

Kristen M. Webb,<sup>1</sup> Ph.D. and Marc W. Allard,<sup>2</sup> Ph.D.

## Identification of Forensically Informative SNPs in the Domestic Dog Mitochondrial Control Region\*

**ABSTRACT:** Dog hair is often found at crime scenes either due to the dog's involvement in the crime or secondary transfer. As little nuclear DNA is present in shed hair, a 1000 base pair fragment of the mitochondrial control region (mtCR) from 552 dogs was assessed for forensically useful sequence variation. Through pairwise alignment to a standard reference sequence, existing haplotypes were further described and 36 new haplotypes and 24 new single nucleotide polymorphisms were identified. The probability of exclusion was found to be 0.957. Breeds were found to have similar sequences, although not identical. No genetic basis was found for grouping dogs by either purebred or mixed or geographic location within the continental United States. Our research demonstrates that the domestic dog mtCR has not been thoroughly surveyed for sequence variation and that a single database comprised of purebred and mixed breed dogs is sufficient for the continental United States.

**KEYWORDS:** forensic science, mitochondrial genome, domestic dog, control region, single nucleotide polymorphism, haplotype

A 2005–2006 survey found that there were *c.* 73 million domestic dogs living in the United States (<http://www.americanpetproducts.org/newsletter/may2005/npos.html>) or one dog for every four people in the country. As demonstrated by several cases, not only is dog hair collected as evidence when a dog is directly involved in a crime (1), dog hair, and other types of canine evidence are frequently found at crime scenes as secondary transfer from the criminal(s) or victim(s) because of the common environments of humans and dogs (2) (State of California vs. David Westerfield, 2002 and State of Iowa vs. Andrew Rich, 2002). Microscopic analyses of hair rarely tell more than species, as hair can vary both between individuals of the same species as well as within an individual (3). There is little or no nuclear DNA in the hair shaft, often leaving mitochondrial DNA (mtDNA) as the only source of DNA that can be recovered from the hair shafts of telogen hairs (4–6). mtDNA from human hair evidence has been used in the United States courts since the case of Tennessee versus Paul William Ware in 1996. The procedures for isolating, analyzing, and presenting human mtDNA data that satisfy the admissibility requirements for scientific or technical evidence are in place and have been accepted by the legal and forensic communities (7).

Mitochondria are organelles that play a role in the body's energy production and are found in numbers as high as 100 per cell with as many as 10 genome copies per mitochondrion (8,9). The high

number of mitochondrial genomes (mtGenomes) per cell is useful for forensic analyses, particularly where the amounts of DNA are small or degraded (10). Additional forensic utility comes from the mtGenome being maternally inherited and not undergoing recombination. The canine mitochondrial genome is a circular, slightly less than 16,800 base pair (bp) genome that codes for 13 protein coding genes, 22 tRNAs, and 2 rRNAs (11). Different regions of the mtGenome accumulate mutations more readily than others. In humans, as well as other mammals, the control region (also known as D-loop or hypervariable region) has the highest mutation rate (12,13), making it a popular region of analysis to search for DNA variation. Relative to humans, dogs have an additional 10 bp variable tandem repeat that is repeated up to 38 times within the control region (14). The number of repeats is known to vary within an individual and between individuals (14,15).

It is well known that other forensic studies have investigated the potential uses of canine mtDNA as evidence and that private databases of canine mtDNA variation exist (2,14,16–19). We plan to use DNA sequencing and analysis to further categorize canine mtDNA haplotypes and develop the first public reference database of canine mtDNA single nucleotide polymorphisms (SNPs) from the control region of the canine mitochondrial genome.

### Materials and Methods

Domestic dog blood, tissue, and buccal swab samples were collected as donations from veterinary practices and private donors across the continental United States. Blood and tissue samples were not collected solely for this study. The collected samples were those that otherwise would have been disposed of. The donor of the sample made the determination of breed type and whether a dog was purebred or mixed. The donor was also asked to indicate any known relationship of a particular sample to other samples donated to this study. As this study focused only on mtDNA, and mtDNA is inherited maternally, siblings would have identical mtDNA sequences. Unrecognized familial relationships could lead to a misinterpretation of individuals of the same breed being

<sup>1</sup>Animal Parasitic Diseases Laboratory, Agricultural Research Service, United States Department of Agriculture Building 1180, Beltsville, MD 20705.

<sup>2</sup>Molecular Methods and Subtyping Branch, Division of Microbiology, Office of Regulatory Science, Center for Food Safety and Applied Nutrition, US Food and Drug Administration, College Park, MD 20740-3835.

\*This work has been presented at The National Institute of Justice Conference 2007 and at The George Washington University Research and Discovery Day, 2007. Both instances were in poster form. A portion of K. M. Webb's support came from a Selective Excellence grant from The George Washington University. This work was supported by the National Institute of Justice through grant 2004-DN-BX-K025 to M. W. Allard.

Received 26 Jan. 2008; and in revised form 16 April 2008; accepted 20 April 2008.

thought to have the same mtDNA and affect estimates of nucleotide diversity. A subset of the blood and tissue samples collected was used for sequencing and analysis.

All blood and tissue samples were stored at  $-20^{\circ}\text{C}$  until needed. Approximately 1 g of tissue was isolated and placed in a culture tube with  $0.1\times$  TAE (Tris-acetate-EDTA [ethylenediamine tetraacetic acid]) for preservation. Each tissue was ground into a single-cell slurry using a Janke and Kunkel Ultra Turrax T25 tissue grinder (Janke and Kunkel, Staufen, Germany). Total genomic DNA was extracted from the blood and tissue samples using the Invitrogen DNA Easy kit following the protocols for "Small Blood Samples and Hair Follicles" or "Small Amounts of Cells, Tissues or Plant Leaves" (Invitrogen Corporation, Carlsbad, CA). Following extraction, DNA samples were stored in  $0.1\times$  TE (Tris-EDTA). DNA was quantified using the Nanodrop ND-1000 spectrophotometer (Nanodrop Technologies, Wilmington, DE).

The oligonucleotide primers used in this study were taken from a recently published previous study (14). PCR primers were redesigned relative to the previous study because those in the previous study yielded double-banded products for some of our samples. The new primers flanked the entire mitochondrial control region (mtCR) and sat slightly further outside of the mtCR than the original primers. The new primers were defined as R51 (5'-TATGTTTATGGAGTCGTGCGA-3') and F15406 (5'-TTTGC-TCCACCATCAGCACC-3'). The previously designed sequencing primers were used for DNA sequencing in this study. As a set, the mtCR primers resulted in bidirectional, overlapping, high quality, 4–6 $\times$  sequence coverage across the entire mtCR excluding the tandem repeat region (Fig. 1). This repeat region is found in both dogs and wolves and is known to vary within and among individuals and thus was not sequenced for the current study (15).

All primers were received lyophilized from Operon Biotechnologies, Inc. (Huntsville, AL) and were resuspended to a concentration of 160 mM in  $0.1\times$  TE. The mtCR was amplified with one primer pair designed to span the entire region. PCR amplifications were performed in 50  $\mu\text{L}$  reactions using 100 ng total DNA,  $1\times$  Buffer (Fisher BioReagents, Fisher Scientific, Pittsburgh, PA), 5 mM  $\text{MgCl}_2$  (Fisher BioReagents, Fisher Scientific), 0.4 mM dNTP mix (Invitrogen Corporation), 0.1  $\mu\text{M}$  of each primer, and 2.5 units of Taq polymerase (Fisher BioReagents, Fisher Scientific). The PCR amplification profile on the thermal cycler consisted of an initial denaturing step of  $96^{\circ}\text{C}$  for 10 min, 39 cycles of amplification which included denaturing at  $94^{\circ}\text{C}$  for 15 sec, annealing at  $56^{\circ}\text{C}$  for 30 sec, extension at  $72^{\circ}\text{C}$  for 1 min, and a final extension at  $72^{\circ}\text{C}$  for 7 min. PCR products were run on a 1% agarose gel at 70 V for 1 h. A 1 kb ladder was used to determine size and a low-mass ladder to determine concentration of each product. Samples were diluted to 10  $\mu\text{L}$  reactions with a concentration of 40–

60 ng/ $\mu\text{L}$  of DNA and cleaned using 2  $\mu\text{L}$  of ExoSAP-IT (USB Corporation, Cleveland, OH) according to the procedure recommended by the vendor. The ExoSAP-IT procedure includes a  $37^{\circ}\text{C}$  incubation step for 15 min followed by an inactivation step at  $80^{\circ}\text{C}$  for 15 min. Samples were then shipped on dry ice overnight to SeqWright DNA Technology Services in Houston, Texas. SeqWright (<http://www.seqwright.com>) completed all DNA sequencing according to their protocols using ABI technology. (Applied Biosystems, Carlsbad, CA)

Representative sequences of previously described haplotypes were downloaded from Genbank (Accession #: AF531654 - AF531741 and AY656703 - AY656710) (16,20). Additionally, 125 domestic dog sequences collected from a previous study (14) were also included in the current dataset (Genbank Accession #: AY240030 - AY240072, AY240074 - AY240093, AY240095 - AY240154, and AY240156 - AY240157).

The forensic version of Sequencher 4.1.4FB19 (Gene Codes Corporation, Ann Arbor, MI) was used to edit and align all mtCR sequences. This version of the software builds alignments according to the previously defined criteria for gap placement and priority for preference of sequence differences in forensic evaluations (21). All alignments were confirmed by eye. Standard International Union of Biochemistry codes were used for polymorphic sites and N's were inserted for positions in which the base could not be determined. As with human forensic studies, a reference control region sequence was used. This has been previously recommended in an effort to standardize canine mitochondrial nucleotide nomenclature (22). As per Pereira's recommendations (22), the reference control region sequence used was the first canine mitochondrial genome to be published (11). Using a reference sequence allows base coordinates to be compared across different studies (22), thus all coordinates mentioned in this research are in terms of the reference sequence.

Within the complete mtGenome, the mtCR begins at position 15,458 and ends at position 16,727. The region spanning from nucleotide 16,663 to 16,676 was removed from the multiple alignment due to sequencing and alignment issues stemming from a polymorphic C/T stretch. However, this region was considered when defining unique haplotypes (see below). The tandem repeat region stretching from 16,130 to 16,430 bps and comprising of a varying number of 10 bp fragments was not sequenced due to variation within an individual (15). To account for this missing region, all mtCR sequences were divided into the region 5' of the repeat (15,458–16,129 bp) and the region 3' of the repeat (16,430–16,727 bp) with respect to the published light strand of the mtCR (Fig. 1). Two multiple alignments of all downloaded and newly sequenced mtCRs were created, one for each region on either side of the repeat. In Winclada (23), the "new matrix merge" command was used to combine the separate alignments based on matching identical taxon names. Arlequin ver 3.11 (24) was then used to search within the dataset for groups of dogs with identical control region sequences or haplotypes, and to calculate the frequency of these haplotypes. Those sequences that were identical to at least one other sequence in the dataset were also identified in Winclada and removed by using the "mark identical taxa" and "delete selected terms" commands. The mtCR sequences, excluding the tandem repeat but including 16,663–16,676 bps of individuals representing unique haplotypes were then aligned to the reference sequence using Sequencher, and the coordinates and base calls of the SNPs were recorded in an Excel spreadsheet. As the 13 bases between 16,663 and 16,676 were not used in the multiple alignment, manual checks were carried out to ensure the uniqueness of all haplotypes. As the majority of the previously published haplotype definitions lacked sequence in the region downstream of the

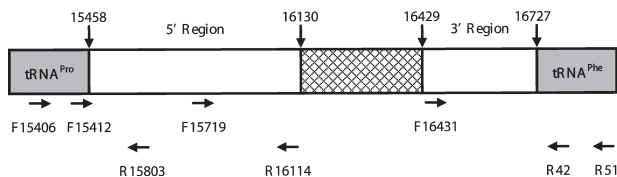


FIG. 1—Canine mitochondrial control region primers. Coordinates and orientation of all canine mitochondrial control region primers. Start and end coordinates of control region are shown as well as coordinates of the unsequenced repeat region, indicated by the checkered box. All coordinates are presented relative to the reference sequence (11). Primers F15406 and R51 were used for PCR amplification of the control region. All eight primers were used to obtain 4–6 $\times$  sequence coverage. All primers except F15406 and R51, which were designed by K.M.W., were designed by Gundry et al. (14).

tandem repeat, our new sequences covered more of the mtCR than many earlier studies (19,20,25,26). New sequences are grouped with previously defined haplotypes based on the SNPs present in the previously sequenced 5' region. If an individual is identical to a previously defined haplotype in the 5' region and new SNPs are found in the 3' region not previously sequenced, then the additional SNPs become the definition of a new haplotype. Also a new haplotype is defined as an individual possessing a unique set of SNPs relative to the reference sequence that do not match the 5' region or complete mtCR sequence of any of the previously published sequences.

In order to determine the relationship relative to previously defined haplotypes of those new mtCR sequences that did not fall within a previously published group, the matrix was transposed from DNA to numeric characters (A = 0, C = 1, G = 2, T = 3). The number "4" was inserted manually to replace any missing data that was truly a gap between the query sequence and the Kim et al. (11) reference sequence so that gaps would be considered as potential informative sites by Winclada as opposed to "missing data" or those regions without base calls due to unobtainable sequence. Winclada was then used to assess the relationships of the different dogs by constructing a phylogenetic tree using a parsimony ratchet search method on the entire dataset. Recommended search strategies for using the parsimony ratchet for large data matrices were followed (27). If multiple equally likely trees were obtained, they were combined to make a consensus tree and the placements of individuals with new mtCR haplotypes were assessed relative to the previously published haplotypes.

Winclada was also used to identify informative SNPs, defined as a nucleotide(s) that supports a group of two or more individuals. This was carried out by using the "mop informative characters/delete selected characters" function and then using the character diagnoser to trace each informative SNP on the tree. The length and retention index (ri) statistics were recorded for each informative SNP. The length is the number of times the nucleotide state at a given position changes on the tree. The ri is a measure of a nucleotide position having the same base in two individuals being the result of shared common ancestry and not convergence. The ri score can range from 100 to 0, with a score of 100 being obtained when the character change arose only once and defines all members of a group. The scores get progressively lower until a score of 0 is reached indicating that all character changes arose independently on that particular reconstruction.

Single nucleotide polymorphisms were classified into three ranks; the first rank was simply the presence of SNP at a nucleotide position. The second rank was assigned to characters found to be phylogenetically informative by Winclada based on character length and ri. The third level of ranking contains informative SNPs that define groups of six or more individuals, or 1% of the total dogs in the dataset. SNPs were evaluated for the third level rank by carefully inspecting the resultant most parsimonious trees.

All population statistics were either calculated in Arlequin or by hand. The dataset was analyzed as a whole with each individual defined as a unique haplotype (ignoring identical taxa). The dataset was also analyzed by separating dogs by their purebred or mixed description to look for suspected evidence of inbreeding in purebred individuals and evaluate if the "purebred" and "mixed" characterizations actually represent two unique populations. The samples were also separated by large regional groupings to look for local substructure and by those breed groups with a high number of purebred individuals ( $n > 6$ ) to look for within-breed structure. Genetic variance was investigated using Analysis of Molecular Variance (AMOVA) with 1023 permutations to assess

the significance of the variation among the various subdivisions of the dataset. The mean number of pairwise differences, nucleotide diversity, and assessment of variation within and between each grouping were calculated through Arlequin. Additional statistics such as the probability of exclusion or  $1 - \sum X_i^2$ , and random match probability or  $\sum X_i^2$  (where  $X_i$  is the frequency of the  $i$ th haplotype) were calculated by hand following the grouping of individuals with identical sequences into haplotypes and removal of those haplotypes missing sequence data.

## Results

Six hundred and ninety-eight domestic dog blood, tissue, and buccal swab samples were collected from various veterinary practices and private donors across the continental United States. As donors were cautioned against sending samples from related individuals, we only received one notification of siblingship between two samples collected. Of the 698 samples collected, 427 blood and tissue samples were used for mtCR sequencing and the results are available on Genbank (Accession #: EU223385 to EU223811). The distribution of these samples across the continental United States was as follows: California = 189, Maryland = 1, Mississippi = 8, New York = 1, Pennsylvania = 100, Nevada = 52, Texas = 14, Vermont = 1, and Virginia = 61. Three hundred and ten of these samples came from purebred individuals and the remaining 116 were mixed breed with individuals of unknown breed type considered mixed. The 427 newly collected samples were combined with the 125 purebred dogs from a previous study (14) for a final dataset of 552 domestic dogs. A complete list of the different breeds and number of each included in this study can be found in Table 1.

The complete mtCR excluding the tandem repeat was sequenced for 417 of the 427 newly collected individuals. The 10 individuals that did not have complete sequence were missing bases immediately after the repeat. The heteroplasmy of the repeat region caused the resultant sequence after this area to be unreadable due to the varying number of repeat units within the same dog. The missing bases in these sequences were coded as missing data and were not considered when looking for SNPs or haplotypes.

Previously defined haplotypes ( $n = 180$ ) were downloaded from Genbank as we planned to continue using the established nomenclature (16,20). Haplotype A15 could not be downloaded from Genbank as it was not found with the other sequences from the publication. Haplotypes labeled A37, A74, A75, A76, A77, A78, and A79 do not appear to exist in previously published datasets.

The sizes of the newly sequenced complete mtCRs ranged from 965 to 975 bp, excluding the tandem repeat. The final dataset of the newly sequenced mtCRs and those from the three previous studies (14,16,20) consisted of 733 taxa including the reference sequence (11). Following the alignments of each unique haplotype to the reference sequence separately, the size of the total matrix was 985 characters. Sixteen deletions were identified when haplotypes were aligned to the Kim et al. (11) reference sequence: 15464.1, 15539.1, 15546.1, 16129.1, 16507.1, 16542.1, 16562.1, 16663.1, 16663.2, 16671.1, 16671.2, 16671.3, 16673.1, 16674.1, 16711.1, 16711.2.

The search for individuals with identical mtCR sequences resulted in 311 unique haplotypes from the starting dataset of 733 domestic dog control region sequences. Tree searches of the unique sequences only resulted in 508 equally parsimonious trees. This means that there were 508 equally likely resolutions of the relationships of the 311 dogs using the control region data and the parsimony ratchet method of grouping. A single consensus tree was made from all resultant trees as a way to summarize nonconflicting

TABLE 1—*Breed list.*

Breed	Purebred	Mixed
Airedale	3	
Airedale Terrier	1	
Akita	2	
Alaskan Husky	1	
Alaskan Malamute	1	
American Cocker Spaniel	1	
American Eskimo dog	2	1
American Spitz	1	
American Staffordshire	1	
Anatolian Shepherd	2	
Australian Shepherd	6	3
Australian Terrier	1	
Basset	1	
Basset Hound	8	
Beagle/Corgi		1
Beagle/Labrador		1
Beagle	5	4
Bearded Collie	1	
Belgian Sheepdog	1	
Bernese Mountain Dog	4	
Bichon Frise	5	4
Blood Hound	1	
Blue Heeler	2	
Bolognese	1	
Border Collie	7	4
Boston Terrier	7	
Boxer	5	1
Brittany Spaniel	2	1
Bulldog	3	
Bull Mastiff	4	
Bull Terrier	2	
Cairn Terrier	2	1
Cardigan Corgi	2	
Catahoula		1
Cavalier King Charles Spaniel	7	
Chesapeake Bay Retriever	3	
Chihuahua	5	9
Chocolate Labrador Retriever	6	1
Chow	1	1
Chow Chow	2	
Cockapoo		2
Cocker Spaniel/Poodle		1
Cocker Spaniel	7	1
Collie	2	1
Corgi	5	1
Coton De Tulear	3	
Cur	1	
Dachshund	8	
Dalmatian	3	1
Doberman	2	
Doberman Pinscher	5	1
Dogue de Bordeaux	1	
English Bulldog	2	
English Mastiff	3	
English Shepherd		1
English Springer Spaniel	2	
English Terrier	1	
Eskimo Dog	1	
Finnish Spitz	1	
Flat Coated Retriever	3	
Fox Terrier	1	1
French Bulldog	1	
German Shepherd	4	1
German Short Haired Pointer	2	
Golden Retriever/Poodle		1
Golden Retriever	39	
Great Dane	6	
Great Pyrenees	1	
Greyhound	1	
Havanese	5	
Hunting Dog	1	
Husky/Retriever		1

TABLE 1—*Continued.*

Breed	Purebred	Mixed
Husky/Shepherd		1
Husky	4	1
Italian Greyhound	1	
Jack Russell/Beagle		1
Jack Russell	7	2
Japanese Chin/Lhasa Apso		1
Keeshond	3	
Kerry Blue Terrier	1	
Labradoodle	3	4
Labrador/Border Collie		1
Labrador/Dane		1
Labrador	2	
Labrador Retriever	31	4
Leonberger	1	
LhasaApso	4	2
Maltese/Shih Tzu		1
Maltese	5	3
Maltipoo	1	
Manchester Terrier	2	
Maremma	2	
Mastiff	2	
Miniature Dachshund	2	
Miniature Pinscher	2	1
Miniature Poodle	4	
Miniature Schnauzer	7	
Mix		3
Munsterlander Pointer	1	
Neapolitan Mastiff	2	
Newfoundland	1	
Norwegian Elkhound	1	
Old English Sheepdog	4	
Papillon/Sheltie		1
Papillon	1	
Pembroke Corgi	1	
Pembroke Welsh Corgi	1	
Pharaoh Hound	1	
Pit Bull	2	
Pit Bull Terrier	5	3
Pointer	1	
Pomeranian	5	3
Poodle	8	6
Portuguese Water Dog	3	1
Pug/Jack Russell		1
Pug/Jug		1
Pug	6	
Rat Terrier	1	
Ridgeback	1	
Rottweiler/Saint Bernard		1
Rottweiler	7	
Rough Collie	1	
Saint Bernard	1	
Samoyed	1	
Sapsarsee*	1	
Schipperke	2	
Schnauzer/Poodle		1
Schnauzer	4	2
Scottish Terrier	2	
Shar Pei	3	
Shar Planinetz	2	
Sheltie	4	1
Shepherd/Chow		2
Shepherd/Labrador		1
Shepherd		8
Shetland Sheepdog	1	
Shiba Inu	6	
Shih Tzu/Lhasa Apso		2
Shih Tzu	5	2
Shilo Shepherd	1	
Siberian Husky	1	
Spitz		1
Springer Spaniel	1	
Stafford Bull Terrier	1	

TABLE 1—Continued.

Breed	Purebred	Mixed
Standard Poodle	1	
Swiss Mountain Dog	1	
Teacup Maltese	1	
Terrier		3
Tibetan Mastiff	1	
Tibetan Spaniel	1	
Tibetan Terrier	1	
Toy Chow	1	
Toy Fox Terrier	1	
Toy Poodle	6	
Unknown	2	1
Vizsla	3	
Walker Hound	1	
Weimaraner	3	
Welsh Corgi	1	
West Highland Terrier	7	
West Highland White Terrier	2	
Wheaton Terrier	2	
Whippet	1	
White Schnauzer		1
Wire-haired Dachshund	1	
Yorkie-Chihuahua		1
Yorkie-Poodle		2
Yorkshire Terrier	6	1

A complete list of all dogs used in the study. Each breed is listed followed by how many purebred and mixed breed members of each breed are contained in the dataset. The sample donor determined breed name and type. \*Sapsarsee is the dog used by Kim et al. (11) and thus the reference sequence.

groupings. These groupings as well as the spreadsheet of all individuals and the variable SNPs they possessed were used to identify haplotypes in the current dataset. Excluding those sequences/haplotypes from previously published studies (16,20) resulted in the identification of 143 unique haplotypes in our dataset of 552 sequences. These 143 haplotypes do not include those sequences that were not sequenced in their entirety.

Canine mtCR sequences had previously been grouped into six main types, A, B, C, D, E, and F (16,20). Review of the parsimony ratchet trees showed that all of the newly sequenced control regions fell within four of these main types, namely, A, B, C, and D.

Hereafter the number and percentage of individuals reported for a haplotype is based only on the 552 dog dataset and does not include previously published individuals. While the majority of the newly sequenced samples aligned with one of the published haplotypes, the additional sequencing of the region downstream of the repeat provided additional variation that allowed the description of a new haplotype. Additionally, 36 haplotypes were identified based on variation in the 5' region only or a combination of 5' and 3' region variation and 60 sequences had ambiguous base calls and could only be classified in terms of a major haplogroup. One of the sequences containing ambiguous base calls also contained missing data. A complete list of all haplotypes found in this study and a list of dogs belonging to each haplogroup are given in Tables 2 and 3.

Haplogroup A was the largest haplogroup in the previously published data and also the haplogroup in which more of the newly sequenced samples clustered relative to the other groups. Previously published studies reported 76 haplotypes within group A, which make up 61.8% of all previously published haplotypes (16,20). Three hundred and sixty-nine of the 552 individuals or 66.8% from the current study fall within haplogroup A. Most of these individuals aligned with one of the 24 previously identified haplotypes, namely, A1, A2, A5, A11, A16, A17, A18, A19, A20, A22, A24, A26, A27, A28, A29, A31, A33, A40, A66, A68, A70, A71, A80,

and A82. Many of these haplotypes were further described as new haplotypes as a result of downstream control region sequence. Expanding on the current naming scheme, these new haplotypes were defined by keeping the 5' haplotype name but adding a lower case letter beginning with "a." Additionally, 24 new A haplotypes were described from the current dataset. In keeping with the previous naming scheme, the new haplotypes are A84–A107. Excluding individuals missing sequence data that clustered into haplogroup A, 77 A haplotypes were found in the current dataset.

Haplogroup B was the second largest set both in terms of previously defined haplotypes and where newly sequenced individuals grouped. Previous studies reported 20 haplotypes within the B haplogroup. These 20 types make up 16.3% of all previously defined groups (16,20). In the current dataset, it was found that 139 or 25.2% of all individuals possessed B haplotypes. New individuals were found to contain 8 of the 20 previously defined haplotypes: B1, B3, B6, B8, B10, B11, B12, and B20. Nine new haplotypes were defined (B21–B29). Excluding individuals with missing sequence data that clustered into haplogroup B, 49 B haplotypes were found in the current dataset. The largest single grouping of individuals ( $n = 59$ ) with the same haplotype occurred in B1a (Table 2).

The third haplogroup described, haplogroup C, was represented by eight haplotypes in the previously published literature (16,20). Again, this distribution of only 6.5% of the total types previously published closely agrees with the distribution of individuals from the current dataset, 7.6% ( $n = 42$ ), grouping within haplogroup C. Of the eight haplotypes, five were represented in the current dataset: C1, C2, C3, C5, and C8. Additionally, three new haplotypes were described, C9, C10, and C11 (while the C11 sequence is unique, it was not sequenced completely and as such is not included in new haplotype counts). Excluding individuals with missing sequence data that clustered into haplogroup C, 16 unique C haplotypes were found in the current dataset.

Haplogroup D was represented by six (4.8%) haplotypes in the literature (16,20) while only one individual (0.2%) from the current dataset fell within haplogroup D possessing haplotype D1a. No individuals from the current dataset matched any types from haplogroups E or F from previously published studies. The number of total haplotypes identified in the current dataset is 143 with a haplotype distribution of A = 77, B = 49, C = 16, D = 1, E = 0, and F = 0.

As can be seen from the distribution chart of haplotypes in Fig. 2, the majority of the haplotypes, 93.7% ( $n = 134$ ), fall into groups 10 or fewer members. These smaller sets of individuals contain only 47.9% ( $n = 259$ ) of the 541 individuals in the dataset for which complete mtCR sequence was obtained. Ambiguous individuals make up 10.5% ( $n = 58$ ) of the total dataset. Fifty-three percent ( $n = 282$ ) of all individuals consist of a total of nine haplotypes shared between 11 and 59 other individuals in the dataset.

Of the 987 characters in the total mtCR dataset, 9.5% ( $n = 94$ ) were found to be SNPs, defined as a different nucleotide character state at a given position possessed by at least one individual relative to the reference sequence (Table 2). Excluding the problematic region between 16,663 and 16,676 bps, 5.6% ( $n = 54$ ) of the characters were found to be informative, meaning the SNP was present in two or more individuals (Table 4). Thirty-three of the 54 (61%) informative SNPs were found to be highly informative, meaning that they defined a group that contains 1% or more of the total dogs in the current dataset. Of the 94 SNPs identified in the current study, 24 had not been previously recognized as variable sites in the published literature (16,20,28) with 6 of these 24 sites found to be informative and 3 highly informative. Of the











TABLE 3—Distribution of haplotypes.

Haplotype	Breed Sample ID	(n) Per Breed	Total (n)	%			
A1a	American Eskimo Dog 988P	1	7	1.27			
	Belgian Sheepdog 968P	1					
	Border Collie 1M	1					
	Catahoula 1M	1					
	Doberman Pinscher 1P	2					
A2b	Rough Collie 1P	1	11	1.99			
	French Bulldog 1P	1					
	Great Dane 2P	5					
	Leonberger 1P	1					
	Saint Bernard 2P	1					
	Schnauzer 4P	1					
	Scottish Terrier 1P	2					
A2a	West Highland Terrier 1P	1	2	0.36			
	Pit Bull Terrier 3M	1					
A5a	Labrador Retriever 2110P	3	3	0.54			
A5b	Jack Russell 7P	1	7	1.27			
	Pug/Jack Russell 1M	1					
	Pug/Jug 1M	1					
A5c	Sheltie 1P	3	1	0.18			
	Shetland Sheepdog 1P	1					
A11f	Labrador Retriever 2148P	1	40	7.25			
	American Staffordshire 1P	1					
	Anatolian Shepherd 532P	4					
	Australian Shepherd 5M	1					
	Border Collie 13P	1					
	Border Collie 5M	2					
	Boston Terrier 2P	1					
	Boxer 8P	1					
	Bulldog 5P	1					
	Chihuahua 439P	1					
	Chihuahua 5M	1					
	Chocolate Labrador Retriever 6P	1					
	Chow Chow 232P	1					
	Cocker Spaniel 1M	1					
	Collie 1P	1					
	English Bulldog 1P	1					
	English Springer Spaniel 2P	1					
	Husky/Shepherd 1M	1					
	Husky 324P	1					
	Jack Russell 1P	2					
	Labrador Retriever 2141P	1					
	Labrador Retriever 2M	1					
	Miniature Dachshund 3P	1					
	Miniature Schnauzer 5P	1					
	Old English Sheepdog 1P	1					
	Pembroke Welsh Corgi 1P	1					
	Pit Bull Terrier 4P	1					
	Rottweiler 1P	2					
	Schnauzer 1M	1					
	Shepherd 3M	3					
	Shih Tzu 1P	1					
	Springer Spaniel 229P	1					
	Yorkshire Terrier 9P	1					
	A11a	Labrador Retriever 9M			1	8	1.45
		Manchester Terrier 1P			2		
		Rottweiler 333P			4		
	A11b	Rottweiler/St. Bernard 1M			1	3	0.54
		Chihuahua 6P			1		
	A11b?	Dachshund 1P			1	1	0.18
		Papillon 1P			1		
	A11c	Airedale 3P			1	1	0.18
	A11d	Terrier 1M			1	1	0.18
	A16a	Greyhound 290P			1	37	6.70
		Brittany Spaniel 1M			1		
		Chesapeake Bay Retriever 974P			2		
		Chocolate Labrador Retriever 1M			1		
		Chow 1M			1		
English Mastiff 2P		2					
Golden Retriever 15P		5					
Italian Greyhound 1P		1					
Labradoodle 2P		2					

TABLE 3—Continued.

Haplotype	Breed Sample ID	(n) Per Breed	Total (n)	%	
A16?	Labradoodle 3M	2	57	10.33	
	Labrador/Border Collie 1M	1			
	Labrador/Dane 1M	1			
	Labrador 994P	15			
	Yorkshire Terrier 1M	1			
	Labrador Retriever 20M	1			
	Labrador Retriever 2108P	1			
	A17a	Beagle 6P			1
		Bichon Frise 1P			2
	A17a?	Bichon Frise 2M			1
		Boston Terrier 4P			3
	A17a?	Boxer 292P			3
		Bull Mastiff 3P			2
	A17a?	Bull Terrier 2P			1
		Cavalier King Charles Spaniel 5P			4
A17a?	Chihuahua 3P	2			
	Chocolate Labrador Retriever 2P	3			
A17a?	Dalmatian 1M	1			
	Dalmatian 2P	1			
A17a?	Dogue de Bordeaux 1P	1			
	English Mastiff 1P	1			
A17a?	Flat Coated Retriever 1P	2			
	Great Dane 5P	1			
A17a?	Jack Russell 2M	1			
	Jack Russell 2P	2			
A17a?	Labrador 980P	2			
	Mastiff 2P	2			
A17a?	Miniature Dachshund 1P	1			
	Miniature Pinscher 3P	1			
A17a?	Pit Bull 227P	3			
	Pomeranian 5P	1			
A17a?	Pug 2P	3			
	Rottweiler 6P	1			
A17a?	Samoyed 1P	1			
	Shar Pei 3P	1			
A17a?	Shepherd/Labrador 1M	1			
	Shepherd 1M	2			
A17a?	Shiba Inu 4P	1			
	Stafford Bull Terrier 1P	1			
A17a?	Toy Fox Terrier 1P	1			
	Unknown 1P	1			
A17a?	Bull Mastiff 1P	1			
	Coton De Tulear 2P	1			
A17a?	Dalmatian 1P	1	1	0.18	
	Yorkshire Terrier 951P	1	1	0.18	
A17c	Pit Bull Terrier 4M	1	1	0.18	
	Bearded Collie 1P	1	44	7.97	
A17d	Chihuahua 10M	3			
	Cockapoo 3M	1			
A18d	Cocker Spaniel 196P	1			
	Dachshund 3P	2			
A18d	English Springer Spaniel 1P	1			
	Fox Terrier 1P	1			
A18d	German Shepherd 3M	1			
	Havanese 1P	4			
A18d	Husky 1M	1			
	Jack Russell 10P	2			
A18d	Lhasa Apso 3P	2			
	Lhasa Apso 3M	1			
A18d	Maltese 2P	1			
	Maltese 3M	2			
A18d	Old English Sheepdog 195P	2			
	Pomeranian 1M	1			
A18d	Poodle 287P	1			
	Pug 4P	3			
A18d	Sheltie 1M	1			
	Shepherd 6M	1			
A18d	Teacup Maltese 1P	1			
	Toy Chow 1P	1			

TABLE 3—Continued.

Haplotype	Breed Sample ID	(n) Per Breed	Total (n)	%
A18a	Toy Poodle 1P	3		
	Vizsla 2P	3		
	Weimaraner 4P	2		
	Whippet 318P	1		
	Miniature Schnauzer 6P	1	4	0.72
	Schnauzer 2P	2		
	White Schnauzer 1M	1		
A18b	American Cocker 1P	1	2	0.36
	Dachshund 5P	1		
A18c	Sheltie 4P	1	1	0.18
A19a	Australian Shepherd 1P	1	13	2.36
	Beagle/Corgi 1M	1		
	Beagle/Labrador 1M	1		
	Dachshund 4P	1		
	English Terrier 982P	1		
	German Shepherd 12P	3		
	German Short Haired Pointer 4P	1		
	Jack Russell/Beagle 1M	1		
	Mix 2M	1		
	Portuguese Water Dog 2P	1		
	Shilo Shepherd 1P	1		
A20c	Chihuahua 11M	1	6	1.09
	Coton De Tulear 1P	1		
	Maremma 533P	1		
	Papillon/Sheltie 1M	1		
	Pharaoh Hound 1P	1		
	Pointer 5P	1		
A20a	Miniature Poodle 2P	1	3	0.54
	Poodle 5P	2		
A20b	English Shepherd 1M	1	1	0.18
A22a	Bernese Mountain Dog 1P	4	7	1.27
	Bull Mastiff 5P	1		
	Neapolitan Mastiff 1P	2		
A24a	Brittany Spaniel 2P	2	3	0.54
	Ridgeback 1P	1		
A26a	Cairn Terrier 1M	1	8	1.45
	Cairn Terrier 4P	1		
	Cavalier King Charles Spaniel 3P	2		
	Newfoundland 1P	1		
	West Highland Terrier 4P	1		
	Wheaton Terrier 1P	2		
A27c	Bichon Frise 4M	1	5	0.91
	Keeshond 1P	3		
	Lhasa Apso 2M	1		
A27a	Corgi 4P	1	1	0.18
A27b	Pit Bull Terrier 3P	1	1	0.18
A28a	Cur 1P	1	2	0.36
	Hunting Dog 1P	1		
A29a	Husky/Retriever 1M	1	4	0.72
	Husky 1P	3		
A31a	Eskimo Dog 168P	1	1	0.18
A33c	Golden Retriever/Poodle 1M	1	16	2.90
	Golden Retriever 1685P	14		
	Labrador Retriever 2113P	1		
A33a	Golden Retriever 1730P	1	1	0.18
A33b	Golden Retriever 1692P	1	1	0.18
A40a	Swiss Mountain Dog 1P	1	1	0.18
A66	Cavalier King Charles Spaniel 4P	1	1	0.18
A68	Shiba Inu 2P	3	3	0.54
A70	Collie 2P	1	1	0.18
A71	Cardigan Corgi 1P	1	5	0.91
	Corgi 2P	2		
	Miniature Pinscher 2P	1		
	Pembroke Corgi 1P	1		
A71a	Akita 1P	1	1	0.18
A80a	Munsterlander Pointer 1P	1	2	0.36
	Yorkshire Terrier 10P	1		
A80b	Yorkshire Terrier 441P	1	1	0.18
A82a	German Shepherd 1P	1	2	0.36
	Terrier 2M	1		
A84*	Poodle 3P	2	2	0.36

TABLE 3—Continued.

Haplotype	Breed Sample ID	(n) Per Breed	Total (n)	%
A85*	Golden Retriever 1696AP	1	5	0.91
	Labrador Retriever 2127P	4		
A86*	Bichon Frise 1M	1	3	0.54
	Beagle 1P	1		
	Boxer 7P	1		
A87*	Miniature Schnauzer 2P	5	5	0.91
A88*	Cocker Spaniel 4P	1	2	0.36
	Shih Tzu 7P	1		
A89*	Maremma 393P	1	1	0.18
A90*	Alaskan Malamute 1P	1	1	0.18
A91*	Miniature Pinscher 1M	1	1	0.18
A92*	Bulldog 4P	1	1	0.18
A93*	Golden Retriever 1701P	1	1	0.18
A94*	Chow Chow 443P	1	1	0.18
A95*	Old English Sheepdog 3P	1	1	0.18
A96*	Beagle 426P	1	1	0.18
A97*	Tibetan Mastiff 1P	1	1	0.18
A98*	Chihuahua 5P	1	1	0.18
A99*	American Spitz 975P	1	1	0.18
A100*	American Eskimo Dog 957P	1	1	0.18
A101*	Mix 1M	1	1	0.18
A102*	Shepherd/Chow 2M	1	1	0.18
A103*	Shar Pei 2P	1	1	0.18
A104*	Finnish Spitz 1P	1	1	0.18
A105*	West Highland Terrier 3P	1	1	0.18
A106*	Alaskan Husky 169P	1	1	0.18
A107*	Doberman 296P	1	1	0.18
AAmbig1?	Akita 2P	1	1	n/a
AAmbig2	Shepherd 5M	1	1	n/a
AAmbig3	Pomeranian 4P	1	1	n/a
AAmbig4	Chihuahua 4M	1	1	n/a
AAmbig5	Beagle 4M	1	1	n/a
AAmbig6	Boxer 2M	1	1	n/a
A2Ambig1	Beagle 4P	1	1	n/a
A11Ambig1	Pit Bull 313P	1	1	n/a
A11Ambig2	Australian Shepherd 5P	1	2	n/a
	Cocker Spaniel 1P	1		n/a
A17Ambig1	Bull Terrier 1P	1	1	n/a
B1a	Airedale 1P	2	59	10.69
	Australian Shepherd 4M	1		
	Basset Hound 233P	5		
	Beagle 1M	1		
	Blue Heeler 2P	1		
	Bolognese 1P	1		
	Border Collie 2P	1		
	Bulldog 2P	1		
	Chihuahua 1M	1		
	Corgi 1M	1		
	Corgi 3P	1		
	Dachshund 326P	1		
	English Bulldog 234P	1		
	Fox Terrier 1M	1		
	German Short Haired Pointer 959P	1		
	Golden Retriever 13P	10		
	Great Pyrenees 1P	1		
	Kerry Blue Terrier 325P	1		
	Labradoodle 1M	1		
	Labradoodle 1P	1		
	Labrador Retriever 2114P	5		
	Lhasa Apso 913P	1		
	Miniature Poodle 1P	1		
	Poodle 6M	2		
	Poodle 9P	1		
	Schnauzer/Poodle 1M	1		
	Schnauzer 3P	1		
	Shar Pei 1P	1		
	Shih Tzu/Lhasa Apso 1M	1		
	Shih Tzu 2M	1		
	Shih Tzu 2P	2		
	Standard Poodle 293P	1		

TABLE 3—Continued.

Haplotype	Breed Sample ID	(n) Per Breed	Total (n)	%
	Terrier 3M	1		
	Tibetan Spaniel 1P	1		
	Tibetan Terrier 1P	1		
	Weimaraner 3P	1		
	Welsh Corgi 1P	1		
	West Highland Terrier 2P	2		
B1b	Beagle 2P	1	7	1.27
	Maltese/Shih Tzu 1M	1		
	Mix 3M	1		
	Poodle 10P	2		
	Rat Terrier 2P	1		
	Shepherd/Chow 1M	1		
B1c	Golden Retriever 1684P	1	1	0.18
B1d	Golden Retriever 1740P	1	1	0.18
B1e	Golden Retriever 2P	1	1	0.18
B1f	Golden Retriever 1699P	1	1	0.18
B3a	Maltipoo 1P	1	6	1.09
	Miniature Poodle 4P	1		
	Poodle 5M	1		
	Toy Poodle 4P	2		
	West Highland White Terrier 2P	1		
B6a	Schipperke 1P	1	2	0.36
	Walker Hound 1P	1		
B6b	Shepherd 7M	1	1	0.18
B8a	Flat Coated Retriever 3P	1	1	0.18
B10a	Cocker Spaniel 8P	1	1	0.18
B10b	Maltese 2M	1	1	0.18
B11a	Cocker Spaniel/Poodle 1M	1	3	0.54
	Dachshund 6P	1		
	Shih Tzu 10P	1		
B12a	Bichon Frise 4P	1	1	0.18
B20a	Portuguese Water Dog 1P	2	2	0.36
B21*	Cocker Spaniel 5P	1	3	0.54
	Labrador Retriever 3M	1		
	Yorkshire Terrier 1P	1		
B22*	Bichon Frise 2P	1	2	0.36
	Maltese 5P	1		
B23*	Maltese 3P	2	3	0.54
	Spitz 1M	1		
B24*	Carrin Terrier 442P	1	1	0.18
B25*	Golden Retriever 11P	1	1	0.18
B26*	Chesapeake Bay Retriever 428P	1	1	0.18
B27*	Unknown 289P	1	1	0.18
B28*	Cockapoo 1M	1	1	0.18
B29*	Japanese Chin/Lhasa Apso 1M	1	1	0.18
BAmbig1	Australian Terrier 1P	1	1	n/a
BAmbig2	American Eskimo Dog 1M	1	1	n/a
BAmbig3	Doberman Pinscher 3P	1	1	n/a
BAmbig4	Doberman Pinscher 4P	1	2	n/a
	Doberman Pinscher 5P	1		n/a
BAmbig5	Basset Hound 1P	1	1	n/a
BAmbig6	Basset Hound 7P	1	1	n/a
BAmbig7	Beagle 6M	1	2	n/a
	Boston Terrier 7P	1		n/a
BAmbig8	Bichon Frise 3M	1	2	n/a
	Bichon Frise 5P	1		n/a
BAmbig9	Chihuahua 2M	1	1	n/a
BAmbig10	Portuguese Water Dog 1M	1	1	n/a
BAmbig11	Unknown 1M	1	1	n/a
BAmbig12	Yorkie-Chihuahua 1M	1	1	n/a
BAmbig13	Schipperke 2P	1	1	n/a
BAmbig14	Wire-haired Dachshund 1P	1	1	n/a
BAmbig15	Jack Russell 1M	1	1	n/a
BAmbig16	Maltese 4P	1	1	n/a
B1Ambig1	Poodle 3M	1	2	n/a
	Poodle 4M	1		n/a
B1Ambig2	Australian Shepherd 4P	1	1	n/a
B1Ambig3	Chocolate Labrador Retriever 3P	1	3	n/a
	Coton De Tulear 3P	1		n/a
	Corgi 1P	1		n/a

TABLE 3—Continued.

Haplotype	Breed Sample ID	(n) Per Breed	Total (n)	%
B1Ambig4	Airedale Terrier 1P	1	7	n/a
	Basset Hound 2P	1		n/a
	Cardigan Corgi 2P	1		n/a
	Chocolate Labrador Retriever 1P	1		n/a
	Labradoodle 2M	1		n/a
	Shih Tzu 4M	1		n/a
	Yorkie-Poodle 4M	1		n/a
B1Ambig5	Basset Hound 5P	1	1	n/a
B1Ambig6	Blood Hound 1P	1	1	n/a
B1Ambig7	Golden Retriever 1P	1	1	n/a
B3Ambig1	Toy Poodle 5P	1	1	n/a
B11Ambig1	Dachshund 2P	1	1	n/a
C1a	Siberian Husky 167P	1	1	0.18
C2a	Dalmatian 993P	1	3	0.54
	West Highland Terrier 6P	2		
C2b	Boston Terrier 3P	1	4	0.72
	Chihuahua 7M	1		
	Lhasa Apso 2P	1		
	Yorkshire Terrier 3P	1		
C3a	Australian Shepherd 3M	1	12	2.17
	Border Collie 11P	4		
	Cocker Spaniel 958P	1		
	Havanese 3P	1		
	Pomeranian 1P	2		
	Pomeranian 2M	1		
	Poodle 2M	1		
	Shiba Inu 6P	1		
C5a	Anatolian Shepherd 531P	1	3	0.54
	Shar Planinetz 179BP	2		
C8a	Border Collie 8M	1	5	0.91
	Doberman 231P	1		
	Doberman Pinscher 2M	1		
	Pit Bull Terrier 1M	1		
	Pomeranian 2P	1		
C9*	Boston Terrier 1P	1	1	0.18
C10*	Pomeranian 4M	1	1	0.18
C11?*	Border Collie 177BP	1	1	0.18
CAmbig1	Blue Heeler 1P	1	1	n/a
CAmbig2	Collie 2M	1	1	n/a
CAmbig3	Chow 1P	1	1	n/a
CAmbig4	Pit Bull Terrier 10P	1	1	n/a
C2Ambig1	Beagle 5M	1	1	n/a
C2Ambig2	West Highland White Terrier 1P	1	1	n/a
C3Ambig1	Cocker Spaniel 3P	1	3	n/a
	Miniature Poodle 3P	1		n/a
	Schnauzer 2M	1		n/a
C3Ambig2	ShibaInu 5P	1	2	n/a
	Yorkie-Poodle 2M	1		n/a
D1a	Norwegian Elkhound 1P	1	1	0.18

mtCR, mitochondrial control region; n/a, not applicable.

The haplotype distribution of 552 domestic dogs relative to the Kim et al. (11) reference sequence. Haplotype name, a representative individual's breed sample ID, number of individuals per breed, number of individuals per haplotype, and frequency (%) of the haplotype observed in the dataset are provided. Asterisks (\*) indicate newly identified complete haplotypes. Question marks (?) indicate haplotypes with missing sequence data. The frequency of the haplotype calculations (%) does not include mtCR sequences with ambiguous sequence data (n/a). For a description of each haplotype refer to Table 2.

with 12 of these 30 SNPs occurring in this 60 bp region. In the current study, 22 SNPs were found in this "hotspot" region, 16 of which were informative. As with the previous study, more SNPs occurred in this 60 bp region than any other comparatively sized region of the mtCR.

Treating all newly collected sequences as a single population, the average pairwise nucleotide difference was  $12.49 \pm 5.65$  and

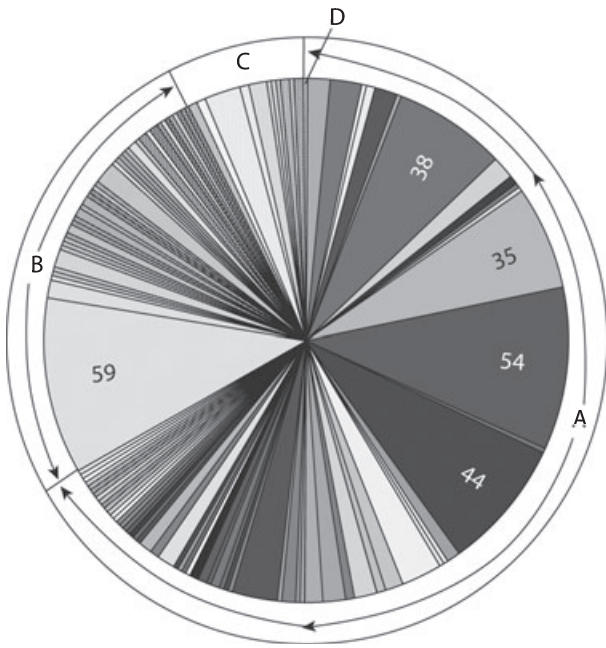


FIG. 2—Distribution of haplotypes. A pie chart showing the distribution of individuals who share an identical mitochondrial control region sequence or haplotype. Individuals with missing data ( $n = 10$ ) were not included. Haplogroup A is the largest, containing 66.8% of all dogs surveyed followed by B with 25.2%, C with 7.6%, and D with 0.2%. The numbers inside of the slices represent the number of individuals found with that particular haplotype. Haplogroup B has the largest single instance of individuals with the same haplotype ( $n = 59$ ). A large portion of the pie is comprised of many individuals sharing a haplotype, indicating the need for identification of additional mitochondrial single nucleotide polymorphisms to break up these large groupings.

the nucleotide diversity was  $0.013 \pm 0.006$ . Excluding the tandem repeat region, the probability of exclusion of the canine mtCR, or  $1 - \sum X_i^2$  (where  $X_i$  is the frequency of the  $i$ th haplotype), was 0.957 and the random match probability,  $\sum X_i^2$ , was 0.043 for all haplotypes in the current dataset. In other words, based on this dataset, the probability of two dogs having the same control region sequence at random is 4.3 out of 100. When the population was split into purebred and mixed breed individuals, the mean number of uncorrected pairwise differences decreased slightly to  $12.36 \pm 5.59$  for purebred and increased for mixed breed to  $12.79 \pm 5.80$ . Rounded to the thousands, the nucleotide diversities of the purebred and mixed separate datasets were identical to the combined dataset:  $0.013 \pm 0.006$ . Accordingly, the AMOVA on the dataset showed that there is an insignificant amount of genetic variation among the purebred and mixed populations ( $p > 0.05$ ) (Table 5). Dogs were also divided based on the large amount of samples from California ( $n = 189$ ), Pennsylvania ( $n = 100$ ), Virginia ( $n = 61$ ), and Nevada ( $n = 52$ ). Again, AMOVA showed that there is no significant difference in genetic variation in dogs sampled from the different geographic regions based on mtCR sequence ( $p > 0.05$ ) (Table 5). The dogs from each state were also evaluated based on how they were distributed among the four major haplogroups. As can be seen from Fig. 3, the distribution of haplogroups is consistent regardless of geographic location. The third AMOVA of large purebred groups ( $n > 6$ ) consisted of Golden Retrievers ( $n = 39$ ), Labrador Retrievers ( $n = 31$ ), Basset Hounds ( $n = 8$ ), Dachshunds ( $n = 8$ ), Poodles ( $n = 8$ ), Border Collies ( $n = 7$ ), Boston Terriers ( $n = 7$ ), Cavalier King Charles Spaniels ( $n = 7$ ), Cocker Spaniels

TABLE 4—Informative sequence variants.

Coordinate	Reference	Observed	L	ri
15,464.1	–	C	6	37
15,475	T	C	1	100
15,483	C	T	1	100
15,508	C	T	1	100
15,513	G	A	1	100
15,526	C	T	2	99
15,553	A	G	13	20
15,595	C	T	7	95
15,611	T	C	1	100
15,612	T	C	1	100
15,620	T	C	68	48
15,621	C	T	3	75
15,622	T	C	1	100
15,625	T	C	5	55
15,627	A	G	85	56
15,628	T	C	2	75
15,632	C	T	2	99
15,635	A	G	2	66
15,639	T	A/C/G	45	84
15,643	A	G	1	100
15,650	T	C	2	97
15,652	G	A	3	98
15,653	A	G	2	66
15,665	T	C	5	60
15,710	C	T	2	95
15,750	C	T	1	100
15,781	C	T	1	100
15,800	T	C	2	99
15,807	C	T	1	100
15,814	C	T	1	100
15,815	T	C	2	99
15,819	T	C	1	100
15,912	C	T	2	99
15,931	A	–	2	92
15,938	G	–	5	90
15,955	C	T	39	85
15,959	C	T	4	25
16,003	A	G	1	100
16,025	T	C	82	38
16,032	A	G	3	71
16,083	A	G	4	97
16,084*	A	G	1	100
16,128	G	A	2	99
16,129.1*	–	G	12	26
16,430*	G	T/–	12	94
16,431	C	–	10	94
16,432*	A	–	8	95
16,433*	C	–	9	95
16,439	T	C	4	97
16,501	T	C	1	100
16,507	T	A	1	100
16,576	A	G	12	21
16,617*	G	A	2	0
16,705	C	T	2	94

The variable nucleotide coordinate relative to the Kim et al. (11) reference sequence base (11), the observed base, the character length (L), and character retention index (ri) are listed. See Materials and Methods for definitions of character length and retention index. Shaded boxes indicate highly informative sites in the current study.

\*Indicate unrecognized sites in previously published literature.

( $n = 7$ ), Jack Russell Terriers ( $n = 7$ ), Miniature Schnauzers ( $n = 7$ ), Rottweilers ( $n = 7$ ), and West Highland Terriers ( $n = 7$ ). As can be seen from Table 4, all dogs from the same breed do not consistently share a haplotype. However, the AMOVA results do show evidence of significant genetic population substructure among dogs grouped according to breed ( $p = 0.00$ ) (Table 5).

TABLE 5—AMOVA results.

Dataset	Source of Variation	Degrees of Freedom	Percentage of Variation
Purebred versus mixed	Among populations	1	1.06
	Within populations	550	98.94
	Total	551	100
	$\Phi_{st}$	0.01057	
By states	Among populations	3	0 (−0.46)
	Within populations	398	100.46
	Total	401	100
	$\Phi_{st}$	0 (−0.00457)	
By breed	Among populations	12	28.14
	Within populations	139	71.86
	Total	151	100
	$\Phi_{st}$	0.28137	

$\Phi_{st}$ , fixation index.

Grouping and results of AMOVA as preformed in Arlequin to assess population structure between purebred and mixed breed dogs, dogs grouped by geographic state of origin, and large breed groups of purebred dogs.

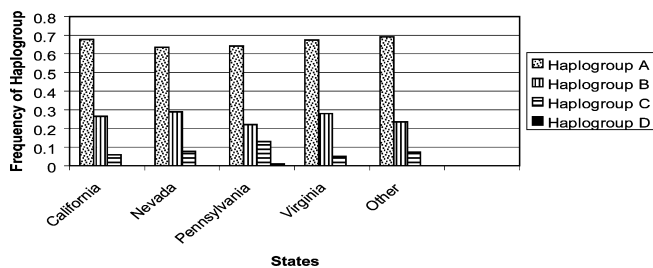


FIG. 3—Distribution of haplogroups by state. A bar graph showing the distribution of domestic dogs based on geographic location and the frequency of each haplogroup within each region. The number of dogs from each location varied: California = 189, Nevada = 52, Pennsylvania = 100, Virginia = 61. The “Other” group comprised the remaining 150 from Maryland (n = 1), Mississippi (n = 8), New York (n = 1), Texas (n = 14), Vermont (n = 1), and Unknown (n = 125). The graph shows that there is no bias towards a specific haplogroup based on geographic region.

## Discussion

This project was intended to survey the largest known sample set of mtCRs isolated from domestic dogs across the continental United States. We sequenced 427 new mtCRs and combined them with 125 previously published mtCRs (14) to search for new SNPs and haplotypes. We evaluated the need to distinguish between purebred and mixed breed dogs and dogs from different geographic regions across the continental United States. We also looked at the necessity of sequencing multiple individuals of the same breed.

When collecting samples, discrepancies were found in breed definition. For example, some samples were received labeled by the donor as “Spitz.” While there is a Finnish Spitz breed and a German Spitz breed, Spitz is not a true breed designation but another name for an American Eskimo Dog. It is unknown whether the donor meant one of the specific Spitz breeds or if the dog was in fact an American Eskimo Dog. Also, samples were received with the breed listed as “unknown,” but some described as purebred and some described as mixed. Descriptions could not be clarified or changed, as this could be error prone without seeing the dog. As a result of each of the above-mentioned problems, the number of distinct breeds collected for this study may be inflated. AMOVA was carried out on various subdivisions of the dataset to assess the severity of this issue.

The tandem repeat region was excluded in this study due to the known possibility of variation within an individual (15). While excluding the tandem repeat region from control region studies has come to be common practice (16,19,20,25,26), it appears that the studies conducted by our lab are the first to have problems obtaining sequence data for the region following the repeat (14). The sequencing problems seem to result from either individuals having a different number of repeats in the tandem repeat region (16,130–16,430 bps), individuals having a different number of C’s and/or T’s at the C/T stretch (16,663–16,676 bps), or a combination of both. This resulted in multiple sequence runs from the same individual being slightly different in length across these regions. There are multiple mitochondria per cell and multiple mitochondrial genomes per mitochondrion. Differences between the genomes caused the DNA sequence reads to be shifted by one or a few bases due to the insertion or deletion of bases in these highly variable regions. This resulted in an increase of ambiguous bases in 3’ region and the region between positions 16,663 and 16,676 being excluded due to variation in the number of repeated nucleotides for this C/T stretch when using a multiple alignment to search for informative SNPs.

The phylogenetic analysis showed that all dogs in the current dataset grouped within previously defined haplogroups A, B, C, and D. The proportions of samples within each group are very similar to the proportions of unique haplotypes previously identified for each group. This is particularly interesting because the samples used in previous studies came from all over the world, while the samples in the current study are from the continental United States alone. It appears that regardless of local origin, more domestic dogs have an A haplotype than any of the other types described followed by B, C, and then D. Additional local studies are needed to confirm this observation. The lack of individuals from groups E and F is most likely due to the fact that the individuals in previous studies that formed groups E and F were collected from Asian and/or Siberian localities (11,20,26,28). Individuals with D, E, and F haplotypes have been found in much lower frequencies compared with individuals with types A, B, and C in world-wide samplings (20). Our dataset is consistent with the conclusion that these haplotypes are relatively rare in the dog population.

One hundred and forty-three haplotypes were found in the current dataset with 53% of dogs possessing 1 of 9 haplotypes shared by between 11 and 59 individuals (Fig. 2). The distribution shows that while there are many canine mtCR haplotypes, the majority of dogs shared a few common haplotypes while the minority had unique or rare haplotypes. These results demonstrate a recurring problem with canine mtCR sequence data: most dogs share identical types. This also indicates a need for the evaluation of the remainder of the canine mitochondrial genome to look for additional SNPs that may distinguish between common mtCR haplotypes.

All of the variable sites identified in the current dataset are listed in Table 2 with the informative and highly informative sites shown in Table 4. Listing informative SNPs is important when attempting to identify the most useful SNPs for assessing population variation. Knowing where these informative SNPs occur in the mtCR allows for the potential development of targeted high throughput SNP genotyping assays wherein one could target the specific sites that define and distinguish between haplotypes, cutting down on resources and DNA necessary for the analysis. Our identification of 24 new SNPs, 6 of which were found to be informative and 3 highly informative, shows that previous studies have not resulted in a complete sampling of dog mtCRs, especially the region downstream of the repeat. All of the newly identified informative and highly informative SNPs were found in this less commonly sequenced region. While this contradicts the other findings of more

informative SNPs upstream of the repeat region, the lack of sequencing and analysis of the region downstream of the repeat most likely explains this finding. How informative a SNP is said to be is relative to the size and variation present in the dataset. As more sequences are added to the dataset, new sites may become phylogenetically informative due to the discovery of shared SNPs. Sites already identified as informative may gain a higher ranking due to their presence in more individuals. Also, the requirement of defining 1% of the total individuals in the dataset as criteria for the third ranking of SNPs is an arbitrary threshold, and changing this requirement may lead to changes in the ranking of SNPs.

As forensic samples are often subjected to conditions that may degrade DNA, the presence of the 60 bp hotspot within the mtCR is particularly useful. While the number of unique haplotypes gleaned from only 60 bases is not going to be as large as those from the entire mtCR, this provides a region of high variability to target when the entire mtCR cannot be sequenced due to DNA degradation.

Conversely, specific SNPs such as position 16439 seem to show higher levels of heteroplasmy relative to the remainder of the dataset and are represented in our dataset as ambiguous base calls. As such, we recommend that future researchers pay close attention to base calls at these sites when editing their raw sequence data, and if possible, clone this region to further investigate these ambiguities.

The probabilities of exclusion and random match probabilities calculated for the dataset are slightly more powerful but similar to those previously reported (17,18). The additional power comes from a larger sampling of dogs, leading to more revealed genetic variation in the target population. These statistics vary depending on the dataset, and ideally, all existing and future control region sequences should be combined and stored in the same database to more closely approximate the population frequencies of individual haplotypes.

The nucleotide diversity and fixation index ( $\Phi_{st}$ ) both identify a lack of genetic structure within dogs when classified as purebred and mixed. This shows that the decision as to how to classify certain breeds (i.e., Labradoodles) is trivial as purebred dogs and mixed breed dogs are not distinct populations based on mtCR sequence (Table 5). The result of the AMOVA when dogs were grouped by state of sample origin was not significant and the distribution of dogs within each major haplogroup was consistent across the different geographic regions. This finding along with the consistent distribution of haplogroups across states (Fig. 3) supports previous studies that there is no need for local canine mtCR SNP databases within the continental United States (17). The significant  $\Phi_{st}$  value when dogs are grouped by breed is most likely due to the strong amount of inbreeding that occurs in purebred dog lineages. While dogs of the same breed do not always share identical mtCR sequences, there are higher  $\Phi_{st}$  values for within breeds than among the population as a whole. This demonstrates why multiple individuals from a single breed should be analyzed and, more importantly, why individuals from a variety of different breed types should comprise a database of domestic dog mtCR SNPs.

The results of the current study are consistent with analyses from earlier studies which showed that domestic dogs can be classified into four previously identified genetic groups as defined by mtCR variation. In conclusion, combining 427 newly sequenced domestic dog mtCRs with a previous study of 125 domestic dog mtCRs (14), we have identified both new haplotypes and new informative SNPs. In the majority of the 552 dogs, 52.2% were classified into 1 of the 9 large haplotype groups with between 11 and 59 individuals per group. The presence of such large groups underscores the need for DNA sequence analysis of the remainder of the domestic

dog mtGenome in the hope of identifying additional discriminatory SNPs to increase the resolution of the analysis. Additionally, 94 SNPs were identified in the current dataset, 54 of which were informative and 33 highly informative. Twenty-four of these SNPs, 6 of which were informative and 3 highly informative, were previously unrecognized in the published literature. In general, population analyses show that domestic dogs consist of one large population. Smaller populations such as "purebred" and "mixed" or geographic populations cannot be distinguished based on mtCR sequences. However, when dogs are grouped by breed, the breeds are found to have less genetic variation than the population as a whole. These population analyses demonstrate the need to sample across a variety of breeds, including multiple individuals of the same breed, and that local mtCR SNP databases are not needed within the continental United States.

#### Acknowledgments

We thank Sheri Church, Gustavo Hormiga, Diana Johnson, and Mark Wilson for careful review of this manuscript. The research was conducted at The George Washington University. Six hundred and ninety-eight new domestic dog samples were collected through donations from veterinary practices and private donors. We would like to thank Adobe Animal Hospital, Austin Vet Hospital, Caring Hands Animal Hospital, Del Paso Vet Clinic, Little River Vet Clinic, The College of Veterinary Medicine at Mississippi State University, Pet Medical Center of Las Vegas, Seneca Hill Animal Hospital, The Animal Clinic of Clifton, West Flamingo Animal Hospital for blood and tissue samples, and 80 private donors who provided buccal swabs to add to our collection. Aisling Kelley, Stephanie Carnation, and Dunia Qutub helped to collect sequence data. Opinions or points of view expressed in this manuscript represent the consensus of the authors and do not necessarily represent the official position or policies of the U.S. Department of Justice or anyone who was acknowledged herein. The products and manufacturers discussed in this document are presented for informational purposes only and do not constitute product approval or enforcement by the U.S. Department of Justice.

#### References

- Schneider PM, Seo Y, Rittner C. Forensic mtDNA hair analysis excludes a dog from having caused a traffic accident. *Int J Legal Med* 1999;112(5):315-6.
- Savolainen P, Lundeberg J. Forensic evidence based on mtDNA from dog and wolf hairs. *J Forensic Sci* 1999;44(1):77-81.
- Deedrick DW. Hairs, fibers, crime, and evidence. *Forensic Sci Commun* 2000;2(3). Available at <http://www.fbi.gov/hq/lab/fsc/backissu/july2000/deedric1.htm>.
- Graham EAM. DNA reviews: hair. *Forensic Sci Med Pathol* 2007; 3(2):133-7.
- Takayanagi K, Asamura H, Tsukada K, Ota M, Saito S, Fukushima H. Investigation of DNA extraction from hair shafts. *Int Congr Ser* 2003;1239:759-64.
- Roberts KA, Calloway C. Mitochondrial DNA amplification success rate as a function of hair morphology. *J Forensic Sci* 2007;52(1):40-7.
- Wilson MR, DiZinno JA, Polansky D, Replogle J, Budowle B. Validation of mitochondrial DNA sequencing for forensic casework analysis. *Int J Legal Med* 1995;108(2):68-74.
- Bogenhagen D, Clayton D. The number of mitochondrial deoxyribonucleic acid genomes in mouse L and human HeLa cells. Quantitative isolation of mitochondrial deoxyribonucleic acid. *J Biol Chem* 1974; 249:7991-5.
- Nass M. Mitochondrial DNA. I. Intramitochondrial distribution and structural relations of single- and double-length circular DNA. *J Mol Biol* 1969;42:521-8.

10. Budowle B, Allard MW, Wilson MR, Chakraborty R. Forensics and mitochondrial DNA: applications, debates, and foundations. *Annu Rev Genomics Hum Genet* 2003;4:119–41.
11. Kim KS, Lee SE, Jeong HW, Ha JH. The complete nucleotide sequence of the domestic dog (*Canis familiaris*) mitochondrial genome. *Mol Phylogenet Evol* 1998;10(2):210–20.
12. Parsons TJ, Coble MD. Increasing the forensic discrimination of mitochondrial DNA testing through analysis of the entire mitochondrial DNA genome. *Croat Med J* 2001;42(3):304–9.
13. Pesole G, Gissi C, De Chirico A, Saccone C. Nucleotide substitution rate of mammalian mitochondrial genomes. *J Mol Evol* 1999;48(4):427–34.
14. Gundry RL, Allard MW, Moretti TR, Honeycutt RL, Wilson MR, Monson KL, et al. Mitochondrial DNA analysis of the domestic dog: control region variation within and among breeds. *J Forensic Sci* 2007;52(3):562–72.
15. Savolainen P, Arvestad L, Lundeberg J. A novel method for forensic DNA investigations: repeat-type sequence analysis of tandemly repeated mtDNA in domestic dogs. *J Forensic Sci* 2000;45(5):990–9.
16. Angleby H, Savolainen P. Forensic informativity of domestic dog mtDNA control region sequences. *Forensic Sci Int* 2005;154(2-3):99–110.
17. Himmelberger AL, Spear TF, Satkoski JA, George DA, Garnica WT, Malladi VS, et al. Forensic utility of the mitochondrial hypervariable region 1 of domestic dogs, in conjunction with breed and geographic information. *J Forensic Sci* 2008;53(1):81–9.
18. Savolainen P, Rosen B, Holmberg A, Leitner T, Uhlen M, Lundeberg J. Sequence analysis of domestic dog mitochondrial DNA for forensic use. *J Forensic Sci* 1997;42(4):593–600.
19. Wetton JH, Higgs JE, Spriggs AC, Roney CA, Tsang CS, Foster AP. Mitochondrial profiling of dog hairs. *Forensic Sci Int* 2003;133(3):235–41.
20. Savolainen P, Zhang YP, Luo J, Lundeberg J, Leitner T. Genetic evidence for an East Asian origin of domestic dogs. *Science* 2002;298(5598):1610–3.
21. Wilson MR, Allard MW, Monson K, Miller KW, Budowle B. Recommendations for consistent treatment of length variants in the human mitochondrial DNA control region. *Forensic Sci Int* 2002;129(1):35–42.
22. Pereira L, Van Asch B, Amorim A. Standardisation of nomenclature for dog mtDNA D-loop: a prerequisite for launching a *Canis familiaris* database. *Forensic Sci Int* 2004;141(2-3):99–108.
23. Nixon KC. Winclada ver. 1.00.08. Ithaca, NY: Published by author, 2002. Available at <http://www.cladistics.org>.
24. Excoffier L, Laval G, Schneider S. Arlequin: ver. 3.0: an integrated software package for population genetics data analysis. *Evolutionary Bioinformatics Online* 2005;1:47–50.
25. Takahasi S, Miyahara K, Ishikawa H, Ishiguro N, Suzuki M. Lineage classification of canine inheritable disorders using mitochondrial DNA haplotypes. *J Vet Med Sci* 2002;64(3):255–9.
26. Tsuda K, Kikkawa Y, Yonekawa H, Tanabe Y. Extensive interbreeding occurred among multiple matriarchal ancestors during the domestication of dogs: evidence from inter- and intraspecies polymorphisms in the D-loop region of mitochondrial DNA between dogs and wolves. *Genes Genet Syst* 1997;72(4):229–38.
27. Nixon KC. The parsimony ratchet, a new method for rapid parsimony analysis. *Cladistics* 1999;15(4):407–14.
28. Okumura N, Ishiguro N, Nakano M, Matsui A, Sahara M. Intra- and interbreed genetic variations of mitochondrial DNA major non-coding regions in Japanese native dog breeds (*Canis familiaris*). *Anim Genet* 1996;27(6):397–405.

Additional information and reprint requests:  
 Marc W. Allard, Ph.D.  
 Food and Drug Administration  
 Office of Regulatory Science  
 Division of Microbiology  
 HFS-712  
 5100 Paint Branch Parkway  
 College Park, MD 20740  
 E-mail: mwallard@gwu.edu