNA: Although the 500K assay uses defined restriction endonucleases for digesting the DNA, we hypothesize that even some proportion of degraded DNA will be ligated and amplified with the assay. We therefore sought to determine what call rates we could expect on 500K assays performed on degraded DNA templates. We generated degraded DNA templates by digestion with DNase I at varying concentrations and times (Figure 4) and then processed the templates further with the WGS assay (Figure 5). The degraded DNA templates were run on the 500K arrays. Call rates were determined using the DM algorithm (Figure 6). The call rates dropped significantly with an increase in DNase I concentration, but leveled off near 50% for the Coriell sample. These data indicate that even highly degraded DNA can yield genotype data on hundreds of thousands of SNPs using this assay. Note that the low intensity filter described in Table 2 was not used in this or subsequent analyses.

Detecting Mixtures: We created quantitative mixtures of a titration of two pure DNAs mixed together in varying ratios: 0:100, 1:99, 5:95, 10:90, 25:75, 50:50, 75:25, 90:10 and 100:0. We also created quantitative mixtures consisting of 3 or 4 components in 1:1:1 and 1:1:1:1 ratios. We then performed the standard 500K assay and genotype analysis. Using % heterozygosity (%AB calls) from the standard genotyping software, we see the expected increase as the ratio of the two samples increases. This simple measurement allows unambiguous detection of mixtures at a ratio of approximately 25:75 (Figure 7-9). With modifications in the analysis algorithm we would likely be able to detect much more dilute mixtures. We computed the accuracy of the genotype calls in each mixture by comparing the resulting genotypes to either of the two components of the mixture alone (Figure 10); as expected, the BRLMM algorithm makes calls which are less accurate as the sample becomes more mixed. Despite the lower accuracy in mixed samples, we tested how downstream determinations, such as ancestry, would be affected. We used the program STRUCTURE to determine ancestry on the mixtures. We found that ancestry could be assigned for each of the two components at a mixture ratio as low as 25:75 (Figure 11). We also

tested pairwise genotype concordance on mixture titrations from 2, 3 and 4 component mixtures and generated colorized heatmaps (Figure 12). The red color indicates the highest genotype concordance (self) and blue colors indicate the lowest genotype concordance (unrelated individuals). Each titration point was run in triplicate to provide data on reproducibility (designated A, B and C). As controls for relatedness, we included the mother and father for each of the components in the mixture. The relationships between the samples in the mixture studies are indicated in the pedigree (Figure 8). Three general trends are visible in this large dataset: 1) the reproducibility of the triplicates is high; 2) each of the two components in the mixtures has high genotype concordance to itself (red, red-orange, yellow), which decreases down to yellow-green, blue-green as the minor component is decreased down to 10-15%; 3) familial relationships are preserved as low as 25% minor component; for example, each of the two children in the mixture show high relatedness (yellow) to each of their parents, but not to the other set of parents (blue). As the minor component is decreased, the color changes from yellow to yellow-green to blue-green, however it is still distinguishable from the dark blue color representative of unrelated individuals.

Theoretical Analysis on simulated mixtures: Ranajit Chakraborty and colleagues are in the process of developing maximum-likelihood-based methods for both detecting mixtures and estimating the number of contributors (Chakraborty et al. 2008). Preliminary simulations are shown in Tables 4-7. Two methods have been developed: the MGP or Multi-locus genotyping profile method; and NHL, or Number of Heterozygous loci method. Each method is used to compute the accuracy of either detection of 1, 2 and 3 component mixtures (Tables 5-6) or estimation of the number of components (Table 4). The accuracy in estimating the number of contributors when the samples are either full-siblings or parent-offspring are also computed (Table 7). As these methods have been developed only for unlinked SNP loci, results from analysis of 500K data are not yet available.

Blinded samples from Bode technologies: We received a total of 43 blinded forensics samples from Bode Technologies. Information on the samples included DNA concentration, source and gender (Table 8). Following the blind experiment, the samples were unblinded and ethnicities were revealed by Bode. We generated data on the samples using the 500K assay. For conciseness, only data for the Sty chip is reported here as the Nsp chip data are highly similar. Standard GTYPE statistics were computed such as %call rate, %AB, %MDR-MCR and Gender (Table 9). Most of the samples, with the exception of hair shafts (no roots), passed the standard 500K assay with high call rates.

Effect of DNA Source on Call Rates: We noted previously that genotype data from blood or from cultured cell lines performs very well with the 500K assay (see previous progress reports). In this project, we analyzed DNA samples from blood stains, hair, semen, buccal and vaginal swabs. As expected, hairs without roots had very low call rates, consistent with the lack of nuclear DNA. While blood, semen and buccal samples show data similar to the controls, the vaginal swabs had call rates that were systematically decreased about 5-10% from the controls (Table 5). This small decrease had no effect on the ethnicity determination by STRUCTURE (see Figure 13 and discussion below) or on correctly assigning gender. Further analysis will be necessary to determine whether possibly other forensically relevant information is compromised in DNA samples taken from vaginal swabs.

Ethnicity Determination: The samples from Bode include Europeans, Africans, Asians and mixed ethnicities. We were blinded to the ethnicities of the samples until after our 500K results were shared with Bode Technologies. We used STRUCTURE to assign ethnicities from the nuclear data (see methods and legend to Figure 11). In all cases, we correctly identified the

ethnicity of the Bode samples, including one sample with a unique ancestry consistent with a Taiwanese aborigine group (Figure 13). We also easily detected several cases of genetic admixture (African-Caucasian and Caucasian-Asian). All ethnicities were confirmed by Bode. We also generated mitochondrial haplogroups on the Bode samples using the Affymetrix mitochondrial v2.0 resequencing chip (Figure 13). The mito haplogroup assignments were 100% consistent with ethnic description provided by Bode.

Allele-sharing determinations: We set out to determine whether there were any duplicates or mixtures in the Bode sample set. We computed pairwise allele-sharing amongst the samples. Duplicates would be expected to have the highest allele-sharing (ie self, indicated by yellow color) and unrelated individuals would expect to randomly share alleles (blue color). As seen in Figure 14, there were many sets of duplicates. Laboratory mixtures, or casework samples coming from more than one contributor, would be expected to also share alleles. Therefore some of the samples showing high levels of allele-sharing (shades of yellow and yellow-green) could indicate mechanical, rather than genetic, mixtures. As predicted, in these cases the allele-sharing value is lower and the yellow color is less intense. Due to the high heterozygosity observed in Bode sample 19, we suspected a mixture. By computing the allele-sharing value for all the Bode mixtures, we concluded that sample 19 was composed of a mixture of samples 07/08 and samples 13/14, each of which were themselves duplicates. Bode later confirmed that Sample 19 is a 2:1 mixture of 13/14 and 07/08. Note that duplicate samples with lower call rates tend to have less allele-sharing due to allele-drop out (e.g. less intense yellow) but are still clearly identified (e.g. Bode samples 03, 04, 05, 06).

FIGURES AND TABLES

Figure 1. Diagram of sample processing workflow. A whole-genome amplification (WGA) reaction is set up with up to 10 ng of input genomic DNA (Qiagen). The resulting product is purified and quantitated according to manufacturer’s instructions. An aliquot of 250 ng of WGA-amplified product is used as input into the 500K Mapping assay, according to standard protocols for Nsp and Sty target preparation (see Figure 2). The targets are hybridized to their respective chips, washed, stained and scanned, and the intensity data interpreted by the Affymetrix software to provide genotype calls for SNP arrays (or base calls for resequencing arrays). Information on each SNP is provided on the NetAffx website.

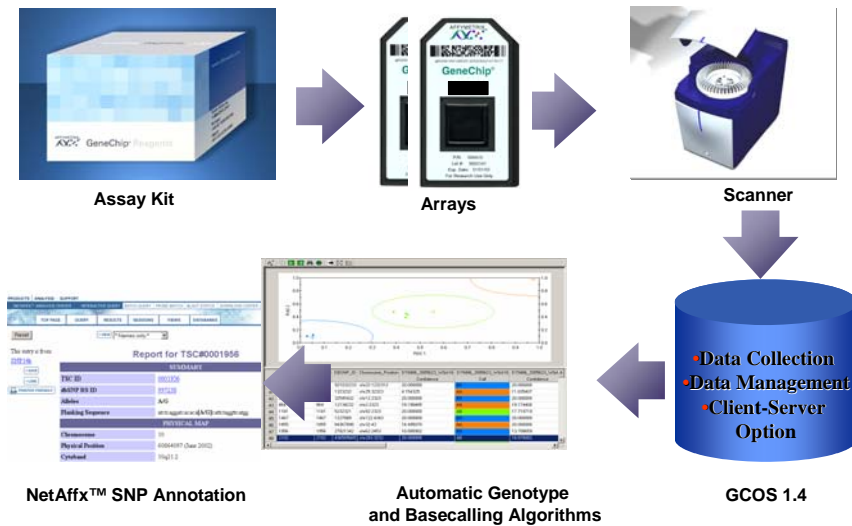


Figure 2. Diagram showing the principle behind the 500K Mapping assay. Genomic DNA is digested with one of two restriction endonucleases (Nsp I or Sty 1) and ligated to the appropriate adaptors. A generic primer is used to amplify fragments within the 200-1100 bp size range, resulting in a genomic fraction of about 500 million basepairs (500Mb) or about 16% of the human genome. The amplification of these fractions is highly reproducible. The fraction is fragmented and labeled and hybridized overnight to the respective array. The array is then washed and scanned using Affymetrix, Inc. GeneChip fluidics stations and scanners, respectively (see Figure 1).

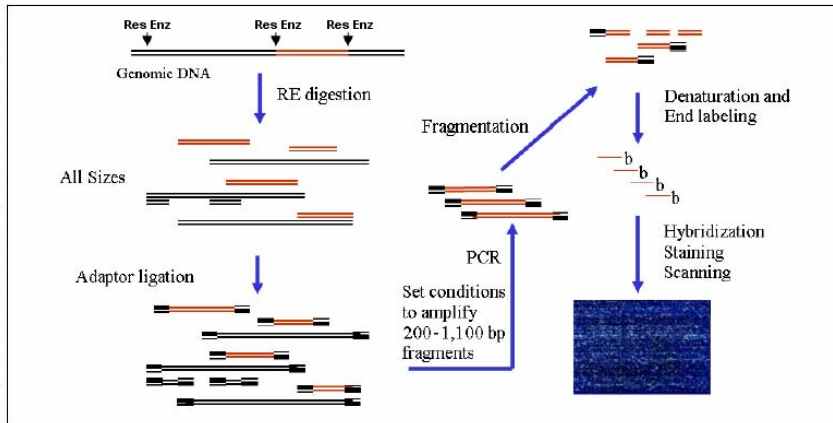


Table 1. Comparison of WGA vs non-WGA pre-amplification on 500K arrays

DNA templates were either pre-amplified using WGA prior to the 500K assay, or directly used for the standard 500K WGS assay (non-WGA). Standard BRLMM genotyping algorithm was used to analyze the data. The WGA call rates range from 80-96%, while the non-WGA call rates are all >98%.

sample	gender	WGA				Non-WGA			
		brimm_call_rate	AB_percent	AA_percent	BB_percent	brimm_call_rate	AB_percent	AA_percent	BB_percent
B020	female	84.06	23.21	32.22	28.64	99.35	26.94	36.88	35.53
B021	male	95.66	29.9	33.82	31.94	98.73	28.85	35.34	34.53
B022	female	88.03	23.31	33.88	30.84	99.69	29.38	35.94	34.37
B023	female	95.74	27.68	35.19	32.87	99.73	27.58	36.8	35.34
B024	male	97.4	27.09	36.37	33.94	99.61	26.61	37.37	35.63
B025	female	96.61	28.06	35.31	33.24	99.41	28.37	36.21	34.84
B026	female	94.47	25.21	35.85	33.41	99.63	24.8	38.11	36.73
B027	female	93.09	27.94	34.18	30.97	99.6	28.46	36.25	34.89
B028	female	94.15	28.79	33.57	31.8	99.75	30.49	35.07	34.19
B029	male	96.17	27.28	35.96	32.93	99.66	26.61	37.37	35.68
B030	female	96.21	27.76	35.35	33.1	99.84	29.4	35.97	34.47
B031	female	94.08	25.8	35.22	33.07	99.71	24.76	38.14	36.81
B032	female	96.5	27.73	35.31	33.46	99.78	27.55	36.84	35.38
B033	male	96.16	26.81	35.85	33.49	99.53	26.51	37.3	35.72
B034	female	96.45	29.87	34	32.59	99.78	30.49	35.08	34.2
B035	female	92.8	26.97	34.4	31.44	99.79	26.96	37.19	35.63
B036	female	96.39	26.76	36.08	33.54	99.68	26.96	37.13	35.6
B037	male	96.72	27.4	35.89	33.43	99.63	26.49	37.36	35.79
B038	female	80.82	24.22	30.01	26.58	98.85	26.97	36.62	35.25
B039	female	93.72	27.34	34.78	31.6	99.81	26.98	37.21	35.63
B040	male	96.77	27.26	35.89	33.62	99.36	26.64	37.11	35.61
B041	male	96.44	28.94	34.41	33.09	99.65	28.69	35.86	35.1
B042	female	80.94	26.18	29.3	25.46	98.59	27.04	36.48	35.07
B043	female	94.25	25.62	35.73	32.9	98.38	24.94	37.43	36.01

Figure 3. Intensity analysis of Ref 103 sensitivity titration. Zero, 1, 10, 100, 1000, 10,000 and 250,000 pg of DNA was amplified by WGA and used in the 500K mapping assay. For simplicity, only 6 of the titration points are shown here. The intensity data is plotted for A-allele probes (x-axis) and B-allele probes (y-axis). The non-template (i.e. negative, water) control is shown in blue (non-WGA) and red (WGA). Each successive titration point is indicated in a separate color, with the negative controls indicated by blue and red in each figure for comparison. The standard assay (250,000 pg) on the lower right clearly shows separation of the AA, AB and BB genotypes in the sample.

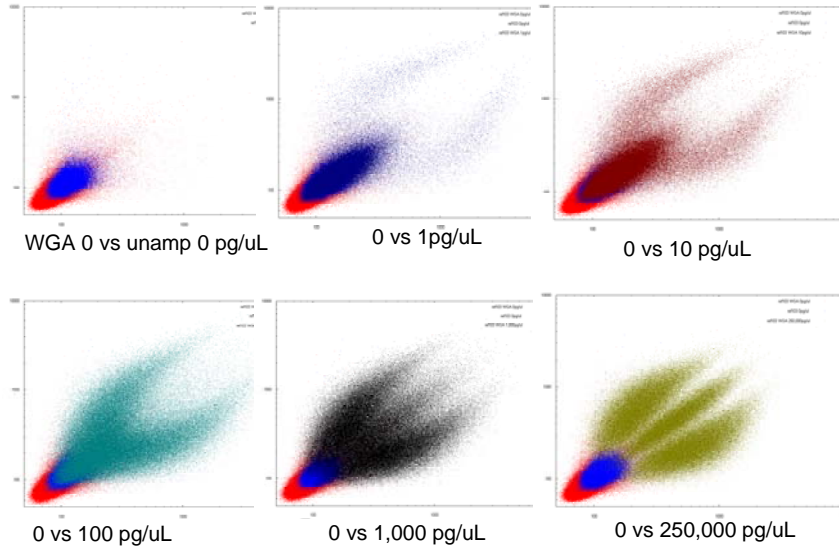
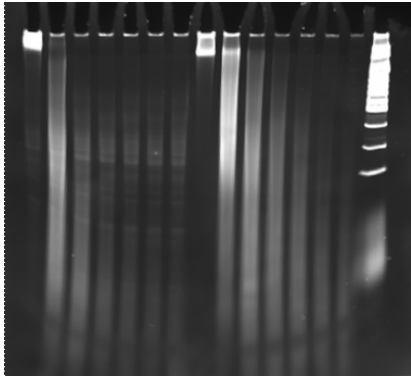


Table 2. Summary statistics on titration data using a low-intensity SNP filter. SNPs with light-unit values of <200 were removed from the analysis. This cut-off was determined from a distribution of intensity values across several experiments (data not shown). Call rate and concordance was then computed for the remaining SNPs. Only ~2000 SNPs remained in the negative control sample following implementation of the SNP filter, resulting in a large decrease in % call rate from 49% to 0.06%. Even at template concentrations of 25 pg, >200,000 SNPs were retained in the analysis, and of these, >80,000 SNPs had 100% accurate genotypes. At 250 pg of template (<100 genome copies) more than 184,000 SNPs were genotyped at 100% accuracy. Thus even low template forensics samples have the capability of generating enough SNP information to aid in multiple levels of identification.

Input DNA amount (pg)	SNPs retained (avg %)	Concordant SNPs (avg % of total comparisons)	Call rate (avg %)	Number of Replicates
0	2651 (1.1%)	525 (35.9%)	0.06	3
2.5	199,658 (83.7%)	28,717 (29.8%)	41.8	3
25	205,138 (86%)	81,488 (64.3%)	55.0	3
250	237,541 (99.7%)	184,715 (94.2%)	85.1	3
2,500	238,218 (100%)	224,359 (99.7%)	97.6	3
25,000	238,258 (100%)	226,213 (99.8%)	98.3	3
250,000	238,169 (99.9%)	230,409 (100 %)	100.0	3

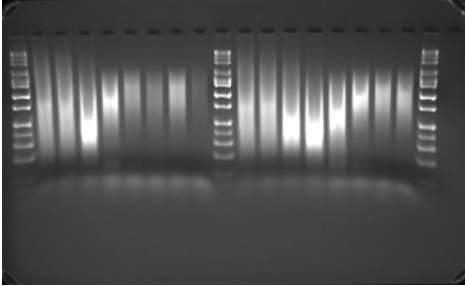
Figure 4. Generation of varying levels of degraded DNA by DNAase I titration. Two DNAs were used in this study. The first, Coriell sample NA10861, is DNA from a single donor. The other, G1471 is a pool of DNA from several sources, which we used in these experiments to reduce reagent costs. The gel below shows the increased fragmentation as the DNAase concentration was increased.



<u>Lane</u>	<u>Source</u>	<u>U DNase/Reaction</u>
1	10861	0.0
2	10861	0.5
3	10861	1.0
4	10861	1.5
5	10861	2.0
6	10861	2.5
7	10861	3.0
8	G1471	0.0
9	G1471	0.5
10	G1471	1.0
11	G1471	1.5
12	G1471	2.0
13	G1471	2.5
14	G1471	3.0
15	HiLow Marker	

Figure 5. WGSa PCR gel results of degraded DNA.

The resulting DNAs shown in Figure 10 then served as templates for the WGSa reaction (see Figure 2). The products were run on an agarose gel, as shown below. As expected, we observed non-standard behavior of the PCR products as the DNA template was increasingly degraded. Each of these WGSa products were further processed and run on 500K arrays (see results in Figure 6)



<u>Lane</u>	<u>Source</u>	<u>U DNase/ Reaction</u>
1	HiLow Marker	
2	10861	0.0
2	10861	0.5
3	10861	1.0
4	10861	1.5
5	10861	2.0
6	10861	2.5
7	10861	3.0
8	G1471	0.0
9	G1471	0.5
10	HiLow Marker	
11	G1471	0.0
12	G1471	0.0
13	G1471	0.5
14	G1471	1.0
15	G1471	1.5
16	G1471	2.0
17	G1471	2.5
18	G1471	3.0
19	HiLow Marker	

Figure 6. Call Rates on degraded DNA samples

The degraded DNA templates as shown in Figures 4 and 5 were further processed and run on the 500K arrays. Call rates were determined using the DM algorithm and plotted below. The call rates dropped significantly with an increase in DNAase concentration, but leveled off near 50% for the Coriell sample. These data indicate that even highly degraded DNA can yield genotype data on hundreds of thousands of SNPs using this assay. The low-intensity filter described in Table 2 was not implemented in this analysis.

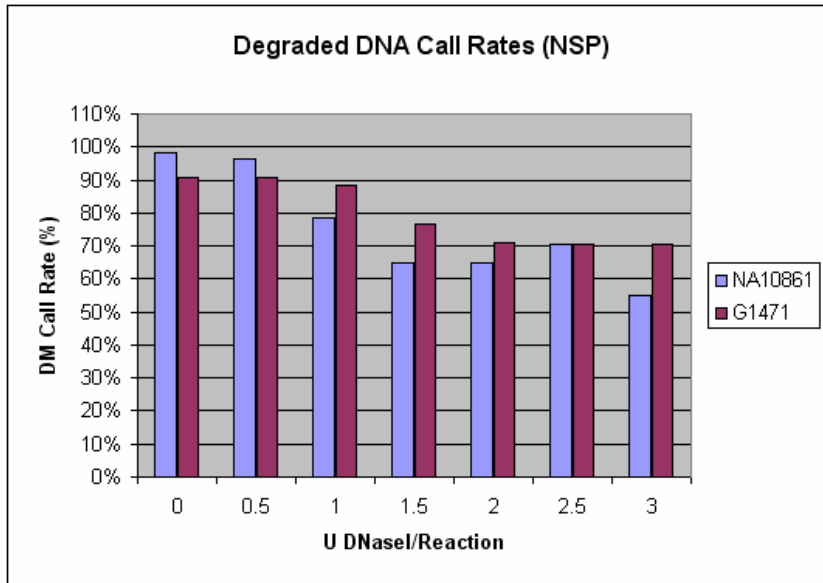


Figure 7. Percent heterozygosity as a function of mixture percentage. Two Coriell samples were mixed in varying proportions: 100:0, 99:1, 95:5, 90:10, 75:25, 50:50 and 0:100. The 500K assay was carried out on each of these mixtures and genotypes called by the standard algorithm. Percent AB was plotted for each mixture proportion. This metric accounts for the vast proportion of variability in the mixture, with high statistical significance.

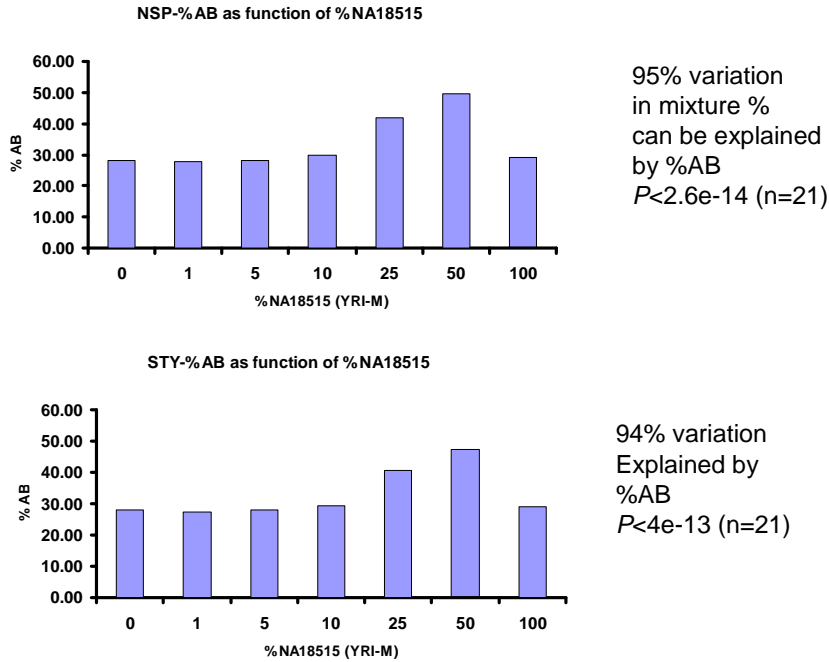


Figure 8. Pedigree of Samples used in the Mixture Study

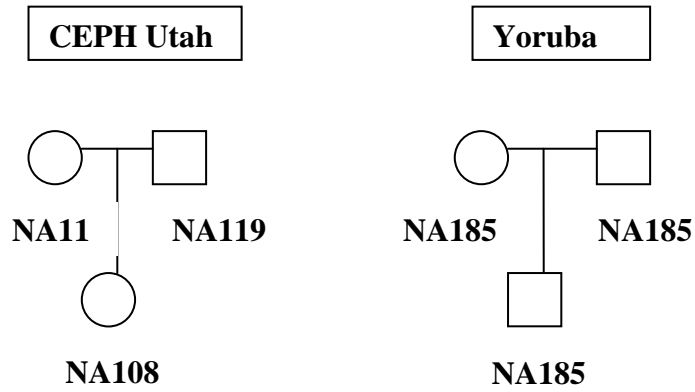


Figure 9. Detecting Mixtures. AB Heterozygosity Call Rates of STY I. Relative DNA mass contributions are indicated on the x-axis. The 3-way mixtures include an African sample, in addition to NA18515 and NA10861, in equal proportions. The 4-way mixture includes 2 Africans and 2 Caucasian samples, in equal proportions.

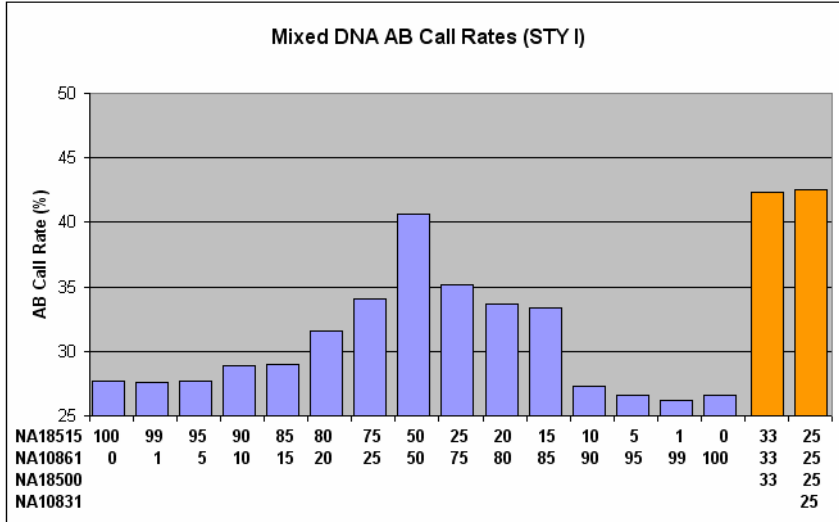
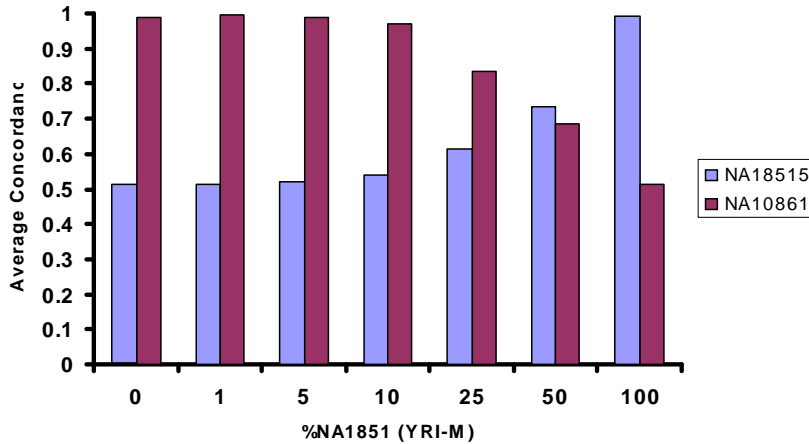


Figure 10. Genotype concordance of both components in mixture (Sty array only). The percent concordance for each of the two components of the mixture was computed at each data point using the standard genotype calling algorithm. The mixtures were composed of varying proportions of two Coriell individuals (NA18515 and NA10861), as explained in the legend to Figure 7.



% NA18515 (M)	%NA10861(F)	Concordance to NA18515	Concordance to NA10861
0	100	0.51	0.99
1	99	0.51	1.00
5	95	0.52	0.99
10	90	0.54	0.97
25	75	0.61	0.84
50	50	0.73	0.69
100	0	0.99	0.51

Figure 11. STRUCTURE analysis of mixture. We used SNPs models and pre-defined quantile normalization from 97 publicly-available Coriell DNA samples as controls. These samples are designated by the Coriell to be Caucasian, Yoruba, and East Asian (Japanese and Chinese). To decrease computation time, we selected a total of 3,700 SNPs, each belonging to an independent genetic distance bin. STRUCTURE was run at 20,000 burnin and 20,000 replications at 10 independent random start points. A k value of 3 was selected to represent the three major ethnic groups in the sample set. The results on the left clearly show 100% correct assignment of individuals as Caucasian (magenta), Yoruba (orange) and Asian (black). The mixtures were composed of one Yoruba (NA18515) and one Caucasian individual (NA10861), at varying proportions, as explained in the legend to Figure 7. STRUCTURE analysis detects the contribution of the Yoruba sample to the mixture starting at the 25% titration point.

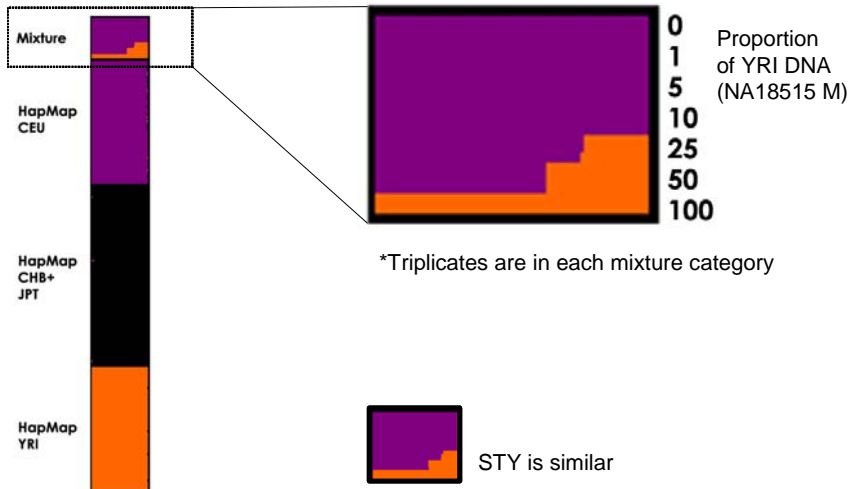


Table 3 Call Rates on 2-, 3- and 4-component Mixture Titrations The low-intensity filter described in Table 2 was not implemented in this analysis.

	Num DNA Sources	Ratios	DM, 0.33				BRLMM			
			NSP		STY		NSP		STY	
			Gender	Call Rate	Gender	Call Rate	Gender	Call Rate	Gender	Call Rate
NA18515_000p_NA10861_100p_A	2	0:100	F	96.6	F	98.7	female	98.1	female	99.0
NA18515_000p_NA10861_100p_B	2	0:100	F	98.0	F	98.3	female	99.3	female	99.0
NA18515_000p_NA10861_100p_C	2	0:100	F	98.7	F	95.5	female	99.1	female	93.7
NA18515_001p_NA10861_099p_A	2	1:100	F	98.4	F	98.2	female	99.3	female	99.0
NA18515_001p_NA10861_099p_B	2	1:100	F	98.0	F	97.3	female	99.1	female	98.3
NA18515_001p_NA10861_099p_C	2	1:100	F	97.9	F	97.7	female	99.0	female	98.4
NA18515_005p_NA10861_095p_A	2	5:100	F	98.2	F	97.4	female	98.8	female	97.8
NA18515_005p_NA10861_095p_B	2	5:100	F	97.3	F	95.6	female	97.6	female	95.4
NA18515_005p_NA10861_095p_C	2	5:100	F	97.4	F	96.9	female	98.2	female	96.7
NA18515_010p_NA10861_090p_A	2	10:90	F	95.4	F	96.3	female	95.3	female	95.9
NA18515_010p_NA10861_090p_B	2	10:90	F	96.4	F	95.9	female	96.1	female	95.3
NA18515_010p_NA10861_090p_C	2	10:90	F	95.4	F	91.9	female	95.2	female	90.8
NA18515_015p_NA10861_085p_A	2	15:85	F	90.5	F	89.5	female	88.1	female	90.6
NA18515_015p_NA10861_085p_B	2	15:85	F	91.2	F	90.3	female	90.6	female	92.0
NA18515_015p_NA10861_085p_C	2	15:85	F	91.5	F	89.8	female	90.5	female	91.6
NA18515_020p_NA10861_080p_A	2	20:80	F	85.9	F	89.4	female	84.4	female	91.3
NA18515_020p_NA10861_080p_B	2	20:80	F	91.3	F	88.2	female	91.0	female	89.5
NA18515_020p_NA10861_080p_C	2	20:80	F	89.6	F	88.0	female	89.4	female	90.8
NA18515_025p_NA10861_075p_A	2	25:75	F	91.5	F	92.4	female	91.0	female	91.1
NA18515_025p_NA10861_075p_B	2	25:75	F	92.3	F	90.8	female	91.2	female	89.5
NA18515_025p_NA10861_075p_C	2	25:75	F	91.3	F	90.0	female	90.8	female	87.8
NA18515_050p_NA10861_050p_A	2	50:50	F	90.3	F	91.6	female	88.6	female	89.6
NA18515_050p_NA10861_050p_B	2	50:50	F	90.2	F	87.7	female	88.8	female	85.2
NA18515_050p_NA10861_050p_C	2	50:50	F	91.0	F	91.9	female	88.6	female	89.4
NA18515_075p_NA10861_025p_A	2	75:25	F	91.8	F	90.3	female	91.1	female	91.6
NA18515_075p_NA10861_025p_B	2	75:25	F	91.3	F	89.4	female	90.4	female	91.8
NA18515_075p_NA10861_025p_C	2	75:25	F	91.0	F	90.1	female	89.7	female	91.1
NA18515_080p_NA10861_020p_A	2	80:20	F	91.2	F	90.8	female	91.3	female	92.8
NA18515_080p_NA10861_020p_B	2	80:20	F	92.1	F	89.6	female	92.0	female	91.7
NA18515_080p_NA10861_020p_C	2	80:20	F	93.7	F	90.1	female	93.2	female	92.2
NA18515_085p_NA10861_015p_A	2	85:15	F	95.1	F	92.4	female	95.7	female	94.9
NA18515_085p_NA10861_015p_B	2	85:15	F	94.8	F	92.0	female	95.8	female	94.8
NA18515_085p_NA10861_015p_C	2	85:15	F	93.8	F	94.7	male	95.2	female	95.7
NA18515_090p_NA10861_010p_A	2	90:10	M	96.3	M	95.1	male	97.5	male	97.6
NA18515_090p_NA10861_010p_B	2	90:10	M	92.0	M	96.0	male	93.3	male	97.9
NA18515_090p_NA10861_010p_C	2	90:10	M	94.3	M	96.6	male	96.0	male	98.1
NA18515_095p_NA10861_005p_A	2	95:5	M	93.4	M	98.2	male	93.8	male	99.4
NA18515_095p_NA10861_005p_B	2	95:5	M	95.2	M	97.7	male	97.5	male	99.3
NA18515_095p_NA10861_005p_C	2	95:5	M	96.5	M	94.7	male	98.6	male	97.0
NA18515_099p_NA10861_001p_A	2	99:1	M	97.0	M	98.1	male	98.0	male	99.3
NA18515_099p_NA10861_001p_B	2	99:1	M	96.8	M	96.4	male	97.8	male	98.6
NA18515_099p_NA10861_001p_C	2	99:1	M	98.2	M	97.6	male	99.4	male	99.3
NA18515_100p_NA10861_000p_A	2	100:0	M	98.8	M	98.4	male	99.5	male	99.3
NA18515_100p_NA10861_000p_B	2	100:0	M	98.3	M	95.9	male	99.4	male	96.1
NA18515_100p_NA10861_000p_C	2	100:0	M	99.0	M	98.3	male	99.6	male	99.3
NA18515_033p_NA10861_033p_NA18500_033p_A	33:33:33	33:33:33	F	89.4	F	90.3	female	85.7	female	86.4
NA18515_033p_NA10861_033p_NA18500_033p_B	33:33:33	33:33:33	F	88.8	F	88.1	female	86.0	female	85.2
NA18515_033p_NA10861_033p_NA18500_033p_C	33:33:33	33:33:33	F	89.5	F	90.5	female	85.4	female	85.7
NA18515_025p_NA10861_025p_NA18500_025p_A	25:25:25:25	25:25:25:25	F	88.5	F	86.4	female	86.4	female	84.3
NA18515_025p_NA10861_025p_NA18500_025p_B	25:25:25:25	25:25:25:25	F	86.9	F	89.5	female	82.9	female	87.2
NA18515_025p_NA10861_025p_NA18500_025p_C	25:25:25:25	25:25:25:25	F	88.4	F	90.0	female	86.2	female	87.1

Figure 12. Genotype Concordance Heatmap on Single, 2-, 3- and 4-component mixtures. Mixtures were made as described in Methods and run in triplicate on the 500K assay (only Sty chip data are shown here). The % genotype concordance was computed for each pair of samples and plotted using Spotfire software. Red color was assigned to the highest % concordance value (ie 100% or self) and blue was assigned to the lowest values of concordance (random, or unrelated). Mother and father for each of the two contributors for the two-way mixture were also genotyped. The pedigree of the samples is shown in Figure 8.

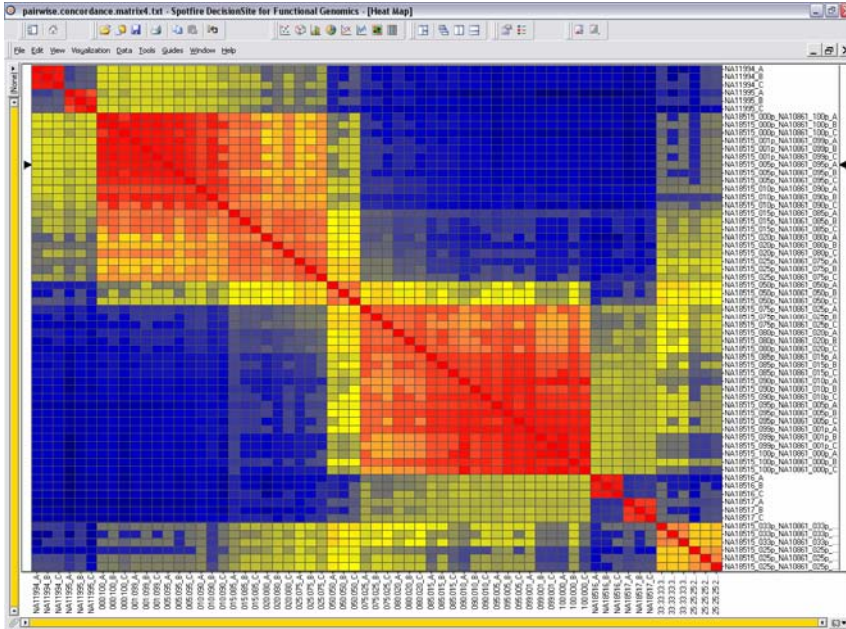


Table 4: Estimation of exact number of contributors in simulated data. The accuracies in estimation of the exact number of contributors for 1-, 2- and 3-component mixtures is shown at several values of θ (inbreeding coefficient) and for a 40- and 100-marker set of unlinked loci. N is the number of contributors in the mixture. Estimation by two maximum-likelihood methods is shown: MGP is the multi-locus genotype profile method and NHL is the number of heterozygous loci method.

θ	Number of loci: 40						Number of loci: 100					
	Single (N = 1)		Pairwise (N = 2)		Triple (N = 3)		Single (N = 1)		Pairwise (N = 2)		Triple (N = 3)	
	MGP	NHL	MGP	NHL	MGP	NHL	MGP	NHL	MGP	NHL	MGP	NHL
0	0.977	0.988	0.803	0.81	0.56	0.473	1	1	0.948	0.973	0.782	0.672
0.02	0.977	0.993	0.803	0.824	0.563	0.457	0.999	0.999	0.929	0.966	0.76	0.657
0.04	0.976	0.985	0.79	0.813	0.518	0.427	0.999	1	0.92	0.962	0.75	0.669
0.06	0.979	0.989	0.772	0.786	0.527	0.451	0.999	1	0.913	0.952	0.699	0.633
0.08	0.974	0.985	0.784	0.815	0.491	0.432	0.998	0.998	0.917	0.955	0.725	0.639
0.10	0.972	0.984	0.736	0.772	0.475	0.424	0.999	1	0.909	0.958	0.714	0.626

Table 5: Detection of sample mixtures. The accuracies in mixture detection of contributors for 1-, 2- and 3-component mixtures is shown at several values of θ (inbreeding coefficient) and for a 40- and 100-marker set of unlinked loci. N is the number of contributors in the mixture. Estimation by two maximum-likelihood methods is shown: MGP is the multi-locus genotype profile method and NHL is the number of heterozygous loci method.

θ	Number of loci: 40				Number of loci: 100			
	Pairwise (N = 2)		Triple (N = 3)		Pairwise (N = 2)		Triple (N = 3)	
	MGP	NHL	MGP	NHL	MGP	NHL	MGP	NHL
0	0.987	0.93	1	1	1	0.99	1	1
0.02	0.983	0.93	1	1	0.999	0.988	1	1
0.04	0.979	0.924	1	0.999	1	0.985	1	1
0.06	0.983	0.918	1	0.999	0.999	0.984	1	1
0.08	0.969	0.908	1	0.999	0.999	0.978	1	1
0.10	0.953	0.883	1	1	0.999	0.978	1	1

Table 6: The accuracies of three simulated datasets ($\theta = 0, 0.05, \text{ and } 0.1$) with different inbreeding coefficient ($\theta = 0, 0.01, \dots, 0.1$). N is the number of contributors in the simulated mixture. Estimation by two maximum-likelihood methods is shown: MGP is the multi-locus genotype profile method and NHL is the number of heterozygous loci method. The peak values for each determination are underlined in bold.

θ	$\theta = 0$ (Simulation)				$\theta = 0.05$ (Simulation)				$\theta = 0.1$ (Simulation)			
	N = 2		N = 3		N = 2		N = 3		N = 2		N = 3	
	MGP	NHL	MGP	NHL	MGP	NHL	MGP	NHL	MGP	NHL	MGP	NHL
<u>0</u>	<u>0.942</u>	<u>0.965</u>	<u>0.781</u>	<u>0.688</u>	<u>0.975</u>	0.945	0.691	0.437	0.954	0.814	0.317	0.154
0.01	0.915	0.964	0.728	0.68	0.969	0.955	0.715	0.485	0.962	0.838	0.382	0.192
0.02	0.871	0.954	0.66	0.677	0.96	<u>0.963</u>	0.762	0.528	0.969	0.866	0.444	0.242
0.03	0.827	0.942	0.593	0.657	0.943	<u>0.963</u>	<u>0.767</u>	0.571	0.977	0.89	0.529	0.3
0.04	0.786	0.924	0.51	0.62	0.93	0.961	0.764	0.605	<u>0.983</u>	0.915	0.592	0.357
<u>0.05</u>	0.73	0.905	0.422	0.566	0.912	0.953	0.741	0.615	0.982	0.929	0.662	0.406
0.06	0.662	0.882	0.322	0.515	0.877	0.943	0.705	0.625	0.977	0.939	0.691	0.458
0.07	0.589	0.851	0.25	0.449	0.831	0.935	0.65	<u>0.626</u>	0.969	0.954	0.727	0.508
0.08	0.521	0.799	0.173	0.398	0.786	0.919	0.583	0.617	0.957	0.962	<u>0.728</u>	0.548
0.09	0.455	0.746	0.127	0.343	0.742	0.894	0.503	0.597	0.943	<u>0.968</u>	0.711	0.579
<u>0.10</u>	0.37	0.703	0.087	0.282	0.674	0.863	0.421	0.577	0.92	<u>0.968</u>	0.683	<u>0.596</u>

Table 7: The accuracies of estimation of exact number of contributors for full-sib and parent-offspring mixtures. θ is the inbreeding coefficient in estimation. Estimation by two maximum-likelihood methods is shown: MGP is the multi-locus genotype profile method and NHL is the number of heterozygous loci method. The peak values for each determination are underlined in bold.

θ	Number of loci: 40				Number of loci: 100			
	Full Sib		Parent Offspring		Full Sib		Parent Offspring	
	MGP	NHL	MGP	NHL	MGP	NHL	MGP	NHL
0	0.552	0.408	0.668	0.49	0.552	0.292	0.742	0.456
0.05	0.683	0.552	0.786	0.634	0.796	0.568	0.923	0.723
0.10	<u>0.797</u>	0.668	<u>0.806</u>	0.721	0.928	0.832	<u>0.967</u>	0.871
0.15	0.785	0.743	0.759	<u>0.776</u>	<u>0.96</u>	0.904	0.917	<u>0.944</u>
0.20	0.712	<u>0.768</u>	0.628	0.752	0.824	<u>0.92</u>	0.693	0.89
0.25	0.522	0.688	0.422	0.618	0.616	0.86	0.457	0.748
0.30	0.344	0.546	0.234	0.444	0.296	0.6	0.146	0.458

Table 8. Sample manifest submitted by Bode with information on DNA concentration, source and gender for two sets of samples. Following the blind experiment, the samples were unblinded and ethnicities were revealed by Bode.

AFFX No.	Bode Set	Sex	Ethnicity	Sample Type	Extraction Method	Quant by RT-PCR	Spec	RT-PCR ng/ul	Spec ng/ul	Nanodrop ng/ul
B001	001	Female	NA	Hair	Qiagen DNA Micro Kit	NA	NA	5.2	NA	NA
B002	001	Female	NA	Hair	Qiagen DNA Micro Kit	NA	NA	1.15	NA	NA
B003	001	Female	NA	Hair	Qiagen DNA Micro Kit	NA	NA	0.00112	NA	NA
B004	001	Female	NA	Hair	Qiagen DNA Micro Kit	NA	NA	0.000749	NA	NA
B005	001	Female	NA	Hair	Qiagen DNA Micro Kit	NA	NA	0	NA	NA
B006	001	Female	NA	Hair	Qiagen DNA Micro Kit	NA	NA	0.000745	NA	NA
B007	001	Male	NA	Blood Stain	Qiagen DNA Micro Kit	NA	NA	21.91	NA	NA
B008	001	Male	NA	Blood Stain	Qiagen DNA Micro Kit	NA	NA	12.72	NA	NA
B009	001	Female	NA	Blood Stain	Qiagen DNA Micro Kit	NA	NA	19.78	NA	NA
B010	001	Female	NA	Blood Stain	Qiagen DNA Micro Kit	NA	NA	9.36	NA	NA
B011	001	Female	NA	Blood Stain	Qiagen DNA Micro Kit	NA	NA	19.39	NA	NA
B012	001	Female	NA	Blood Stain	Qiagen DNA Micro Kit	NA	NA	7.91	NA	NA
B013	001	Female	NA	Blood Stain	Qiagen DNA Micro Kit	NA	NA	22.81	NA	NA
B014	001	Female	NA	Blood Stain	Qiagen DNA Micro Kit	NA	NA	10.12	NA	NA
B015	001	Female	NA	Blood Stain	Qiagen DNA Micro Kit	NA	NA	9.42	NA	NA
B016	001	Female	NA	Blood Stain	Qiagen DNA Micro Kit	NA	NA	10.69	NA	NA
B017	001	Female	NA	Blood Stain	Qiagen DNA Micro Kit	NA	NA	5.88	NA	NA
B018	001	Female	NA	Blood Stain	Qiagen DNA Micro Kit	NA	NA	3.21	NA	NA
B019	001	Female	NA	Blood Stain	Qiagen DNA Micro Kit	NA	NA	10.73	NA	NA
B020	002	Female	European	Vaginal	Qiagen Micro	9.49	0.047	474.5	235	98.9
B021	002	Male	African American	Buccal	Qiagen Micro	1.37	0.014	68.5	70	25.5
B022	002	Female	Mixed-Caucasian and Asian	Buccal	Qiagen Micro	2.66	0.0727	133	363.5	27.6
B023	002	Female	Hispanic	Buccal	Qiagen Micro	4.79	0.0156	239.5	78	60.5
B024	002	Male	European	Semen	Qiagen Micro	4.7	0.04	235	200	26.7
B025	002	Female	Hispanic	Buccal	Qiagen Micro	1.76	0.014	88	70	47.1
B026	002	Female	Asian	Vaginal	Qiagen Micro	6.82	0.03	341	150	89.4
B027	002	Female	Hispanic	Buccal	Qiagen Micro	8.92	0.157	446	785	41.2
B028	002	Female	African American	Buccal	Qiagen Micro	6.36	0.018	318	90	47.1
B029	002	Male	European	Semen	Qiagen Micro	2.2	0.007	110	35	13.7
B030	002	Female	Mixed-	Buccal	Qiagen Micro	1.47	0.028	73.5	140	68.9

			Caucasian and Asian							
B031	002	Female	Asian	Vaginal	Qiagen Micro	15.14	0.185	757	925	128.2
B032	002	Female	Hispanic	Buccal	Qiagen Micro	0.77	0.04	38.5	200	57.5
B033	002	Male	European	Semen	Qiagen Micro	3.52	0.074	176	370	26.3
B034	002	Female	African American	Buccal	Qiagen Micro	5.36	0.033	268	165	48.8
B035	002	Female	European	Vaginal	Qiagen Micro	39.63	0.284	1981.5	1420	251.2
B036	002	Female	European	Vaginal	Qiagen Micro	53.64	0.298	2682	1490	384.5
B037	002	Male	European	Semen	Qiagen Micro	4.6	0.038	230	190	32.5
B038	002	Female	European	Vaginal	Qiagen Micro	11.61	0.192	580.5	960	88.0
B039	002	Female	European	Vaginal	Qiagen Micro	15.06	0.109	753	545	117.6
B040	002	Male	Unknown	Semen	Qiagen Micro	4.04	0.039	202	195	28.8
B041	002	Male	African American	Buccal	Qiagen Micro	1.91	0.017	95.5	85	33.7
B042	002	Female	European	Vaginal	Qiagen Micro	20.24	0.202	1012	1010	121.1
B043	002	Female	Asian	Vaginal	Qiagen Micro	3.34	0.036	167	180	73.4
Ref103	001	Male	NA	NA	NA	NA	NA	5 ng/uL	NA	48.9
Ref103	002	Male	NA	NA	NA	NA	NA	5 ng/uL	NA	48.9

Table 9. Summary statistics on analysis of 43 Bode samples on 500K. Only data for the Sty chip is shown here. Standard GTYPE statistics were computed such as %call rate, %AB, %MDR-MCR and Gender. The low-intensity filter described in Table 2 was not implemented in this analysis.

DM 500K STY WGA DNA							
CellID	Call Rate	PCR Gel	Called Gender	AB Call	MCR	MDR	MDR - MCR
B001_WGA_STY_01	90.7%	Pass	F	16.0%	75.5%	92.8%	17.3%
B002_WGA_STY_01	84.9%	Pass	F	10.5%	65.3%	82.2%	16.9%
B003_WGA_STY_01	50.4%	Fail	F	12.5%	2.4%	5.7%	3.3%
B004_WGA_STY_01	51.6%	Fail	F	14.7%	3.6%	8.8%	5.3%
B005_WGA_STY_01	51.7%	Fail	F	13.7%	3.4%	7.9%	4.5%
B006_WGA_STY_01	56.0%	Fail	F	8.8%	3.0%	6.2%	3.2%
B007_WGA_STY_01	93.1%	Pass	M	23.6%	84.7%	95.7%	10.9%
B008_WGA_STY_01	94.3%	Pass	M	24.4%	87.0%	96.8%	9.8%
B009_WGA_STY_01	93.3%	Pass	F	23.3%	85.2%	95.4%	10.2%
B010_WGA_STY_01	93.4%	Pass	F	22.1%	81.4%	96.1%	14.7%
B011_WGA_STY_01	90.8%	Pass	F	22.1%	80.1%	92.7%	12.7%
B012_WGA_STY_02	96.0%	Pass	F	24.7%	89.5%	98.4%	8.9%
B013_WGA_STY_01	93.3%	Pass	M	22.3%	83.2%	95.4%	12.2%
B014_WGA_STY_01	93.9%	Pass	M	23.6%	86.4%	95.7%	9.2%
B015_WGA_STY_01	92.2%	Pass	F	22.5%	81.2%	94.8%	13.7%
B016_WGA_STY_01	93.4%	Pass	F	23.2%	83.5%	95.3%	11.8%
B017_WGA_STY_01	93.3%	Pass	F	23.1%	84.0%	95.4%	11.3%
B018_WGA_STY_01	50.8%	Fail	F	10.7%	2.1%	5.3%	3.1%
B019_WGA_STY_01	90.3%	Pass	F	27.0%	72.3%	95.9%	23.6%
B020_WGA_STY_01	80.3%	Pass	F	16.3%	58.9%	80.7%	21.8%
B021_WGA_STY_01	89.8%	Pass	M	25.6%	82.3%	93.8%	11.5%
B022_WGA_STY_01	84.1%	Pass	F	16.9%	64.7%	85.9%	21.2%
B023_WGA_STY_01	90.7%	Pass	F	23.7%	80.1%	93.8%	13.7%
B024_WGA_STY_01	93.9%	Pass	M	25.3%	88.5%	96.1%	7.6%
B025_WGA_STY_01	92.7%	Pass	F	25.5%	84.5%	95.3%	10.9%
B026_WGA_STY_01	89.4%	Pass	F	20.8%	77.9%	92.8%	14.9%
B027_WGA_STY_01	87.7%	Pass	F	23.3%	74.2%	92.7%	18.5%
B028_WGA_STY_01	89.2%	Pass	F	25.1%	77.2%	93.5%	16.3%
B029_WGA_STY_01	90.9%	Pass	M	23.7%	82.1%	92.8%	10.7%
B030_WGA_STY_01	93.3%	Pass	F	26.4%	83.3%	96.2%	12.9%
B031_WGA_STY_01	87.9%	Pass	F	20.6%	75.6%	90.7%	15.1%
B032_WGA_STY_01	92.3%	Pass	F	24.4%	84.4%	94.9%	10.5%
B033_WGA_STY_01	90.9%	Pass	M	23.2%	81.2%	93.9%	12.8%
B034_WGA_STY_01	92.3%	Pass	F	27.2%	85.0%	95.7%	10.7%
B035_WGA_STY_01	86.9%	Pass	F	20.9%	72.7%	90.4%	17.7%
B036_WGA_STY_01	93.8%	Pass	F	25.5%	87.2%	96.3%	9.1%
B037_WGA_STY_01	92.5%	Pass	M	25.1%	86.9%	95.4%	8.5%

B038_WGA_STY_01	76.9%	Fail	F	14.1%	52.1%	73.4%	21.3%
B039_WGA_STY_01	88.6%	Pass	F	22.9%	75.1%	90.9%	15.8%
B040_WGA_STY_01	92.9%	Pass	M	25.2%	87.9%	95.7%	7.8%
B041_WGA_STY_01	91.7%	Pass	M	26.2%	84.5%	95.1%	10.6%
B042_WGA_STY_01	74.3%	Fail	F	15.0%	49.4%	72.2%	22.8%
B043_WGA_STY_01	89.9%	Fail	F	22.5%	78.9%	93.1%	14.3%
Ref103_WGA_STY_01	93.5%	Pass	M	21.8%	83.7%	96.7%	13.0%
NA	NA	NA	NA	NA	NA	NA	NA

Figure 13. STRUCTURE analysis on Bode samples. Methods are described in legend to Figure 11. The samples include Europeans, Africans, Asians and mixed ethnicities. Mito refers to haplogroup data generated on the Affymetrix mitochondrial resequencing chip. The mito haplogroup assignments were 100% consistent with ethnic description provided by Bode.

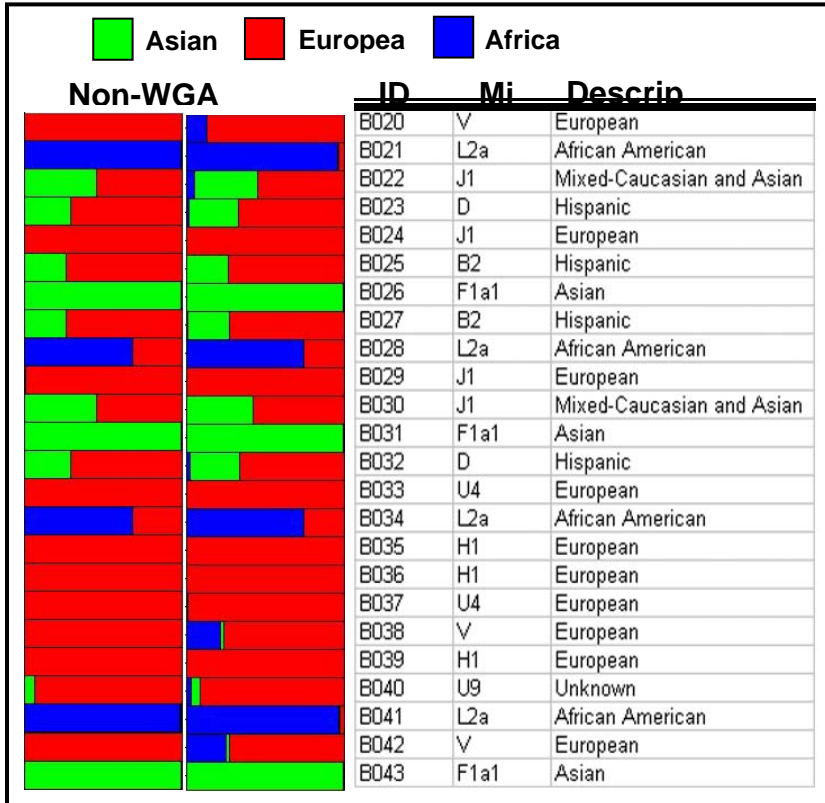
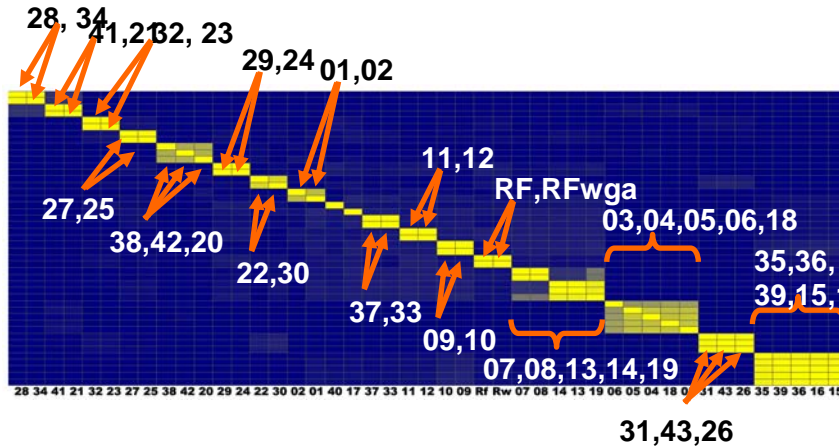


Figure 14. Pairwise Allele-sharing heatmap on Bode samples sets 1 and 2. Bright yellow indicates the highest level of allele sharing (ie self) and dark blue indicates random allele sharing (ie non-related). This analysis reveals which samples are duplicates (bright yellow) and also reveals a mixture. Sample 19 has a high degree of allele sharing with samples 13 and 14 (which are duplicates of each other) and also lesser sharing with samples 07 and 08 (which are duplicates of each other), shown by the less intense yellow color. This analysis suggests that samples 13/14 are a larger contributor to the mixture than samples 07/08. Bode later confirmed that Sample 19 is a 2:1 mixture of 13/14 and 07/08. Note that duplicate samples with lower call rates tend to have lower allele-sharing due to allele-drop out (e.g. less intense yellow) but are still clearly identified.



IV. Conclusions

Overall Conclusions

- The goal of this work is to use high-density SNP assays to extract accurate genotypic information from reduced quantity and quality forensics samples. We began by laying a solid foundation of assay and algorithm development work which then allowed us to successfully test blinded forensics samples supplied by Bode Technologies.
- The assay development successfully introduced a whole genome amplification (WGA) step to adapt smaller sample quantities to the standard 500K protocol. While call rates are generally lower with pre-amplified samples, nonetheless sufficient data are collected on hundreds of thousands of SNPs to provide critical genotype information such as relatedness and ancestry.
- We generated DNase I-degraded samples and obtained useful genotype information on a subset of SNPs, using the standard algorithm.
- We generated panels of defined mixtures of two, three and four-component DNAs at varying proportions. We processed them on the standard 500K assay and determined that the simple metric of % heterozygosity (ie %AB calls) is statistically sufficient to detect mixtures at the 25:75 ratio. Genotype concordance measurements can detect mixtures down to 5% contribution of the minor component.

Discussion of findings

Low Copy Template

Forensics samples rarely yield large quantities of nucleic acids, therefore methods which pre-amplify genomic DNA are desirable for downstream molecular assays. We used isothermal whole-genome amplification (WGA) with phi29 polymerase to amplify 10 ng of starting DNA many thousand-fold to microgram quantities. We tested whether this pre-amplified material could then be used downstream in the 500K (WGS) assay and compared its performance to non-amplified genomic DNA. Table I shows a comparison of WGA vs non-WGA pre-amplification on 500K arrays. The WGA call rates range from 80-96%, while the non-WGA call rates are all >98%. This shows that pre-amplification of template results in 2-20% marker loss in the 500K assay. This is most likely due to incomplete pre-amplification of the starting genomic DNA by the phi29 polymerase. Due to the large number of markers captured, it is likely that downstream analyses and positive identification would not be compromised with such an overabundance of markers.

We also tested whether smaller amounts of starting material could be analyzed by WGA/WGS. We titrated 0, 1, 10, 100, 1000, 10,000 and 250,000 pg of DNA in the 500K mapping assay. Figure 3 shows a raw intensity plot of the signal on the arrays. While the standard BRLMM algorithm is designed to be used on the full amount of 250 ng, the arrays themselves are capable of detecting signal down to pg levels of DNA. The standard assay (250,000 pg) on the lower right clearly shows separation of the AA, AB and BB genotypes in the sample. As the starting DNA quantity is decreased, the number of AB genotypes called clearly decreases. This is most likely due to stochastic drop-out of diploid DNA as lower starting amounts of DNA are used. Nonetheless, novel algorithms designed to analyze incomplete genomes (ie a mixture of haploid and diploid SNPs) will be necessary to accurately capture the genomic variation in these ultra-low

copy templates. By employing a simple low-signal cutoff (ie low intensity filter) we were able to remove some of the background noise in the low template samples (Table 2)

Degraded Samples

Forensics samples are not always composed of intact genomic DNA. The nucleic acids can be subject to varying degrees of degradation. To mimic degradation in forensics samples, we digested DNAs in a controlled fashion with DNase I (Figure 4). When we use these degraded templates for the WGS reaction, we observe non-standard behavior of the PCR products as the DNA template is increasingly degraded (Figure 5). Each of these WGS products were further processed and run on 500K arrays. Call rates dropped significantly with an increase in DNase I concentration, but leveled off to ~50% for the Coriell sample (Figure 6). These data indicate that even highly degraded DNA can yield genotype data on hundreds of thousands of SNPs using this assay.

Mixtures

Mixtures of two or more contributors are frequently encountered in forensics casework. The standard 500K assay was designed to genotype only one sample at a time, therefore we set out to determine whether the arrays could detect mixtures. We created a set of controlled mixtures of two Coriell samples mixed in varying proportions: 100:0, 99:1, 95:5, 90:10, 75:25, 50:50 and 0:100. The 500K assay was carried out on each of these mixtures and genotypes called by the standard algorithm. Percent AB was plotted for each mixture proportion (Figure 7); clearly this measurement alone accounts for the vast proportion of variability in the mixture, with high statistical significance ($10e-13$). This measurement was informative even with mixtures of 3 and 4 DNAs (Figure 9) indicating that %AB would be an adequate metric for detecting putative mixtures, the composition of which would then need to be further explored. As DNAs are mixed, it is expected that the BRLMM genotype algorithm (which was developed for single donors only) would make erroneous genotype calls due to mixed signal from multiple SNPs. To measure the extent of error, we computed the genotype concordance to the known genotypes in the individual contributors as a function of mixture percentage. As expected, the error rate increases as the minor contributor is increased from a few percent to 50% (Figure 10). Given the large number of SNPs that are measured in the 500K assay, we set out to determine whether the increase in genotype error would affect a downstream analysis such as ancestry determination.

STRUCTURE results clearly show 100% correct assignment of individuals as Caucasian, Yoruba or Asian. The mixtures were composed of one Yoruba and one Caucasian individual, at varying proportions. STRUCTURE analysis detects the contribution of the Yoruba sample to the mixture as low as at the 25% titration point (Figure 11). This indicates that even if a subset of SNP markers is genotyped in mixtures, it is still possible to generate information about that sample. This may be useful in cases where one or more STR markers fail.

We also tested pairwise genotype concordance on mixture titrations from 2, 3 and 4 component mixtures and generated colorized heatmaps from the results (Figure 12). The red color indicates the highest genotype concordance (100%, ie self) and blue colors indicate the lowest genotype concordance (~50%, ie unrelated individuals). Each titration point was run in triplicate to provide data on reproducibility (designated A, B and C). As controls for relatedness, we included the mother and father for each of the components in the mixture. The relationships between the samples in the mixture studies are indicated in the pedigree (Figure 8). Three general trends are visible in this large dataset: 1) the reproducibility of the triplicates is high; 2) each of the two components in the mixtures has high genotype concordance to itself at high ratios (red, red-orange, yellow), which decreases down to yellow-green, blue-green as the minor component is decreased down to 10-15%; 3) familial relationships are preserved as low as 25% minor component; for example, each of the two children in the mixture show high relatedness (yellow)

to each of their parents, but not to the other set of parents (blue). As the minor component is decreased, the genotype concordance also decreases, as can be observed as the color changes from yellow to yellow-green to blue-green. Nonetheless, concordance values are still distinguishable from the dark blue color representative of unrelated individuals.

It is likely that a smaller set of SNP markers can be used to detect and evaluate mixtures. The optimal number and characteristics of this subset of markers has not yet been determined. Our collaborators from the University of Cincinnati, Ranajit Chakraborty and his laboratory, have begun to develop theoretical methods for detecting the presence of mixtures and for determining the number of contributors in simulated data. They have developed two methods: MGP (multi-locus genotype profile) and NHL (number of heterozygous loci). Using these two maximum-likelihood based methods, Dr. Chakraborty can theoretically detect mixtures with 2 and 3 contributors at >98% accuracy across a wide range of inbreeding coefficients (Table 5). Furthermore, these methods predict the exact number of contributors with >95% accuracy in two-way mixtures, and >78% accuracy for 3-way mixtures (Figures 4-7) (Chakraborty et al. 2008). At the present time, applications for linked loci has not been developed, and hence no mixture analyses for 500K sets is as yet available.

Blinded Samples From Bode Technologies

We received a total of 43 blinded forensics samples from Bode Technologies. Information on the samples included DNA concentration, source and gender (Table 8). Following the blind experiment, the samples were unblinded and ethnicities were revealed by Bode. We generated data on the samples using the 500K assay. For conciseness, only data for the Sty chip is reported here as the Nsp chip data are highly similar. Standard GTYPE statistics were computed such as %call rate, %AB, %MDR-MCR and Gender (Table 9). Most of the samples, with the exception of hair shafts (no roots), passed the standard 500K assay with high call rates.

Effect of DNA Source on Call Rates: We noted previously that genotype data from blood or from cultured cell lines performs very well with the 500K assay (see previous progress reports). In this project, we analyzed DNA samples from blood stains, hair, semen, buccal and vaginal swabs. As expected, hairs without roots had very low call rates, consistent with the lack of nuclear DNA. While blood, semen and buccal samples show data similar to the controls, the vaginal swabs had call rates that were systematically decreased about 5-10% from the controls (Table 5). This small decrease had no effect on the ethnicity determination by STRUCTURE (see Figure 13 and discussion below) or on correctly assigning gender. Further analysis will be necessary to determine whether possibly other forensically relevant information is compromised in DNA samples taken from vaginal swabs.

Ethnicity Determination: The samples from Bode include Europeans, Africans, Asians and mixed ethnicities. We were blinded to the ethnicities of the samples until after our 500K results were shared with Bode Technologies. We used STRUCTURE to assign ethnicities from the nuclear data (see methods and legend to Figure 11). In all cases, we correctly identified the ethnicity of the Bode samples, including one sample with a unique ancestry consistent with a Taiwanese aborigine group (Figure 13). We also easily detected several cases of genetic admixture (African-Caucasian and Caucasian-Asian). All ethnicities were confirmed by Bode. We also generated mitochondrial haplogroups on the Bode samples using the Affymetrix mitochondrial v2.0 resequencing chip (Figure 13). The mito haplogroup assignments were 100% consistent with ethnic description provided by Bode.

Allele-sharing determinations: We set out to determine whether there were any duplicates or mixtures in the Bode sample set. We computed pairwise allele-sharing amongst the samples. Duplicates would be expected to have the highest allele-sharing (ie self, indicated by yellow color) and unrelated individuals would expect to randomly share alleles (blue color). As seen in Figure 14, there were many sets of duplicates. Laboratory mixtures, or casework samples coming from more than one contributor, would be expected to also share alleles. Therefore some of the samples showing high levels of allele-sharing (shades of yellow and yellow-green) could indicate mechanical, rather than genetic, mixtures. In these cases, the allele-sharing value is lower and the yellow color is less intense. Due to the high heterozygosity observed in Bode sample 19, we suspected a mixture. By computing the allele-sharing value for all the Bode mixtures, we concluded that sample 19 was composed of a mixture of samples 07/08 and samples 13/14, each of which were duplicates. Bode later confirmed that Sample 19 is a 2:1 mixture of 13/14 and 07/08. Note that duplicate samples with lower call rates tend to have less allele-sharing due to allele-drop out (e.g. less intense yellow) but are still clearly identified (e.g. Bode samples 03, 04, 05, 06).

Implications for policy and practice

As feasibility is demonstrated for SNP-based platforms in forensics, practitioners may be interested in testing the robustness in their own laboratories. The ability to access samples and genetic information that is not amenable to standard STR-based analysis may prove to be an incentive for exploring other genetic marker systems such as high-density SNP marker sets.

Implications for further research

Much of this work can be expanded more fully to meet the needs of the forensics community. We have developed an ongoing dialog with the California Department of Justice to identify and refine system requirements for a SNP-based platform. In an effort to build upon the work accomplished in this grant period, investigators such as Martin Buoncristiani and Eva Steinberg at the CalDOJ will be writing new NIJ grant proposals based on this work. Using real forensic casework and databank samples will be very important for the next phase of this project and this can only be done in an accredited and licensed forensics laboratory with access to test development expertise, such as the CalDOJ. Here, authentic degraded samples from casework that fail STR analysis can be tested with a high-density SNP panel to determine whether an exclusion or inclusion can be made with high statistical probability. Similarly, complex mixtures that give confusing STR results may be amenable to deconvolution using high-density SNP panels. Our mixture titrations and allele-sharing data suggest that large numbers of SNPs may be useful for identification as sample complexity increases.

V. References Cited

- Affymetrix (2006). Product Update: BRLMM Analysis Tool (Affymetrix, Inc.).
URL: http://www.affymetrix.com/support/technical/product_updates/brlmm_algorithm.affx
- Budowle, B., Allard, M. W., Wilson, M. R., and Chakraborty, R. (2003). Forensics and mitochondrial DNA: applications, debates, and foundations. *Annu Rev Genomics Hum Genet* 4, 119-141.
- Chakraborty, R. et al. (2008) Detection of DNA mixtures and estimation of the number of contributors. (manuscript in preparation)
- Chakraborty R, and Jin L (1992) Heterozygote deficiency, population substructure and their implications in DNA fingerprinting. *Hum Genet* 88:267-272
- Chakraborty R, Stivers DN, Su B, Zhong Y, Budowle B (1999) The utility of short tandem repeat loci beyond human identification: Implications for development of new DNA typing systems. *Electrophoresis* 20: 1682-1696
- Chee, M., Yang, R., Hubbell, E., Berno, A., Huang, X. C., Stern, D., Winkler, J., Lockhart, D. J., Morris, M. S., and Fodor, S. P. (1996). Accessing genetic information with high-density DNA arrays. *Science* 274, 610-614.
- Chicago, University of. (2007). STRUCTURE.
URL: <http://pritch.bsd.uchicago.edu/structure.html>
- Curran JM, Triggs CM, Buckleton J, Weir BS (1999) Interpreting DNA mixtures in structured populations. *J Forensic Sci* 44:987-995
- Egeland T, Dalen I, Mostad PF (2003) Estimating the number of contributors to a DNA profile. *Int J Legal Med* 117(5): 271-275
- Evelt IW, Buffery C, Willott G, Stoney D (1991) A guide to interpreting single locus profiles of DNA mixtures in forensic cases. *J Forensic Sci Soc* 31:41-47
- Evelt IW, Weir BS (1998) Interpreting DNA Evidence. Sinauer Associates Inc., Massachusetts.
- Falush, D., Stephens, M., and Pritchard, J. K. (2003). Inference of population structure using multilocus genotype data: linked loci and correlated allele frequencies. *Genetics* 164, 1567-1587.
- Fodor, S. P., Read, J. L., Pirrung, M. C., Stryer, L., Lu, A. T., and Solas, D. (1991). Light-directed, spatially addressable parallel chemical synthesis. *Science* 251, 767-773.
- Frudakis, T., Venkateswarlu, K., Thomas, M. J., Gaskin, Z., Ginjupalli, S., Gunturi, S., Ponnuswamy, V., Natarajan, S., and Nachimuthu, P. K. (2003). A classifier for the SNP-based inference of ancestry. *J Forensic Sci* 48, 771-782.
- Fung WK, Hu YQ (2000a) Interpreting forensic DNA mixtures: allowing for uncertainty in population substructure and dependence. *J R Statist Soc A* 163:241-254

Fung WK, Hu YQ (2000b) Interpreting DNA mixtures based on the NRC-II recommendation 4.1. *Forensic Sci Commun* 2(4)

Fung WK, Hu YQ (2002) The statistical evaluation of DNA mixtures with contributors from different ethnic groups. *Int J Legal Med* 116:79–86

Fung WK, Hu YQ (2004) Interpreting DNA Mixtures with Related Contributors in Subdivided Populations. *Scand J Stat* 31 (1): 115-130

Gill P (2001) An assessment of the utility of single nucleotide polymorphisms (SNPs) for forensic purposes. *Int J Legal Med* 114:204–210

Gill, P. (2002). Role of short tandem repeat DNA in forensic casework in the UK--past, present, and future perspectives. *Biotechniques* 32, 366-368, 370, 372, passim.

Ginther, C., Issel-Tarver, L., and King, M. C. (1992). Identifying individuals by sequencing mitochondrial DNA from teeth. *Nat Genet* 2, 135-138.

Holland, M. M., Fisher, D. L., Mitchell, L. G., Rodriguez, W. C., Canik, J. J., Merrill, C. R., and Weedn, V. W. (1993). Mitochondrial DNA sequence analysis of human skeletal remains: identification of remains from the Vietnam War. *J Forensic Sci* 38, 542-553.

Holland, M. M., and Parsons, T. J. (1999). Mitochondrial DNA Sequence Analysis: Validation and Use for Forensic Casework. *Forensic Sci Rev* 11, 21-50.

Higuchi, R. G., von Beroldingen, C. H., Sensabaugh, G. F., and Erlich, H. A. (1988). DNA typing from single hairs. *Nature* 332, 543-546.

Hu YQ, Fung WK (2003) Interpreting DNA mixtures with the presence of relatives. *Int J Legal Med* 117:39–45

Kennedy, G. C., Matsuzaki, H., Dong, S., Liu, W. M., Huang, J., Liu, G., Su, X., Cao, M., Chen, W., Zhang, J., *et al.* (2003). Large-scale genotyping of complex DNA. *Nat Biotechnol* 21, 1233-1237.

Kwok, P. Y. (2001). Methods for genotyping single nucleotide polymorphisms. *Annu Rev Genomics Hum Genet* 2, 235-258.

Lauritzen SL, Mortera J (2002) Bounding the number of contributors to mixed DNA stains. *Forensic Sci Int* 130:125–126

Pritchard, J. K., Stephens, M., and Donnelly, P. (2000). Inference of population structure using multilocus genotype data. *Genetics* 155, 945-959.

Rabbee, N., and Speed, T. P. (2005). A genotype calling algorithm for Affymetrix SNP arrays. UC Berkeley Statistics Online Tech Reports.

Rabbee, N., and Speed, T. P. (2006). A genotype calling algorithm for affymetrix SNP arrays. *Bioinformatics* 22, 7-12.

Stockmarr A (2000) The choice of hypotheses in the evaluation of DNA profile evidence. In: Gastwirth JL (ed) *Statistical science in the courtroom*. Springer, Berlin Heidelberg New York, pp 143–160

Stoneking, M., Hedgecock, D., Higuchi, R. G., Vigilant, L., and Erlich, H. A. (1991). Population variation of human mtDNA control region sequences detected by enzymatic amplification and sequence-specific oligonucleotide probes. *Am J Hum Genet* 48, 370-382.

Syvanen, A. C. (2001). Accessing genetic variation: genotyping single nucleotide polymorphisms. *Nat Rev Genet* 2, 930-942.

United States Department of Justice, Office of Justice Programs (2002). *Using DNA to Solve Cold Cases*.

United States Department of Justice, Office of Justice Programs (2003). *Report to the Attorney General on Delays in Forensic DNA Analysis*.

Weir, B. S. (1996). *Genetic Data Analysis II* (Sunderland, MA, Sinauer Associates, Inc.).

Weir BS, Triggs C, Starling L, Stowell L, Walsh K, Buckleton J (1997) Interpreting DNA mixtures. *J Forensic Sci* 47:213–222

Xiong, M., and Jin, L. (1999). Comparison of the power and accuracy of biallelic and microsatellite markers in population-based gene-mapping methods. *Am J Hum Genet* 64, 629-640.

VI. Dissemination of Research Findings

The results of our studies will be published in peer-reviewed journals. Several manuscripts are in preparation.