

Standardisation of nomenclature for dog mtDNA D-loop: a prerequisite for launching a *Canis familiaris* database

Luísa Pereira^{a,*}, Bárbara Van Asch^a, António Amorim^{a,b}

^a*Instituto de Patologia e Imunologia Molecular da Universidade do Porto (IPATIMUP),*

R. Dr. Roberto Frias s/n, 4200-465 Porto, Portugal

^b*Faculdade de Ciências da Universidade do Porto, Pr. Gomes Teixeira, 4050 Porto, Portugal*

Received 23 August 2003; received in revised form 28 November 2003; accepted 1 December 2003

Abstract

Domestic dogs are increasingly involved, often as protagonists, in the forensic scene. Acknowledging this fact and benefiting from the accumulated experience on human mitochondrial DNA (mtDNA) analyses, we propose a standard for *Canis familiaris* mtDNA sequences as a prerequisite for the launching of the corresponding database.

© 2004 Elsevier Ireland Ltd. All rights reserved.

Keywords: *Canis familiaris*; Dogs; Database; mtDNA

1. Introduction

Information obtained from sequencing hypervariable segments of the control region (D-loop) of mitochondrial DNA (mtDNA) has been systematically used in the study of human populations for more than two decades [1]. Standardisation of data reporting has been immediately recognised to be a key issue, namely in forensic genetics [2]. Valuable lessons from this accumulated experience can and may be applied to other species. One of such species is *Canis familiaris*, on which mtDNA studies are being undertaken, not only as population genetics worldwide comparisons for unravelling breed evolution and conservation [3], but also, and increasingly, on the forensic field [4–6] and even on ancient DNA [7]. Furthermore, apart from its growing importance in forensics, the species is also of great economic interest [8,9]. The purpose of this work is to: (1) standardise the nomenclature for dog mtDNA sequences; and (2) to launch a database that includes most of the dog published sequences, according to that procedure.

2. A reference sequence

The mtDNA polymorphisms are usually reported in comparison with a reference sequence, which must be the same for all the scientific community, in order to make the analysis of results easy and straightforward. In humans, the chosen reference was the firstly reported complete mtDNA sequence [10], denominated, accordingly, as Cambridge Reference Sequence, or CRS (later corrected by [11]). In dogs, some studies were performed [4,12] prior to the publication of a complete dog mtDNA sequence [13]. So, some authors employed an observed sequence as reference: Savolainen et al. [4] used one and maintained it as such in later works (the A2 haplotype of [3]); whereas Vila et al. [12] referred to a wolf sequence (denominated W12); Takahasi et al. [14] used a Shiba 1 sequence; and Kim et al. [15] a Sapsare A (KS1). Therefore, authors are currently converting previously published data for comparison to their own reference sequence. This tedious and error-prone work can be avoided if from now on, the first complete mtDNA dog sequence published by Kim et al. [13] (GeneBank accession number: NC_002008) is used as reference. Using this complete sequence as reference and not, for instance, the partial one used in the database collected by Savolainen et al. [3], will prevent future problems of enlargement of the mtDNA region studied, as for the screening of coding regions and

* Corresponding author. Tel.: +351-22-5570700;

fax: +351-22-5570799.

E-mail address: lpereira@ipatimup.pt (L. Pereira).

even complete sequencing (as it happened in the last years for humans).

3. The numbering of polymorphisms

At the moment, to check if a sequence was already published, we need to (1) download from GeneBank each deposited sequences, and to (2) use sequence alignment software to compare the pattern of the segment under analysis with the downloaded ones. We could also search in the tables reported in published works, but in most cases, the polymorphic positions are numbered in relation to the beginning of the segment analysed in each study [4,12,14,15]. And some reports [3] even do not indicate the relative positions of the polymorphic bases described.

The easiest solution (and the one adopted for studies on humans) would be to have a catalogue of sequences just denoting the differences relatively to the reference sequence, numbering unambiguously the corresponding positions, so that search matches for a specific sequence or position could then be comfortably performed.

The Kim et al. [13] reference sequence consists in 16727 bp, and the D-loop is located between positions 15458–16727, where a 10 bp imperfect minisatellite is repeated 30 times, with a polymorphic transition (A/G) in its ninth base (the repeats are located between 16130 and 16429). However, Savolainen et al. [16] reported large size variations in this minisatellite that could confound numbering of its upstream D-loop region. In humans, the correction of the complete sequence by Andrews et al. [11] revealed some missing bases in the reported sequence by Anderson et al. [10], but in order to avoid errors in conversions between different numbering systems, the bases are still numbered according to Anderson, although incorrect. So far, only Okumura et al. [17] surveyed the dog D-loop region upstream the minisatellite, and did not report the variation of the minisatellite region. We suggest the following procedure for the numbering of the complete dog D-loop: (1) refer how many times the minisatellite is repeated in the typical VNTR way, as for instance (GTACAGTNC)₂₀, and then (2) number the following sequence variation as beginning in base 16430 (regardless of the number of repeat units).

Concerning the recording of substitutions, for simplicity reasons, a haplotype can be described as being 15627 15639^{T/A} 15814, where numbers without superscript denote transitions (15627 refers the A to G transition and 15814 the C to T), and other base changes (as in the case of 15639) being explicitly indicated. For indels, the numbering of the reference should be maintained: new positions must be considered insertions, and referred as 15534.1C if the base inserted is a C (or X.2C if there is insertion of 2 Cs); and missing ones coded as deletions (e.g., 15938^{del}). In cases where addition or deletion occurs in a homopolymeric tract (sequence stretch of the same base), the gaps are placed in

the 3' end of the tract [2] (e.g., a C insertion into a 15461–15464 homopolymeric C tract is recorded as 15464.1C instead of 15461.1C).

In some cases the alignment of a certain gap can be interpreted in different ways, conducting to potentially miscoded variation. In this dog mtDNA database there is such a case with haplotype F1, below compared with the reference:

```
F1      AAACCCTCCCCCTATG
Ref     AAACCCTTCTCCCCCTCCCCTATG
```

Wilson et al. [18] recommend an alignment approach that is based on a phylogenetic context using differential weighting of transitions, transversions and indels. Basically, they proved that most variants could be characterised if the following three recommendations are followed:

- (1) Characterise profiles using the least number of differences from the reference sequence.
- (2) If there is more than one way to maintain the same number of differences relatively to the reference, differences should be prioritised in the following manner:
 - (a) indels
 - (b) transitions
 - (c) transversions
- (3) Indels should be placed 3' with respect to the light strand. Insertions and deletions should be combined in situations where the same number of differences to the reference sequence is maintained.

The alignment proposed in [16] is:

```
F1      AAACCCT-----CCCCCTATG
Ref     AAACCCTTCTCCCCTCCCC-TATG
```

that is: 15523^{del} 15524^{del} 15525^{del} 15526^{del} 15527^{del} 15528^{del} 15529^{del} 15530^{del} 15534.1C, where the combination of the insertion with the deletions is supported by phylogeny, since all the remaining F haplotypes have this insertion in comparison with the reference. But according to the first of the above rules, the following alignment must be considered

```
F1      AAACCCTC-----CCCCTATG
Ref     AAACCCTTCTCCCCTCCCCTATG
```

that is: 15523 15524^{del} 15525^{del} 15526^{del} 15527^{del} 15528^{del} 15529^{del} 15530^{del}, being the transition at position 15523 also supported by phylogeny. There are many examples in [18,19] that can be helpful in deciding further alignments in newly discovered confusing gaps.

4. Towards the organisation of a *Canis familiaris* mtDNA database

We provide a catalogue of the most significant datasets published so far reorganised under the standardised procedures proposed here. We indicate the polymorphism

Table 1 (Continued)

Reference	1	1	11
	A-TTCAACGCCGATTCTCCCT-CCATACCTCATCTATCAGATTTTAAATCGATAACACCACCCTTCTCACCTTCCTCGGAACCAGCCCTCATAAAGACTAAT--TCC		
A64	D83629(1)	?-.....-.....T.....A.....C.....T.T.T.....T....-C..T.....-C..
	AB007400(2)	?-.....-.....T.....A.....C.....T.T.T.....???????????T.....-C..
A65	D83635(1)	?-.....-.....A.....T.....T.....T.....T.....-C..T.....-C..
	I-22(4)	?-.....-.....A.....T.....T.....T.....T.....-C..T.....-C..
A66	D83602(1)	?-.....-.....C.....T.....T.....T.....T.....-C..T.....-C..
	I-41(4)	?-.....-.....C.....T.....T.....T.....T.....-C..T.....-C..
A67	D83608(1)	?-.....-.....C.....T.....T.....T.....T.....-C..T.....-C..
A68	D83628(1)	?-.....-.....G.....T.....T.....T.....T.....-C..T.....-C..
	AB007389(2)	?-.....-.....G.....T.....T.....T.....T.....-C..T.....-C..
	UK26(6)-.....G.....T.....T.....T.....T.....-C..T.....-C..
A69	D83631(1)	?-.....-.....T.....G.....T.....T.....T.....T.....-C..T.....-C..
A70	D83626(1)	?-.....-.....C.CG.....A.....T.....T.....T.....T.....-C.T.T.....-C.T.
	AB007395(2)	?-.....-.....C.CG.....A.....T.....T.....T.....T.....-C.T.T.....-C.T.
	I-1(4)	?-.....-.....C.CG.....A.....T.....T.....T.....T.....-C.T.T.....-C.T.
A71	D83612(1)	?-.....-.....C.G.....A.....C.....C.....C.....C.....C.T.T.....C.T.
	UK2(6)-.....C.G.....A.....C.....C.....C.....C.....C.T.T.....C.T.
A72	D83621(1)	?-.....-.....G.....C.A.....A.C.....T.....T.....T.....T.....-C..T.....-C..
A73	D83613(1)	?-.....-.....G.....A.....A.C.....T.....T.....T.....T.....-C..T.....-C..
	AB007384(2)	?-.....-.....G.....A.....A.C.....T.....T.....T.....T.....-C..T.....-C..
	A73(03)	?-.....-.....G.....A.....A.C.....T.....T.....T.....T.....-C..T.....-C..
	KS3/KS5/KJ4(05)	??G.....A.....A.C.....T.....T.....T.....T.....-C..T.....-C..
A74	AB007386(2)	?-.....-.....G.....A.....T.....C.....T.....T.....T.....T.....-C..T.....-C..
	UK5(6)-.....G.....A.....T.....C.....T.....T.....T.....T.....-C..T.....-C..
A75	AB007398(2)	?-.....-.....G.....A.G.....A.....T.....T.....T.....T.....-C..T.....-C..
	A75(3)	?-.....-.....G.....A.G.....A.....T.....T.....T.....T.....-C..T.....-C..
NewA	I-7(4)	?-.....-.....C.G.....A.....T.....T.....T.....T.....C.....???????????????T.....C.....???????????????
NewA	I-27(4)	?-.....-.....G.....A.....T.....T.....T.....T.....C.....???????????????T.....T.....C.....???????????????
NewA	I-33(4)	?-.....G.....T.....T.....T.....T.....T.....T.....T.....T.....C.....???????????????T.....T.....T.....T.....C.....???????????????
NewA	I-34(4)	?-.....GT.....T.....T.....T.....T.....T.....T.....T.....T.....C.....???????????????T.....T.....T.....T.....C.....???????????????
NewA	I-36(4)	?-.....-.....C.....T.....T.....T.....T.....T.....T.....T.....T.....C.....???????????????T.....T.....T.....T.....C.....???????????????
NewA	I-40(4)	?-.....-.....T.....A.....G.....A.....A.....A.....T.....T.....T.....T.....C.....???????????????T.....T.....T.....T.....C.....???????????????
NewA	I-13(4)	?-.....-.....A.....G.....A.....A.....A.....T.....T.....T.....T.....C.....???????????????T.....T.....T.....T.....C.....???????????????
NewA	WGD3(5)	??G.....A.....A.C.....T.....T.....T.....T.....C.....???????????????T.....T.....T.....T.....C.....???????????????
NewA	WGD1(5)	??G.....A.....A.C.....T.....T.....T.....T.....C.....???????????????T.....T.....T.....T.....C.....???????????????
NewA	KC1/KS7(5)	??G.....A.....A.C.....T.....T.....T.....T.....C.....???????????????T.....T.....T.....T.....C.....???????????????
NewA	Apu(5)	??G.....A.....A.C.....T.....T.....T.....T.....C.....???????????????T.....T.....T.....T.....C.....???????????????
NewA	KJ2(5)	??G.....A.....A.C.....T.....T.....T.....T.....C.....???????????????T.....T.....T.....T.....C.....???????????????
NewA	UK11(6)-.....A.....T.....T.....T.....T.....T.....T.....T.....T.....C.....???????????????T.....T.....T.....T.....C.....???????????????
NewA	UK14(6)-.....A.....C.....T.....T.....T.....T.....T.....T.....T.....T.....C.....???????????????T.....T.....T.....T.....C.....???????????????
NewA	UK15(6)-.....A.....G.....T.....T.....T.....T.....T.....T.....T.....T.....C.....???????????????T.....T.....T.....T.....C.....???????????????
NewA	UK1(6)-.....C.G.....A.....T.....T.....T.....T.....T.....T.....T.....T.....C.....???????????????T.....T.....T.....T.....C.....???????????????
NewA	UK12(6)-.....G.....A.....T.....T.....T.....T.....T.....T.....T.....T.....C.....???????????????T.....T.....T.....T.....C.....???????????????

positioning of sequences relatively to the Kim's reference sequence in two alternative ways: (1) in Table 1 as a dot list of all the different haplotypes observed; and (2) in supplementary material as a database containing information of the number of individuals for each haplotype, and, when available, the breed, place of origin and GeneBank accession number for the sequence.

The size of the database presented here is of 1146 individuals, most of them with indication of breed (amounting to 143) and place of origin. Most of the sequences reported survey the D-loop region between positions 15458 and 16039, a segment approximately 582 bp long, for which 96 were polymorphic (72 transitions, 13 transversions and 14 indels), defining a total of 139 haplotypes.

5. Haplogrouping

We further want to stress another possible source of trouble in the report of mtDNA results, namely the attribution of different names/symbols to the same clades/haplogroups (groups of related sequences sharing specific polymorphisms). This situation is already arising in dogs. Initially, clade names were designed by Roman numerals (from I to IV), followed by an Arabic number, denoting haplotypes inside the clade [4,12], whereas recently a system using Latin alphabet was employed (from A to F; [3]). The correspondence for both nomenclatures is: I to A; II to D; III to B; and IV to C. At the moment, only notoriously different clades were named in dogs, but, in the future, the accumulation of information on position heterogeneity, leading to sub-classification based on less recurrent positions will originate a pressing need for a unitary haplogroup classification. Networks reported in Savolainen et al. [3] already show sub-hierarquistion especially within clade A, and the further indication of the positions will turn it possible to easily classify sub-clades. This information can be used to sub-classify the sequences reported as in list and tables presented here, without the need to build networks from those samples. A possible solution for a functional nomenclature could be the one adopted by the Y Chromosome Consortium [20]. This nomenclature is open, that is, uses capital letters, followed by numbers, which can still be subdivided by lowercase letters (and additional numbers), in order to allow sub-hierarquistion inside a clade (e.g., A1a and A1b3).

6. Conclusions and future considerations

We think that the organisation of the database for dog mtDNA sequences, which is by now considerably extended, will be of extreme importance in the forensic field, where sequence matches are already being searched in order to solve forensic cases [21]. Furthermore, it can also be of considerable value for all researchers on *Canis familiaris* as well as for its breeders.

Acknowledgements

We wish to thank the considerable improvement contributed to the manuscript by two anonymous referees. This work was partially supported by a research grant to LP (SFRH/BPD/7121/2001) from Fundação para a Ciência e a Tecnologia and IPATIMUP by Programa Operacional Ciência, Tecnologia e Inovação (POCTI), Quadro Comunitário de Apoio III. This study is included in the Project "GENCERT", financed by POCTI Medida 2.3 and POSI Medida 1.3 (PO-QCA III).

References

- [1] M. Richards, V. Macaulay, The mitochondrial gene tree comes of age, *Am. J. Hum. Genet.* 68 (2001) 1315–1320.
- [2] A. Carracedo, W. Bar, P. Lincoln, W. Mayr, N. Morling, B. Olaisen, P. Schneider, B. Budowle, B. Brinkmann, P. Gill, M. Holland, G. Tully, M. Wilson, DNA commission of the international society for forensic genetics: guidelines for mitochondrial DNA typing, *Forensic Sci. Int.* 110 (2000) 79–85.
- [3] P. Savolainen, Y.P. Zhang, J. Luo, J. Lundeberg, T. Leitner, Genetic evidence for an East Asian origin of domestic dogs, *Science* 298 (2002) 1610–1613.
- [4] P. Savolainen, B. Rosen, A. Holmberg, T. Leitner, M. Uhlen, J. Lundeberg, Sequence analysis of domestic dog mitochondrial DNA for forensic use, *J. Forensic Sci.* 42 (1997) 593–600.
- [5] P. Savolainen, J. Lundeberg, Forensic evidence based on mtDNA from dog and wolf hairs, *J. Forensic Sci.* 44 (1999) 77–81.
- [6] J.H. Wetton, J.E. Higgs, A.C. Spriggs, C.A. Roney, C.S. Tsang, A.P. Foster, Mitochondrial profiling of dog hairs, *Forensic Sci. Int.* 133 (2003) 235–241.
- [7] J.A. Leonard, R.K. Wayne, J. Wheeler, R. Valadez, S. Guillen, C. Vila, Ancient DNA evidence for old world origin of new world dogs, *Science* 298 (2002) 1613–1616.
- [8] B.M. Kuehn, Breed discrimination bites homeowners: insurance companies dropping home insurance coverage for owners of large dog breeds, *J. Am. Vet. Med. Assoc.* 222 (2003) 1337–1338.
- [9] J.M. Dobson, S. Samuel, H. Milstein, K. Rogers, J.L. Wood, Canine neoplasia in the UK: estimates of incidence rates from a population of insured dogs, *J. Small Anim. Pract.* 43 (2002) 240–246.
- [10] S. Anderson, A.T. Bankier, B.G. Barrell, M.H. de Bruijn, A.R. Coulson, J. Drouin, I.C. Eperon, D.P. Nierlich, B.A. Roe, F. Sanger, P.H. Schreier, A.J. Smith, R. Staden, I.G. Young, Sequence and organization of the human mitochondrial genome, *Nature* 290 (1981) 457–465.
- [11] R.M. Andrews, I. Kubacka, P.F. Chinnery, R.N. Lightowlers, D.M. Turnbull, N. Howell, Reanalysis and revision of the Cambridge reference sequence for human mitochondrial DNA, *Nat. Genet.* 23 (1999) 147.
- [12] C. Vila, P. Savolainen, J.E. Maldonado, I.R. Amorim, J.E. Rice, R.L. Honeycutt, K.A. Crandall, J. Lundeberg, R.K. Wayne, Multiple and ancient origins of the domestic dog, *Science* 276 (1997) 1687–1689.

- [13] K.S. Kim, S.E. Lee, H.W. Jeong, J.H. Ha, The complete nucleotide sequence of the domestic dog (*Canis familiaris*) mitochondrial genome, *Mol. Phylogenet. E* 10 (1998) 210–220.
- [14] S. Takahasi, K. Miyahara, H. Ishikawa, N. Ishiguro, M. Suzuki, Lineage classification of canine inheritable disorders using mitochondrial DNA haplotypes, *J. Vet. Med. Sci.* 64 (2002) 255–259.
- [15] K.S. Kim, H.W. Jeong, C.K. Park, J.H. Ha, Suitability of AFLP markers for the study of genetic relationships among Korean native dogs, *Genes Genet. Syst.* 76 (2001) 243–250.
- [16] P. Savolainen, L. Arvestad, J. Lundeberg, A novel method for forensic DNA investigations: repeat-type sequence analysis of tandemly repeated mtDNA in domestic dogs, *J. Forensic Sci.* 45 (2000) 990–999.
- [17] N. Okumura, N. Ishiguro, M. Nakano, A. Matsui, M. Saharal, Intra- and interbreed genetic variations of mitochondrial DNA major non-coding regions in Japanese native dog breeds (*Canis familiaris*), *Anim. Genet.* 27 (1996) 397–405.
- [18] M.R. Wilson, M.W. Allard, K. Monson, K.W. Miller, B. Budowle, Recommendations for consistent treatment of length variants in the human mitochondrial DNA control region, *Forensic Sci. Int.* 129 (2002) 35–42.
- [19] M.R. Wilson, M.W. Allard, K. Monson, K.W. Miller, B. Budowle, Further discussion of the consistent treatment of length variants in the human mitochondrial DNA control region, *Forensic Sci. Commun.* 4 (2002) 4.
- [20] Y Chromosome Consortium, A nomenclature system for the tree of human Y-chromosomal binary haplogroups, *Genome Res.* 12 (2002) 339–348.
- [21] P.M. Schneider, Y. Seo, C. Rittner, Forensic mtDNA hair analysis excludes a dog from having caused a traffic accident, *Int. J. Legal Med.* 112 (1999) 315–316.
- [22] K. Tsuda, Y. Kikkawa, H. Yonekawa, Y. Tanabe, 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.* 72 (1997) 229–238.