

Forensic mtDNA analysis of highly degraded materials, or samples lacking sufficient quantity of nuclear DNA (e.g. shed hairs), has found an important niche in DNA testing. Recent mtDNA research has focused on two important limitations for mtDNA testing: (a) the cost of generating a mtDNA sequence profile, and (b) the low power of discrimination associated with common mtDNA types. To overcome the cost prohibition of mtDNA testing, Linear Arrays have been evaluated as a screening tool (Kline et al. 2005). To increase the power of discrimination for individuals sharing one of the few common mtDNA types, strategies to identify resolving polymorphisms in the coding region through sequencing short (~100 bp) fragments of (Allen and Andersson 2005) or through the identification of SNPs (Coble et al. 2004; Vallone et al. 2004) have been proposed.

To assess the amount of variation gained by entire control region sequencing compared to Linear Array mitotyping, a comparison of discrimination from Linear Array – HV1 – HV1/HV2 – Control Region was determined among 666 population samples. Further discrimination for the most common haplotype, MCH (A263G; 315.1C), was determined by coding region SNP analysis (Coble et al. 2004; Vallone et al. 2004). In addition, a survey of the underlying source of null alleles (blanks) in Linear Arrays, as determined by sequence information was performed.

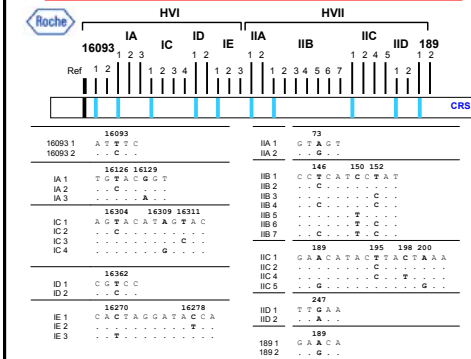
Finally, an evaluation of coding region variation in both a global dataset of mtDNA genomes, and among the most common HV1/HV2 haplotype in Caucasians (A263G; 315.1C) was also examined to determine if sequencing strategies for identifying mtDNA variation are more effective than targeting synonymous SNPs (Budowle et al. 2005; Coble et al. 2006).

Evaluation of mtDNA Sequence Data to the Linear Array

Materials and Methods

Linear Array mitotypes for 666 samples representing African American, Caucasian, and Hispanic populations were generated using the manufacturer's protocol and the results were published in Kline et al. (2005). Complete control region sequence data was generated at the Armed Forces DNA Identification Laboratory (AFDIL) using the established protocol and strategy described in Brandstadter et al. (2004). Haplogroup associations were determined from Kivisild et al. (2006) and references within.

Schematic representation of mtDNA Linear Array HV1/HV2 probe regions



Breakdown of Discrimination among 666 population samples

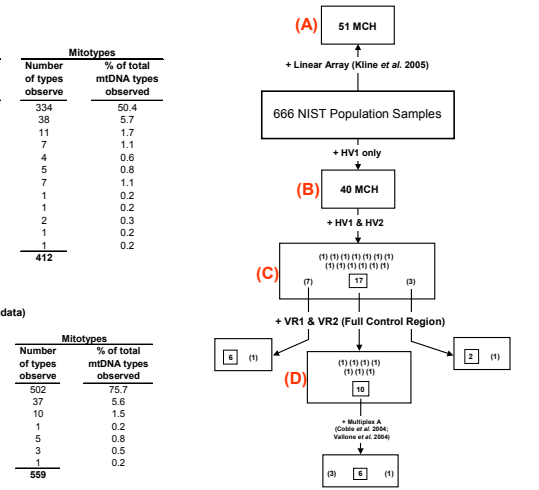
A Linear Array (HV1/HV2) Kline et al. (2005) J. Forensic Sci. 50(2): 377-385.				B HV1 only (Sequence data)							
Times type observed	Frequency	% of individuals tested	Number of types observed	Mitotypes	% of total mtDNA types observed	Times type observed	Frequency	% of individuals tested	Number of types observed	Mitotypes	% of total mtDNA types observed
1	27.9	165	27.9	1	50.4	1	75.7	1	50.4	1	50.4
2	15.8	45	6.8	2	33.4	2	11.5	2	33.4	2	11.5
3	8.1	18	2.7	3	5.0	3	5.0	3	5.0	3	5.0
4	2.4	4	0.6	4	4.2	4	1.1	4	4.2	4	1.1
5	3.0	4	0.6	5	3.0	4	0.6	5	3.0	4	0.6
6	2.7	3	0.5	6	2.7	3	0.5	6	2.7	3	0.5
7	1.1	1	0.2	7	4.5	5	0.8	7	4.5	5	0.8
8	10.9	9	1.4	8	7.4	7	1.1	8	14.3	10	2.2
9	2.7	2	0.3	9	1.5	1	0.2	9	1.5	1	0.2
10	6.0	4	0.6	10	3.3	2	0.3	10	6.0	4	0.6
11	1.7	1	0.2	11	1.5	1	0.2	11	1.7	1	0.2
12	1.8	1	0.2	12	1.8	1	0.2	12	1.8	1	0.2
13	2.7	1	0.2	13	2.7	1	0.2	13	2.7	1	0.2
14	1.8	1	0.2	14	1.8	1	0.2	14	1.8	1	0.2
15	4.2	1	0.2	15	4.2	1	0.2	15	4.2	1	0.2
16	7.7	1	0.2	16	7.7	1	0.2	16	7.7	1	0.2
17	281	1	0.2	17	281	1	0.2	17	281	1	0.2
18	518	1	0.2	18	518	1	0.2	18	518	1	0.2

An increase in discrimination was observed when comparing the number of types assayed with the Linear Array (A) compared to sequence data generated from HV1 (B), HV1/HV2 (C), and the entire control region (D). Using the Linear Array (A) resolved the population samples into 281 different haplotypes. About 28% of the samples were observed once (185 unique individuals among the 666 population samples), with one common haplotype shared among 51 individuals (7.7% of the individuals tested). Sequence information from HV1 (B), HV1/HV2 (C), and the entire control region (D) gave an increasing number of haplotypes and a decreasing number of common types (see also the schematic representation to the right).

Analysis of Null Alleles Comparing Sequence Information

A	HV1 Array Locus (CRS position)	African American	Hispanic	Caucasian	Haplogroup Associated Polymorphisms	Percentage of Null Alleles from Haplogroup Polymorphisms
HV1A (16126; 16124)	24	2	7	16124C - L3b and L3d	27/33 (82%)	
HV1C (16304; 16309; 16311)	28	38	10	16320T - L3e2 16290T and 16319A - A (Asian)	69/76 (91%)	
HV1D (16362)	29	3	2	16360T - L1c	29/34 (85%)	
HV1E (16270; 16278)	40	11	8	16270T and 16277T - L1b 16264T - L3e4 16265T - L3e3 16265G - L1c2	33/39 (85%)	
	129	62	34	225	158/225 (70%)	
B	HV1 Array Locus (CRS position)	African American	Hispanic	Caucasian	Haplogroup Associated Polymorphisms	Percentage of Null Alleles from Haplogroup Polymorphisms
HV1A (73)	1	0	2	721C - p1ev	1/3 (33%)	
HV1B (146; 150; 152)	46	36	15	151C - L1c 1430A - L2a 153A - G - A2, X*	88/97 (91%)	
HV1C (188; 195; 198; 200)	66	19	37	186 C - L1c 189A - L1b 189A - L2bc 185G - A, 189A - G, and 200A - G - L3e* 194C - T - D194b2 199T - C and 204T - C - I	78/122 (64%)	
HV1D (247)	5	27	7	248A - del - CZ 242C - J1* 250T - C - L2*	32/39 (82%)	
HV1E (189)	103	16	32	182C - T and 195T - C - L1L12* 185G - T - L1b 195T - C and 196C - T - L2a* 189A - G - L2bc 185G - A and 200A - G - L3e* 185G - A - J1* 194C - T - D194b2	123/153 (80%)	
	221	100	93	414	322/414 (78%)	

Null alleles (or "blanks") are produced by a failure of PCR product to hybridize to the immobilized SSO probe of the Linear Array. We compared the Linear Array mitotypes for 666 samples (Kline et al. 2005) to the underlying polymorphisms from sequence data (A, HV1; B, HV2) to characterize the reasons for null alleles. We determined that a considerably large number of null alleles produced by the Linear Array were caused by the existence of haplogroup-associated polymorphisms (70% in HV1 and 78% in HV2). Polymorphisms noted with an asterisk indicate instances where most, but not necessarily all, individuals within the haplogroup carry the polymorphism.



Schematic representation of the most common haplotype observed at each level of discrimination. (A), Fifty-one individuals shared the most common haplotype (MCH) using the Linear Array (Kline et al. 2005). (B) HV1 sequence data produced 40 unique individuals having the MCH. (C) HV2 Sequence information produced 13 types with individuals still unresolved. The sequencing strategy produced 4 types with 21 individuals still unresolved. Additional results were observed with a set of 54 haplogroup H sequences, indicating the overall usefulness of SNPs discovered by Coble et al. (2004) even among sequences outside of the most common types (compared to a random sequencing strategy).

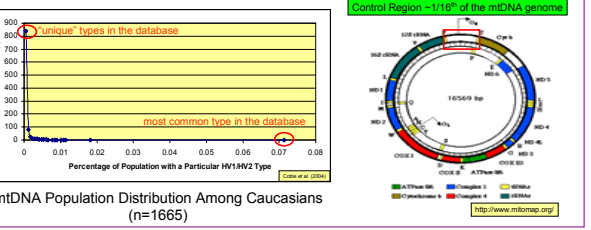
Acknowledgments
*Funding from the National Institute of Justice through the NIST Office of Law Enforcement Standards
*Margaret Kline, Jan Redman, Peter Valtos, and John Butler (NIST) for assistance in Linear Array analysis
*Jodi Irwin, Rebecca Just, Toni Diegoli, Bronn Smith, Tom Parsons, and members of the AFDIL Research Section for the high quality sequence data used in this study.

An Evaluation of Coding Region mtDNA Sequence Data to Increase Forensic Discrimination – Effective Strategies

The Problem: One major disadvantage of mtDNA testing lies with the low power of discrimination for individuals that share one of the few common types. In Caucasians, the most common haplotype (MCH) occurs in ~7% of the population (H1 – A263G; 315.1C).

A Possible Solution: Forensic Scientists have started to examine variation in the mtDNA coding region to increase the discrimination among the common mtDNA haplotypes. Coble et al. (2004) and Vallone et al. (2004) proposed a targeted SNP approach to increase discrimination among common types. This approach considers the possibility that a limited amount of highly degraded template may only facilitate a small number of additional mtDNA tests.

Objective: Recently, Budowle et al. (2005) offered several criticisms to this approach (below). To evaluate these criticisms, we examined the performance of selected SNPs to a published method that sequences short mtDNA coding region fragments for increasing mtDNA discrimination (Coble et al. 2006).



Criticisms of Budowle et al. (2005) using mtDNA coding region SNPs rather than sequencing short fragments of "informative" regions.

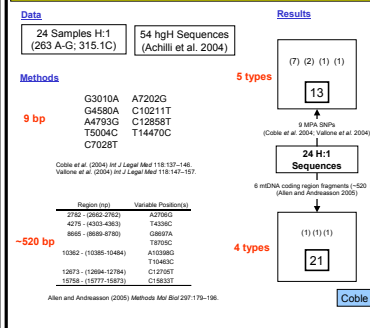
[Coble and Vallone] have proposed that forensic analyses of the coding region [should] be restricted to synonymous substitutions [and] suggest that sequencing strategies for forensic analyses of the coding region of the mtDNA genome should be avoided [and] that only SNP-based systems should be employed."

"We disagree with this proposition as applying such a strict criterion is *not well thoughtout* and would severely hamper the use of mtDNA in forensic testing."

"By limiting the analysis only to synonymous polymorphisms that cannot have any phenotypic effect, a large part of the polymorphic positions (and thus forensically informative) would be excluded."

"Even if only synonymous changes are being used, since the mtDNA effectively is one locus, such polymorphisms would still be linked to functional polymorphisms in other regions of the molecule."

An evaluation of the performance of selected SNPs compared to the sequencing of short fragments of "highly variable" mtDNA coding region DNA for increased discrimination



We sequenced 17 MCH Samples from the NIST population samples and included 7 MCH samples from Achilli et al. (2004). We also analyzed the haplogroup H sequences from Achilli et al. (2004) to determine how well the 30 haplogroup H SNPs discovered by Coble et al. (2004) performed among individuals not selected for the MCH.

We found that the 9 coding region SNPs performed better as discrimination than sequencing ~520 additional bp in the coding region. Of the 24 individuals sharing the MCH, the 9 SNPs produced 5 types with 15 individuals still unresolved. The sequencing strategy produced 4 types with 21 individuals still unresolved. Similar results were observed with a set of 54 haplogroup H sequences, indicating the overall usefulness of SNPs discovered by Coble et al. (2004) even among sequences outside of the most common types (compared to a random sequencing strategy).

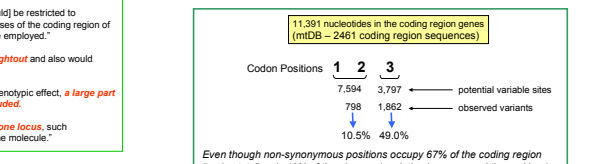
Why Did 9 SNPs Outperform Sequencing?

Region (np)	Variable Position*	mtDNA Frequency	Haplogroup Association	Comments
2782 - (2662-2762)	A2706G	1631/2063 (79.1%)	Diagnostic H1	Uninformative
4275 - (4300-4363)	T4336C	24/2064 (1.2%)	"H1S"	Redundant for C456T
8665 - (8688-8788)	G8697A	109/2064 (5.3%)	Diagnostic T2	Uninformative
	T8705C	15/2064 (0.7%)	Diagnostic X2	Uninformative
10362 - (10385-10484)	A10398G	890/2063 (43.1%)	Superhaplogroup N	Uninformative
	T10206C	110/2064 (5.3%)	Diagnostic S1	Uninformative
12673 - (12694-12784)	C12705T	902/2063 (43.7%)	Diagnostic Superhaplogroup R	Uninformative
15759 - (15777-15833)	C15833T	20/2064 (1.0%)	"H1S"	Uninformative

*As determined by dispersion order (P=3) from Table 4 of Allen and Andersson (2005)
*Dispersion order for this fragment was not listed in Table 3 of Allen and Andersson (2005), 100 bp were used as default
*Dispersion order in Table 3 of Allen and Andersson (2005) did not cover the 4336 polymorphism, 89 bp were used as default

We found that the 6 "highly variable" coding regions identified by Allen and Andersson (2005) contained several SNPs that, though observed at high frequencies, were uninformative for resolving individuals sharing the most common haplotype for haplogroup H. For example, 12705 C (CRS) variant is shared among all Caucasians, and would offer no additional information from the sequence data generated by HV1/HV2 alone.

What is the Effect of Avoiding Non-Synonymous Polymorphisms?



Even though non-synonymous positions occupy 67% of the coding region 'real estate', only 10% of the observed variation has occurred there. Nearly half of all potential synonymous positions in a global database of 2461 coding region sequences have been observed, making synonymous positions a better repository of genetic variation.

Using the information from mDB (<http://www.genpat.uu.se/mDB/>), we found that the overall percentage of variants observed at non-synonymous 1st and 2nd codon positions (among 2461 coding region sequences) to be significantly less than synonymous 3rd position codons – despite the fact that synonymous codons are outnumbered 2:1. Of the 198 synonymous polymorphisms occurring at a frequency of 1% or greater in the global database, all could be readily assigned to a haplogroup.

A mtDNA genome from an African-derived sequence (Hausa) belonging to haplogroup L0a1 (Ingman et al. 2000)

Region	Variable Position*	mtDNA Frequency	Haplogroup Association	Comments
750	G750A	5/60	A	10688 A
769	T769A	5/60	T	10810 T
825	A825G	1/60	C	10813 C
858	G858A	1/60	G	10815 A
868	T868C	1/60	T	10816 C
872	G872A	1/60	G	10817 A
873	A873G	1/60	A	10818 G
874	G874A	1/60	G	10819 A
875	T875C	1/60	T	10820 C
876	A876G	1/60	A	10821 G
877	G877A	1/60	G	10822 A
878	T878C	1/60	T	10823 C
879	A879G	1/60	A	10824 G
880	G880A	1/60	G	10825 A
881	T881C	1/60	T	10826 C
882	A882G	1/60	A	10827 G
883	G883A	1/60	G	10828 A
884	T884C	1/60	T	10829 C
885	A885G	1/60	A	10830 G
886	G886A	1/60	G	10831 A
887	T887C	1/60	T	10832 C
888	A888G	1/60	A	10833 G
889	G889A	1/60	G	10834 A
890	T890C	1/60	T	10835 C
891	A891G	1/60	A	10836 G
892	G892A	1/60	G	10837 A
893	T893C	1/60	T	10838 C
894	A894G	1/60	A	10839 G
895	G895A	1/60	G	10840 A
896	T896C	1/60	T	10841 C
897	A897G	1/60	A	10842 G
898	G898A	1/60	G	10843 A
899	T899C	1/60	T	10844 C
900	A900G	1/60	A	10845 G
901	G901A	1/60	G	10846 A
902	T902C	1/60	T	10847 C
903	A903G	1/60	A	10848 G
904	G904A	1/60	G	10849 A
905	T905C	1/60	T	10850 C
906	A906G	1/60	A	10851 G
907	G907A	1/60	G	10852 A
908	T908C	1/60	T	10853 C
909	A909G	1/60	A	10854 G
910	G910A	1/60	G	10855 A
911	T911C	1/60	T	10856 C
912	A912G	1/60	A	10857 G
913	G913A	1/60	G	10858 A
914	T914C	1/60	T	10859 C
915	A915G	1/60	A	10860 G
916	G916A	1/60	G	10861 A
917	T917C	1/60	T	10862 C
918	A918G	1/60	A	10863 G
919	G919A	1/60	G	10864 A
920	T920C	1/60	T	10865 C
921	A921G	1/60	A	10866 G
922	G922A	1/60	G	10867 A
923	T923C	1/60	T	10868 C
924	A924G	1/60	A	10869 G
925	G925A	1/60	G	10870 A
926	T926C	1/60	T	10871 C
927	A927G	1/60	A	10872 G
928	G928A	1/60	G	10873 A
929	T929C	1/60	T	10874 C
930	A930G	1/60	A	10875 G
931	G931A	1/60	G	10876 A
932	T932C	1/60	T	10877 C
933	A933G	1/60	A	10878 G
934	G934A	1/60	G	10879 A
935	T935C	1/60	T	10880 C
936	A936G	1/60	A	10881 G
937	G937A	1/60	G	10882 A
938	T938C	1/60	T	10883 C
939	A939G	1/60	A	10884 G
940	G940A	1/60	G	10885 A
941	T941C	1/60	T	10886 C
942	A942G	1/60	A	10887 G