

1 **Contrasting patterns of coding and flanking region evolution in mammalian keratin**  
2 **associated protein-1 genes**

3

4 Huitong Zhou<sup>\*,†,1</sup>, Tina Visnovska<sup>‡,1</sup>, Hua Gong<sup>\*,†</sup>, Sebastian Schmeier<sup>‡</sup>, Jon Hickford<sup>\*,†</sup>, and  
5 Austen R.D. Ganley<sup>‡,§</sup>

6

7 \* State Key Laboratory of Sheep Genetic Improvement and Healthy Production,  
8 Shihezi 832000, China

9 † Faculty of Agricultural and Life Sciences, Lincoln University, Lincoln 7647, New  
10 Zealand

11 ‡ Institute of Natural and Mathematical Sciences, Massey University Auckland,  
12 Auckland 0632, New Zealand

13 § School of Biological Sciences, University of Auckland, Auckland 1142, New  
14 Zealand

15

16 1 These authors contributed equally to this work

17

18 **Running title:** Contrasting *KRTAPI* evolutionary patterns

19

20 **Key words:**

21 concerted evolution, gene conversion, keratin associated protein, *krtp1*, tandem repeat

22

23 **To whom correspondence should be addressed:**

24 **Austen Ganley:** School of Biological Sciences, University of Auckland, 3A Symonds St,

25 Building 110N, Auckland 1142, New Zealand; +64 9 923 2906; [a.ganley@auckland.ac.nz](mailto:a.ganley@auckland.ac.nz)

26

27 **Huitong Zhou:** Faculty of Agriculture and Life Sciences, Lincoln University,

28 Cnr Springs Road & Ellesmere Junction Road, Lincoln 7647, New Zealand; +64 3 423 0684;

29 [zhouh@lincoln.ac.nz](mailto:zhouh@lincoln.ac.nz)

30

31

32 **Abstract**

33 DNA repeats are common elements in eukaryotic genomes, and their multi-copy nature  
34 provides the opportunity for genetic exchange. This exchange can produce altered  
35 evolutionary patterns, including concerted evolution where within genome repeat copies are  
36 more similar to each other than to orthologous repeats in related species. Here we  
37 investigated the genetic architecture of the keratin-associated protein (KAP) gene family,  
38 *KRTAPI*. This family encodes proteins that are important components of hair and wool in  
39 mammals, and the genes are present in tandem copies. Comparison of *KRTAPI* gene repeats  
40 from species across the mammalian phylogeny shows strongly contrasting evolutionary  
41 patterns between the coding regions, which have a concerted evolution pattern, and the  
42 flanking regions, which have a normal, radiating pattern of evolution. This dichotomy in  
43 evolutionary pattern transitions abruptly at the start and stop codons, and we show it is not  
44 the result of purifying selection acting to maintain species-specific protein sequences, nor of  
45 codon adaptation or reverse transcription of *KRTAPI-n* mRNA. Instead, the results are  
46 consistent with short-tract gene conversion events coupled with selection for these events in  
47 the coding region driving the contrasting evolutionary patterns found in the *KRTAPI* repeats.  
48 Our work shows the power that repeat recombination has to complement selection and finely  
49 tune the sequences of repetitive genes. Interplay between selection and recombination may be  
50 a more common mechanism than currently appreciated for achieving specific adaptive  
51 outcomes in the many eukaryotic multi-gene families, and our work argues for greater  
52 emphasis on exploring the sequence structures of these families.

53

54

## INTRODUCTION

55 Repetitive DNA is widespread in most eukaryote genomes (BRITTEN AND KOHNE 1968;  
56 RICHARD *et al.* 2008; LOPEZ-FLORES AND GARRIDO-RAMOS 2012). There are two basic repeat  
57 DNA types: tandem repeats that are typically arranged in head-to-tail arrays; and dispersed  
58 repeats, and these can occur in either coding or non-coding DNA. Repeats are thought to  
59 arise from recombination-based duplication/amplification events (STEPHAN 1989). Sequence  
60 identity between duplicates will then decay through the diversifying force of mutation, unless  
61 counteracting processes operate (BROWN *et al.* 1972; DOVER 1982). The balance between  
62 duplication, diversification, selection, and counteracting forces thus dictate the evolutionary  
63 dynamics of repeats. Two main paradigms have been proposed to account for the long-term  
64 maintenance of repeat identity: concerted evolution and birth-and-death evolution. Concerted  
65 evolution describes a pattern of evolution where the repeats within a genome show greater  
66 sequence identity to each other than to orthologous repeats in related genomes (ELDER AND  
67 TURNER 1995). The pattern of concerted evolution is proposed to result from recombination-  
68 based processes, such as gene conversion and unequal cross-over events, that replace the  
69 DNA sequence from one repeat with that from another repeat (LIAO 1999). In so doing, these  
70 recombination processes maintain sequence identity between repeat copies in the face of  
71 mutation, and thus homogenize the repeats (DOVER 1982). ‘Birth-and-death’ evolution  
72 involves purifying selection maintaining sequence identity between repeats that are generated  
73 by occasional duplication events (i.e. birth), as well as death, which results from repeat loss  
74 or pseudogenization (NEI *et al.* 1997; NEI *et al.* 2000). While there has been debate as to  
75 which of these processes best describes the evolutionary dynamics of repetitive DNA (NEI  
76 AND ROONEY 2005; ROONEY AND WARD 2005; EIRIN-LOPEZ *et al.* 2012), a basic  
77 characterization of the evolutionary dynamics of most repeat families is lacking.

78 The keratin-associated proteins (KAPs) are a diverse group of proteins, and are rich in either  
79 sulphur, or glycine and tyrosine. They are important structural components of hair and wool  
80 fibres, and form a matrix that cross-links the keratin intermediate filaments. The genes  
81 encoding the KAPs are called *KRTAPs* (GONG *et al.* 2012), and can be classified into 27  
82 families, with each family comprising 1-12 members that are usually tandemly arranged  
83 (ROGERS AND SCHWEIZER 2005; ROGERS *et al.* 2006; GONG *et al.* 2016). The *KRTAPs* are  
84 single exon (intron-less) genes, with small coding sequences (less than 1 kb) (ROGERS AND  
85 SCHWEIZER 2005), and they have low numbers of pseudogenes. For example, in humans the  
86 pseudogene:gene *KRTAP* ratio is approximately 1:5 (GONG *et al.* 2016), while across all  
87 human genes the ratio is close to 1:1 (TORRENTS *et al.* 2003; STEIN 2004). In addition, the  
88 *KRTAPs* show high levels of population variation, with all known *KRTAP* genes being  
89 polymorphic in sheep (GONG *et al.* 2010b; GONG *et al.* 2016; ZHOU *et al.* 2016), where they  
90 are well studied because of their roles in determining wool phenotypes (ZHOU *et al.* 2015; LI  
91 *et al.* 2017a; LI *et al.* 2017b; LI *et al.* 2017c; TAO *et al.* 2017a; TAO *et al.* 2017b). Despite this  
92 variation, it has been reported that at least some *KRTAP* genes show a pattern of concerted  
93 evolution between the paralogous gene copies (ROGERS *et al.* 1994; WU *et al.* 2008; KHAN  
94 *et al.* 2014).

95 The KAP1 proteins form the best characterised KAP family, and they show a high degree of  
96 sequence heterogeneity compared to other KAP families. These KAP1 proteins appear to be  
97 restricted in expression to the middle to upper cortex region of the hair and wool follicle, and  
98 are absent in the cuticle (POWELL AND ROGERS 1997; SHIMOMURA *et al.* 2002). Their precise  
99 role in hair and wool function, has yet to be determined. The genes encoding the KAP1  
100 proteins (*KRTAPI-n*) have been characterized in a number of mammalian species, where they  
101 are usually arranged as four tandem copies (**Figure 1**) (KHAN *et al.* 2014). The coding  
102 regions of the *KRTAPI-n* genes vary in length within species, predominantly as a

103 consequence of variation in the number of imperfect tandem decapeptide repeat units (GONG  
104 *et al.* 2016) (**Figure 1**).

105 Here we analyse the *KRTAPI* genes from a number of mammalian species, including four  
106 species for which the *KRTAPI-n* loci have not been described. Together with the existing  
107 *KRTAPI-n* sequences, we reveal that the *KRTAPI-n* coding regions display a pattern of  
108 concerted evolution. In stark contrast to the coding region though, we find that the repeat  
109 flanking regions display no evidence of concerted evolution, and instead appear to be  
110 evolving by normal vertical or radiating evolution. Surprisingly, we find that this pattern of  
111 coding region restricted concerted evolution is not the result of purifying selection, nor does  
112 it result from codon adaptation or reverse transcription/reintegration of *KRTAPI-n* mRNA  
113 sequences. Instead, the results are best explained by a combination of on-going short-tract  
114 gene conversion events between the *KRTAPI-n* copies, and negative selection. We argue that  
115 these gene conversion events act as an unusual mechanism of purifying selection to prevent  
116 excessive intra-genomic divergence between the four gene copies, while also allowing inter-  
117 species diversity. This unusual mode of evolution may apply to other multicopy genes that  
118 encode products subject to diversifying selection.

119

120

## 121 **MATERIALS AND METHODS**

122 **Sequence Resources and Gene Identification:** All genome sequences were sourced from  
123 the NCBI GenBank. Previously identified *KRTAPI-n* sequences (ITENGE-MWEZA *et al.* 2007;  
124 WU *et al.* 2008; GONG *et al.* 2010a; GONG *et al.* 2011) were used to search the genomes of  
125 cattle, horses, rabbits and African elephants using BLAST with default parameters, and the

126 genes retrieved were identified by sequence identity within both the coding and flanking  
127 regions (**Table S1**).

128 **Sequence Alignments:** *KRTAPI* nucleotide sequences (**Table S1**) for all four paralogs from  
129 the ten species (sheep, cattle, dog, elephant, horse, human, macaque, mouse, rat and rabbit)  
130 were separated into 5' flanking regions, coding sequences, and 3' flanking regions. The  
131 multiple sequence alignment tool *mafft* (v7.123b) (KATOH AND STANDLEY 2013) was used to  
132 separately align the 5' and 3' flanking regions as nucleotide sequences, using the arguments '-  
133 -nuc --localpair --maxiterate 1000'. To align the coding sequences at the predicted amino acid  
134 level, *mafft* with the arguments '--amino --localpair --maxiterate 1000' was run.

135 The coding sequence alignment was subsequently reverse translated using *revTrans* (v1.4)  
136 (WERNERSSON AND PEDERSEN 2003) with two input files: the sequences of all the coding  
137 regions, and the amino acid sequence alignments. The sequences in the two files were paired  
138 by name using the '-match name' parameter, and default values were used for all other  
139 parameters. A number of regions align poorly and have many indels, therefore we used the  
140 longest continuous coding sequence block (198 nucleotides; covers on average around 40%  
141 of the coding region) where none of the 40 sequences had indels. For the flanking region  
142 alignments, we used *Gblocks* (v0.91b) (TALAVERA AND CASTRESANA 2007) to select blocks  
143 that cover approximately 40% of the flanking regions having the best alignment. We also  
144 used *Gblocks* with less stringent criteria to create multiple sequence alignments of the coding  
145 and flanking regions that included more poorly aligning regions.

146 **Phylogenetic Trees:** *PhyML* (v3.1) (GUINDON *et al.* 2010) was used to construct phylogenies  
147 based on the coding and flanking region sequences. The number of resampled bootstrap data  
148 sets was set to 1000 (parameter '-b 1000'), and the additional arguments '-q -s BEST -o tlr'

149 were employed. The Bioconductor package *ggtree* (v1.9.4) (YU *et al.* 2017) was used to plot  
150 the phylogenies.

151 **Codon Adaptation Index:** The CAIcal server (<http://genomes.urv.es/CAIcal>) (PUIGBO *et al.*  
152 2008) was used to calculate CAI values for the *KRTAP1s*, as well as expected CAI values  
153 from permuted sequences using default parameters and published codon usage data  
154 (NAKAMURA *et al.* 2000).

155 **Motifs in the Coding Sequences:** We used MEME motif finder (v4.12.0) (BAILEY *et al.*  
156 2006) to explore repetitive elements in the coding sequences. The repetitive structure of the  
157 coding regions reported in the Results was obtained with parameters ‘-dna -oc . -nostatus -  
158 time 18000 -maxsize 60000 -mod anr -nmotifs 6 -minw 6 -maxw 30 -minsites 20 -maxsites  
159 600 -revcomp’ and all the other parameters set to the default values.

160 ***KRTAP1-n* Polymorphism in Sheep:** Intra-specific variation was assessed using three  
161 sequences for *KRTAP1-1* (ITENGE-MWEZA *et al.* 2007), eleven sequences for *KRTAP1-2*  
162 (GONG *et al.* 2011; GONG *et al.* 2015), nine sequences for *KRTAP1-3* (ITENGE-MWEZA *et al.*  
163 2007), and nine sequences for *KRTAP1-4* (GONG *et al.* 2010a). These were aligned using  
164 DNAMAN (v5.2.10; Lynnon BioSoft, Canada) with default parameters, and polymorphic  
165 sites were identified manually.

166 **Data Availability Statement:** Sequence data are available at GenBank and the accession  
167 numbers and positions are listed in the **Materials and Methods** (sheep polymorphism data)  
168 and Table S1 (*KRTAP1* sequences). Descriptions of the supplemental material are given in  
169 **File S1**, and the supplemental material are available on figshare.

170

171



172

## RESULTS

### 173 **Mammalian *KRTAPI-n* repeats show a concerted evolution pattern in the coding but** 174 **not the flanking regions**

175 To better understand the genetic architecture of the mammalian *KRTAPI* cluster, we selected  
176 the *KRTAPI* genomic region from key members of the mammalian phylogeny for analysis.  
177 The Basic Local Alignment Search Tool (BLAST) was used to search GenBank with known  
178 *KRTAPI-n* sequences to identify and retrieve the *KRTAPI* clusters from the genomes of four  
179 species (cattle, horses, rabbits and African elephants) for whom *KRTAPI-n* sequence  
180 information has not been reported (**Figure S1**). We then combined these with previously-  
181 identified *KRTAPI-n* sequences from other mammalian species to obtain sampling across the  
182 mammalian phylogeny (**Figure 2**).

183 Previously, the *KRTAPI* genes of sheep were shown to contain a variable number of  
184 occurrences of a QTSCCQPXXX decapeptide tandem repeat in the N-terminal region of the  
185 protein (ROGERS *et al.* 1994; GONG *et al.* 2011; GONG *et al.* 2016). We used a motif finding  
186 tool (MEME; (BAILEY *et al.* 2006) to search for repetitive motifs in the coding regions of all  
187 the mammalian *KRTAPI-n* sequences. This revealed that the decapeptide repeat is present at  
188 the N-terminus in all mammalian *KRTAPI-n* genes we obtained (**Figure S2**), albeit with less  
189 amino acid conservation than that observed in sheep. MEME also identified nucleotide level  
190 tandem copies of this repeat at the C-terminus of the protein. Furthermore, both the N- and C-  
191 terminal repeats vary in copy number, within and between genomes. This copy number  
192 variation is responsible for much of the length variation between *KRTAPI-n* sequences.

193 To determine the genetic relationships between of the mammalian *KRTAPI-n* genes, we  
194 generated a *KRTAPI* phylogenetic tree from an alignment of our mammalian *KRTAPI-n*  
195 coding region sequences. This revealed that, in most cases, the *KRTAPI* genes are more

196 related to each other within a species than to their orthologs in other species, thus exhibiting a  
197 concerted evolution pattern. This manifests as clades that group by species, rather than by  
198 repeat, in the phylogenetic tree (**Figure 3**). This concerted evolution pattern breaks down  
199 between the most closely-related species pairs (cattle/sheep, rat/mouse, human/monkey),  
200 presumably because the signal is confounded by these species having more recent shared  
201 ancestry. Nevertheless, for most species there is a clear pattern of concerted evolution.

202 For concertedly evolving tandem repeat sequences such as the ribosomal RNA gene repeats,  
203 homogenization occurs for the complete repeat unit, including the non-coding regions  
204 (GANLEY AND KOBAYASHI 2007). To test whether the *KRTAPI* clusters display a ‘whole-  
205 unit’ pattern of concerted evolution, we generated *KRTAPI* phylogenetic trees from multiple  
206 alignments of the 5’ and 3’ flanking sequences of the mammalian *KRTAPI* genes.

207 Surprisingly, the phylogenies derived from these flanking sequences did not show any pattern  
208 of concerted evolution, and in contrast to the coding region phylogeny, the clades in these  
209 phylogenetic trees were group by *KRTAPI* repeat number, not by species (**Figure 3**). We  
210 note that bootstrap support is not strong for all the clades in these phylogenetic trees, but the  
211 contrast between the coding region concerted versus flanking region radiating evolutionary  
212 patterns is unmistakable. Furthermore, the topology within many of the *KRTAPI* flanking  
213 region clades is consistent with the reported mammalian phylogeny (refer to **Figures 2** and  
214 **3**). These phylogenies were generated from multiple sequence alignments that encompass the  
215 regions that align well, but phylogenies derived from sequence alignments that include poorly  
216 aligned regions give qualitatively similar results (**Figure S3**). Overall, in stark contrast to the  
217 coding region, the flanking regions show a phylogenetic pattern expected for normal  
218 radiating evolution, and exhibit no evidence of concerted evolution.

219

220 **What is responsible for the different evolutionary patterns of the *KRTAP1* coding and**  
221 **flanking regions?**

222 The difference in evolutionary pattern between the coding and flanking regions is striking,  
223 hence we sought to identify the mechanism(s) responsible.

224 **Purifying selection:** Previous studies have shown that multi-gene loci undergoing birth-and-  
225 death evolution can show high levels of identity within the coding region due to strong  
226 purifying selection (NEI *et al.* 2000; PIONTKIVSKA *et al.* 2002). It is possible that purifying  
227 selection maintains sequence identity between *KRTAP1-n* copies within a species, whilst  
228 diversifying selection results in differences between species. If so, we would predict that  
229 while the non-synonymous sites would show a concerted evolution pattern, the synonymous  
230 sites would instead show a normal radiating pattern of evolution (resembling the flanking  
231 regions).

232 To investigate this, we looked at the pattern of evolution of the synonymous sites in the  
233 coding sequences compared to the non-synonymous sites. The number of KAP1 amino acid  
234 changes present within and between species makes it difficult to consistently call sites as  
235 synonymous or non-synonymous, so third codon positions were used as a proxy for  
236 synonymous sites, and first and second codon positions were used as a proxy for non-  
237 synonymous sites. We generated phylogenetic trees from multiple sequence alignments of the  
238 first-second (which we refer to as “non-synonymous”), and third (which we refer to as  
239 “synonymous”) codon sites of the *KRTAP1-n* coding regions to test for different evolutionary  
240 patterns. Surprisingly, while the non-synonymous sites displayed a pattern of concerted  
241 evolution as was expected (**Figure 4A**), the synonymous sites also revealed the same pattern  
242 of concerted evolution (**Figure 4B**). The concerted evolution pattern for the synonymous  
243 sites seems to be stronger than that of the non-synonymous sites, as they separate sheep and

244 cattle into separate clades, and also resolve dog, elephant, and rat/mouse into separate clades  
245 (**Figure 4**).

246 **Codon adaptation:** We considered whether this pattern of concerted evolution amongst the  
247 synonymous sites might result from codon adaptation (LIN *et al.* 2006), as a result of  
248 synonymous mutations being selected to follow changes in the favoured codons between  
249 species. The *KRTAPI-n* genes display strong evidence for codon adaptation (the degree to  
250 which the favoured codons for that species are used in a gene). For example, the human  
251 *KRTAPI-n* genes collectively show a codon adaptation index (CAI) of 0.91 (out of a  
252 maximum of 1), higher than the CAI of randomly permuted human *KRTAPI* sequences  
253 (CAI=0.78). Using the *KRTAPI* coding sequence alignment used for the phylogenies  
254 presented in **Figure 3**, we identified nine synonymous differences between human and mouse  
255 that exhibit a concerted evolution pattern (similarity within species versus difference between  
256 species). If codon adaptation can explain this pattern, these synonymous mutations should  
257 change in a manner consistent with a change in codon usage preference for that amino acid.  
258 Five of these mutations show the pattern expected, given the change in codon usage between  
259 human and mouse (synonymous change creates the more favoured codon in the species it is  
260 found in). However, four of these mutations show the opposite pattern, and most of the codon  
261 usage preference changes between human and mouse are small (**Table S2**). These results  
262 provide no evidence for adaptation to different codon usage preferences driving the pattern of  
263 *KRTAPI* concerted evolution.

264 **Reverse transcription of *KRTAPI* mRNA:** Another potential explanation for the  
265 incongruence in evolutionary pattern between the *KRTAPI* coding and flanking regions is  
266 reverse transcription of *KRTAPI-n* mRNAs, followed by homologous recombination-  
267 mediated replacement of a genomic *KRTAPI-n* with the reverse transcribed copy

268 (COULOMBE-HUNTINGTON AND MAJEWSKI 2007). This is feasible given that *KRTAPI-n* are  
269 single-exon genes. If reverse transcription events occur, the 5' and particularly 3' flanking  
270 regions should show a concerted evolution pattern that is similar to the coding region.  
271 Inspection of the 5' and 3' flanking regions revealed that sequence similarity between  
272 *KRTAPI-n* sequences within a genome tends to decay immediately upstream of the ATG  
273 codon and downstream of the stop codon (**Figure 5**). This suggests that reverse  
274 transcription/integration of *KRTAPI-n* mRNA is unlikely to explain the pattern of *KRTAPI*  
275 concerted evolution, as the transcribed flanking regions of the gene would be expected to  
276 'hitch-hike' with the coding regions through such a mechanism.

277 We also considered whether the *KRTAPI-n* sequences might have arisen through a pure  
278 birth-and-death process by independent gene duplication events. However, we think this is  
279 improbable as it would require the same number of duplications to occur in at least seven of  
280 the species, and, independently, that each of these duplications would not involve any  
281 flanking sequence (including promoter and terminator sequences) and have inserted into the  
282 same site in each species.

283 **Gene conversion:** Finally, we considered whether gene conversion could explain the pattern  
284 of *KRTAPI* repeat evolution. Gene conversion events within a genome that convert a section  
285 of one repeat to the sequence of another can create homogeneity (CHEN *et al.* 2007), and the  
286 degree of homogeneity depends on the relative rates of gene conversion and mutation  
287 (TESHIMA AND INNAN 2004; HARPAK *et al.* 2017). Our results imply that if gene conversion  
288 does occur, it is somehow restricted to the coding region. This pattern could occur if there is  
289 selective pressure to maintain a degree of intra-genome homogeneity between the repeat  
290 copies. If so, under the assumption that gene conversion occurs in both the coding and  
291 flanking regions, those events occurring in the flanking region will not have a selective

292 advantage, while those occurring in the coding region will. Therefore, the probability of gene  
293 conversion events becoming fixed in the population will be greater for events that involve the  
294 coding region. There is considerable intra-genomic variation between *KRTAPI* repeats  
295 (**Figure 3**), but this incomplete level of homogenization can be explained by relatively  
296 infrequent gene conversion events and/or relative infrequent fixation of these events.  
297 Therefore, the sequence features of the *KRTAPI* repeats that we document here can all be  
298 accounted for by gene conversion coupled with selection.

299

### 300 **Evidence for gene conversion events in the *KRTAPI-n* repeats**

301 Inspection of the *KRTAPI* coding region multiple sequence alignment provides evidence for  
302 tracts of gene conversion. Specifically, sites where there are mutations that are shared  
303 between copies within a species, but that differ between species, are frequently clustered  
304 together rather than scattered throughout the gene (**Figure 6**). Such patches of homogeneity  
305 are expected if there has been occasional, short-tract gene conversion events. The patches we  
306 observe are small, but are within the expected range for mammalian gene conversion events  
307 (CHEN *et al.* 2007). In addition, we collected population polymorphism data for *KRTAPI-n*  
308 sequences in sheep, as comprehensive sequence variation data are scarce in other species. For  
309 many of the sites that are polymorphic, the polymorphism is shared across some, or all, of the  
310 *KRTAPI-n* sequences (**Figure 7**). While we cannot rule out independent mutation events in  
311 each *KRTAPI* copy, we think that gene conversion is a more parsimonious explanation for  
312 this observation, particularly for the polymorphisms at synonymous sites. Gene conversion  
313 has also previously been suggested as an explanation for the pattern of polymorphism in the  
314 ovine *KRTAPI* genes (ROGERS *et al.* 1994). Collectively, our results suggest that the unusual  
315 evolutionary pattern of the *KRTAPI* repeats, where the coding region evolutionary dynamics

316 are uncoupled from those of the flanking region, is the result of occasional short-tract gene  
317 conversion events that are selected for in the coding region but not the flanking regions, and  
318 that drive partial homogenization.

319

320

## 321 DISCUSSION

322 Here we have shown that *KRTAPI-n* genes are conserved as a block of four tandem repeats in  
323 mammalian species, and this suggests they derive from a relatively ancient gene-  
324 amplification event or events that probably pre-date mammalian speciation. These four  
325 tandem copies display a strong pattern of concerted evolution in the coding regions, yet the  
326 regions flanking show a normal radiating pattern of evolution. We suggest that this  
327 dichotomous pattern of evolution is not the result of purifying selection acting to retard  
328 changes to the amino acid sequence, but instead results from short gene conversion tracts that  
329 periodically homogenize sequences between the four *KRTAPI* genes.

330 The role of gene conversion is supported by two key pieces of evidence: 1) unique amino  
331 acid tracts that are shared by *KAP1* copies within a species, but are unique to that  
332 species/group of related species; and 2) the possession of shared nucleotide variants between  
333 *KRTAPI* gene copies in sheep populations. These results extend previous reports of  
334 homogenization via ongoing short-tract gene conversion events in other protein coding genes  
335 (NOONAN *et al.* 2004; LAMPING *et al.* 2017).

336 We propose that gene conversion is being utilized as an unusual form of purifying selection  
337 that prevents accumulation of too much divergence between *KRTAPI* gene copies. We  
338 speculate that homogeneity of the *KRTAPI* coding sequences is beneficial as it enables the

339 production of more homogenous components of the hair and wool fibre matrix, and thus  
340 potentially facilitates better associations with the keratin intermediate filaments. We cannot,  
341 however, rule out the possibility that individual *KRTAPI* repeats might have functional  
342 differences, the signal of which is overwhelmed by the concerted evolution signal from the  
343 majority of the gene. However, we note that, particularly in dogs, some of the *KRTAPI-n*  
344 genes are very similar in sequence. Therefore, we favour the explanation that *KRTAPI*  
345 concerted evolution results from ongoing, stochastic gene conversion events coupled with  
346 selection within the coding region against inter-repeat heterogeneity.

347 Purifying selection is evident in the *KRTAPI-n* coding regions, as the rate of synonymous  
348 change is about twice that of the non-synonymous rate (**Figure 4**). While this may seem to  
349 contradict the similarity in the synonymous and non-synonymous concerted evolution tree  
350 topologies, it can be simply explained by purifying selection acting on residues that are  
351 conserved between species, and thus not contributing to the synapomorphies that influence  
352 the tree topologies. Any gene conversion events that homogenize unfavourable amino acids  
353 will be selected against, thereby preventing deleterious mutations from spreading between  
354 copies. However, this same process also allows tolerable and advantageous amino acid  
355 changes to sweep through the copies (DOVER 1982). The *KRTAPI-n* sequences from closely  
356 related species (i.e. human and macaque, rat and mouse, sheep and cattle) were not separated  
357 into different clades for most of the phylogenetic trees we generated (**Figures 3 and 4**). This  
358 suggests that the rate of homogenization is relatively slow, and insufficient to drive  
359 substantial homogeneity over the evolutionary time frames separating these species pairs. In  
360 this context, the shared polymorphisms that we observe in sheep (that are evidence for gene  
361 conversion events) are likely intermediate stages in the accumulation of homogenized  
362 *KRTAPI-n* sequences.



363 The sharp border between a concerted evolution pattern in the coding region and a radiating  
364 evolution pattern in the immediate flanking regions is striking. This can partially be explained  
365 by the selection for gene conversion events within the coding region, as we have proposed.  
366 However, it is intriguing to speculate that this may also be a consequence of differential  
367 expression between the *KRTAPI* genes that is mediated by copy-specific differences in the  
368 regulatory regions. Although not direct, some evidence for differential regulation of  
369 *KRTAPI-n* gene expression was found in two transcriptome studies looking for differentially  
370 expressed genes (FAN *et al.* 2013; CHANG *et al.* 2014). If the *KRTAPI-n* genes do have  
371 functionally distinct roles, gene conversion events in the *KRTAPI* regulatory regions that  
372 perturb their differential regulation may be maladaptive and therefore selected against. Thus,  
373 selective pressure for coding region homogeneity versus regulatory region diversity, coupled  
374 with ongoing gene conversion, may be a powerful way to achieve the dichotomy in  
375 evolutionary patterns we observe. Clearly, a better understanding of the transcriptional  
376 regulation of the *KRTAPI* genes is required to address this hypothesis.

377 Gene conversion is frequently viewed through the lens of impeding sub-functionalization of  
378 gene duplicates. This view is consistent with the well characterized case of the opsin gene  
379 duplicates in primates, where there is a much stronger signal of gene conversion/concerted  
380 evolution in the introns, than in the exons (SHYUE *et al.* 1994; HIWATASHI *et al.* 2011). The  
381 interpretation is that selection has largely rejected gene conversion events that include the  
382 coding (exon) regions, whilst allowing those occurring in the non-coding (intron) regions to  
383 spread in the population (SHYUE *et al.* 1994). This is the opposite of what we observe, and  
384 illustrates how gene conversion and selection can intersect to produce a constellation of  
385 evolutionary patterns: homogenization of the non-coding but not the coding regions in the  
386 opsin paralogs (SHYUE *et al.* 1994); homogenization of the coding but not the non-coding

387 regions in the *KRTAPI* genes (this study); and homogenisation of both coding and non-  
388 coding regions equally in the ribosomal RNA gene repeats (GANLEY AND KOBAYASHI 2007).

389 The extent to which gene conversion acts to homogenize gene duplicates remains  
390 controversial (GAO AND INNAN 2004; CASOLA *et al.* 2012; HARPAK *et al.* 2017). Furthermore,  
391 even in examples where recurrent gene conversion events can be detected, they are often not  
392 sufficient to produce a strong concerted evolution pattern (PETRONELLA AND DROUIN 2011;  
393 PETRONELLA AND DROUIN 2014). There are two potential explanations for why such a strong  
394 pattern of concerted evolution is observed in the case the *KRTAPI* genes, despite the  
395 relatively high levels of divergence between copies. First, unlike many of the examples that  
396 have aroused controversy (GAO AND INNAN 2004; CASOLA *et al.* 2012; HARPAK *et al.* 2017),  
397 the *KRTAPI-n* repeats are tandemly-arranged. Proximity effects as a consequence of tandem  
398 arrangement may increase the chances of unequal alignment of the repeats during DNA  
399 repair-based homologous recombination compared to dispersed repeats, and thus may  
400 increase the chances of conversion events. However, this does not explain examples where  
401 tandemly repeated paralogs do not show a strong concerted evolution pattern (NEI *et al.* 2000;  
402 PERINA *et al.* 2011). A second explanation relates to the imperfect decapeptide tandem repeat  
403 motif found in the coding region. Variation in the copy number of decapeptide repeats  
404 between *KRTAPI* genes is possibly the result of unequal recombination (LIAO AND WEINER  
405 1995; GANLEY AND SCOTT 1998; MORRILL *et al.* 2016). If so, the *KRTAPI* genes may  
406 harbour a recombination hotspot that drives both decapeptide repeat copy number variation  
407 and gene conversion at higher than average levels.

408 Repeats are ubiquitous denizens of eukaryote genomes, where they exist in different forms  
409 (coding, non-coding) and organizations (tandem, dispersed). Our results add to the growing  
410 list of examples that illustrate how different molecular and evolutionary processes can

411 impinge on repeats to structure their sequences and create distinctive patterns of evolution  
412 (SHYUE *et al.* 1994; NOONAN *et al.* 2004; GANLEY AND KOBAYASHI 2007; STORZ *et al.* 2007;  
413 HIWATASHI *et al.* 2011; LAMPING *et al.* 2017). However, it is unclear how widespread these  
414 sorts of evolutionary dynamics are for eukaryotic gene repeats, largely because the patterns of  
415 evolution have not been investigated for the vast majority of multi-gene families. The  
416 increasing availability of high quality genome sequences for a wide range of eukaryotes puts  
417 us in an excellent position to determine, on a much more systematic and wide-ranging basis,  
418 the patterns of repeat sequence dynamics and evolution. This will, in turn, make it clear  
419 whether the impact of recombination on the *KRTAP*s is unusual, or highlights a common  
420 mechanism to finely scale patterns of homogeneity and divergence between repeat copies  
421 over time.

422

423

424

## ACKNOWLEDGEMENTS

425 This work was supported by a Marsden Fund award (14-MAU-053) to ARDG, an  
426 AGMARDT Postdoctoral Fellowship to HG, and a Vernon Willey Trust Fellowship to HZ.

427

428

429

## REFERENCES

- 430 Bailey, T. L., N. Williams, C. Misleh and W. W. Li, 2006 MEME: discovering and analyzing DNA  
431 and protein sequence motifs. *Nucleic Acids Res.* 34: W369-373.
- 432 Britten, R. J., and D. E. Kohne, 1968 Repeated sequences in DNA. *Science* 161: 529-540.

- 433 Brown, D. D., P. C. Wensink and E. Jordan, 1972 A comparison of the ribosomal DNA's of  
434 *Xenopus laevis* and *Xenopus mulleri*: the evolution of tandem genes. J. Mol. Biol. 63:  
435 57-73.
- 436 Casola, C., G. C. Conant and M. W. Hahn, 2012 Very low rate of gene conversion in the yeast  
437 genome. Mol. Biol. Evol. 29: 3817-3826.
- 438 Chang, T. H., H. D. Huang, W. K. Ong, Y. J. Fu, O. K. Lee *et al.*, 2014 The effects of actin  
439 cytoskeleton perturbation on keratin intermediate filament formation in  
440 mesenchymal stem/stromal cells. Biomaterials 35: 3934-3944.
- 441 Chen, J. M., D. N. Cooper, N. Chuzhanova, C. Ferec and G. P. Patrinos, 2007 Gene  
442 conversion: mechanisms, evolution and human disease. Nature Reviews Genetics 8:  
443 762-775.
- 444 Coulombe-Huntington, J., and J. Majewski, 2007 Characterization of intron loss events in  
445 mammals. Genome Res. 17: 23-32.
- 446 Dover, G. A., 1982 Molecular drive: a cohesive mode of species evolution. Nature 299: 111-  
447 117.
- 448 Eirin-Lopez, J. M., L. Rebordinos, A. P. Rooney and J. Rozas, 2012 The birth-and-death  
449 evolution of multigene families revisited. Genome Dynamics 7: 170-196.
- 450 Elder, J. F., Jr., and B. J. Turner, 1995 Concerted evolution of repetitive DNA sequences in  
451 eukaryotes. Q. Rev. Biol. 70: 297-320.
- 452 Fan, R., J. Xie, J. Bai, H. Wang, X. Tian *et al.*, 2013 Skin transcriptome profiles associated with  
453 coat color in sheep. BMC Genomics 14: 389.
- 454 Ganley, A. R. D., and T. Kobayashi, 2007 Highly efficient concerted evolution in the  
455 ribosomal DNA repeats: total rDNA repeat variation revealed by whole-genome  
456 shotgun sequence data. Genome Res. 17: 184-191.

- 457 Ganley, A. R. D., and B. Scott, 1998 Extraordinary ribosomal spacer length heterogeneity in a  
458 *Neotyphodium* endophyte hybrid: implications for concerted evolution. *Genetics*  
459 150: 1625-1637.
- 460 Gao, L.-Z., and H. Innan, 2004 Very low gene duplication rate in the yeast genome. *Science*  
461 306: 1367-1370.
- 462 Gong, H., H. Zhou, R. H. J. Forrest, S. Li, J. Wang *et al.*, 2016 Wool keratin-associated protein  
463 genes in sheep—a review. *Genes* 7: 24.
- 464 Gong, H., H. Zhou and J. G. H. Hickford, 2010a Polymorphism of the ovine keratin-associated  
465 protein 1-4 gene (KRTAP1-4). *Mol. Biol. Rep.* 37: 3377-3380.
- 466 Gong, H., H. Zhou, S. Hodge, J. M. Dyer and J. G. H. Hickford, 2015 Association of wool traits  
467 with variation in the ovine KAP1-2 gene in Merino cross lambs. *Small Rumin. Res.*  
468 124: 24-29.
- 469 Gong, H., H. Zhou, G. W. McKenzie, J. G. Hickford, Z. Yu *et al.*, 2010b Emerging issues with  
470 the current keratin-associated protein nomenclature. *International Journal of*  
471 *Trichology* 2: 104-105.
- 472 Gong, H., H. Zhou, G. W. McKenzie, Z. Yu, S. Clerens *et al.*, 2012 An updated nomenclature  
473 for keratin-associated proteins (KAPs). *Int. J. Biol. Sci.* 8: 258-264.
- 474 Gong, H., H. Zhou, Z. Yu, J. Dyer, J. E. Plowman *et al.*, 2011 Identification of the ovine  
475 keratin-associated protein KAP1-2 gene (KRTAP1-2). *Exp. Dermatol.* 20: 815-819.
- 476 Guindon, S., J. F. Dufayard, V. Lefort, M. Anisimova, W. Hordijk *et al.*, 2010 New algorithms  
477 and methods to estimate maximum-likelihood phylogenies: assessing the  
478 performance of PhyML 3.0. *Syst. Biol.* 59: 307-321.

- 479 Harpak, A., X. Lan, Z. Gao and J. K. Pritchard, 2017 Frequent nonallelic gene conversion on  
480 the human lineage and its effect on the divergence of gene duplicates. PNAS 114:  
481 12779-12784.
- 482 Hiwatashi, T., A. Mikami, T. Katsumura, B. Suryobroto, D. Perwitasari-Farajallah *et al.*, 2011  
483 Gene conversion and purifying selection shape nucleotide variation in gibbon L/M  
484 opsin genes. BMC Evol. Biol. 11: 312.
- 485 Itenge-Mweza, T. O., R. H. Forrest, G. W. McKenzie, A. Hogan, J. Abbott *et al.*, 2007  
486 Polymorphism of the KAP1.1, KAP1.3 and K33 genes in Merino sheep. Mol. Cell.  
487 Probes 21: 338-342.
- 488 Katoh, K., and D. M. Standley, 2013 MAFFT multiple sequence alignment software version 7:  
489 improvements in performance and usability. Mol. Biol. Evol. 30: 772-780.
- 490 Khan, I., E. Maldonado, V. Vasconcelos, J. O. Stephen, W. E. Johnson *et al.*, 2014 Mammalian  
491 keratin associated proteins (KRTAPs) subgenomes: disentangling hair diversity and  
492 adaptation to terrestrial and aquatic environments. BMC Genomics 15: 779.
- 493 Lamping, E., J. Y. Zhu, M. Niimi and R. D. Cannon, 2017 Role of ectopic gene conversion in  
494 the evolution of a *Candida krusei* pleiotropic drug resistance transporter family.  
495 Genetics 205: 1619-1639.
- 496 Li, S., H. Zhou, H. Gong, F. Zhao, J. Hu *et al.*, 2017a Identification of the ovine keratin-  
497 associated protein 26-1 gene and its association with variation in wool traits. Genes  
498 8: 225.
- 499 Li, S., H. Zhou, H. Gong, F. Zhao, J. Wang *et al.*, 2017b Identification of the ovine keratin-  
500 associated protein 22-1 (KAP22-1) gene and its effect on wool traits. Genes 8: 27.

- 501 Li, S., H. Zhou, H. Gong, F. Zhao, J. Wang *et al.*, 2017c Variation in the ovine KAP6-3 gene  
502 (KRTAP6-3) is associated with variation in mean fibre diameter-associated wool  
503 traits. *Genes* 8: 204.
- 504 Liao, D., 1999 Concerted evolution: molecular mechanism and biological implications. *Am. J.*  
505 *Hum. Genet.* 64: 24-30.
- 506 Liao, D., and A. M. Weiner, 1995 Concerted evolution of the tandemly repeated genes  
507 encoding primate U2 small nuclear RNA (the RNU2 locus) does not prevent rapid  
508 diversification of the (CT)<sub>n</sub>.(GA)<sub>n</sub> microsatellite embedded within the U2 repeat unit.  
509 *Genomics* 30: 583-593.
- 510 Lin, Y.-S., J. K. Byrnes, J.-K. Hwang and W.-H. Li, 2006 Codon-usage bias versus gene  
511 conversions in the evolution of yeast duplicate genes. *PNAS* 103: 14412-14416.
- 512 Lopez-Flores, I., and M. A. Garrido-Ramos, 2012 The repetitive DNA content of eukaryotic  
513 genomes. *Genome Dynamics* 7: 1-28.
- 514 Morrill, S. A., A. E. Exner, M. Babokhov, B. I. Reinfeld and S. M. Fuchs, 2016 DNA instability  
515 maintains the repeat length of the yeast RNA polymerase II C-terminal domain. *J.*  
516 *Biol. Chem.* 291: 11540-11550.
- 517 Nakamura, Y., T. Gojobori and T. Ikemura, 2000 Codon usage tabulated from international  
518 DNA sequence databases: status for the year 2000. *Nucleic Acids Res.* 28: 292.
- 519 Nei, M., X. Gu and T. Sitnikova, 1997 Evolution by the birth-and-death process in multigene  
520 families of the vertebrate immune system. *PNAS* 94: 7799-7806.
- 521 Nei, M., I. B. Rogozin and H. Piontkivska, 2000 Purifying selection and birth-and-death  
522 evolution in the ubiquitin gene family. *PNAS* 97: 10866-10871.
- 523 Nei, M., and A. P. Rooney, 2005 Concerted and birth-and-death evolution of multigene  
524 families. *Annu. Rev. Genet.* 39: 121-152.

- 525 Noonan, J. P., J. Grimwood, J. Schmutz, M. Dickson and R. M. Myers, 2004 Gene conversion  
526 and the evolution of protocadherin gene cluster diversity. *Genome Res.* 14: 354-366.
- 527 Perina, A., D. Seoane, A. M. Gonzalez-Tizon, F. Rodriguez-Farina and A. Martinez-Lage, 2011  
528 Molecular organization and phylogenetic analysis of 5S rDNA in crustaceans of the  
529 genus *Pollicipes* reveal birth-and-death evolution and strong purifying selection.  
530 *BMC Evol. Biol.* 11: 304.
- 531 Petronella, N., and G. Drouin, 2011 Gene conversions in the growth hormone gene family of  
532 primates: stronger homogenizing effects in the Hominidae lineage. *Genomics* 98:  
533 173-181.
- 534 Petronella, N., and G. Drouin, 2014 Purifying selection against gene conversions in the folate  
535 receptor genes of primates. *Genomics* 103: 40-47.
- 536 Piontkivska, H., A. P. Rooney and M. Nei, 2002 Purifying selection and birth-and-death  
537 evolution in the histone H4 gene family. *Mol. Biol. Evol.* 19: 689-697.
- 538 Powell, B. C., and G. E. Rogers, 1997 The role of keratin proteins and their genes in the  
539 growth, structure and properties of hair, pp. 59-148 in *Formation and structure of*  
540 *human hair*. Birkhäuser Verlag.
- 541 Puigbo, P., I. G. Bravo and S. Garcia-Vallve, 2008 CAIcal: a combined set of tools to assess  
542 codon usage adaptation. *Biol. Direct* 3: 38.
- 543 Richard, G. F., A. Kerrest and B. Dujon, 2008 Comparative genomics and molecular dynamics  
544 of DNA repeats in eukaryotes. *Microbiol. Mol. Biol. Rev.* 72: 686-727.
- 545 Rogers, G. R., J. G. H. Hickford and R. Bickerstaffe, 1994 Polymorphism in two genes for B2  
546 high sulfur proteins of wool. *Anim. Genet.* 25: 407-415.
- 547 Rogers, M. A., L. Langbein, S. Praetzel-Wunder, H. Winter and J. Schweizer, 2006 Human hair  
548 keratin-associated proteins (KAPs). *Int. Rev. Cytol.* 251: 209-263.

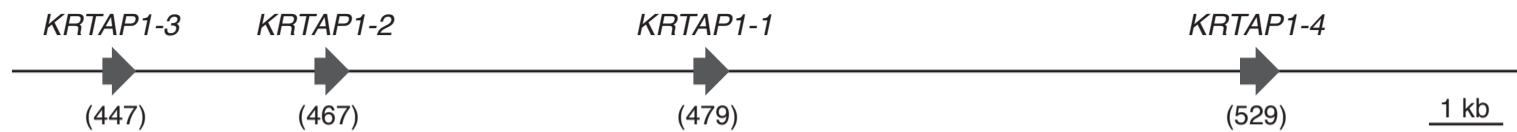


- 549 Rogers, M. A., and J. Schweizer, 2005 Human KAP genes, only the half of it? Extensive size  
550 polymorphisms in hair keratin-associated protein genes. *J. Invest. Dermatol.* 124: vii-  
551 ix.
- 552 Rooney, A. P., and T. J. Ward, 2005 Evolution of a large ribosomal RNA multigene family in  
553 filamentous fungi: birth and death of a concerted evolution paradigm. *PNAS* 102:  
554 5084-5089.
- 555 Shimomura, Y., N. Aoki, J. Schweizer, L. Langbein, M. A. Rogers *et al.*, 2002 Polymorphisms in  
556 the human high sulfur hair keratin-associated protein 1, KAP1, gene family. *J. Biol.*  
557 *Chem.* 277: 45493.
- 558 Shyue, S. K., L. Li, B. H. Chang and W.-H. Li, 1994 Intronic gene conversion in the evolution of  
559 human X-linked color vision genes. *Mol. Biol. Evol.* 11: 548-551.
- 560 Stein, L. D., 2004 End of the beginning. *Nature* 431: 915-916.
- 561 Stephan, W., 1989 Tandem-repetitive noncoding DNA: Forms and forces. *Mol. Biol. Evol.* 6:  
562 198-212.
- 563 Storz, J. F., M. Baze, J. L. Waite, F. G. Hoffmann, J. C. Opazo *et al.*, 2007 Complex signatures  
564 of selection and gene conversion in the duplicated globin genes of house mice.  
565 *Genetics* 177: 481-500.
- 566 Talavera, G., and J. Castresana, 2007 Improvement of phylogenies after removing divergent  
567 and ambiguously aligned blocks from protein sequence alignments. *Syst. Biol.* 56:  
568 564-577.
- 569 Tao, J., H. Zhou, H. Gong, Z. Yang, Q. Ma *et al.*, 2017a Variation in the KAP6-1 gene in  
570 Chinese Tan sheep and associations with variation in wool traits. *Small Rumin. Res.*  
571 154: 129-132.

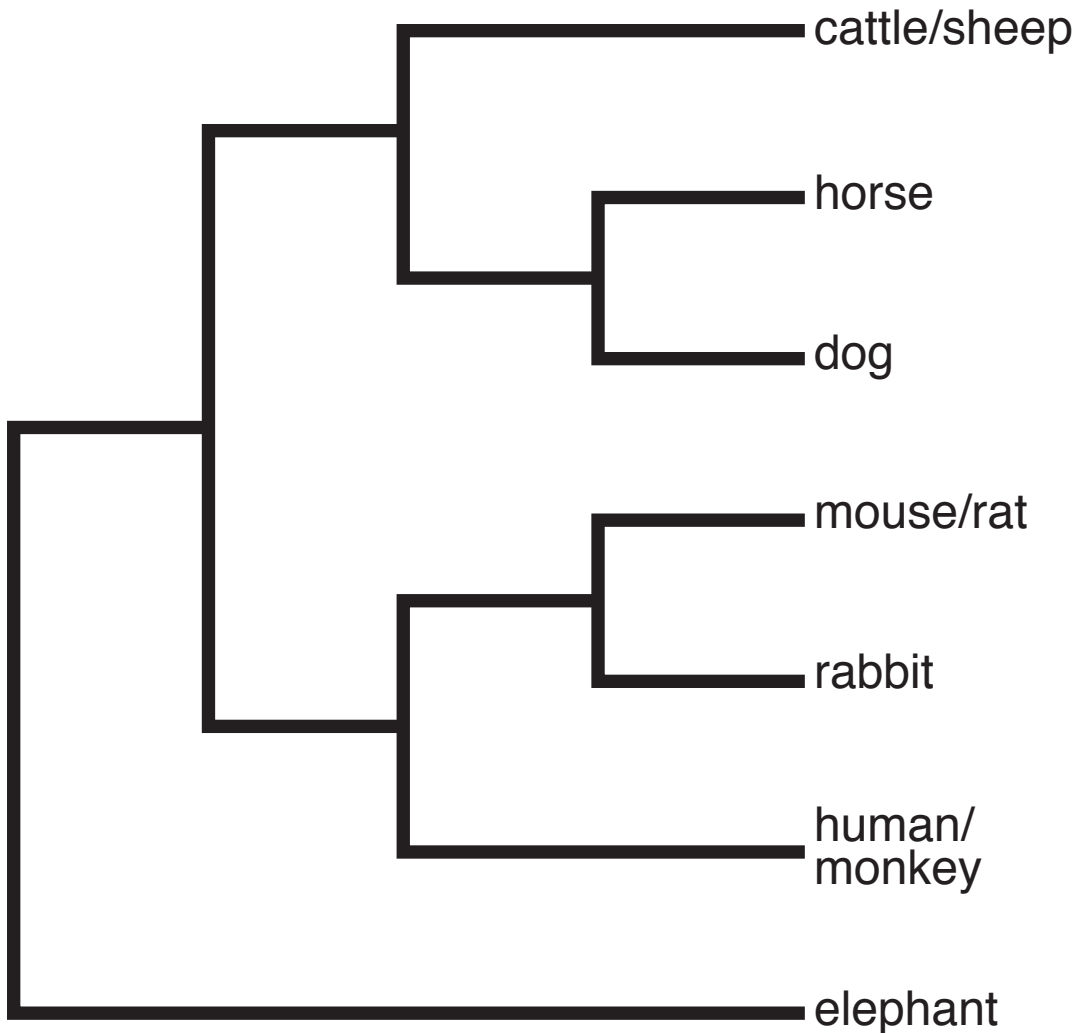
- 572 Tao, J., H. Zhou, Z. Yang, H. Gong, Q. Ma *et al.*, 2017b Variation in the KAP8-2 gene affects  
573 wool crimp and growth in Chinese Tan sheep. *Small Rumin. Res.* 149: 77-80.
- 574 Teshima, K. M., and H. Innan, 2004 The effect of gene conversion on the divergence  
575 between duplicated genes. *Genetics* 166: 1553-1560.
- 576 Torrents, D., M. Suyama, E. Zdobnov and P. Bork, 2003 A genome-wide survey of human  
577 pseudogenes. *Genome Res.* 13: 2559-2567.
- 578 Wernersson, R., and A. G. Pedersen, 2003 RevTrans: Multiple alignment of coding DNA from  
579 aligned amino acid sequences. *Nucleic Acids Res.* 31: 3537-3539.
- 580 Wu, D. D., D. Irwin and Y. P. Zhang, 2008 Molecular evolution of the keratin associated  
581 protein gene family in mammals, role in the evolution of mammalian hair. *BMC Evol.*  
582 *Biol.* 8: 241.
- 583 Yu, G., D. K. Smith, H. Zhu, Y. Guan and T. T.-Y. Lam, 2017 GGTREE: an R package for  
584 visualization and annotation of phylogenetic trees with their covariates and other  
585 associated data. *Methods Ecol. Evol.* 8: 28–36.
- 586 Zhou, H., H. Gong, S. Li, Y. Luo and J. Hickford, 2015 A 57-bp deletion in the ovine KAP6-1  
587 gene affects wool fibre diameter. *Journal of Animal Breeding and Genetics.*
- 588 Zhou, H., H. Gong, J. Wang, J. M. Dyer, Y. Luo *et al.*, 2016 Identification of four new gene  
589 members of the KAP6 gene family in sheep. *Sci. Rep.* 6: 24074.

590

591



**Figure 1. Tandem repeat organization of the keratin associated protein-1 (*KRTAP1*) genes**  
The organization of mammalian *KRTAP1* genes is illustrated by the arrangement found in sheep. The four *KRTAP1-n* paralogs are represented by arrows that indicate the direction of transcription. Diagram is drawn to scale, with *KRTAP1-n* lengths bracketed below the genes. These repeats are numbered *KRTAP1-1*, 3, 4, and 5 in human.

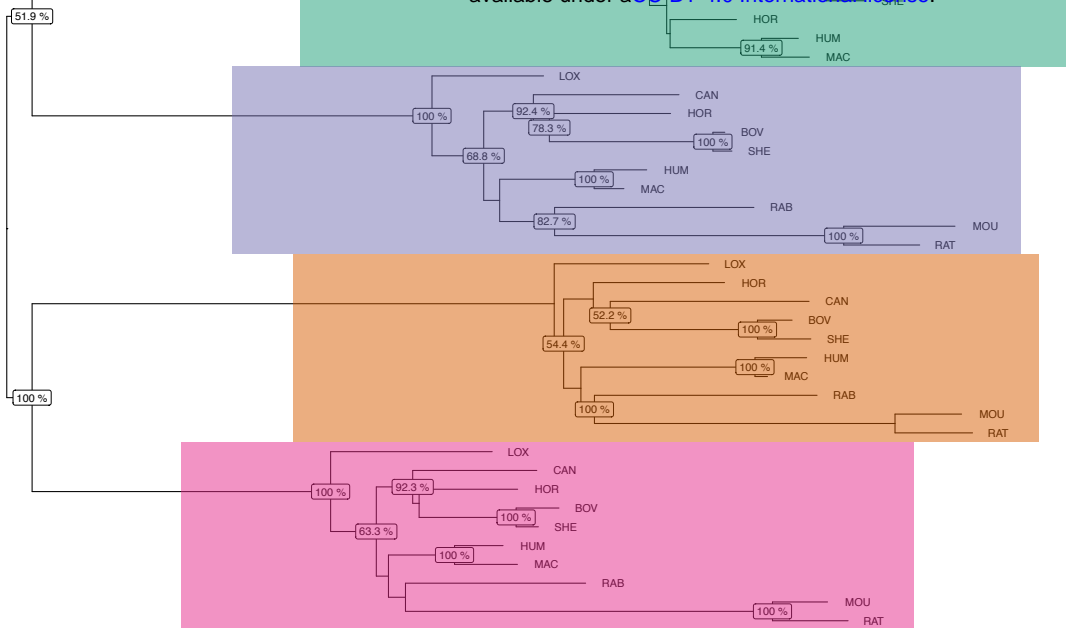


**Figure 2. Mammalian *KRTAP1-n* gene phylogenetic relationships**  
Representative phylogenetic tree illustrating the relationships between the *KRTAP1-n* genes in the species used in this study. Branch lengths are not to scale. The phylogeny is adapted from that presented in McCormack et al. (2012).

0.1

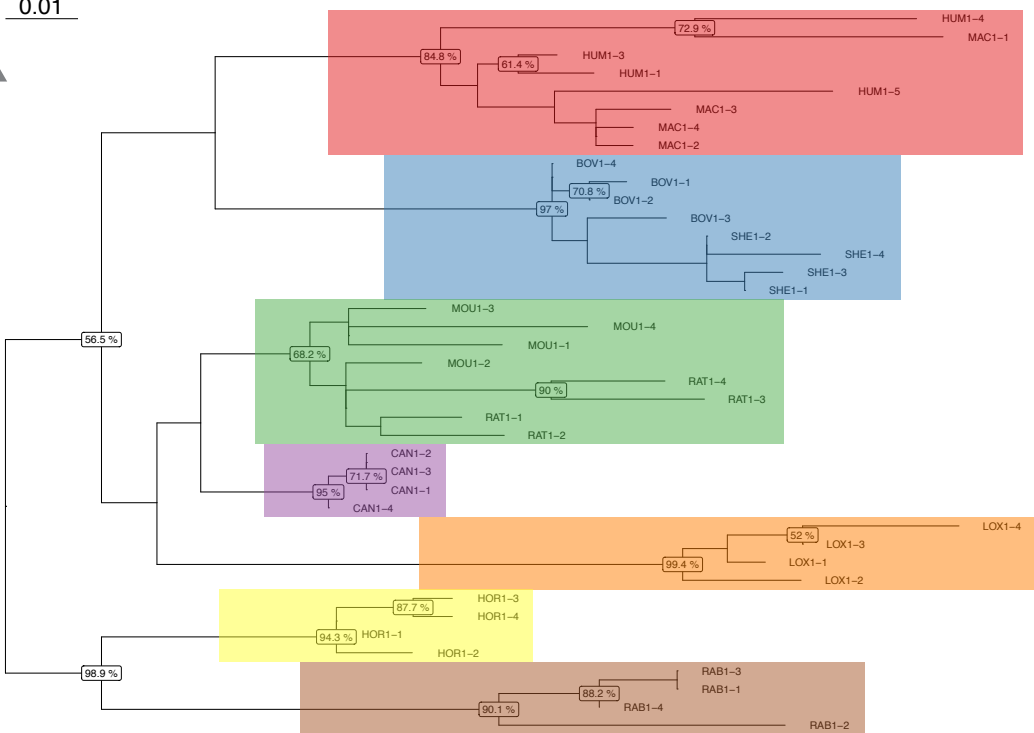
bioRxiv preprint doi: <https://doi.org/10.1101/282418>; this version posted March 14, 2018. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY 4.0 International license.

3-prime flanking



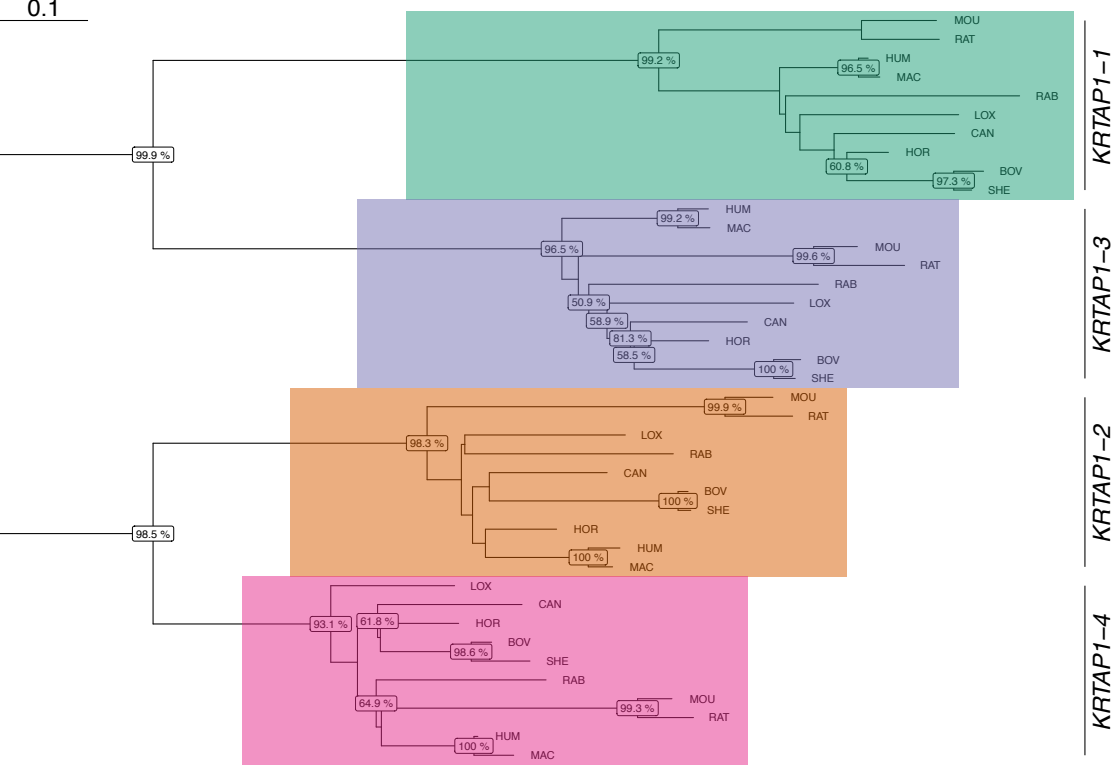
0.01

coding region



0.1

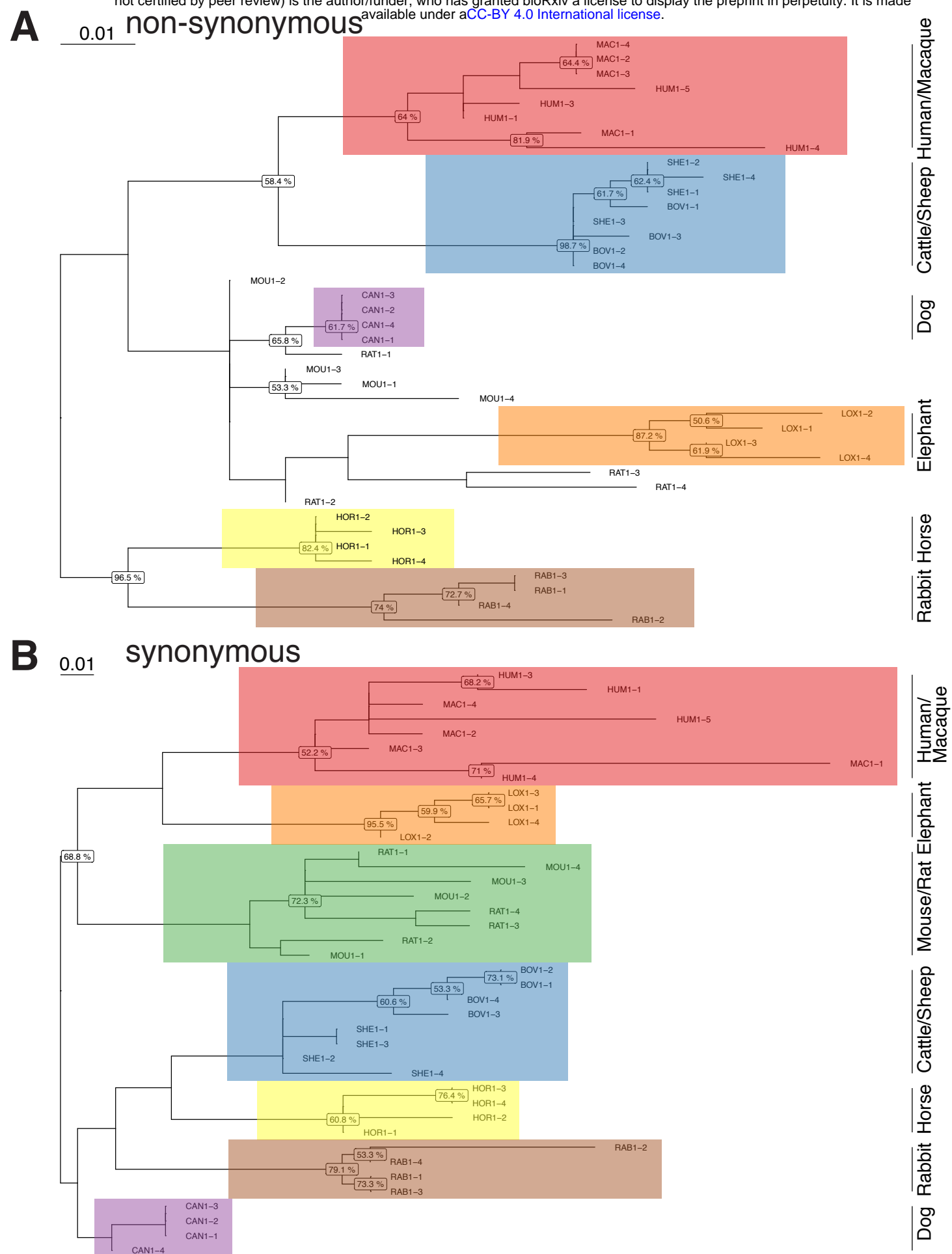
5-prime flanking



### Figure 3. Phylogenetic trees of *KRTAP1-n* coding and flanking region sequences

Phylogenetic trees were constructed for the mammalian *KRTAP1-n* 5' flanking region, coding region, and 3' flanking region using PhyML. The species are indicated by three-letter abbreviations. The number following this for the coding regions indicates the *KRTAP1-n* gene name. The major clades within the trees are indicated by coloured boxes. The 5' and 3' flanking region phylogenies group by repeat number, while the coding region phylogeny tends to group by species. Numbers on nodes indicate bootstrap supports over 50%, and substitution rates are indicated at the top left. Human *KRTAP1-n* gene names have been altered for consistency with other species.

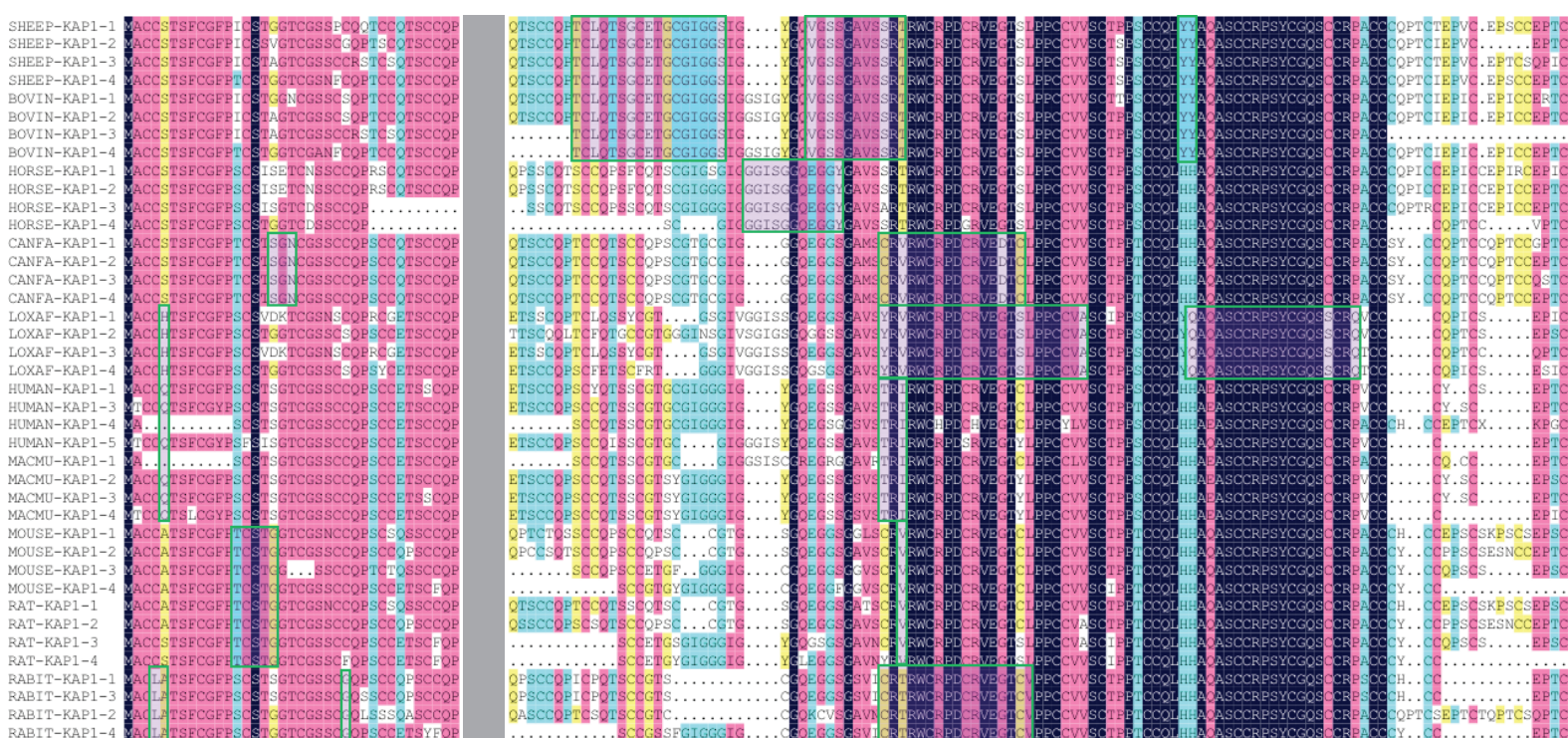
Rabbit Horse Elephant Dog Mouse/Rat Cattle/Sheep Human/Macaque



**Figure 4. The *KRTAP1-n* concerted evolution pattern is not explained by purifying selection**

Phylogenetic trees were constructed for the 1st and 2nd codon positions ("non-synonymous"; **A**), and the 3rd codon position ("synonymous"; **B**), as per **Figure 3**. The major clades in both phylogenies tend to group by species, with this concerted evolution pattern being stronger for the synonymous phylogeny. Numbers on nodes indicate bootstrap supports with values over 50%, and substitution rates are indicated at the top left.





**Figure 6. Evidence for short gene conversion tracts within species**  
 Alignment of KAP1 amino acid sequences from the ten mammalian species. Amino acid tracts boxed in green represent sequences unique to a species or related species pairs. The grey vertical box represents the conserved decapeptide repeat sequences (which have been removed). Dots represent gaps in the alignment.



*KRTAP1-1* ATGGCCTGCTGTTCCACCAGCTTCTGTGGATTTCATCTGT**Y**CCTACTGGTGGGACCTGTGGCTCCAGTCCCTGCCAGCMGACCTGCTGC <30 bp repeat>  
*KRTAP1-2* -----T-----T--T--T--A-----CTG-G-----CA---C----- <30 bp repeat>  
*KRTAP1-3* -----**Y**-----C-----CTG-----GATCA-----A-T <30 bp repeat>  
*KRTAP1-4* -----T--C-----CT--CT-----T--ACTTT-----CA----- <30 bp repeat>

*KRTAP1-1* CAGACCAGTGGCTGTGAGAC**S**GGCTGTGGCATTGGTGGCAGCATTGGYTATGGCCAGGTGGGTAGCAGCGGAGCTGTGAGCAGCCGCACCAGGTGGTGCCGCC  
*KRTAP1-2* -----**V**-----**W**-----C--R-----Y-----  
*KRTAP1-3* -----Y--R--**S**-----**Y**--C-----R--  
*KRTAP1-4* -----**B**-----**W**-----**Y**--RC-----

*KRTAP1-1* CTGACTGCCGCGTGGAGGGCACCAGCCTGCCWCCCTGCTG**Y**GTGGTGGAGCTGCAC**A****Y**CCCCGTCCTGCTGCCAGCTGTACTATGCCAGGCCTCCTGCTGCCG  
*KRTAP1-2* -----T-----**Y**-----T-----R-----  
*KRTAP1-3* -----R--T-----**Y**-----**Y**-----  
*KRTAP1-4* -----T-----**Y**-----T-----

*KRTAP1-1* CCCATCCTACTGTGGACAGTCCTGCTGCCGCCAGCCTGCTGCT**K**CCAGCCCACCTGC**A****Y**TGAGCCC**R**TCTGTGAGCCCAGCTGCTGTGAGCCCACCTGCTGA  
*KRTAP1-2* -----G-----T-----G-----**S**-----.....  
*KRTAP1-3* -----G-----**Y**-----G-----C-----CCC-A---T-M-T-A-  
*KRTAP1-4* -----**K**-----T-----**R**-----**S**-----R-----

### Figure 7. Shared polymorphisms between *KRTAP1-n* sequences in sheep

Alignment of the four sheep *KRTAP1-n* coding region sequences. Dashes represent nucleotides identical to the top sequence, and dots represent gaps. The 30 bp repeats are not shown, as the insertion/deletion positions cannot be precisely determined. Shared nucleotide substitutions between repeat copies are highlighted in red.