



OPEN

Genome sequence and comparative analysis of reindeer (*Rangifer tarandus*) in northern Eurasia

Melak Weldenogdguad^{1,2,6}, Kisun Pokharel^{1,6}, Yao Ming^{3,6}, Mervi Honkatukia^{1,4}, Jaana Peippo¹, Tiina Reilas¹, Knut H. Røed⁵ & Juha Kantanen¹✉

Reindeer are semi-domesticated ruminants that have adapted to the challenging northern Eurasian environment characterized by long winters and marked annual fluctuations in daylight. We explored the genetic makeup behind their unique characteristics by *de novo* sequencing the genome of a male reindeer and conducted gene family analyses with nine other mammalian species. We performed a population genomics study of 23 additional reindeer representing both domestic and wild populations and several ecotypes from various geographic locations. We assembled 2.66 Gb (N50 scaffold of 5 Mb) of the estimated 2.92 Gb reindeer genome, comprising 27,332 genes. The results from the demographic history analysis suggested marked changes in the effective population size of reindeer during the Pleistocene period. We detected 160 reindeer-specific and expanded genes, of which zinc finger proteins (n = 42) and olfactory receptors (n = 13) were the most abundant. Comparative genome analyses revealed several genes that may have promoted the adaptation of reindeer, such as those involved in recombination and speciation (*PRDM9*), vitamin D metabolism (*TRPV5*, *TRPV6*), retinal development (*PRDM1*, *OPN4B*), circadian rhythm (*GRIA1*), immunity (*CXCR1*, *CXCR2*, *CXCR4*, *IFNW1*), tolerance to cold-triggered pain (*SCN11A*) and antler development (*SILT2*). The majority of these characteristic reindeer genes have been reported for the first time here. Moreover, our population genomics analysis suggested at least two independent reindeer domestication events with genetic lineages originating from different refugial regions after the Last Glacial Maximum. Taken together, our study has provided new insights into the domestication, evolution and adaptation of reindeer and has promoted novel genomic research of reindeer.

Reindeer (*Rangifer tarandus*) have pivotal economic, societal, cultural and ecological values for indigenous people and pastoralists in northern and subarctic regions of Eurasia. Reindeer are a source of meat, hide and occasionally milk and have been used for transportation. Reindeer were crucial for the colonization of the northernmost parts of Eurasia and have a central symbolic role for the indigenous Sami, Nenets, and Evenki cultures and several other North Eurasian cultures^{1,2}.

During the evolutionary time scale and in recent demographic history, reindeer were distributed across the northern Eurasian regions, except during the glacial periods³, and adapted to a challenging northern environment characterized by short daylight and limited vegetation for feeding during the long winter and prolonged daylight during the short summer period^{4,5}. The marked fluctuations in daylight may have led to weakened circadian rhythms^{6–8}. Similarly, reindeer exhibit retinal structural adaptation to the extreme seasonal light change to increase retinal sensitivity in dim light⁹. Moreover, reindeer have developed unique features, such as fat metabolism processes, a low resting metabolic rate⁵ and the annual growth and loss of antlers both in males and females¹⁰.

The genetic diversity, population structure, and prehistorical demographic events of *Rangifer* populations have been investigated, and inferences of domestication have been drawn using mitochondrial DNA (mtDNA) D-loop polymorphisms and autosomal microsatellites as genetic markers^{3,11–14}. For example, mtDNA and microsatellite

¹Natural Resources Institute Finland, FI-31600, Jokioinen, Finland. ²Department of Environmental and Biological Sciences, University of Eastern Finland, FI-70201, Kuopio, Finland. ³BGI-Genomics, BGI-Shenzhen, Shenzhen, Guangdong, 518083, China. ⁴Nordic Genetic Resource Centre – NordGen, c/o NMBU – Biovit Box 5003, Ås, NO-1432, Norway. ⁵Department of Basic Sciences and Aquatic Medicine, Norwegian University of Life Sciences, P.O. Box 369 Centrum, 0102, Oslo, Norway. ⁶These authors contributed equally: Melak Weldenogdguad, Kisun Pokharel and Yao Ming. ✉e-mail: juha.kantanen@luke.fi

Sequencing	Insert size	Clean data (Gb)	Genome coverage (×)
Paired-end libraries	170, 500, 800 bp	181.39	60.46
	2, 5, 10, 20 kb	119.16	39.72
	Total	300.55	100.18
Assembly	N50 (Kb)	Longest (Kb)	Size (Gb)
Contig	48.8	441.2	2.62
Scaffold	5,023.6	28,190.2	2.66
Annotation	Number	Total length (Mb)	Percentage of genome
Genes	27,332	787.3	—
Repeats	—	1,012.6	38.02

Table 1. Assembly and annotation statistics.

diversity patterns between domestic and wild Eurasian tundra reindeer (*R. t. tarandus*), domestic and wild boreal forest reindeer (*R. t. fennicus*), and Arctic Svalbard reindeer (*R. t. platyrhynchus*), wild North American Caribou (*R. t. caribou*) and several other populations appear to reflect their refugial origins and colonization of circumpolar regions after the Last Glacial Maximum rather than following the classic division of *Rangifer* populations into three main ecological groups adapted to a particular set of environmental conditions (forest, tundra and arctic environments) and sedentary or migratory life-history strategies^{3,13,15}. Moreover, mtDNA analyses in ancient and modern *Rangifer* populations indicate at least three different domestication events in the prehistory of reindeer and suggest that one of the domestication sites previously assumed may not have been located in northern Fennoscandia^{13,15}.

For most domesticated animal species (e.g., chicken¹⁶, taurine cattle¹⁷, pig¹⁸, sheep¹⁹ and horse²⁰), several genomic tools, such as annotated *de novo* genome assemblies, resequencing data sets, and whole-genome single nucleotide polymorphism (SNP) panels, are available for genomics studies. These whole-genome research approaches have provided critical, new knowledge of the evolution, domestication, genetic diversity, selection patterns and adaptation of domestic animal species. For reindeer genetics studies, however, these genomics tools have become available only recently: Li *et al.* (2017)²¹ reported the *de novo* sequencing of one female reindeer (*R. tarandus*) originating from Inner Mongolia, PR of China. Subsequently, the reindeer reference genome was compared with the genomes of 50 additional ruminant species, and the genomic data from three Chinese domestic reindeer and three Norwegian wild tundra reindeer were resequenced in the study by Lin *et al.* (2019)⁴. The genetic basis of specific *R. tarandus* characteristics, such as adaptation to marked annual regulation in daylight, exceptional vitamin D metabolism and behavioural traits, was identified. In the present study, we provide an alternative draft reference genome for reindeer studies and particularly those representing arctic and sub-arctic regions which show phylogenetic distinctiveness compared to existing reference assemblies^{21–23}. Here, we deep-sequenced and *de novo* assembled the genome of a male reindeer originating from Sodankylä, northern Finland by using the Illumina HiSeq platform. In addition, we analysed the whole transcriptome of six tissues of the reference animal to improve the annotation. Moreover, we resequenced 23 *Rangifer* from Norway and Russia (wild and domestic), together with wild reindeer from Svalbard and Alaska, representing different ecotypes, and report here the largest reindeer genome data sequenced thus far. With these data, we examined the genetic diversity, demographic evolution, adaptation and domestication of reindeer and compared the genome structure of reindeer with that of several other mammalian species, e.g., with several domesticated ruminant species having a *de novo* reference genome sequenced.

Results

Genome assembly and annotation. Genomic DNA extracted from the blood sample of a male reindeer was sequenced using the Illumina HiSeq 2500 and 4000 platforms. Libraries with insert sizes ranging from 170 base pairs (bp) to 20 kilobases (kb) were sequenced and generated a total of 513 gigabases (Gb) (170×) raw sequence data (Supplementary Table S1). After filtering low-quality reads, 300.55 Gb (100×) clean data was retained for assembly (Table 1 and Supplementary Table S1). The genome size of reindeer was estimated to be 2.92 Gb using k-mer analysis²⁴ (Supplementary Table S2 and Supplementary Fig. S1). High-quality reads were assembled using SOAPdenovo V2.04²⁵, and an assembly spanning 2.66 Gb of an estimated 2.92 Gb was generated, with N50 scaffolds and contig sizes of 5 megabases (Mb) and 48.8 kb, respectively (Table 1, Supplementary Table S3). The final assembly comprised 23,450 scaffolds (170 kb–28.2 Mb), which represented 2.62 Gb and 107,910 contigs (200 bp–22.5 kb) that in turn represented an additional 44.9 Mb of the assembled genome. The assemblies covered ~91% of the estimated reindeer genome size. The GC content of the assembled genome was approximately 41.4% (Supplementary Fig. S2 and Fig. S3), which is similar to the GC content of the genomes of other related species^{26,27}.

We evaluated the quality of the draft assembly by aligning high-quality reads from the short insert-size libraries against the assembly using Burrows-Wheeler Alignment (BWA)²⁸ and found that the mapping rate was 98.3%. In addition, assembly quality was assessed by aligning the *de novo* assembled transcriptome against the assembly using BLAT²⁹. Approximately 99% of the assembled transcripts were mapped to the corresponding assembly

with higher than 50% sequence identity, and approximately 97% of the transcripts mapped to the assembly with 90% sequence identity (Supplementary Table S4). Furthermore, to assess the assembly quality, we also used Benchmarking Universal Single-Copy Orthologs (BUSCO)³⁰ genes and found that the assembly contained 94.9% (3,895 of 4,104) complete BUSCOs of which 7.1% were duplicated (Supplementary Table S5).

We predicted a total of 27,332 protein-coding genes with an average of 7.4 exons, 1,305 bp coding sequences (CDS) and 28,808 bp transcripts per gene in the reindeer genome by employing a homology-based and an RNA-assisted approach as described by Curwen *et al.* (2004)³¹ (Table 1, Supplementary Table S7). In total, 26,838 (98.19%) of the total predicted genes were functionally annotated according to public databases (InterPro, Kyoto Encyclopedia of Genes and Genomes (KEGG), Swiss-Prot, TrEMBL), while a total of 494 (1.81%) genes remained unannotated (Supplementary Table S8). To assess the quality of gene prediction, the gene length distribution, the length of the CDS, exons and introns, and the distribution of exon number per gene were compared with those of the five mammalian genomes (see Methods) in a homology-based prediction. No significant differences in exon length and intron length distribution were observed among the six species (Supplementary Fig. S4). Among the compared species, the gene models derived for reindeer were similar to those for dog with respect to all of these key parameters (Supplementary Fig. S4). The annotated gene set was evaluated by BUSCO, and our gene set mapped 96.8% (3,973 of 4,104) complete BUSCOs (Supplementary Table S6), indicating the high quality of our draft genome annotation. We also annotated 6,580 non-coding RNAs (ncRNAs), of which 551 were microRNAs (miRNAs), 329 were rRNAs, 3,994 were transfer RNAs (tRNAs), and 1,706 were small nuclear RNAs (snRNAs) (Supplementary Table S9). Compared with the analyses of cattle, yak, sheep, goat, horse, human and the first reindeer genome assembly, with less than 900 tRNAs predicted, our analysis revealed an exceptionally high number of tRNAs (the analysis was repeated for confirmation)²¹. This unique finding may be due to some unknown bioinformatics bias or could be related to unique characteristics of the Fennoscandian reindeer.

In total, 38.02% of the assembled reindeer genome comprised repetitive sequences (transposable elements (TEs) and tandem repeats; Table 1, Supplementary Table S10), most of which were transposon-derived repeats. All TEs obtained by different methods (Repbase library³², alignment with known transposable-element-related genes, and *de novo* RepeatModeller³³ identification) were combined, and the result indicated that 36.78% of the reindeer genome consisted of TE sequences (Supplementary Table S11), which is similar to the results obtained for the panda (36.2%)²⁶ and dog genomes (36.1%)³⁴, while lower than the results obtained for the cattle (46.5%)¹⁷, goat (42.2%)³⁵ and sheep genomes (42.7%)¹⁹.

Evolution of the reindeer genome. We compared the novel reindeer reference genome to nine different mammalian species (*Camelus dromedarius*, *Capra hircus*, *Ovis aries*, *Bos taurus*, *Bos grunniens*, *Equus caballus*, *Canis familiaris*, *Ursus maritimus* and *Homo sapiens*) to examine the evolution of genes and gene families. The *R. tarandus* genome revealed 17,709 orthologous gene families that were shared by at least two species (pairwise comparison to reindeer) and a total of 11,640 orthologous gene families that were shared by all ten species (Fig. 1a, Supplementary Table S12).

Gene families specific to reindeer. A total of 14,520 orthologous gene families were shared between the domestic ruminant species included in this study, while 547 gene families were only found in reindeer, hereinafter referred to as 'reindeer-specific' gene families (Fig. 1b, Supplementary Data 1). The reindeer-specific gene families contained 1,090 genes, of which 530 were associated with 945 known InterPro domains (Supplementary Data 2). A number of olfactory receptors ($n = 29$ genes), zinc finger proteins ($n = 71$ genes) and ribosomal proteins ($n = 29$ genes) were present in the reindeer-specific gene families (Supplementary Data 1). Out of 1,090 reindeer-specific genes, 691 genes lacked Gene Ontology (GO) annotations. Analysis of the remaining 399 genes revealed 39 significant (false discovery rate (FDR) < 0.05) GO terms, of which 23 were associated with biological processes (BP) and the rest were categorized into molecular function (MF) (Supplementary Table S13). The reindeer-specific gene families enriched in biological processes revealed several GO terms associated with biosynthetic processes ($n = 10$ GO terms) and metabolic processes ($n = 7$ GO terms) (Supplementary Table S13). The biological processes associated with reindeer-specific gene families also included "GO:0097659, nucleic acid-templated transcription, $P = 6 \times 10^{-6}$ ", "GO:0010468, regulation of gene expression, $P = 6.2 \times 10^{-6}$ " and "GO:0007186, G protein-coupled receptor signaling pathway, $P = 0.00028$ " (Supplementary Table S13). Similarly, the molecular function category of the GO terms revealed a number of receptor activities ($n = 7$ GO terms), such as "GO:0099600, transmembrane receptor activity, $P = 3.3 \times 10^{-6}$ ", "GO:0004888, transmembrane signaling receptor activity, $P = 1.8 \times 10^{-5}$ ", "GO:0004984, olfactory receptor activity, $P = 5.9 \times 10^{-5}$ ", "GO:0038023, signaling receptor activity, $P = 6.5 \times 10^{-5}$ " and "GO:0001637, G protein-coupled chemoattractant receptor activity, $P = 2.9 \times 10^{-5}$ " (Supplementary Table S13). Moreover, three molecular function categories associated with NADH dehydrogenase activity were significantly enriched in reindeer-specific gene families (Supplementary Table S13).

We further examined reindeer-specific gene families to search for genes related to adaptation to subarctic climatic and other environmental factors, such as cold tolerance, marked fluctuations in daylight and feed shortage during winter. Several genes (*AKAP5*, *CACNB2*, *CEBPB*, *CEBPE*, *COL1A2*, *CYSLTR1*, *DNAJB6*, *EIF2B2*, *FOXE3*, *GATA3*, *GNAO1*, *MICB*, *PKD1*, *PRDM7*, *PRDM9*, *RPS6*, *TRPV5* and *TRPV6*) (Supplementary Data 1) from our list of reindeer-specific gene families have been reported to be associated with cold adaptation in Siberian human populations³⁶. We also identified genes related to unique characteristics of the reindeer species⁴, such as vitamin D metabolism (*TRPV5* and *TRPV6*) and antler development (*SILT2*). In addition, we found that the melanopsin gene (*OPN4B*), which is a photoreceptor, has an important function in the retina^{37–39}.

Expanded and contracted gene families. Based on the comparison of orthologous gene families among the 10 species, the reindeer genome revealed 3,312 and 1,840 expanded and contracted gene families, respectively (Fig. 1c). The significantly ($P < 0.01$) expanded ($n = 368$) and contracted ($n = 15$) gene families contained 2,683

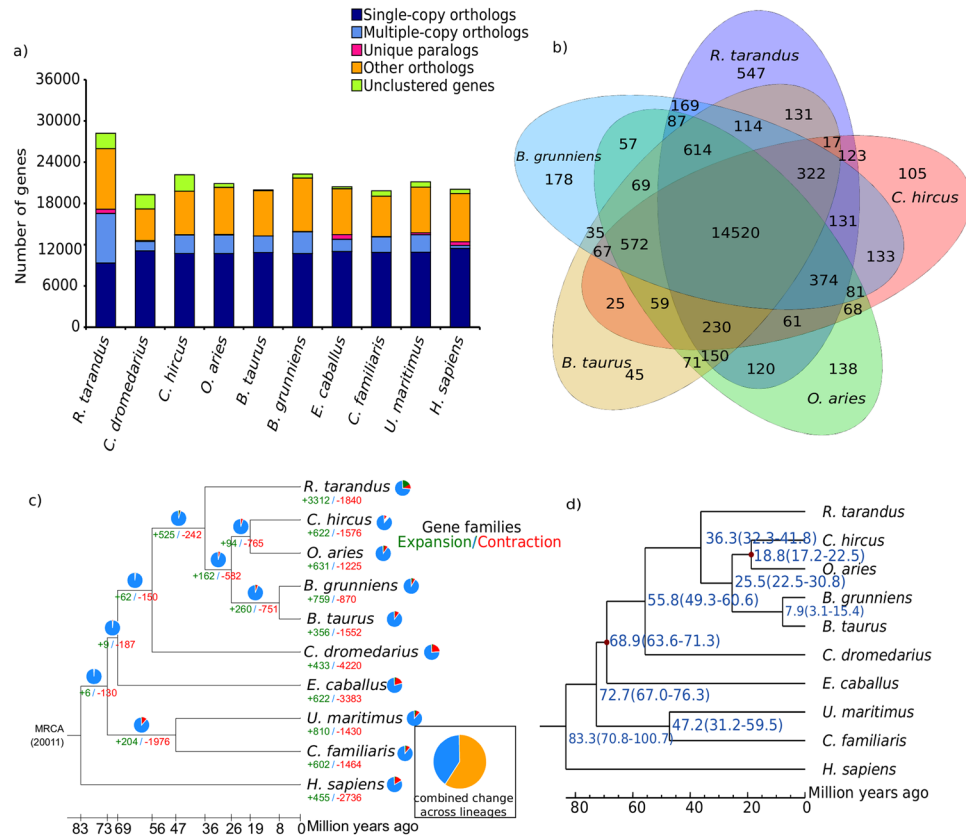


Figure 1. Gene ontology comparison among genomes of reindeer and nine mammalian species. **(a)** Distribution of gene clusters among the compared mammalian species. **(b)** Venn diagram showing unique and shared gene families among five ruminants (cattle, goat, reindeer, sheep and yak). **(c)** Gene expansion and contraction in the reindeer genome. The numbers of gene families that have expanded (green, +) and contracted (red, -) after splitting are shown on the corresponding branch and pie charts. MRCA, most recent common ancestor. **(d)** Phylogenetic tree based on 4-fold degenerate sites of 7,951 single-copy orthologous genes. Estimates of divergence time with 95% confidence interval (CI) are shown at each node.

and 19 genes, respectively (Supplementary Data 3 and Data 4). A high number of expanded genes were associated with ribosomal proteins ($n = 658$ genes), zinc finger proteins ($n = 155$ genes) and olfactory receptors ($n = 35$ genes) (Supplementary Data 3). We performed functional enrichment analysis for only the 368 expanded (2,688 genes) gene families (Supplementary Table S14). The expanded genes revealed 163 significant ($FDR < 0.05$) GO terms, of which 122 were categorized into biological processes and 41 were associated with molecular function (Supplementary Table S14 and Data 5). The expanded gene families enriched in biological processes were mainly associated with metabolic processes ($n = 33$ GO terms), biosynthetic processes ($n = 23$ GO terms), ion homeostasis ($n = 14$ GO terms) and transport ($n = 14$ GO terms) (Supplementary Table S14 and Data 5). Similarly, the expanded gene families enriched in the molecular function category were dominated by binding activities ($n = 16$ GO terms), such as “GO:0008199, ferric iron binding, $P = 4.2 \times 10^{-57}$ ”, “GO:0005506, iron ion binding, $P = 2.5 \times 10^{-17}$ ”, “GO:0031492, nucleosomal DNA binding, $P = 9.2 \times 10^{-13}$ ” and “GO:0043566, structure-specific DNA binding, $P = 9.2 \times 10^{-13}$ ”.

We further explored expanded gene families to search for genes related to arctic and subarctic adaptation and unique reindeer characteristics. Several genes (*ADRA2A*, *ADRA2B*, *ADRA2C*, *CIDEA*, *CRYAB*, *CYCS*, *DNAJA1*, *DRD3*, *EPHA3*, *EPHB1*, *FOXC1*, *FOXC2*, *GNG5*, *HMG3*, *HOXA5*, *HSP1*, *HSPE1*, *HTR1B*, *HTR2A*, *PARK2*, *PPP2R1A*, *PRDM1*, *PRDM7*, *PRDM9*, *RAP1B*, *RHOC*, *RPS6*, *SLC25A5*, *SLC8A1*, *TRPV5* and *TWIST1*) (Supplementary Data 3) that were predicted to be in expanded gene families have been reported to be associated with cold adaptation in Siberian human populations³⁶. Moreover, we also identified genes related to vitamin D metabolism (*TRPV5*) among the list of expanded gene families⁴.

Genes under positive selection in reindeer. We conducted an analysis of the signature of positive selection by assessing dN/dS ratios of 7,951 single-copy orthologous gene sets found in the genomes of *R. tarandus* and the other nine mammalian species. A total of 418 of the 7,951 (5.26%) single-copy orthologous genes in reindeer revealed a significant signature of positive selection ($P < 0.01$) (Supplementary Data 6). Analysis of the GO enrichment for the genes under positive selection in reindeer revealed seven statistically significant ($FDR \leq 0.05$) enriched GO terms particularly related to channel activities ($n = 6$ GO terms) (Supplementary Table S15).

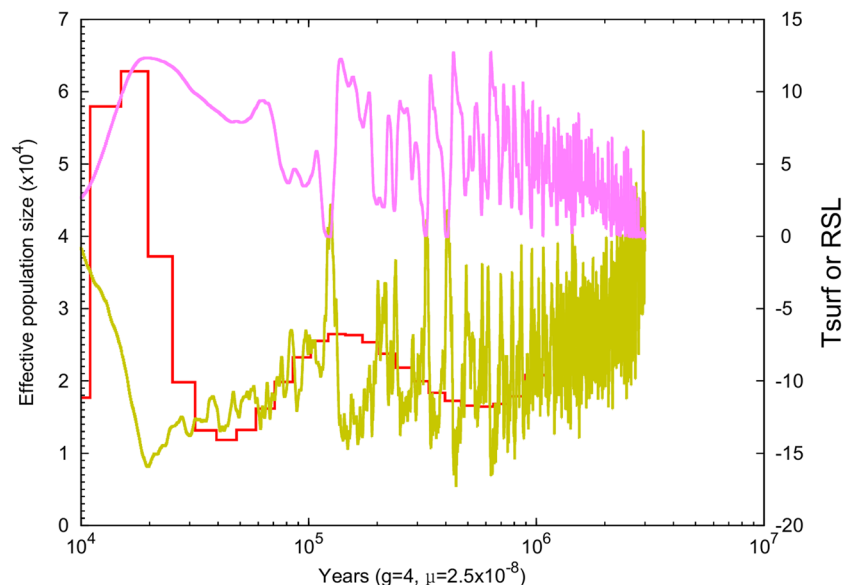


Figure 2. Demographic history of *R. tarandus* reconstructed from the draft reference genome by using PSMC. The X axis shows the time in thousands of years and the Y axis shows the effective population size. In the figure, the red line describes the fluctuation of effective population size, the pink line the relative sea level (RSL, 10 m sea level equivalent) and the yellow line atmospheric surface air temperature °C (Tsurf).

A number of genes (*BDKRB2*, *CORIN*, *GCG*, *GNB1*, *IL6*, *KCNMA1*, *KCNMB2*, *MED1*, *MLXIPL*, *NPY*, *NTS*, *PRDM10*, *PRKCQ*, *RNPEP*, *SHC1*, *SORBS3*, *STRADA*, *TCF7L2*, *TRPC1* and *VEGFA*) (Supplementary Data 6) showing positive selection in reindeer have been reported to be associated with cold adaptation in Siberian human populations³⁶. Moreover, we noticed that a gene regulating circadian rhythm (*GRIA1*)⁴ showed signs of positive selection.

Evolutionary divergence analysis. We also estimated the evolutionary divergence of reindeer by comparing the assembled reindeer genome with four other domesticated ruminant species (goat, sheep, yak and taurine cattle), dromedary camel, horse, polar bear, dog and human. A phylogenetic tree was constructed based on 4-fold degenerate codon sites extracted from 7,951 single-copy orthologous genes identified by TreeFam⁴⁰. The estimated divergence time showed that reindeer shared a common ancestor with other domesticated ruminant species approximately 36 million years ago (95% confidence interval (CI) 32–42 million years ago) (Fig. 1d). Our divergence estimate is longer than the divergence estimate from a previous reindeer study²¹. Our analysis suggests that the divergence between the domesticated ruminant species and the dromedary camel occurred in the early phase of *Artiodactyla* evolution⁴¹, 56 million years ago.

Demographic history. We applied the pairwise sequentially Markovian coalescent (PSMC)⁴² method to the reindeer reference genome to examine the changes in the effective population size (N_e) of the ancestral population of the Fennoscandian reindeer in the course of the Pleistocene period in the world's history (Fig. 2). The N_e of the ancestor of reindeer gradually decreased between 1 million years ago and 500 thousand years ago (Kya). The ancestral N_e of reindeer showed peaks at 150 Kya and 20 Kya, while the population underwent three major bottlenecks at 600 Kya, 40 Kya and 11 Kya.

Heterozygosity estimation. Using the assembled draft reindeer genome sequence as a reference, we mapped the reads of short insert-size libraries to the genome assembly of reindeer and found high coverage, ensuring a high level of accuracy at the nucleotide level. Approximately 98.34% of the reads were successfully mapped to the genome, and the sequencing coverage was approximately $60\times$. We used the uniquely mapped reads to call variants using the Genome Analysis Toolkit (GATK)⁴³ best practice pipeline. We identified a total of 5.46 million (M) heterozygous SNPs in the genome. The estimated heterozygosity rate of reindeer was 2.05×10^{-3} , which is 3.48 and 2.3 times higher than that of cattle (0.59×10^{-3}) and yak (0.89×10^{-3})²⁷.

Whole-genome sequence analysis. The raw sequences for the 23 reindeer samples were generated at Beijing Genomics Institute (BGI) using the Illumina HiSeq 4000 platform. The raw reads were preprocessed; adapters and low-quality data were removed. After the data were processed, we generated a total of 680 Gb clean paired-end resequence data (Supplementary Table S16). On average per individual, we achieved 197 M and 29.57 Gb clean reads and bases, respectively (Supplementary Table S16). On average per individual, 98.78% of the clean reads were successfully mapped to the reindeer genome assembly (Supplementary Table S17) and represented 9.46-fold coverage.

A total of 28.6 M high-quality SNPs were detected in the mapped reads across all 23 individuals. The average number of SNPs detected per individual was 7.24 M (Supplementary Table S18). In the wild tundra reindeer

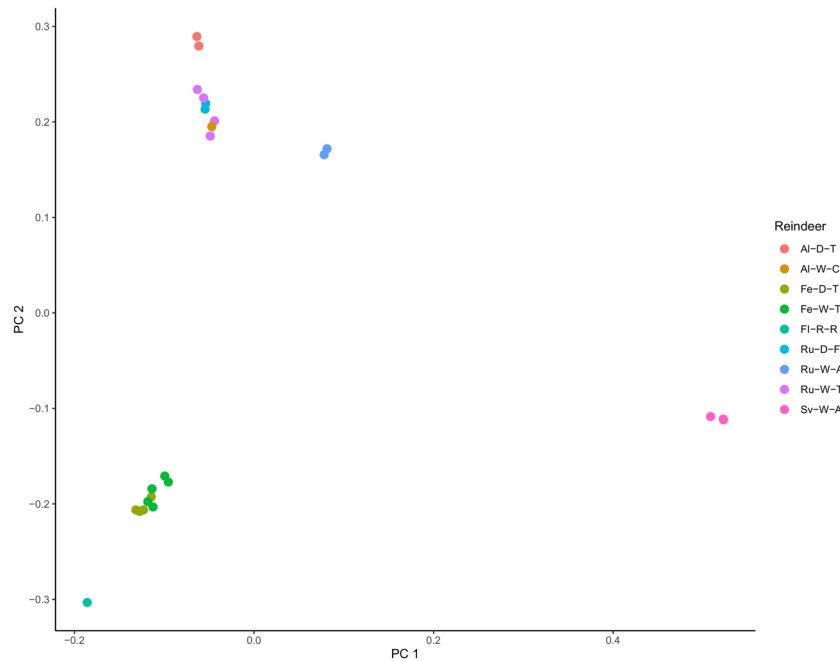


Figure 3. Principal component analysis (PCA) of 23 resequenced animals and the reference individual. (FI-R-R, Finnish reindeer (i.e. the reference animal); Fe-W-T, Fennoscandian wild tundra reindeer; Fe-D-T, Fennoscandian domestic tundra reindeer; Ru-W-T, Russian wild tundra reindeer; Ru-D-F, Russian domestic forest reindeer; Ru-W-A, Russian wild arctic reindeer (i.e. Novaya Zemlya); Sv-W-A, Svalbard wild arctic reindeer; Al-D-T, Alaska domestic tundra reindeer; Al-W-C, Alaska wild caribou).

in Russia (four samples collected at two sites), 8.00 M SNPs were identified on average per individual; in the wild Norwegian tundra reindeer, 7.09 M (five animals) SNPs were identified; and in the domestic tundra reindeer, 6.79 M (four animals from two sites) SNPs were identified. Three arctic reindeer samples from Svalbard Norway showed low diversity, with an average of 1.37 M (20.42%) heterozygous and 5.35 M (79.58%) homozygous SNPs per individual, while in the rest of the samples, the respective estimates were 3.81 M (52.14%) and 3.50 M (47.86%). We also detected a total of 3.93 M indels across all 23 samples, and on average, we detected approximately 895,386 indels per individual.

Functional annotation of SNPs. SnpEff⁴⁴ was used to annotate all filtered SNPs identified across all samples, and the proportions of SNPs in the intergenic, intronic, downstream and upstream, and transcript regions were 58.27%, 14.0%, 12.52% and 14.61%, respectively, while a small number of variants were annotated in exonic (0.57%) regions (Supplementary Table S19 A and B). The overall estimated missense to silent ratio for all samples was 0.7963. The average SNP-variant rate was 1 variant every 90 bases for all samples. The transition to transversion (Ts/Tv) ratio was calculated for each sample, and the average Ts/Tv ratio was 1.97 for all reindeer samples, which is slightly lower than that of human (2.1)⁴⁵ and lower than that of bovine (2.2)^{46–48}.

Population genetic analysis. Genetic relationships between the 23 resequenced animals and the reference animal were investigated with a principal component analysis (PCA). The PCA plot revealed two major groups: Fennoscandian reindeer versus forest and tundra reindeer/caribou from Russia and Alaska, as shown in Fig. 3. The reindeer from Svalbard and Novaya Zemlya, both wild arctic ecotype, grouped separately from the main clusters. We computed genetic diversity parameters for the two main clusters. The overall genome-wide genetic diversity, as measured by Watterson's (θ) and pairwise nucleotide diversity (π), was smaller in the Fennoscandian group (0.002022733 and 0.002096655, respectively) than in the Russian – Alaskan group (0.002441424 and 0.002265819, respectively). The genetic differentiation between the two main groups, as measured by fixation index (F_{ST}), was 0.039786, suggesting that these two cluster populations exhibit low differentiation.

Assembly and annotation of the reindeer mitochondrial genome. The reindeer mitochondrial genome was assembled using the mitochondrial baiting and iterative mapping (MITObim)⁴⁹ approach. A total of 16,962 reads were utilized in the assembly, corresponding to 0.036% of the reads generated from the 170 bp insert-size library. The reads used in the assembly mapped with 84% of bases mapping with quality greater than Phred quality score of 20 (Q20). The assembled mitochondrial genome is 16,357 bp in length with an average conversion of 131.84 \times and a GC content of 36.28%. The reindeer mitogenome consists of 22 tRNAs, 13 protein-coding genes and two rRNAs (Supplementary Table S20).

Discussion

With our present study and the recent publications by Li *et al.* (2017)²¹ and Lin *et al.* (2019)⁴, a new genomic era in the study of reindeer has been established. Here, we have generated and described an alternative draft genome assembly for reindeer and published the largest resequenced *Rangifer* sp. dataset for the investigation of genetic diversity, domestication, demographic aspects and genomic adaptation.

Our reindeer genome assembly is 2.66 Gb in size and represents ~91% of the estimated 2.92 Gb genome size according to the k-mer analysis. The estimated genome size of reindeer in our study is slightly higher than that in two recent studies, i.e., 2.86⁵⁰ Gb and 2.76 Gb²¹. Species belonging to the *Cervidae* family⁵⁰, such as Indian muntjac (2.92 Gb), Chinese muntjac (2.99 Gb), Milu (3.09 Gb), Black muntjac (3.24 Gb) and White-lipped deer (3.23 Gb), appear to have larger genome sizes than domesticated mammalian species in general (~2.6–2.7 Gb), such as pig (*S. scrofa*¹⁸), yak (*B. grunniens*²⁷), dromedary camel (*C. dromedarius*⁵¹) and sheep (*O. aries*¹⁹). In addition, the number of protein-coding annotated genes (>27,000) in our draft genome assembly is higher than that typically predicted in the latest assemblies of several other domesticated mammalian species (~21,000)^{18–20,27,35} and in the first draft genome assembly of reindeer (21,555)²¹. We assume that the differences in the gene numbers may be due to the annotation approach (homology-based approach) used in this study and will require further validation. It has been shown that homology-based annotation approach commonly used in draft reference genomes may also include split or fragmented genes⁵². However, the BUSCO assessment revealed 97% orthologous genes in the present assembly, which is better than the BUSCO assessment results (92.6%) in a previous study²¹, thus indicating the high quality of our reindeer gene annotation. In addition, the key gene parameters (see Supplementary Table S7 and Fig. S4) also revealed no significant differences between the reindeer genome and the five species that were used in homology prediction related to gene and exon length distribution.

We applied the present draft assembly to estimate the evolutionary divergence times between reindeer and nine other mammalian species and to reveal recent demographic events in the history of reindeer that occurred during the last 1 million years. Our phylogeny suggests that reindeer diverged from four other domestic ruminants ca. 36 million years ago, which is earlier than previously estimated^{21,50}. The demographic analysis (Fig. 2) showed that there have been marked fluctuations in the N_e of the reindeer species: two population expansions and three bottlenecks. The decline in the N_e occurred during the mid-Pleistocene period, and the peak during Last Glacial Maximum period may have been associated with global temporal changes in climate^{53–55}. The pattern of the N_e of the reindeer genome during the Pleistocene period is consistent with the previously suggested hypotheses of the responses of other domestic mammalian species, such as pig (*S. scrofa*¹⁸), horse (*E. caballus*⁵⁶), sheep (*O. aries*⁵⁷) and cattle (*B. taurus*²³), to temporal climate changes in the past. However, the heterozygosity rate (2.05×10^{-3}) of the reindeer genome estimated in our analysis was 3.48 and 2.3 times higher than that of *B. taurus*-cattle and *B. grunniens*-yak, respectively²⁷, suggesting a larger founder population size of the contemporary semi-domesticated reindeer. In addition, possible gene flow from wild reindeer populations, a less intensive artificial selection history and the early phase of reindeer domestication history in general may have contributed to higher diversity estimates for reindeer than for the two other ruminant species¹⁵.

To explore the genetic distinctiveness of reindeer in the present research context, we focused on reindeer-specific and expanded orthologous gene families and their genes and gene functions. The changes in genetic makeup (i.e., gene copy and/or protein domain number changes) are potentially linked to unique phenotypic and functional changes, suggesting a mechanism for adaptive divergence of closely related species^{26,27,58}. Gene family analysis revealed that 160 genes were shared between 547 reindeer-specific gene families and 368 expanded gene families. Therefore, the 160 genes that we identified as unique to reindeer compared to four other ruminants are also rapidly evolving (among nine other mammals) and may thus represent characteristic features of reindeer. Zinc finger proteins ($n = 42$) dominated the list of uniquely expanding genes followed by olfactory receptors ($n = 13$). Zinc finger proteins have multi-functional roles, including lipid binding and the differentiation of adipose tissues^{59,60}. Interestingly, among the list of zinc finger proteins was PR domain zinc finger protein 9 (*PRDM9*), which is considered one of the fastest evolving genes in mammalian species. *PRDM9*, also referred to as the speciation gene, has an important function in enhancing recombination and epigenetic modification^{61,62}. We speculate that the *PRDM9* gene has played a vital role in the evolution and adaptation of *R. tarandus* to challenging northernmost biogeographic and climatic conditions. Another member of the zinc finger protein, PR domain zinc finger protein 1 (*PRDM1*), is associated with retinal development^{63–65} and is a candidate gene promoting adaptation to extreme seasonal light changes. Furthermore, the improved sensory systems of reindeer were also reflected in our study, where we observed several genes associated with olfactory receptors and G protein-coupled receptors in specific and expanded gene families (Supplementary Data 1 and Data 3). The mammalian genome contains 4–5% of olfactory receptor genes, which are sensors to the extracellular environment and essential for the survival of animals⁶⁶. G protein-coupled receptors are also known to be involved in sensing of the extracellular environment⁶⁶. During the long winter season, the reindeer sense of smell or olfaction has an important role in finding feed covered with snow.

Moreover, the examination of genes in the reindeer-specific and expanded orthologous gene families provided information on genes associated with specific phenotypic characteristics of reindeer. We identified two genes belonging to the transient receptor potential (TRP) cation channel subfamily, *TRPV5* and *TRPV6*, among the reindeer-specific gene families. *TRPV5* and *TRPV6* are involved in calcium reabsorption and maintenance of blood calcium levels in higher organisms^{67,68}. Among these, *TRPV5* has been found to be associated with vitamin D metabolism⁴. Previous studies have reported that reindeer require efficient calcium metabolism and reabsorption for antler growth during the period of continuous winter darkness and low solar energy^{69,70}. Studies also indicated that efficient vitamin D metabolism is needed to sustain calcium metabolism and reabsorption^{69,70}. The nucleotide sequences of both the *TRPV5* and *TRPV6* genes harbour 15 exons, encode proteins of approximately 730 amino acids and have 75% sequence similarity^{67,68}. Both genes have similar functional properties and regulation mechanisms⁶⁷. We speculate that in addition to *TRPV5*, *TRPV6* also has a significant role in vitamin D

metabolism in reindeer. In addition, three adaptive immune-related chemokines (*CXCR1*, *CXCR2*, *CXCR4*)^{71–73} and three interferons (*IFNT1*, *IFNT3*, *IFNW1*) that have antiviral functions^{74,75} may have promoted the adaptation of reindeer to the given environment. Finally, insulin-like growth factor 1 receptor (*IGF1R*) promotes the activity of *IGF1*, which has a pivotal role in antler formation (*IGF1* regulates *RUNX1* expression via *IRS1/2* signalling for antler growth)^{76,77}.

We also observed large numbers of ribosomal proteins present in the reindeer-specific and expanded gene families. A number of reports have suggested that ribosomal proteins are essential for survival and adaptation to environmental stress^{78–80}. For example, one of the ribosomal proteins we found was the ribosomal large subunit protein 7 (*RPL7*), which has been previously reported in Russian cattle breeds⁸¹ to be related to adaptation to harsh (cold) environments. The ribosomal large subunit protein 7 (*RPL7*) gene was found to be upregulated in the skin of freeze-tolerant wood frogs compared to non-tolerant species, showing that this gene is cold-responsive and that the gene was upregulated in skeletal muscle and the brain under freezing stress⁸⁰. Moreover, in the list of reindeer-specific gene families, we found opsin 4B (*OPN4B*), which is a photoreceptor that appears to be involved in the regulation of the circadian clock^{37–39}. Several studies^{39,82–84} have reported the existence of a melanopsin-associated photoreceptive system in the mammalian retina that helps to transmit photic information and is also involved in the regulation of photoentrainment of the circadian clock. A previous study⁹ reported that reindeer possess a *tapetum lucidum* (reflective surface behind the central retina) to cope with winter darkness, which may help in changing wavelength reflection and increase the sensitivity of the animal's vision in dim light. In winter, the reindeer *tapetum lucidum* changes to blue, which may help scatter light through photoreceptors, and less light is reflected⁹. We speculate that this melanopsin-related gene might have important functions for the adaptation to the light system in circumpolar environments. Furthermore, we found Slit homolog 2 protein (*SLIT2*) in reindeer-specific gene families, which has been reported to be associated with antler development in a recent reindeer study⁴.

In our study, signs of positive selection were found in genes in the genome of *R. tarandus* mainly enriched in GO terms related to channel activity (Supplementary Table S15), such as sodium channel activity and ion channel activity, which are thought to be relevant in cold sensing mechanisms⁸⁵. The ability to detect and adapt to cold temperature is crucial for the survival of an organism^{86,87}. The process of sensory transduction and a large array of ion channels, such as Na⁺ and K⁺ channels, are involved in cold temperature detection⁸⁵. Among the genes categorized into these GO terms, interestingly, *SCN11A1* has an essential role in pain tolerance caused by cold stress^{88,89}, suggesting the gene's role in survival in extremely cold environments. One additional gene showing signatures of positive selection, glutamate ionotropic receptor AMPA type subunit 1 (*GRIA1*), appears to have a role in the circadian clock⁴.

Finally, we investigated the genetic relationships between the *de novo* reference animal and the 23 resequenced individuals, and two main genetic clusters (Fennoscandia and Russia-Alaska) were identified in our data (Fig. 3). We suggest that this observation reflects the historical spread of reindeer populations in northern Eurasia and reveals the domestication history of reindeer. The Russian/northern American cluster probably reflects the Euro-Beringian lineage evolved from Pleistocene population in northern Eurasia^{3,90} while the Fennoscandian cluster may have descended from some refugia populations in southern Europe partly isolated from the Euro-Beringian lineage⁹¹. In these two main genetic clusters, we found both domestic and wild reindeer from respective Fennoscandia and Russia, which gives support to previously indicated independent domestication origin for the breeds in these regions^{13,92}. Our findings are consistent with the previous study¹⁵. Our data suggest that the semi-domesticated tundra reindeer show lower genetic diversity in terms of the number of SNPs per individual than the wild tundra reindeer, implying a genetic bottleneck effect caused by domestication as also seen in temporal change in mitochondrial DNA^{13,93}. Moreover, our data clearly indicated the effects of geographic isolation and genetic bottleneck in the Svalbard reindeer: the subspecies (or an arctic ecotype) displayed exceptionally low diversity compared to the other resequenced Rangifer populations.

Here, we have reported a 2.66 Gb genome assembly of a Finnish male reindeer originating from Sodankylä (67.34°N, 26.83°E), where the annual mean temperature is 0 °C (the mean temperature in January and February –13 °C), the snow cover is on average 202 days per year, the daylight is less than 3 h from December until mid-January (during a short period in December, the sun is down all day) and the sun is up all day from June until mid-July (<https://www.timeanddate.com/>). The first draft assembly of the reindeer genome was obtained from a female animal²¹ from Inner Mongolia Autonomous Region, China (50.77°N, 121.47°E), where the climate is continental, but the annual daylight rhythm (at least 8 h of daylight and no periods of total darkness or “nightless nights” exist) is different from that typically found in traditional reindeer herding regions in northern Eurasia. In addition to these climatic differences in their geographic origins, we assume that the Finnish reindeer and Inner Mongolian reindeer do not share their genetic origins in the same refugial, ancestral populations. The Finnish reference reindeer was found to belong to the cluster of northern Fennoscandian reindeer distinct from the “eastern” and North American genetic cluster. It is recommended that phylogenetically different animal populations have their own reference assemblies²³. With our GO family-based comparisons with nine other mammalian genomes, we identified hundreds of genes and gene families that are unique, under positive selection and rapidly expanding in reindeer. Exploring these genes and gene families revealed several reindeer-specific characteristics that have helped these animals survive in the arctic and subarctic conditions. We identified genes that are important for vitamin D metabolism (*TRPV5*, *TRPV6*), antler formation (*SLIT2*), and circadian rhythm (*OPN4B*). While many of the findings from our study complement reports from other species, the functions of *TRPV6* and *OPN4B* have been reported here for the first time. Population genomics analyses based on 23 individuals representing domestic and wild reindeer suggested two domestication events in *R. tarandus* subspecies, but additional larger scale studies are required for validating this finding. Collectively, our study provides new insight into the genomic and evolutionary characteristics of reindeer, and the resequencing data (the most comprehensive in reindeer thus far) will serve as an invaluable resource for future studies involving reindeer and other cervids.

Materials and methods

Sampling and genome sequencing. All protocols and sample collections were performed in accordance with the legislations approved by the Animal Experiment Board in Finland (ESAVI/7034/04.10.07.2015).

A one-year-old male reindeer (*R. tarandus*) from Sodankylä, Finland, was selected, and blood samples were collected for sequencing. Genomic DNA was extracted from the blood using a standard phenol/chloroform method and sequenced at high coverage on the Illumina HiSeq 2500 and 4000 platforms using a shotgun-sequencing approach. We constructed seven paired-end DNA libraries with insert sizes ranging from 170 bp to 20 kb. Summary statistics of the generated reads are shown in Table 1 and Supplementary Table S1. The short reads with low-quality bases were removed before assembly by filtering as follows:

1. Reads from short insert-size libraries having 'N' over 2% of its length, and the reads from large insert-size libraries having 'N' over 5% of its length.
2. Reads from short insert-size libraries having more than 40% bases with Q20 less than or equal to 7, and the reads from large insert-size libraries having more than 30% bases with Q20 less than or equal to 7.
3. Reads with adapter sequences.
4. Paired-end reads of read 1 and read 2 that were completely identical.

Estimation of genome size using k-mer analysis. For the short reads with an insert size of 500 bp, a total of 64.7 Gb (approximately $21.6\times$) was used to estimate the genome size using the k-mer method²⁴, as described in previous studies²⁶. We used the following formula to estimate the genome size (G): $G = K_num/K_depth$, where K_num is the total number of k-mers, and K_depth is the frequency occurring more frequently than other frequencies. In this study, the size of k-mer for reindeer was set to 17, K_num was 55,451,706,382 and K_depth was 19. Hence, the genome size was estimated to be approximately 2.92 Gb (Supplementary Table S2).

Genome assembly. After the data were filtered, clean data were retained for assembly (Table 1, Supplementary Table S1). The high-quality clean reads were assembled first into contigs and subsequently into scaffolds using SOAPdenovo V2.04 software²⁵.

The quality of the genome assembly was evaluated by aligning high-quality reads from short insert-size libraries to the *de novo* genome assembly using BWA²⁸ (version. 0.7.15) with no more than five mismatches and then by determining the percentage of total aligned reads. Reads from the transcriptome were assembled by trinityrnaseq-2.0.6⁹⁴; all sequences were mapped to the assembly by BLAT²⁹ with the default parameters. We also employed BUSCO (version 3.0.2)³⁰ to assess the completeness of the genome assembly and annotated gene set using the conserved 4,104 mammalian genes.

RNA sequencing. To improve the reference genome annotation, we extracted RNA using the RNeasy Plus Universal Kit (Qiagen, Valencia, CA, USA) from six tissues (scapular adipose, metacarpal adipose, tailhead adipose, perirenal adipose, liver and *Musculus gluteobiceps*) from the same male reindeer subjected to *de novo* sequencing. The tissue samples were collected at slaughter and stored in RNAlater[®] Solution (Ambion/Qiagen, Valencia, CA, USA) according to the manufacturer's instructions. mRNA libraries were sequenced using Illumina HiSeq 2500 with a paired-end (2×150 bp) strategy. The raw data were preprocessed, and adapters and low-quality reads were filtered out using the cutadapt tool⁹⁵.

Genome annotation. Repeat annotation. Tandem repeats were detected in the generated genome assembly using Tandem Repeats Finder (TRF) software⁹⁶. Transposable elements (TEs) were predicted in the genome using the combination of both known and *de novo* strategies. For the known-based strategy, we identified known TEs in the genome using RepeatMasker⁹⁷ and RepeatProteinMask⁹⁸ by homology search against the Repbase³² library with default parameters. For the *de novo* way, we constructed a *de novo* repeat library using RepeatModeler³³, with the default parameters. RepeatMasker was used again with the *de novo* libraries to identify new TEs in the genome. All repeats obtained by the different methods were combined according to the coordinate in the genome and overlapping TEs belonging to the same repeat class to form a non-redundant list of reindeer repeats.

Gene prediction and functional annotation. Gene annotation consists of structure annotation and function annotation. We used homology-based and RNA-seq data to annotate coding genes of the genome. Firstly, protein sequences of *H. sapiens* (GRCh38.p13, Ensembl release 84), *M. musculus* (GRCm.p4 Ensembl release 84), *B. taurus* (UMD3.1, Ensembl release 84), *C. dromedarius* (PRJNA234474_Ca_dromedarius_V1.0, NCBI GeneBank accession GCA_000767585.1) and *C. familiaris* (CanFam3.1, Ensembl release 84) were downloaded and mapped onto the *R. tarandus* genome using software tblastn²⁹. Secondly, high-score segment pairs were concatenated between the same pair of proteins by solar (in-house software, version 0.9.6), the results were then filtered by the overlap cutoff of 50% to get the non-redundant result. Thirdly, the Genewise⁹⁹ was employed to define accurate gene models according to the sequences extracted from the non-redundant result. Finally, we filtered the alignment with coverage <30% and filtered redundancy from multiple species based on the alignment score of the Genewise. For the RNA prediction, RNA-seq data of six samples representing six tissues were mapped to the genome with Tophat2 (v2.08)¹⁰⁰ with default parameters and assembled into transcripts using Cufflinks (v2.2.1)¹⁰¹. These transcripts were used to extend the homology gene models to refine the open reading frame (ORF) and untranslated exons (UTR) prediction. The rules for adding UTRs to genewise predictions was described in Ensembl pipeline³¹. Functions of genes were assigned based on the best match derived from the alignments to proteins annotated in four protein databases: InterPro¹⁰², KEGG¹⁰³, Swiss-Prot¹⁰⁴ and TrEMBL¹⁰⁵.

ncRNA annotation. We used INFERNAL¹⁰⁶ and tRNAscan-SE¹⁰⁷ to predict ncRNA. Four types of ncRNAs were annotated in our analysis: tRNA, rRNA, miRNA and snRNA. tRNA genes were predicted by tRNAscan-SE with eukaryotic parameters. rRNA fragments were identified by aligning the rRNA template sequences from human genomes using BLASTN with an E-value cutoff of 1.0×10^{-5} . miRNA and snRNA genes were inferred by the INFERNAL software against the Rfam database (v11.0)¹⁰⁸.

Gene family construction. To define gene family evolution in the reindeer genome, we used the TreeFam methodology⁴⁰ as follows: BLAST was used to compare all protein sequences from 10 species, *R. tarandus*, *C. dromedarius*, *C. hircus*, *O. aries*, *B. Taurus*, *B. grunniens*, *E. caballus*, *C. familiaris*, *U. maritimus* and *H. sapiens*, with the E-value threshold set to 1.0×10^{-7} . After that, HSPs of each protein pair were concatenated by Solar software. In this analysis, the protein datasets were used from Ensembl release 84 for *H. sapiens*, *B. taurus*, *O. aries* and *C. familiaris*. For *C. dromedarius* we used NCBI version NW_011591059.1. For *C. hircus*³⁵, *B. grunniens*²⁷ and *U. maritimus*¹⁰⁹ we used BGI internal version from <http://gigadb.org>. H-scores were computed based on Bit-scores, and these were taken to evaluate the similarity among genes. Finally, gene families were obtained by clustering homologous gene sequences using Hcluster_sg (version 0.5.0). If these genes had functional motifs, they were annotated by GO¹¹⁰.

Phylogenetic tree construction and divergence time estimation. The phylogenetic tree was constructed based on the single-copy orthologous genes from the ten species identified by gene family analysis. CDS from each single-copy gene cluster were aligned using MUSCLE¹¹¹, and their protein sequences were concatenated to form one super gene for every species. Four-fold degenerate sites (4d sites) of aligned CDS were extracted for subsequent analysis. The phylogenetic relationships were inferred using PhyML software^{112,113}, with the GTR substitution model and gamma distribution rates model as the chosen parameters. We constructed a phylogenetic framework including six *Artiodactyla* and four other vertebrates, with *H. sapiens* as the outgroup. Divergence times were estimated using PAML mcmctree (PAML version 4.5)^{114–116} by implementing the approximate likelihood calculation method.

Expansion and contraction of gene families. We analysed the expansion and contraction of gene families using the CAFE program¹¹⁷, which employs a random birth and death model across a user-specified phylogeny. The global parameter λ , which describes both the gene birth (λ) and death ($\mu = -\lambda$) rate across all branches in the tree for all gene families, was estimated through maximum likelihood. A conditional P-value was calculated for each gene family, and families with conditional P-values less than the threshold (0.01) were considered to have notable gains or losses. Finally, the branches responsible for low overall P-values were defined as significant families.

Detection of positively selected genes (PSGs). To estimate positive selection, dN/dS ratios were calculated for all single-copy orthologues of *R. tarandus* and nine other vertebrates (*C. dromedarius*, *C. hircus*, *O. aries*, *B. taurus*, *B. grunniens*, *E. caballus*, *C. familiaris*, *U. maritimus* and *H. sapiens*). Orthologous genes were first aligned by PRANK¹¹⁸, then Gblocks 0.91b was used to remove ambiguously aligned blocks within PRANK alignments and 'codeml' in the PAML package¹¹⁶ was employed with the free-ratio model to estimate Ka, Ks, and Ka/Ks ratios of different branches. The difference in mean Ka/Ks ratios for single-copy genes between *R. tarandus* and each of the other species were compared with paired Wilcoxon rank-sum tests. The genes that showed values of Ka/Ks higher than 1 along the branch leading to *R. tarandus* were reanalysed using the codon-based branch-site tests implemented in PAML. The branch-site model, which allowed ω to vary both at sites in the protein and across branches, was used to detect episodic positive selection.

Functional enrichment analysis. To gain insight into the biological functions and relevance of the detected specific, expanded and positively selected genes in reindeer, functional enrichment analysis was conducted using GO Analysis Toolkit and Database for Agricultural Community (AgriGO V2.0)¹¹⁹ with default parameters. The significantly enriched GO terms were detected by Fisher's exact test with the Bonferroni correction using default parameters (significance level threshold 0.05 (FDR < 0.05) and minimum number of mapping entries 5). In this analysis, GO annotation file of the assembled genome was used as a background reference. Two types of GO categories, molecular function and biological processes were examined in more detail in this study.

SNP calling and heterozygosity estimation. We used high-quality reads from the short insert-size libraries (170, 500 and 800 bp) to call SNPs. The high quality reads were mapped against the assembly using BWA²⁸ (version. 0.7.15) software with the default parameters. After mapping, we employed Picard tools (<http://broadinstitute.github.io/picard/>) (version.2.5.0) to preprocess alignment reads and remove PCR duplicates, and then the uniquely mapped reads were used for SNP calling. Further, we used GATK (version. 3.6)⁴³ to identify poorly mapped regions nearby indels, realigned and performed SNP calling according to the GATK Best Practices pipeline¹²⁰. We estimated the heterozygosity rate of the reindeer genome; first, we counted the identified heterozygous SNPs and then divided the total heterozygous SNPs by the effective genome size.

Population history estimation. We inferred a marked population bottleneck in the demographic history of *R. tarandus* using the pairwise sequentially Markovian coalescent (PSMC) model⁴² with generation times and mutation rates.

Assembly and annotation of the reindeer mitochondrial genome. Complete mitochondrial genomes are useful materials for phylogenetic studies. The complete mitochondrial genome for reindeer was assembled using MITObim v 1.9 software⁴⁹. Among the reads generated from the seven libraries on the reference

individual, for the mitochondrial genome assembly, we used only reads from the 170 bp insert-size library (10% of 66.2 Gb, Table S1) for MITObim input. The short DNA sequence (506 bp), cytochrome oxidase subunit 1 (COI) (GenBank accession number COI: KX085230.1)¹² from *R. tarandus* was used to seed the initial baiting of mitochondrial reads. Following the assembly, the *de novo* assembly was annotated using the MITOs web server¹²¹ with the default parameters.

Whole-genome resequencing and variant calling of 23 domestic and wild reindeer. To investigate genome variation and perform population analyses, we resequenced 23 domesticated and semi-domesticated reindeer individuals from Russia, Norway, Alaska and USA. Resequenced samples originated from the following populations and subspecies: domesticated forest reindeer from Russia (n = 2), domesticated tundra reindeer from Norway (n = 4), wild island tundra-mountain reindeer from Russia (n = 2), wild tundra reindeer from Russia (n = 4), wild tundra reindeer from Norway (n = 5), wild arctic reindeer from Svalbard Norway (n = 3), Alaskan domestic reindeer from Alaska-USA (n = 2) and Alaskan wild caribou from Alaska-USA (n = 1). Genomic DNA was extracted from 23 blood samples (obtained from the previous study by Flagstad and Røed 2003³), and sequence libraries with an average fragment size of 500 bp were constructed for each individual. Paired-end reads were generated using Illumina HiSeq 4000 at BGI. After sequencing, the raw reads were filtered to remove adapter sequences, contamination and low-quality reads. Following the data treatment, we achieved an average of 197 M and 29.57 Gb clean reads and bases, respectively (Supplementary Table S16).

The clean reads were mapped against the present draft reindeer assembly genome using BWA with the default parameters. Using Picard tools, we preprocessed and filtered alignments for SNP calling, including the removal of low-quality alignments, the sorting of alignments and the removal of PCR duplicates. Further, using GATK, we identified poorly mapped regions nearby indels from the alignments and performed a realignment and base quality recalibration on the bases that were disrupted by indel sites. Finally, using the GATK Best Practices pipeline, we performed SNP discovery and genotyping across the 23 samples.

Functional annotation of SNPs. SnpEff V4.3T⁴⁴ was used to annotate the SNPs and categorize them into coding (synonymous and non-synonymous), upstream/downstream and intronic/intergenic classes. For this analysis, we constructed the SnpEff databases using the draft reindeer reference genome and the corresponding genome annotation (.gtf) file.

Population genetics statistics. First, to examine the genetic relationships across the 23 reindeer samples, we conducted PCA with smartpca in EIGENSOFT3.0 software¹²² using the detected SNPs in each individual. Second, using the PCA plot result, we divided the 23 samples according to their grouping in the PCA plot and computed the average pairwise nucleotide diversity within a population π and the proportion of polymorphic sites Watterson's θ for the main clusters using the Bio::PopGen::Statistics package in BioPerl (v1.6.924)¹²³. The same program was used to compute the population differentiation estimate, with pairwise F_{st} between the two population groups (i.e., the PCA cluster result).

Data availability

Raw sequence data of Finnish male reindeer used for genome assembly is available under NCBI BioProject accession code PRJNA609893 and the whole genome sequence data of 23 reindeer is available under ENA (European Nucleotide Archive) study accession code PRJEB37216.

Received: 11 September 2019; Accepted: 5 May 2020;

Published online: 02 June 2020

References

- Helskog, K. & Indrelid, S. Humans and reindeer. *Quat. Int.* **238**, 1–3 (2011).
- Bjørklund, I. Domestication, Reindeer Husbandry and the Development of Sámi Pastoralism. *Acta Borealis*. **30**, 174–189 (2013).
- Flagstad, O. & Røed, K. H. Refugial origins of reindeer (*Rangifer tarandus* L.) inferred from mitochondrial DNA sequences. *Evolution* **57**, 658–670 (2003).
- Lin, Z. *et al.* Biological adaptations in the Arctic cervid, the reindeer (*Rangifer tarandus*). *Science* (80-). **364**, eaav6312 (2019).
- Blix, A. S. Adaptations to polar life in mammals and birds. *J. Exp. Biol.* **219**(1093), LP-1105 (2016).
- Van Oort, B. E. H. *et al.* Circadian organization in reindeer. *Nature* **438**, 1095–1096 (2005).
- Lu, W., Meng, Q.-J., Tyler, N. J. C., Stokkan, K.-A. & Loudon, A. S. I. A Circadian Clock Is Not Required in an Arctic Mammal. *Curr. Biol.* **20**, 533–537 (2010).
- Stokkan, K.-A., Van Oort, B. E. H., Tyler, N. J. C. & Loudon, A. S. I. Adaptations for life in the Arctic: evidence that melatonin rhythms in reindeer are not driven by a circadian oscillator but remain acutely sensitive to environmental photoperiod. *J. Pineal Res.* **43**, 289–293 (2007).
- Karl-Arne, S. *et al.* Shifting mirrors: adaptive changes in retinal reflections to winter darkness in Arctic reindeer. *Proc. R. Soc. B Biol. Sci.* **280**, 20132451 (2013).
- Miller, J., Patrick, D. & Volker, B. Antlers on the Arctic Refuge: capturing multi-generational patterns of calving ground use from bones on the landscape. *Proc. R. Soc. B Biol. Sci.* **280**, 20130275 (2013).
- Anderson, D. G., Kvie, K. S., Davydov, V. N. & Røed, K. H. Maintaining genetic integrity of coexisting wild and domestic populations: Genetic differentiation between wild and domestic *Rangifer* with long traditions of intentional interbreeding. *Ecol. Evol.* **7**, 6790–6802 (2017).
- Kvie, K. S., Heggnes, J. & RÅ, ed H., K. Merging and comparing three mitochondrial markers for phylogenetic studies of Eurasian reindeer (*Rangifer tarandus*). *Ecol. Evol.* **6**, 4347–4358 (2016).
- Røed, K. H., Bjørklund, I. & Olsen, B. J. From wild to domestic reindeer – Genetic evidence of a non-native origin of reindeer pastoralism in northern Fennoscandia. *J. Archaeol. Sci. Reports* **19**, 279–286 (2018).
- Weckworth, B. V., Musiani, M., McDevitt, A. D., Hebblewhite, M. & Mariani, S. Reconstruction of caribou evolutionary history in Western North America and its implications for conservation. *Mol. Ecol.* **21**, 3610–3624 (2012).

15. Røed, K. H. *et al.* Genetic analyses reveal independent domestication origins of Eurasian reindeer. *Proceedings. Biol. Sci.* **275**, 1849–55 (2008).
16. Hillier, L. W. *et al.* Sequence and comparative analysis of the chicken genome provide unique perspectives on vertebrate evolution. *Nature* **432**, 695–716 (2004).
17. Bovine Genome Sequencing and Analysis Consortium. C. G. *et al.* The genome sequence of taurine cattle: a window to ruminant biology and evolution. *Science* **324**, 522–8 (2009).
18. Groenen, M. A. M. *et al.* Analyses of pig genomes provide insight into porcine demography and evolution. *Nature* **491**, 393–398 (2012).
19. Jiang, Y. *et al.* The sheep genome illuminates biology of the rumen and lipid metabolism. *Science* **344**, 1168–1173 (2014).
20. Wade, C. M. *et al.* Genome sequence, comparative analysis, and population genetics of the domestic horse. *Science* **326**, 865–7 (2009).
21. Li, Z. *et al.* Draft genome of the reindeer (*Rangifer tarandus*). *Gigascience* **6**, 1–5 (2017).
22. Taylor, S. R. *et al.* The Caribou (*Rangifer tarandus*) Genome. *Genes* **10**, (2019).
23. Weldenegodguad, M. *et al.* Whole-Genome Sequencing of Three Native Cattle Breeds Originating From the Northernmost Cattle Farming Regions. *Frontiers in Genetics* **9**, 728 (2019).
24. Liu, B. *et al.* Estimation of genomic characteristics by analyzing k-mer frequency in de novo genome projects. *ArXiv e-prints* **1308**, arXiv **1308**, 2012 (2013).
25. Luo, R. *et al.* SOAPdenovo2: an empirically improved memory-efficient short-read de novo assembler. *Gigascience* **1**, 18 (2012).
26. Li, R. *et al.* The sequence and de novo assembly of the giant panda genome. *Nature* **463**, 311–7 (2010).
27. Qiu, Q. *et al.* The yak genome and adaptation to life at high altitude. *Nat. Genet.* **44**, 946–949 (2012).
28. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics* **25**, 1754–60 (2009).
29. Kent, W. J. BLAT—the BLAST-like alignment tool. *Genome Res.* **12**, 656–64 (2002).
30. Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs.
31. Curwen, V. *et al.* The Ensembl Automatic Gene Annotation System. *Genome Res.* **14**, 942–950 (2004).
32. Jurka, J. *et al.* Repbase Update, a database of eukaryotic repetitive elements. *Cytogenet. Genome Res.* **110**, 462–7 (2005).
33. Smit, A. F. A. & Hubley, R. RepeatModeler Open-1.0. 2008–2010.
34. Lindblad-Toh, K. *et al.* Genome sequence, comparative analysis and haplotype structure of the domestic dog. *Nature* **438**, 803–819 (2005).
35. Dong, Y. *et al.* Sequencing and automated whole-genome optical mapping of the genome of a domestic goat (*Capra hircus*). *Nat. Biotechnol.* **31**, 135–141 (2012).
36. Cardona, A. *et al.* Genome-Wide Analysis of Cold Adaptation in Indigenous Siberian Populations. *PLoS One* **9**, e98076 (2014).
37. Drivenes, Ø. *et al.* Isolation and characterization of two teleost melanopsin genes and their differential expression within the inner retina and brain. *J. Comp. Neurol.* **456**, 84–93 (2003).
38. Berson, D. M., Dunn, F. A. & Takao, M. Phototransduction by Retinal Ganglion Cells That Set the Circadian Clock. *Science (80-)*. **295**(1070), LP-1073 (2002).
39. Hannibal, J., Hindersson, P., Knudsen, S. M., Georg, B. & Fahrenkrug, J. The Photopigment Melanopsin Is Exclusively Present in Pituitary Adenylate Cyclase-Activating Polypeptide-Containing Retinal Ganglion Cells of the Retinohypothalamic Tract. *J. Neurosci.* **22**(RC191), LP-RC191 (2002).
40. Li, H. *et al.* TreeFam: a curated database of phylogenetic trees of animal gene families. *Nucleic Acids Res.* **34**, D572–D580 (2005).
41. Consortium, T. B. C. G. S. and A. *et al.* Genome sequences of wild and domestic bactrian camels. *Nat. Commun.* **3**, 1202 (2012).
42. Li, H. & Durbin, R. Inference of human population history from individual whole-genome sequences. *Nature* **475**, 493–496 (2011).
43. McKenna, A. *et al.* The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303 (2010).
44. Cingolani, P. *et al.* A program for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster* strain w(1118); iso-2; iso-3. *Fly (Austin)*. **6**, 80–92 (2012).
45. Lachance, J. *et al.* Evolutionary history and adaptation from high-coverage whole-genome sequences of diverse African hunter-gatherers. *Cell* **150**, 457–469 (2012).
46. Choi, J.-W. *et al.* Whole-Genome Analyses of Korean Native and Holstein Cattle Breeds by Massively Parallel Sequencing. *PLoS One* **9**, e101127 (2014).
47. Choi, J.-W. *et al.* Whole-Genome Resequencing Analysis of Hanwoo and Yanbian Cattle to Identify Genome-Wide SNPs and Signatures of Selection. *Mol. Cells* **38**, 466–473 (2015).
48. Stothard, P. *et al.* Whole genome resequencing of black Angus and Holstein cattle for SNP and CNV discovery. *BMC Genomics* **12**, 559 (2011).
49. Hahn, C., Bachmann, L. & Chevreaux, B. Reconstructing mitochondrial genomes directly from genomic next-generation sequencing reads—a baiting and iterative mapping approach. *Nucleic Acids Res.* **41**, e129–e129 (2013).
50. Chen, L. *et al.* Large-scale ruminant genome sequencing provides insights into their evolution and distinct traits. *Science (80-)*. **364**, eaav6202 (2019).
51. Wu, H. *et al.* Erratum: Camelid genomes reveal evolution and adaptation to desert environments. *Nat. Commun.* **6**, 6107 (2015).
52. Denton, J. F. *et al.* Extensive Error in the Number of Genes Inferred from Draft Genome Assemblies. *PLoS Comput. Biol.* **10**, e1003998 (2014).
53. Barnosky, A. D., Koch, P. L., Feranec, R. S., Wing, S. L. & Shabel, A. B. Assessing the Causes of Late Pleistocene Extinctions on the Continents. *Science (80-)*. **306**, 70 (2004).
54. Hughes, P. D., Woodward, J. C. & Gibbard, P. L. Middle Pleistocene cold stage climates in the Mediterranean: New evidence from the glacial record. *Earth and Planetary Science Letters* **253**, 50–56 (2007).
55. Yokoyama, Y., Lambeck, K., De Deckker, P., Johnston, P. & Fifield, L. K. Timing of the Last Glacial Maximum from observed sea-level minima. *Nature* **406**, 713 (2000).
56. Librado, P. *et al.* The Evolutionary Origin and Genetic Makeup of Domestic Horses. *Genetics* **204**, 423 (2016).
57. Yang, J. *et al.* Whole-Genome Sequencing of Native Sheep Provides Insights into Rapid Adaptations to Extreme Environments. *Mol. Biol. Evol.* **33**, 2576–2592 (2016).
58. Sudmant, P. H. *et al.* Diversity of human copy number variation and multicopy genes. *Science* **330**, 641–6 (2010).
59. Laity, J. H., Lee, B. M. & Wright, P. E. Zinc finger proteins: new insights into structural and functional diversity. *Curr. Opin. Struct. Biol.* **11**, 39–46 (2001).
60. Wei, S. *et al.* Emerging roles of zinc finger proteins in regulating adipogenesis. *Cell. Mol. Life Sci.* **70**, 4569–4584 (2013).
61. Grey, C., Baudat, F. & de Massy, B. PRDM9, a driver of the genetic map. *PLoS Genet.* **14**, e1007479 (2018).
62. Hohenauer, T. & Moore, A. W. The Prdm family: expanding roles in stem cells and development. *Development* **139**(2267), LP-2282 (2012).
63. Brzezinski, J. A., Lamba, D. A. & Reh, T. A. Blimp1 controls photoreceptor versus bipolar cell fate choice during retinal development. *Development* **137**(619), LP-629 (2010).
64. Brzezinski, J. A., Uoon Park, K. & Reh, T. A. Blimp1 (Prdm1) prevents re-specification of photoreceptors into retinal bipolar cells by restricting competence. *Dev. Biol.* **384**, 194–204 (2013).

65. Katoh, K. *et al.* Blimp1 Suppresses Chx10 Expression in Differentiating Retinal Photoreceptor Precursors to Ensure Proper Photoreceptor Development. *J. Neurosci.* **30**(6515), LP-6526 (2010).
66. Niimura, Y. Olfactory receptor multigene family in vertebrates: from the viewpoint of evolutionary genomics. *Curr. Genomics* **13**, 103–14 (2012).
67. Peng, J.-B., Suzuki, Y., Gyimesi, G. & Hediger, M. A. TRPV5 and TRPV6 Calcium-Selective Channels. <https://doi.org/10.1201/9781315152592-13> (2018).
68. Van Abel, M., Hoenderop, J. G. J. & Bindels, R. J. M. The epithelial calcium channels TRPV5 and TRPV6: regulation and implications for disease. *Naunyn. Schmiedebergs. Arch. Pharmacol.* **371**, 295–306 (2005).
69. Van der Eems, K. L., Brown, R. D. & Gundberg, C. M. Circulating levels of 1,25 dihydroxyvitamin D, alkaline phosphatase, hydroxyproline, and osteocalcin associated with antler growth in white-tailed deer. *Acta Endocrinol. (Copenh).* **118**, 407–414 (1988).
70. Sempere, A. J., Grimberg, R., Silve, C., Tau, C. & Garabedian, M. Evidence for Extrarenal Production of 1,25-Dihydroxyvitamin During Physiological Bone Growth: *In Vivo* and *In Vitro* Production by Deer Antler Cells. *Endocrinology* **125**, 2312–2319 (1989).
71. Susek, K. H., Karvouni, M., Alici, E. & Lundqvist, A. The Role of CXC Chemokine Receptors 1–4 on Immune Cells in the Tumor Microenvironment. *Frontiers in Immunology* **9**, 2159 (2018).
72. Zarbock, A. & Stadtmann, A. CXCR2: From Bench to Bedside. *Frontiers in Immunology* **3**, 263 (2012).
73. Costa, M. J. *et al.* Optimal design, anti-tumour efficacy and tolerability of anti-CXCR4 antibody drug conjugates. *Sci. Rep.* **9**, 2443 (2019).
74. Samuel, C. E. Antiviral actions of interferons. *Clin. Microbiol. Rev.* **14**, 778–809 (2001).
75. Goodbourn, S., Didcock, L. & Randall, R. E. Interferons: cell signalling, immune modulation, antiviral response and virus countermeasures. *J. Gen. Virol.* **81**, 2341–2364 (2000).
76. Yang, Z.-Q. *et al.* IGF1 regulates RUNX1 expression via IRS1/2: Implications for antler chondrocyte differentiation. *Cell Cycle* **16**, 522–532 (2017).
77. Hu, W. *et al.* MicroRNA let-7a and let-7f as novel regulatory factors of the sika deer (*Cervus nippon*) IGF-1R gene. *Growth Factors* **32**, 27–33 (2014).
78. Jones, P. G. & Inouye, M. The cold-shock response — a hot topic. *Mol. Microbiol.* **11**, 811–818 (1994).
79. Thieringer, H. A., Jones, P. G. & Inouye, M. Cold shock and adaptation. *BioEssays* **20**, 49–57 (1998).
80. Wu, S., De Croos, J. N. A. & Storey, K. B. Cold acclimation-induced up-regulation of the ribosomal protein L7 gene in the freeze tolerant wood frog, *Rana sylvatica*. *Gene* **424**, 48–55 (2008).
81. Yurchenko, A. A. *et al.* Scans for signatures of selection in Russian cattle breed genomes reveal new candidate genes for environmental adaptation and acclimation. *Sci. Rep.* **8**, 12984 (2018).
82. Provencio, I., Jiang, G., De Grip, W. J., Hayes, W. P. & Rollag, M. D. Melanopsin: An opsin in melanophores, brain, and eye. *Proc. Natl. Acad. Sci.* **95**(340), LP-345 (1998).
83. Provencio, I. *et al.* A Novel Human Opsin in the Inner Retina. *J. Neurosci.* **20**(600), LP-605 (2000).
84. Gooley, J. J., Lu, J., Chou, T. C., Scammell, T. E. & Saper, C. B. Melanopsin in cells of origin of the retinohypothalamic tract. *Nat. Neurosci.* **4**, 1165 (2001).
85. Lolignier, S. *et al.* New Insight in Cold Pain: Role of Ion Channels, Modulation, and Clinical Perspectives. *J. Neurosci.* **36**(11435), LP-11439 (2016).
86. Luiz, A. P. *et al.* Cold sensing by Nav1.8-positive and Nav1.8-negative sensory neurons. *Proc. Natl. Acad. Sci.* **116**(3811), LP-3816 (2019).
87. Zimmermann, K. *et al.* Transient receptor potential cation channel, subfamily C, member 5 (TRPC5) is a cold-transducer in the peripheral nervous system. *Proc. Natl. Acad. Sci.* **108**(18114), LP-18119 (2011).
88. Lolignier, S. *et al.* The Nav1.9 Channel Is a Key Determinant of Cold Pain Sensation and Cold Allodynia. *Cell Rep.* **11**, 1067–1078 (2015).
89. Leipold, E. *et al.* Cold-aggravated pain in humans caused by a hyperactive Nav1.9 channel mutant. *Nat. Commun.* **6**, 10049 (2015).
90. Yannic, G. *et al.* Genetic diversity in caribou linked to past and future climate change. *Nat. Clim. Chang.* **4**, 132 (2013).
91. Roed Knut H. A4 - Bjørnstad, Gro A4 - Flagstad, Øystein A4 - Haanes, Hallvard A4 - Hufthammer, Anne K. A4 - Jordhøy, Per A4 - Rosvold, Jørgen, K. H. A.-R. Ancient DNA reveals prehistoric habitat fragmentation and recent domestic introgression into native wild reindeer. *Conserv. Genet.* v. 15, 1137-1149-2014 v.15 no.5 (2014).
92. Roed, K. H. *et al.* Male phenotypic quality influences offspring sex ratio in a polygynous ungulate. *Proc. R. Soc. B Biol. Sci.* **274**, 727–733 (2007).
93. Bjørnstad, G., Flagstad, Ø., Hufthammer, A. K. & Roed, K. H. Ancient DNA reveals a major genetic change during the transition from hunting economy to reindeer husbandry in northern Scandinavia. *J. Archaeol. Sci.* **39**, 102–108 (2012).
94. Grabherr, M. G. *et al.* Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nat Biotech* **29**, 644–652 (2011).
95. Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet journal* **17**, 10–12 (2011).
96. Benson, G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic Acids Res.* **27**, 573–580 (1999).
97. Tempel, S. Using and understanding repeatMasker. *Methods Mol. Biol.* **859**, 29–51 (2012).
98. Maja, T. & Nansheng, C. Using RepeatMasker to Identify Repetitive Elements in Genomic Sequences. *Curr. Protoc. Bioinforma.* **25**, 4.10.1–4.10.14 (2009).
99. Birney, E., Clamp, M. & Durbin, R. GeneWise and Genomewise. *Genome Res.* **14**, 988–95 (2004).
100. Kim, D. *et al.* TopHat2: accurate alignment of transcriptomes in the presence of insertions, deletions and gene fusions. *Genome Biol.* **14**, R36 (2013).
101. Trapnell, C. *et al.* Differential gene and transcript expression analysis of RNA-seq experiments with TopHat and Cufflinks. *Nat. Protoc.* **7**, 562–578 (2012).
102. Mulder, N. J. *et al.* New developments in the InterPro database. *Nucleic Acids Res.* **35**, D224–8 (2007).
103. Kanehisa, M. & Goto, S. KEGG: kyoto encyclopedia of genes and genomes. *Nucleic Acids Res.* **28**, 27–30 (2000).
104. Bairoch, A. & Apweiler, R. The SWISS-PROT protein sequence database and its supplement TrEMBL in 2000. *Nucleic Acids Res.* **28**, 45–48 (1999).
105. The UniProt Consortium. UniProt: the universal protein knowledgebase. *Nucleic Acids Res.* **45**, D158–D169 (2016).
106. Nawrocki, E. P., Kolbe, D. L. & Eddy, S. R. Infernal 1.0: inference of RNA alignments. *Bioinformatics* **25**, 1335–1337 (2009).
107. Lowe, T. M. & Eddy, S. R. tRNAscan-SE: a program for improved detection of transfer RNA genes in genomic sequence. *Nucleic Acids Res.* **25**, 955–964 (1997).
108. Gardner, P. P. *et al.* Rfam: updates to the RNA families database. *Nucleic Acids Res.* **37**, D136–D140 (2008).
109. Liu, S. *et al.* Population Genomics Reveal Recent Speciation and Rapid Evolutionary Adaptation in Polar Bears. *Cell* **157**, 785–794 (2014).
110. The Gene, O. C. *et al.* Gene Ontology: tool for the unification of biology. *Nat. Genet.* **25**, 25–29 (2000).
111. Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high throughput. *Nucleic Acids Res.* **32**, 1792–1797 (2004).
112. Guindon, S. & Gascuel, O. A Simple, Fast, and Accurate Algorithm to Estimate Large Phylogenies by Maximum Likelihood. *Syst. Biol.* **52**, 696–704 (2003).

113. Guindon, S. *et al.* New Algorithms and Methods to Estimate Maximum-Likelihood Phylogenies: Assessing the Performance of PhyML 3.0. *Syst. Biol.* **59**, 307–321 (2010).
114. Rannala, B., Yang, Z. & Anderson, F. Inferring Speciation Times under an Episodic Molecular Clock. *Syst. Biol.* **56**, 453–466 (2007).
115. Yang, Z. & Rannala, B. Bayesian Estimation of Species Divergence Times Under a Molecular Clock Using Multiple Fossil Calibrations with Soft Bounds. *Mol. Biol. Evol.* **23**, 212–226 (2006).
116. Yang, Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* **24**, 1586–91 (2007).
117. De Bie, T., Cristianini, N., Demuth, J. P. & Hahn, M. W. CAFE: a computational tool for the study of gene family evolution. *Bioinformatics* **22**, 1269–71 (2006).
118. Löytynoja, A. & Goldman, N. An algorithm for progressive multiple alignment of sequences with insertions. *Proc. Natl. Acad. Sci. USA* **102**, 10557–62 (2005).
119. Tian, T. *et al.* agriGO v2.0: a GO analysis toolkit for the agricultural community, 2017 update. *Nucleic Acids Res.* **45**, W122–W129 (2017).
120. DePristo, M. A. *et al.* A framework for variation discovery and genotyping using next-generation DNA sequencing data. *Nat. Genet.* **43**, 491–498 (2011).
121. Bernt, M. *et al.* MITOS: Improved de novo metazoan mitochondrial genome annotation. *Molecular Phylogenetics and Evolution* **69**, 313–319 (2013).
122. Patterson, N., Price, A. L. & Reich, D. Population Structure and Eigenanalysis. *PLoS Genet* **2**, e190 (2006).
123. Stajich, J. E. *et al.* The Bioperl Toolkit: Perl Modules for the Life Sciences. *Genome Res.* **12**, 1611–1618 (2002).

Acknowledgements

This study was financially supported by the Academy of Finland in the Arctic Research Programme ARKTIKO (decision number 286040). M.W. acknowledges a study grant from the Finnish Cultural Foundation. We are grateful to Juhani Majjala for providing the reindeer that was used as a reference genome. The authors thank Nuccio Mazzullo, Päivi Soppela and Anna Stammler-Gossmann from the Arctic Centre, University of Lapland, Rovaniemi, Finland for collaboration in the sampling of the Finnish reference reindeer. The authors wish to acknowledge the CSC – IT Center for Science, Finland, for computational resources. The owners of the animals included in the study are acknowledged for providing samples for this study.

Author contributions

J.K. conceived and designed the project. J.K., T.R., M.H., J.P. and K.R. collected the samples. M.W., K.P. and M.Y. performed bioinformatics analyses. M.W. prepared the manuscript draft with substantial contribution from J.K. and K.P. All authors read, reviewed and approved the final manuscript.

Competing interests

The authors declare no competing interests.

Additional information

Supplementary information is available for this paper at <https://doi.org/10.1038/s41598-020-65487-y>.

Correspondence and requests for materials should be addressed to J.K.

Reprints and permissions information is available at www.nature.com/reprints.

Publisher's note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons license, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons license, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons license and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this license, visit <http://creativecommons.org/licenses/by/4.0/>.

© The Author(s) 2020