

GSD_1.0 (canFam4): A novel canine reference genome resolves genomic architecture and uncovers transcript complexity

Chao Wang (✉ chao.wang@imbim.uu.se)

Uppsala University <https://orcid.org/0000-0003-3936-4023>

Ola Wallerman

Uppsala University

Maja-Louise Arendt

Uppsala University; University of Copenhagen

Elisabeth Sundström

Uppsala University <https://orcid.org/0000-0001-9526-8541>

Åsa Karlsson

Uppsala University

Jessika Nordin

Uppsala University <https://orcid.org/0000-0002-8414-2190>

Suvi Mäkeläinen

Swedish University of Agricultural Sciences <https://orcid.org/0000-0001-7378-4991>

Gerli Rosengren Pielberg

Uppsala University

Jeanette Hanson

Swedish University of Agricultural Sciences

Åsa Ohlsson

Swedish University of Agricultural Sciences

Sara Saellström

Swedish University of Agricultural Sciences <https://orcid.org/0000-0001-5253-5830>

Henrik Rönnberg

Swedish University of Agricultural Sciences <https://orcid.org/0000-0001-6098-5364>

Ingrid Ljungvall

Swedish University of Agricultural Sciences

Jens Häggström

Swedish University of Agricultural Sciences

Tomas Bergström

Swedish University of Agricultural Sciences (SLU) <https://orcid.org/0000-0002-7480-2669>

Åke Hedhammar

Swedish University of Agricultural Sciences

Jennifer Meadows

Uppsala University

Kerstin Lindblad-Toh (✉ kersli@broadinstitute.org)

The Broad Institute

Article

Keywords: GSD_1.0, RNA-seq, miRNA-seq, ATAC-seq, Dog Leukocyte Antigen (DLA), T Cell Receptor (TCR), 366 COSMIC

Posted Date: August 25th, 2020

DOI: <https://doi.org/10.21203/rs.3.rs-57997/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Communications Biology on February 10th, 2021. See the published version at <https://doi.org/10.1038/s42003-021-01698-x>.

Abstract

We present GSD_1.0, a novel high-quality domestic dog reference genome with chromosome length scaffolds and gap number reduced 41-fold, from 23,836 to 585. Annotation with novel and existing long and short read RNA-seq, miRNA-seq and ATAC-seq, revealed that 32.1% of closed gaps harboured previously hidden functional elements, including promoters, genes and miRNAs. A catalogue of canine “dark” regions was made to facilitate mapping rescue. Alignment in these regions is difficult, but we demonstrate that they harbour trait-associated variation. Key genomic regions were completed, including the Dog Leukocyte Antigen (DLA), T Cell Receptor (TCR) and 366 COSMIC cancer genes. The sequencing of 27 dogs from 19 breeds with linked read technology uncovered 22.1 million SNPs, indels and larger structural variants. Intersection with protein coding genes showed that 1.4% could directly influence gene products, and so provide a source of normal or aberrant phenotypic modifications.

Introduction

Domestic dogs have lived alongside humans for at least 10,000 years^{1,2}, and during this time, they have adapted to a shared environment and diet, whilst being selectively bred for traits such as morphology³ and behavior⁴. Man and dog also share orthologous genes, genomic architecture and disease sets, placing the dog as an important comparative species for human genetics and genomics. Taking advantage of pet dog medical records, within breed homogeneity and disease risk enrichment, it has been possible to provide insights into both rare and common spontaneous disease. The Online Mendelian Inheritance in Animals website (OMIA, June 2020, omia.org) currently catalogues 774 canine traits with linked genetic associations, 234 of which are likely causative in the canine models for human disease. The types of canine variants implicated in disease range from SNPs (e.g. a missense variation in *SOD1* leading to degenerative myelopathy⁵) through to complex genomic rearrangements (e.g. a deletion in the repetitive interferon alpha gene cluster associated with hypothyroidism⁶), and were identified with canine SNP chips, e.g. 170K Canine HD Array (Illumina), genotyping complemented with imputation⁷ or genome and transcriptome sequencing of individuals, families⁸ or large populations³. Clearly, genome contiguity, as well as gene and regulatory element annotation from a range of diverse breeds and tissues are all required to translate association to causation.

The current canine reference genome, CanFam3.1, is based on a 2005 7.4X Sanger sequencing framework⁹, improved in 2014 with multiple methods to better resolve euchromatic regions and annotate transcripts from gross tissues¹⁰. However, it still contains 23,876 gaps, with 19.6% of these within gene bodies, and a further 9.8% located a mere 5 kb upstream of predicted gene start sites. These gaps result from the accumulation regions which are difficult to sequence, and are in part due to the loss of *PRDM9* which leads to genomic sections with very high GC content¹¹. The consequence of this is loss of promoters, CpG islands and other regulatory elements from the reference; sequences which may hold the key to deciphering complex traits^{12,13}.

To drive canine comparative genomics forward, we generated a high-quality canine reference assembly using a combination of PacBio long read sequencing, 10x Genomics Chromium Linked-Reads (henceforth called 10x) and Hi-C proximity ligation. The new reference, GSD_1.0 (canFam4), with 41-fold gap reduction, was subsequently annotated with both novel and published WGS, ATAC- and RNA sequencing to enhance gene models and variant annotation. The result was 159 thousand transcripts across 29,583 genes. This novel data opens the door to the identification of functional variants underlying complex traits, especially in difficult to sequence, and often biologically important, regions.

Results And Discussion

De novo assembly. Mischka, a 12-year-old female German Shepherd, was selected as the source for our high-quality reference genome assembly (UU_CFam_GSD_1.0/canFam4; henceforth called GSD_1.0). Mischka was free of known genetic disorders, and when compared with additional German Shepherd sourced from within Sweden, was found to be genetically representative of the breed (Supplementary Fig. 1). We sequenced the genome using ~ 100X coverage PacBio long reads and assembled these in contigs with the standard FALCON method¹⁴. Further scaffolding using 94X of 10x and 48X of Hi-C linked reads resulted in 39 single-scaffold chromosomes (total 2.35 Gb) and 2,159 unplaced scaffolds (total 128.5 Mb; Fig. 1a). The latter contigs predominantly contain segmental duplications (58.1%) and centromeric repeats (30.1%, Supplementary Fig. 2).

Reference benchmarking. Compared to CanFam3.1, the contiguity of GSD_1.0 has been improved 41-fold, reaching a contig N50 of 14.8 Mb (Supplementary Fig. 3), with only 367 gaps in the chromosome (chr) scaffolds (Table 1; Fig. 1a). The quantity of sequence with extreme GC content (> 90% in 50 bp-windows) has doubled to 1.7 Mb (Fig. 1b), leading to a 14% increase in the average length of CpG islands (1,056 bp versus 926 bp, $P = 8.4 \times 10^{-4}$, t-test). Filled gaps were found to have either high GC or high repeat content (Fig. 1c).

Table 1
Assembly statistics of GSD_1.0 compared to CanFam3.1

| | GSD_1.0 | CanFam3.1 |
|-----------------------|--------------------|--------------------|
| Number of contigs | 2,783 | 27,104 |
| N50 (L50) contig | 14,840,767 bp (57) | 267,478 bp (2,436) |
| Number of scaffolds | 2,198 | 3,268 |
| N50 (L50) scaffolds | 64,299,765 bp (15) | 63,241,923 bp (15) |
| Number of Gaps | 585 | 23,876 |
| Gap density (gaps/Mb) | 0.24 | 9.9 |
| Total bases | 2,482,000,080 bp | 2,410,976,875 bp |
| Total ungapped bases | 2,481,941,580 bp | 2,392,715,236 bp |

Repeat structure. Approximately 42.7% of the genome is repetitive sequence, with the three major categories being LINEs (504 Mb), SINEs (253 Mb) and LTRs (120 Mb) (Supplementary Fig. 4, Supplementary Table 1). Long read technology allowed for the further resolution of centromeric repeats, and based on their positions, the orientation of chr 27 and 32 were reversed compared to CanFam3.1. In addition, the q-arms of 21 autosomes now begin with centromeric repeats, and 17 autosomes end in telomeric repeats (Fig. 1a). As expected, the sub-metacentric chr X has telomeric repeats at each end, and a clear centromeric signal at 49.4–49.9 Mb. Throughout the genome we found 10 internal centromeric, and 7 internal telomeric repeats. These may indicate ancient centromere and telomere positions prior to chromosomal rearrangements and most were also present in the previous reference genome assembly.

Functional annotation. To resolve transcript complexity and account for the gap closures in GSD_1.0, we generated more than 70 M nanopore and PacBio full-length cDNA reads from 40 tissues (including 15 brain regions; Supplementary Table 2), and combined this with 24 billion public RNA-seq paired reads (Supplementary Table 3). The annotation consisted of 159 thousand transcripts in 29,583 genes; of which 20,654 had an open reading frame (ORF) of at least 100 amino acids, and 19,691 genes had a significant BLAST hit against proteins in Swissprot or ENSEMBL. Further, 7,725 were defined as long noncoding genes. Compared to proteins extracted from CanFam3.1, our new GSD-1.0 annotation has a higher number of genes with BLAST hits and the number of genes with a full-length match has increased by 11% (Supplementary Fig. 5). Gene predictions and non-dog refSeq alignments were used to identify potentially missed genes that did not overlap with our annotation, yielding an additional 874 protein coding genes with BLAST evidence. Using a combination of new miRNA-seq reads and public data we identified a conservative set of 719 miRNAs, similar to the set found for CanFam3.1¹⁵. Among the novel miRNAs, a copy of the highly expressed Mirlet-7i was identified in a filled gap region (Supplementary Fig. 6). This miRNA has been implicated in several human diseases, including multiple sclerosis¹⁶, gastric cancer¹⁷ and breast cancer¹⁸, but has yet to be extensively studied in dogs.

We identified 7,468 closed gaps containing either an exon or promoter sequence as defined by ATAC-seq peaks, accounting for 5,743 unique coding exons which were missing in CanFam3.1 (Fig. 1e). Notably, eight genes with expression across multiple tissues were completely absent in CanFam3.1 but were now available for interrogation (*PSMA4*, *CDHR5*, *SCT*, *PAOX*, *UTF1*, *EFNA2*, *GPX4* and *SLC25A22*). These genes have diverse functions ranging from embryonic stem cell co-activator (*UTF1*), through to osmoregulation (*SCT*). Both *CDHR5* and *SLC25A22* (Fig. 2a) have been investigated as biomarkers for either renal¹⁹ or colorectal²⁰ cancers.

Implications for research. We assessed the chromosomal order and contiguity of regions essential to the study of cancer and immunological disease. Using the human COSMIC²¹ gene list as a baseline, we affirmed that 282 tier1 and 78 tier2 genes are now completely captured, including *HOXD13* and *KLF4* (Supplementary Table 4). Both have been implicated in human breast cancer; *HOXD13* methylation status functions as a prognostic indicator²² and deubiquitination of *KLF4* promotes metastasis²³ (Supplementary Fig. 7). The Dog Leukocyte Antigen (DLA) regions positioned on chr 12 (Fig. 1f) and 35 (Supplementary Fig. 8a) are contiguous in GSD_1.0 (covering 2.58 Mb and 0.61 Mb, respectively) and contain new coding and potential regulatory sequences in filled gaps. Contiguous sequence was also reported for both the T-cell receptor alpha (TRA) and T-cell receptor beta (TRB) loci on chr 8 and 16 respectively (Supplementary Fig. 8b&c).

Comparison to canine assemblies. Five additional canine genome assemblies have recently been deposited in NCBI (Supplementary Table 5). For each assembly, we compared BUSCO²⁴ scores and mappability using in-house Iso-Seq cDNA alignments generated above from a beagle dog (Supplementary Table 2). With GSD_1.0 it was possible to map > 5% more bases from 25,609 of Iso-Seq reads compared to CanFam3.1 (4.8% of total reads; Supplementary Fig. 9). This was a higher fraction than for the other assemblies (Supplementary Table 6). GSD_1.0 had the second highest BUSCO score for complete genes (95.5%), but each canine assembly, including the dingo, is of value to the community and may serve different experimental goals (Supplementary Fig. 10).

Genome variation. Polymorphisms detected in 27 dogs (19 breeds) were extracted from 10x sequencing data to facilitate the investigation of genome features and across-breed variant segregation (Supplementary Table 7). We identified 14,953,199 SNPs, 6,958,645 indels and 217,951 structural variants (SV, average 2.4 kb; Fig. 3b). Of these, 42.1% were private, 57.9% polymorphic across multiple individuals, and 1.4% overlapped with protein-coding regions (295,112 SNPs and 16,654 SVs). Intersection with existing SV catalogues based on either SNP or aCGH arrays²⁵⁻²⁷ showed between 12.6–39.0% agreement, but these numbers are likely a reflection of within project breed and detection technology. 10x sequencing allowed for the detection of many novel SVs with small to medium size (\geq 30 kb) with accurate breakpoints (Test set of 9 validated with PCR and Sanger sequencing, data not shown).

Genome “dark” regions unmasked. The majority of publicly available dog WGSs were generated with short read technologies. To facilitate the reanalysis of these resources with GSD_1.0 we aimed to identify

the genome's "dark" regions²⁸; those sections either not adequately covered by sequencing method ("covered", COV) or to which unique alignment is not possible ("camouflaged", CAM). We defined GSD_1.0 dark regions for Illumina short reads (ISR), 10x, and PacBio (PB) sequencing (see Methods). COV comprised 5.8 Mb, 5.7 Mb and 6.4 Mb respectively, while CAM comprised 15.9 Mb, 6.4 Mb and 1.0 Mb (Fig. 2c). Intersection showed that while 10x could rescue 11.3 Mb dark regions not seen with ISR (9.73 + 1.56 Mb), more than half of this again (5.9 Mb) could be further recovered by PacBio (Fig. 2d). We noted six tier 1 & 2 COSMIC genes that contained COV/CAM regions (*EPHA3*, *RALGDS*, *LRP1B*, *CSMD3*, *ZMYM2*, *PTEN*; 0.8–6.6% of coding region dark), potentially masking drivers of disease. Due to the nature of COV/CAM regions, default practices will not allow for the mapping of IRS reads to, and subsequent variant extraction from, these regions. Instead, we extracted variants overlapping annotated dark regions from our "healthy" 10x data set, and in doing so, identified 51,994 SNPs and indels, including 19,340 intronic and 2,074 exonic variants. Many of these variants were embedded in genes that may be important for morphology or associated with disease. For example, 14 variants were found within seven intronic *TYRP1* ISR CAM dark regions (Supplementary Fig. 11a): a gene linked to brown color in dogs²⁹ and melanoma in humans^{30,31}. Likewise, 76 variants were found in *ADCY2* ISR CAM regions (Supplementary Fig. 11b). Polymorphisms in this gene have previously been associated with neurological disorders (Bipolar disorder³²; and Alzheimer's disease³³), and response to associated drug therapies of schizophrenia³⁴ in humans.

Chromosome mis-assembly resolved. A direct comparison of CanFam3.1 and GSD_1.0 revealed a complex ~ 10 Mb inverted region on chr 9 that harboured *SOX9* and was previously implicated in canine XX disorder of sex development (DSD)^{35–37}. Three polymorphic regions homologous to parts of *MAGI2* on chr 18 (M1, M2, M3), have been inserted upstream of *SOX9* (Fig. 3a,b). In DSD, having multiple copies of a CNV overlapping M2³⁶, was shown to be associated with altered *SOX9* function during gonadal development. Using HiC and BAC end sequencing data, we confirmed the inverted GSD_1.0 orientation was correct and refined the placement of regions M1, M2 and M3 (Fig. 3a). These chr 9 insertions are missing from GSD_1.0, but allelic depth analysis revealed that most 10x dogs (26/27) carry between 2–6 chr 9 copies (Fig. 3c,d), similar to the estimates reported for non-DSD dogs³⁷. Recently it was shown that the DSD phenotype presents in a breed specific manner, and is influenced by the combination of a SNP and CNVs in this region^{35,37}. However, as this inversion contains numerous genes and regulatory elements, this rearrangement including multiple CNV expansions, has the potential to impact additional canine traits.

CYP1A2 locus variation. To further investigate the impact of SVs on coding genes, we examined the 16.2 kb copy number locus which encompassed *CYP1A2* (Fig. 4a). Dogs are used as comparative models for human xenobiotic metabolism, and while a *CYP1A2* premature stop codon (rs852922442 C > T) has been reported^{38,39}, the locus expansion has not. The homozygous T mutant is polymorphic across breeds⁴⁰ and results in an array of pharmacokinetic effects, including reduced hepatic drug metabolism⁴¹. This T variant was observed in 4/27 10x dogs, but in heterozygous form and not segregating with CNV count (2–5 copies; Fig. 4b). Differential gene expression analyses for this and

neighboring genes outside the locus were performed using either liver or spleen tissue from additional individuals (Supplementary Tables 8 and 2). After accounting for *CYP1A2* SNP rs852922442-T, no significant relative gene expression difference was observed, leaving the phenotypic consequence of this expansion unresolved (CNV 3 vs > 3; Supplementary Table 9). It may be that the effect in this region is subtle, and so not detectable with qPCR, however, *CYP1A2* is an inducible gene and so the true outcome may only be observed after a drug challenge⁴².

Conclusion

Through the combination of sequencing technologies, PacBio (> 90X) long read, 10x and Hi-C proximity ligation, we have generated the contiguous, chromosome length scaffolded GSD_1.0 canine reference genome. GSD_1.0 has a 41-fold gap reduction (from 23,836 to 585) compared with its predecessor CanFam3.1. This brings the canine reference genome quality in line with other key mammalian species e.g. human⁴³, mouse⁴⁴, and gorilla⁴⁵. For both human and mouse projects, the *de novo* sequence assembly of multiple individuals from different population backgrounds has revealed novel sequence not found in the single (hybrid in the case of human) species reference, and facilitated the search for population specific variants which likely contribute to traits of interest, including within the highly polymorphic immune gene clusters^{43,44}. While this type of *de novo* collection is on-going within the canine community, GSD_1.0 is the first genome of reference quality that is further annotated with novel long read RNA-sequencing data, allowing for the resolution of transcript complexity through regions with high GC context, or “dark” regions²⁸.

The resolution and placement of repeats in GSD_1.0, including non-LTR retrotransposons, will facilitate the study of gene and genome evolution and the process of neofunctionalization across mammalian lineages to an extent not possible previously. Over more recent timespans, these mobile elements can allow for genome slippage, and to the accumulation of within and across population SVs. In human clinical genomics, SVs spanning coding and/or non-coding sequence have been responsible for a range of maladies including cardiac anomalies (OMIM 192430) and intellectual delay and autism (OMIM 608 636). Accordingly, this source of variation is of keen interest in canine genetics, and should facilitate this line of investigation. The technology applied to GSD_1.0 to read across repeats, was also successful in reading into regions of constitutive heterochromatin, allowing for the correction of chromosomal direction (chr 27 and 32) and revealing novel centromeric and telomeric sequences.

Perhaps the largest gain offered by the contiguity of GSD_1.0 is to the accelerating field of low pass genotyping and imputation for trait mapping⁷. The completion of key regions to the investigation of immunological disease and cancer, e.g. DLA and TCR, when combined with large reference populations, will facilitate the more accurate genotyping of these regions and hopefully fast track the process from association to causation. We believe that the catalogues generated here (extended gene models, “dark” regions, within and across breed variation), based on the GSD_1.0 framework, will propel the comparison of canine and human genetic disease forward by leaps and bounds.

Methods

Detailed protocols can be found in the Supplementary Methods.

Reference individual. Mischka, a 12-year-old female German Shepherd, was born and raised in Sweden with known ancestral background and no medical history of genetic disease. Mischka was genotyped with the CanineHD BeadChip (Illumina) and compared to a population of 260 German Shepherds from a previous study⁴⁶. Mischka was assessed to be representative of the population via expected inbreeding value ($F = 0.037$) and MDS genetic distance measures (PLINK v1.9) and selected for the genome assembly. High molecular weight (HMW) DNA was extracted from blood with MagAttract HMW DNA Kit (QIAGEN).

Genome sequencing. The assembly used multiple sequencing technologies. Long read libraries were prepared with SMRTbell Template Prep Kit 1.0 and 70 SMRT cells were sequenced on the PacBio Sequel system with v2.1 chemistry (Pacific Biosciences; 276.86 Gb data). Linked reads were sequenced from HMW DNA with Chromium libraries (10x Genomics) on Illumina HiSeq X (2×150 bp; 269.75 Gb of data). Dovetail Genomics prepared three HiC libraries which were sequenced on Illumina HiSeq X (2×150 bp paired-end reads; 121.47 Gb data, Supplementary Table 10).

Assembly construction. *de novo* assembly used PacBio subreads (> 8 kb) with the standard FALCON⁴⁷ v0.5.0 method. After Arrow⁴⁷ v2.3.3 polishing, the assembly yielded 3,656 contigs with an N50 and mean length of 4.66 Mb and 677 kb respectively. ARCS⁴⁸ v1.05 and LINKS⁴⁹ v1.8.6, with the recommended link ratio (-a) 0.9, were used to scaffold contigs with 10x reads. 1,170 FALCON contigs were joined in this step, increasing the scaffold N50 to 18.5 Mb.

Conflict resolution. The assembled scaffold sequences were aligned onto the high-density canine linkage map⁵⁰. 21,278 of 22,362 markers (95%) were unambiguously mapped to the assembly by BLAT⁵¹ v36. Synteny of genetic and physical location of markers was further compared with Chromonomer⁵² v1.0, which showed 207 scaffolds were anchored correctly, but that four had conflicting markers. These four scaffolds were split after careful sequence review confirmed that each discrepancy arose from incorrect inter-chromosomal joining.

Gap filling and assembly polishing. PBjelly from PBSuite⁵³ v15.8.24 was used with PacBio subreads to close 648 gaps. An initial QC scan showed no putative wrong joins, and so long-distance interaction information from HiC (HiRise, Dovetail Genomics) was used to successfully extend scaffolds to chromosome-level (scaffold N50: 64.3 Mb). These results were evaluated with the JUICER⁵⁴ pipeline; HiC reads were mapped back to the HiRise assembly and HiC map with intra- and inter-chromosomal interactions visualised. We identified and manually adjusted contigs placed in either the wrong order or orientation (chr 6, 14, 17, 26 and X), and joined separated contigs from the same chromosome (chr 8 and 18). A second round of PBjelly gap filling closed another 110 gaps. The assembly was polished with Arrow (PacBio subreads) and Pilon⁵⁵ v1.22(10x Genomics reads, BWA⁵⁶ v0.7.15 mem mapping). A

FreeBayes-based method was applied to further correct the indel errors⁵⁷. SNPs and indels were called from short reads aligned to the polished assembly (FreeBayes⁵⁸ v1.1.0). The reference base was replaced with the variant allele at 149,264 positions where 10x sequencing depth was at least 30X and the variant allele ratio was > 90% using FastaAlternateReferenceMaker from GATK⁵⁹ v4.1.1.0. A final round of Pilon short read polishing was completed prior to the removal of 68 unplaced contigs with suspected bacterial contamination (Kraken2⁶⁰ v2.0.8).

The correctness of a large rearranged region on chr 9 of GSD1.0 was confirmed through comparison to end sequences from original CanFam BAC clones (CH82 library; NCBI TraceDB). BAC sequences were mapped as paired-reads (BWA⁵⁶ mem default setting), to GSD_1.0 and CanFam3.1. End pairs that mapped to both assemblies were compared and defined as concordant when they aligned in forward and reverse direction with a distance < 500 kb.

GC-content and repetitive elements. GC-content (%) was assessed in 50 bp windows (NUC from BEDTools⁶¹ v2.29.2). CpG islands were detected with the “cpg_lh” script, modified from Gardiner-Garden⁶². The unique mappability of GSD_1.0 was tested with different k-mers (50/150/250 bp in GEM-Tools⁶³ v1.71). Repetitive elements were annotated with Repeat Masker v4.0.8 sensitive mode (<http://www.repeatmasker.org>) and a combined library (dc20171107-rb20181026). Telomere repeats, “TTAGGG”, were highlighted on both strands with fuzznuc (EMBOSS⁶⁴ v6.6.0). Putative telomere sequences were defined as \geq [TTAGGG]₁₂ repeats, with less than 11 variant bases between each, and multiple sequences were merged if within 100 bp. Centromeric regions were defined based on satellite repeat⁶⁵ (CarSat1/Carsat2/SAT1_CF) content in 5 kb windows. Putative centromere sequences were annotated if the repeat content was > 80%.

RNA preparation and long read cDNA sequencing. Multiple RNA samples were used for sequencing (Supplementary Table 2). First, commercial hypothalamus RNA (Zyagen) was used in the PacBio Iso-Seq express protocol. Libraries were run on two separate SMRT-cells using the Sequel system (~ 500,000 reads; mean read length/library 2,452 bp and 3,451 bp). TRIzol (Invitrogen) extracted total RNA from a further 40 tissues (including 15 brain regions; Supplementary Table 2) was used for nanopore cDNA and Illumina miRNA-sequencing.

Gene annotation. Public canine Illumina RNA-seq samples (Supplementary Table 3) from diverse tissues and breeds was downloaded (<https://sra-explorer.info/>) and Stringtie2⁶⁶ used to assemble and merge transcripts with both PacBio iso-seq and nanopore full-length cDNA alignments. Coding regions were identified with TAMA⁶⁷ tool and BLAST to the curated Uniprot_Swissprot and ENSEMBL dog annotation v100. QC and settings for each stage are detailed in extended methods.

miRNA identification. Public micro-RNA-seq samples (Supplementary Table 3) were combined with the above brain micro-RNA-seq reads (Total reads, 1.3 billion). Reads were included if they were between 20–30 bases after adaptor trimming. Bowtie alignments of unique sequences were used for MiRDeep2⁶⁸

analysis and compared to known dog and human miRNAs (miRBase) in order to identify the position of both known and novel miRNAs.

ATAC-seq analysis. Reads from BARKbase⁶⁹ (Supplementary Table 3) were aligned with BWA mem and peaks called with Genrich (<https://github.com/jsh58/Genrich>). BedGraph files were produced with BEDTools.

GSD_1.0 gap closure. CanFam3.1 gaps, continuous "N" bases, and 1 kb of flanking sequences were extracted and mapped as pairs to GSD_1.0 (BWA mem). CanFam3.1 gaps were considered closed when, 1) flanking sequence pairs could be mapped properly in the same scaffold with mapping quality > 20; 2) the distance between pairs was less than 100 kb; and 3) no GSD_1.0 gap was present in the sequence between pairs. Closed gaps were intersected with annotations (BEDTools), and novel genes defined if it 1) had at least 80% of the gene body identified from the filled gaps; 2) was not a pseudogene; 3) had not been annotated in the unplaced scaffolds of CanFam3.1; 4) did not have the duplicated/homologous fragment in another region of the genome. Synteny was compared to hg38.

Assembly benchmark with Busco and Iso-Seq data. BUSCO²⁴ v3.0.2b was run with the mammalia_odb9 dataset. Mappability was assessed with Iso-Seq data using only PacBio CCS reads supported by > 10 subreads (483,702 reads). CCS reads were mapped with minimap2 v2.17, and the percentage of mapped bases per read calculated according to the "difference string" in cs tag. With these methods, GSD_1.0, CanFam3.1 and five newly released canine assemblies, Luka (Basenji), Nala (German Shepherd), Zoey (Great Dane), Scarlet (Golden Retriever), and Sandy (Dingo) were benchmarked (Supplementary Table 5).

10x and standard Illumina short read (ISR) mapping. HMW DNA was extracted from the blood of 27 additional dogs (19 breeds), and Chromium library preparation and sequencing completed as per "Genome sequencing". Sequencing depth ranged between 30-93X (Supplementary Table 7). Unplaced GSD_1.0 scaffolds were concatenated into a single scaffold with 500 "N" base spacers and 10x reads were mapped to each with the Long Ranger v2.2.2 WGS pipeline (10x Genomics). 10x breed-matched ISR data was downloaded for 25 individuals (Supplementary Table 11) and mapped to GSD_1.0 (BWA mem, default settings). SNPs and short indels were detected in 10x and ISR datasets (GATK4) and called (HaplotypeCaller), prior to merging (CombineGVCFs and GentoypesGVCFs) and filtering (SelectVariants).

Dark region detection. Both depth and mapping quality were calculated for each sample in each 10x or ISR dataset. For coverage, bamCoverage (Deeptools⁷⁰ v3.3.2) with a 25 bp window was used, with unmapped reads and secondary alignments excluded from the analysis. For the same windows, the proportion of reads with mapping quality > 10 was also assessed. Dark regions by coverage (COV), were defined as windows with coverage $\leq 5X$, with threshold adjusted for sequencing depth. A lower cutoff was applied in low coverage samples to select a maximum of 60 Mb (Supplementary Table 12). The individual COV dark regions were merged, and the COV fraction for each window was assessed for both ISR and 10x datasets: windows with $F_{COV} > 0.9$ (90% individuals, in at least 23 ISR dogs or 25 10x dogs) retained as the candidate COV dark regions. Regions were defined camouflaged (CAM) if the coverage

was $\geq 10X$ and the proportion of high mapping quality reads was less than 10%. We searched for and merged the genomic windows that reached the threshold from each dog. As the CAM regions detected in one individual could have been assigned as COV in others, we excluded those COV dark dogs before we calculated the fraction of CAM for each window. Any window with $F_{CAM} > 0.9$ was selected as a candidate.

Structural variation (SV) detection. Four SV callers were used to call SVs from 10x sequences. The first, Long Ranger, was used to call medium (50 bp to 30 kb) and larger-scale (> 30 kb) SVs. SVs were categorized as deletion (DEL), copy number variant (CNV) or inversion (INV). Callers GridSS⁷¹ and Manta⁷² also detected medium SVs (50 bp – 30 kb), and high-quality SVs from each, plus Long Ranger, were marked as “PASS” and kept for analysis. CNVnator⁷³ predicted CNVs by a read-depth approach in 150 bp window size. For each 10x sample, the filtered median SVs from all four callers were merged with SURVIVOR⁷⁴, and combined with the large size SVs called from Long Ranger. Chr X SVs that were only supported by CNVnator were pruned as the algorithm lacks the right model for chr X SV detection in females. SVs were further merged across individuals into a nonredundant SVs set.

SV validation and genotyping. Four DELs and four CNVs which overlapped protein-coding genes that were polymorphic within the 10x data set ($> 3/27$ individuals) were selected (Supplementary Table 8). SV breakpoints were confirmed with Sanger sequencing where possible. PCR was performed with either PrimeSTAR GXL DNA Polymerase (Takara) or AmpliTaq Gold DNA Polymerase (Applied Biosystems) according to manufacturer's recommendations. PCR fragments were cloned using either Zero Blunt or TOPO TA Cloning Kit (Invitrogen) depending on PCR overhang. Plasmid DNA was extracted using QIAprep Spin Miniprep Kit (QIAGEN), PCR products and plasmids sequenced using the Mix2Seq service (Eurofins Genomics) and analysed using CodonCode Aligner v6.0.2 (CodonCode). For *CYP1A2* CNV genotyping, ddPCR absolute quantification (BioRad) was performed and quantified as before⁷⁵. *CYP1A2* C1117T was genotyped according to published method⁷⁶. New Primers and probes were designed using Primer3 v0.4.0 (<http://bioinfo.ut.ee/primer3-0.4.0/>) and collated in Supplementary Table 8.

Gene expression. Total RNA was extracted from liver and spleen tissues using the AllPrep DNA/RNA/miRNA Universal Kit (QIAGEN) according to manufacturer's specification and including on-column *DNaseI* treatment (Supplementary Table 13). 1000 ng of total RNA was reverse transcribed using the Advantage RT-for-PCR Kit (Takara) and qPCR performed in quadruplet using SYBR Green PCR Master Mix (Thermo Fisher Scientific) and 900 nM primers in a QuantStudio 6 Real-Time system (Thermo Fisher Scientific) with standard cycling and dissociation curve analysis. Two housekeeper primer sets (RPS19 and RPS5) were assessed for stability (*Normfinder*⁷⁷ R package) and used in combination to calculate relative gene expression⁷⁸. These calculations included primer specific efficiencies and used the average Ct from all control samples for initial delta Ct normalisation. *wilcox.test* in R was used to assess the significance of between genotypic class gene expression changes.

Supplementary Information

Supplementary Figure 1. Genetic position of Mischka using the 170K SNP chip

Supplementary Figure 2. Sequence characteristics of 2,159 unplaced scaffolds in GSD_1.0.

Supplementary Figure 3. Tree map of contig sizes of CanFam3.1 and GSD1.0.

Supplementary Figure 4. Content of repetitive element on each chromosome of GSD_1.0.

Supplementary Figure 5. Comparison between length of the best BLAST hit for genes in GSD_1.0 and CanFam3.1 annotation.

Supplementary Figure 6. Mirlet7i identified from GSD_1.0.

Supplementary Figure 7. Closure of gaps in the cancer genes from COSMIC.

Supplementary Figure 8. Sequence comparison of immunity loci between GSD_1.0 and CanFam3.1.

Supplementary Figure 9. Iso-Seq data mapped to GSD_1.0 and CanFam3.1

Supplementary Figure 10. Mapping ISO-seq data to different canine assemblies.

Supplementary Figure 11. Illumina short reads (ISR) dark regions rescued by 10x sequencing

Supplementary Figure 12. RNA sequencing of different tissues.

Supplementary Table1. Summary of repetitive elements in GSD_1.0

Supplementary Table2. Tissues used in RNA experiments

Supplementary Table3. Datasets used for annotation

Supplementary Table4. Filled gaps within the cancer genes from COSMIC

Supplementary Table5. Summary of available canine assemblies from public resources

Supplementary Table6. ISO-seq reads mapping in different canine assemblies

Supplementary Table7. 10x sequencing of Mischka and 27 dogs

Supplementary Table8. Primers and probes used for validation and genotyping

Supplementary Table9. Gene expression summary for structural variant loci.

Supplementary Table10. Sequencing data generated from three HiC libraries

Supplementary Table11. Public resources of illumina short-read data of 25 dogs

Supplementary Table12. Summary of dark regions detected in each dog

Declarations

Ethics approval and consent to participate

Approval was obtained from dog owners before collecting the biological samples at veterinary clinics. Ethical approvals for sampling were granted by Uppsala Animal Ethical Committee and Swedish Board of Agriculture (C139/9, C2/12, C12/15). Importation of canine tissues was approved by Jordbruksverket (6.7.18-14513/17).

Competing interests

The authors declare no competing financial or non-financial interests.

Data availability

The PacBio-long reads, HiC, and Illumina 10x data of Mischka are available in SRA under BioProject PRJNA587469. The Illumina 10x data of 27 dogs are available in SRA under BioProject PRJNA588624. Scripts used in the study are available at the GitHub repository (<https://github.com/Chao912/Mischka/>). The canFam_GSD_1.0 assembly is deposited in DDBJ/ENA/GenBank under JAAHUQ000000000.

Author contributions

KLT, JRSM and MLA conceived the study and designed the experiments. GP and MLA collected the samples with the help of JH1, ÅO, SS, HR, IL, SM, JH2 and ÅH. MLA, ÅK and ÅO performed the DNA/RNA extractions. ÅK, ES and JRSM performed the validation of structural variation, genotyping and expression analyses. CW, OW, MLA and KLT contributed to the data analysis of the genome assembly. OW performed the gene annotation with the help of TB and SM. JRSM and KLT oversaw and interpreted the results together with CW, OW and MLA and ES. CW, OW, JRSM and KLT wrote the manuscript with input from all authors.

Acknowledgements

We thank Mischka's owners who kindly allowed us to collect blood and tissues for scientific purposes, Susanne Gustafsson from the SLU Canine Biobank for the management of these and other canine samples used throughout the project and Anna Darlene van der Heiden for generating retina data. We would like to acknowledge Mats Pettersson, Olga Vinnere Pettersson and Ignas Bunikis for helpful suggestions. Next generation sequencing was made possible with assistance from the Uppsala Genome Center (PacBio) and the SNP&SEQ Technology Platform (10x Chromium). Both fall under the umbrella of National Genomics Infrastructure (NGI) Sweden and Science for Life Laboratory, Sweden and themselves are supported by RFI/VR and the Swedish Research Council and the Knut and Alice Wallenberg Foundation respectively. We thank SNIC through Uppsala Multidisciplinary Center for Advanced

Computational Science (UPPMAX) for providing computation resources under Projects SNIC 2017/7-384, 2017/7-385 and 2020/5-190.

KLT is a Distinguished Professor at the Swedish Research Council. Research reported in this publication was supported by the National Cancer Institute of the National Institutes of Health under Award Number R01CA225755, The Knut and Alice Wallenberg Foundation and Agria och Svenska Kennelklubben Forskningsfond" (<https://www.sk.se/sv/Agria-SKK-Forskningsfond/>, grant numbers: P2012-0015, N2013-0020, P2014-0018, P2015-0012).

References

1. Axelsson, E. *et al.* The genomic signature of dog domestication reveals adaptation to a starch-rich diet. *Nature* **495**, 360–364 (2013).
2. Freedman, A. H. *et al.* Genome sequencing highlights the dynamic early history of dogs. *PLoS Genet.* **10**, e1004016 (2014).
3. Plassais, J. *et al.* Whole genome sequencing of canids reveals genomic regions under selection and variants influencing morphology. *Nat. Commun.* **10**, 1489 (2019).
4. Friedrich, J. *et al.* Genetic dissection of complex behaviour traits in German Shepherd dogs. *Heredity* **123**, 746–758 (2019).
5. Awano, T. *et al.* Genome-wide association analysis reveals a SOD1 mutation in canine degenerative myelopathy that resembles amyotrophic lateral sclerosis. *Proc. Natl. Acad. Sci. U. S. A.* **106**, 2794–2799 (2009).
6. Bianchi, M. *et al.* Whole-genome genotyping and resequencing reveal the association of a deletion in the complex interferon alpha gene cluster with hypothyroidism in dogs. *BMC Genomics* **21**, 307 (2020).
7. Friedenber, S. G. & Meurs, K. M. Genotype imputation in the domestic dog. *Mamm. Genome Off. J. Int. Mamm. Genome Soc.* **27**, 485–494 (2016).
8. Oliver, J. A. C., Ricketts, S. L., Kuehn, M. H. & Mellersh, C. S. Primary closed angle glaucoma in the Basset Hound: Genetic investigations using genome-wide association and RNA sequencing strategies. *Mol. Vis.* **25**, 93–105 (2019).
9. Lindblad-Toh, K. *et al.* Genome sequence, comparative analysis and haplotype structure of the domestic dog. *Nature* **438**, 803–819 (2005).
10. Hoepfner, M. P. *et al.* An Improved Canine Genome and a Comprehensive Catalogue of Coding Genes and Non-Coding Transcripts. *PLoS ONE* **9**, (2014).
11. Axelsson, E. *et al.* Death of PRDM9 coincides with stabilization of the recombination landscape in the dog genome. *Genome Res.* **22**, 51–63 (2012).
12. Datlinger, P. *et al.* Pooled CRISPR screening with single-cell transcriptome readout. *Nat. Methods* **14**, 297–301 (2017).

13. Boyle, E. A., Li, Y. I. & Pritchard, J. K. An Expanded View of Complex Traits: From Polygenic to Omnigenic. *Cell* **169**, 1177–1186 (2017).
14. Chin, C.-S. *et al.* Phased diploid genome assembly with single-molecule real-time sequencing. *Nat. Methods* **13**, 1050–1054 (2016).
15. Penso-Dolfin, L. *et al.* An Improved microRNA Annotation of the Canine Genome. *PLoS One* **11**, e0153453 (2016).
16. Kimura, K. *et al.* Circulating exosomes suppress the induction of regulatory T cells via let-7i in multiple sclerosis. *Nat. Commun.* **9**, 17 (2018).
17. Shi, Y. *et al.* Down-regulation of the let-7i facilitates gastric cancer invasion and metastasis by targeting COL1A1. *Protein Cell* **10**, 143–148 (2019).
18. de Anda-Jáuregui, G., Espinal-Enríquez, J., Drago-García, D. & Hernández-Lemus, E. Nonredundant, Highly Connected MicroRNAs Control Functionality in Breast Cancer Networks. *Int. J. Genomics* 2018, 9585383 (2018).
19. Bläsius, F. M. *et al.* Loss of cadherin related family member 5 (CDHR5) expression in clear cell renal cell carcinoma is a prognostic marker of disease progression. *Oncotarget* **8**, 75076–75086 (2017).
20. Wong, C. C. *et al.* SLC25A22 Promotes Proliferation and Survival of Colorectal Cancer Cells With KRAS Mutations and Xenograft Tumor Progression in Mice via Intracellular Synthesis of Aspartate. *Gastroenterology* **151**, 945–960.e6 (2016).
21. Sondka, Z. *et al.* The COSMIC Cancer Gene Census: describing genetic dysfunction across all human cancers. *Nat. Rev. Cancer* **18**, 696–705 (2018).
22. Zhong, Z. *et al.* HOXD13 methylation status is a prognostic indicator in breast cancer. *Int. J. Clin. Exp. Pathol.* **8**, 10716–10724 (2015).
23. Zou, H., Chen, H., Zhou, Z., Wan, Y. & Liu, Z. ATXN3 promotes breast cancer metastasis by deubiquitinating KLF4. *Cancer Lett.* **467**, 19–28 (2019).
24. Seppey, M., Manni, M. & Zdobnov, E. M. BUSCO: Assessing Genome Assembly and Annotation Completeness. *Methods Mol. Biol. Clifton NJ* **1962**, 227–245 (2019).
25. Berglund, J. *et al.* Novel origins of copy number variation in the dog genome. *Genome Biol.* **13**, R73 (2012).
26. Molin, A.-M., Berglund, J., Webster, M. T. & Lindblad-Toh, K. Genome-wide copy number variant discovery in dogs using the CanineHD genotyping array. *BMC Genomics* **15**, 210 (2014).
27. Nicholas, T. J. *et al.* The genomic architecture of segmental duplications and associated copy number variants in dogs. *Genome Res.* **19**, 491–499 (2009).
28. Ebbert, M. T. W. *et al.* Systematic analysis of dark and camouflaged genes reveals disease-relevant genes hiding in plain sight. *Genome Biol.* **20**, 97 (2019).
29. Schmutz, S. M., Berryere, T. G. & Goldfinch, A. D. TYRP1 and MC1R genotypes and their effects on coat color in dogs. *Mamm. Genome* **13**, 380–387 (2002).

30. Gilot, D. *et al.* A non-coding function of TYRP1 mRNA promotes melanoma growth. *Nat. Cell Biol.* **19**, 1348–1357 (2017).
31. Goldstein, A. M. *et al.* Rare germline variants in known melanoma susceptibility genes in familial melanoma. *Hum. Mol. Genet.* **26**, 4886–4895 (2017).
32. Mühleisen, T. W. *et al.* Genome-wide association study reveals two new risk loci for bipolar disorder. *Nat. Commun.* **5**, 3339 (2014).
33. Silver, M. *et al.* Identification of gene pathways implicated in Alzheimer's disease using longitudinal imaging phenotypes with sparse regression. *NeuroImage* **63**, 1681–1694 (2012).
34. Jajodia, A. *et al.* Evaluation of genetic association of neurodevelopment and neuroimmunological genes with antipsychotic treatment response in schizophrenia in Indian populations. *Mol. Genet. Genomic Med.* **4**, 18–27 (2016).
35. Meyers-Wallen, V. N. *et al.* XX Disorder of Sex Development is associated with an insertion on chromosome 9 and downregulation of RSP01 in dogs (*Canis lupus familiaris*). *PloS One* **12**, e0186331 (2017).
36. Nowacka-Woszuik, J. *et al.* Deep sequencing of a candidate region harboring the SOX9 gene for the canine XX disorder of sex development. *Anim. Genet.* **48**, 330–337 (2017).
37. Nowacka-Woszuik, J. *et al.* Association between polymorphisms in the SOX9 region and canine disorder of sex development (78,XX; SRY-negative) revisited in a multibreed case-control study. *PloS One* **14**, e0218565 (2019).
38. Tenmizu, D., Endo, Y., Noguchi, K. & Kamimura, H. Identification of the novel canine CYP1A2 1117 C > T SNP causing protein deletion. *Xenobiotica Fate Foreign Compd. Biol. Syst.* **34**, 835–846 (2004).
39. Mise, M. *et al.* Polymorphic expression of CYP1A2 leading to interindividual variability in metabolism of a novel benzodiazepine receptor partial inverse agonist in dogs. *Drug Metab. Dispos. Biol. Fate Chem.* **32**, 240–245 (2004).
40. Court, M. H. Canine cytochrome P-450 pharmacogenetics. *Vet. Clin. North Am. Small Anim. Pract.* **43**, 1027–1038 (2013).
41. Mise, M., Hashizume, T. & Komuro, S. Characterization of substrate specificity of dog CYP1A2 using CYP1A2-deficient and wild-type dog liver microsomes. *Drug Metab. Dispos. Biol. Fate Chem.* **36**, 1903–1908 (2008).
42. Graham, R. A. *et al.* In vivo and in vitro induction of cytochrome P450 enzymes in beagle dogs. *Drug Metab. Dispos. Biol. Fate Chem.* **30**, 1206–1213 (2002).
43. Aneur, A. *et al.* De Novo Assembly of Two Swedish Genomes Reveals Missing Segments from the Human GRCh38 Reference and Improves Variant Calling of Population-Scale Sequencing Data. *Genes* **9**, (2018).
44. Lilue, J. *et al.* Sixteen diverse laboratory mouse reference genomes define strain-specific haplotypes and novel functional loci. *Nat. Genet.* **50**, 1574–1583 (2018).
45. Gordon, D. *et al.* Long-read sequence assembly of the gorilla genome. *Science* **352**, aae0344 (2016).

46. Tengvall, K. *et al.* Genome-Wide Analysis in German Shepherd Dogs Reveals Association of a Locus on CFA 27 with Atopic Dermatitis. *PLoS Genet.* **9**, e1003475 (2013).
47. Chin, C.-S. *et al.* Nonhybrid, finished microbial genome assemblies from long-read SMRT sequencing data. *Nat. Methods* **10**, 563–569 (2013).
48. Yeo, S., Coombe, L., Warren, R. L., Chu, J. & Birol, I. ARCS: scaffolding genome drafts with linked reads. *Bioinforma. Oxf. Engl.* **34**, 725–731 (2018).
49. Warren, R. L. *et al.* LINKS: Scalable, alignment-free scaffolding of draft genomes with long reads. *GigaScience* **4**, 35 (2015).
50. Wong, A. K. *et al.* A comprehensive linkage map of the dog genome. *Genetics* **184**, 595–605 (2010).
51. Kent, W. J. BLAT—the BLAST-like alignment tool. *Genome Res.* **12**, 656–664 (2002).
52. Catchen, J., Amores, A. & Bassham, S. Chromonomer: a tool set for repairing and enhancing assembled genomes through integration of genetic maps and conserved synteny. *bioRxiv* 2020.02.04.934711 (2020) doi:10.1101/2020.02.04.934711.
53. English, A. C. *et al.* Mind the Gap: Upgrading Genomes with Pacific Biosciences RS Long-Read Sequencing Technology. *PLOS ONE* **7**, e47768 (2012).
54. Durand, N. C. *et al.* Juicer Provides a One-Click System for Analyzing Loop-Resolution Hi-C Experiments. *Cell Syst.* **3**, 95–98 (2016).
55. Walker, B. J. *et al.* Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PloS One* **9**, e112963 (2014).
56. Li, H. & Durbin, R. Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinforma. Oxf. Engl.* **25**, 1754–1760 (2009).
57. Kronenberg, Z. N. *et al.* High-resolution comparative analysis of great ape genomes. *Science* **360**, (2018).
58. Garrison, E. & Marth, G. Haplotype-based variant detection from short-read sequencing. *ArXiv12073907 Q-Bio* (2012).
59. Van der Auwera, G. A. *et al.* From FastQ data to high confidence variant calls: the Genome Analysis Toolkit best practices pipeline. *Curr. Protoc. Bioinforma.* **43**, 11.10.1–11.10.33 (2013).
60. Wood, D. E., Lu, J. & Langmead, B. Improved metagenomic analysis with Kraken 2. *Genome Biol.* **20**, 257 (2019).
61. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features. *Bioinforma. Oxf. Engl.* **26**, 841–842 (2010).
62. Gardiner-Garden, M. & Frommer, M. CpG islands in vertebrate genomes. *J. Mol. Biol.* **196**, 261–282 (1987).
63. Derrien, T. *et al.* Fast Computation and Applications of Genome Mappability. *PLOS ONE* **7**, e30377 (2012).
64. Rice, P., Longden, I. & Bleasby, A. EMBOSS: the European Molecular Biology Open Software Suite. *Trends Genet. TIG* **16**, 276–277 (2000).

65. Hayden, K. E. & Willard, H. F. Composition and organization of active centromere sequences in complex genomes. *BMC Genomics* **13**, 324 (2012).
66. Kovaka, S. *et al.* Transcriptome assembly from long-read RNA-seq alignments with StringTie2. *Genome Biol.* **20**, 278 (2019).
67. Kuo, R. I., Cheng, Y., Smith, J., Archibald, A. L. & Burt, D. W. Illuminating the dark side of the human transcriptome with TAMA Iso-Seq analysis. *bioRxiv* 780015 (2019) doi:10.1101/780015.
68. Friedländer, M. R., Mackowiak, S. D., Li, N., Chen, W. & Rajewsky, N. miRDeep2 accurately identifies known and hundreds of novel microRNA genes in seven animal clades. *Nucleic Acids Res.* **40**, 37–52 (2012).
69. Megquier, K. *et al.* BarkBase: Epigenomic Annotation of Canine Genomes. *Genes* **10**, (2019).
70. Ramírez, F. *et al.* deepTools2: a next generation web server for deep-sequencing data analysis. *Nucleic Acids Res.* **44**, W160-165 (2016).
71. Cameron, D. L. *et al.* GRIDSS: sensitive and specific genomic rearrangement detection using positional de Bruijn graph assembly. *Genome Res.* **27**, 2050–2060 (2017).
72. Chen, X. *et al.* Manta: rapid detection of structural variants and indels for germline and cancer sequencing applications. *Bioinforma. Oxf. Engl.* **32**, 1220–1222 (2016).
73. Abyzov, A., Urban, A. E., Snyder, M. & Gerstein, M. CNVnator: an approach to discover, genotype, and characterize typical and atypical CNVs from family and population genome sequencing. *Genome Res.* **21**, 974–984 (2011).
74. Jeffares, D. C. *et al.* Transient structural variations have strong effects on quantitative traits and reproductive isolation in fission yeast. *Nat. Commun.* **8**, 14061 (2017).
75. Olsson, M. *et al.* Absolute quantification reveals the stable transmission of a high copy number variant linked to autoinflammatory disease. *BMC Genomics* **17**, 299 (2016).
76. Mise, M., Hashizume, T., Matsumoto, S., Terauchi, Y. & Fujii, T. Identification of non-functional allelic variant of CYP1A2 in dogs. *Pharmacogenetics* **14**, 769–773 (2004).
77. Andersen, C. L., Jensen, J. L. & Ørntoft, T. F. Normalization of real-time quantitative reverse transcription-PCR data: a model-based variance estimation approach to identify genes suited for normalization, applied to bladder and colon cancer data sets. *Cancer Res.* **64**, 5245–5250 (2004).
78. Vandesompele, J. *et al.* Accurate normalization of real-time quantitative RT-PCR data by geometric averaging of multiple internal control genes. *Genome Biol.* **3**, RESEARCH0034 (2002).

Figures

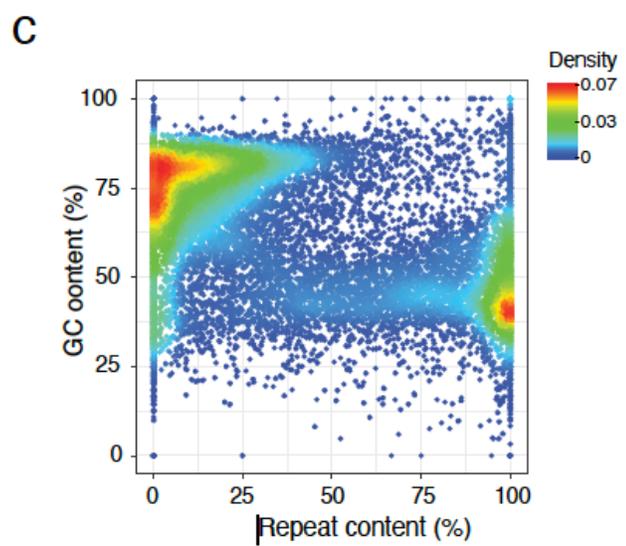
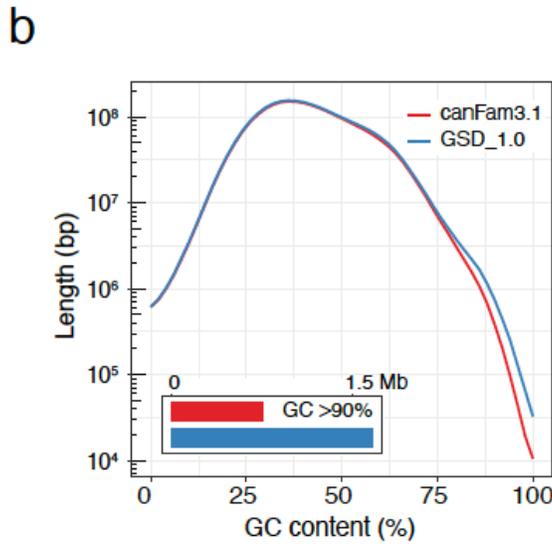
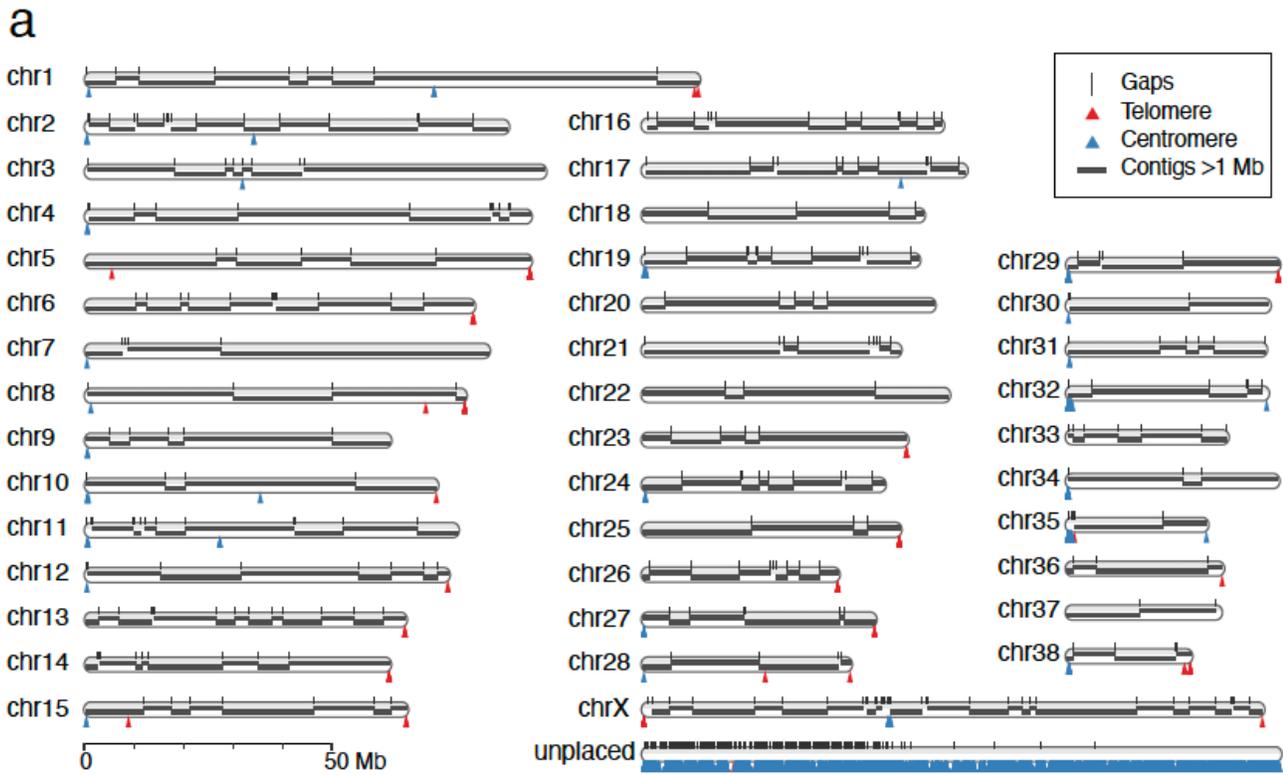


Figure 1

Features of the novel canine assembly. a) GSD_1.0 ideogram showing chromosomes, contigs, gaps, centromere and telomere repeats. All unplaced sequences were concatenated into a single scaffold (segmental duplications, 58.1% and centromeric repeats, 30.1%). b) Comparison of GC content (50 bp-window) between GSD_1.0 and CanFam3.1. c) Sequence characteristics of closed gaps in GSD_1.0. These are predominately high in GC or repeat content.

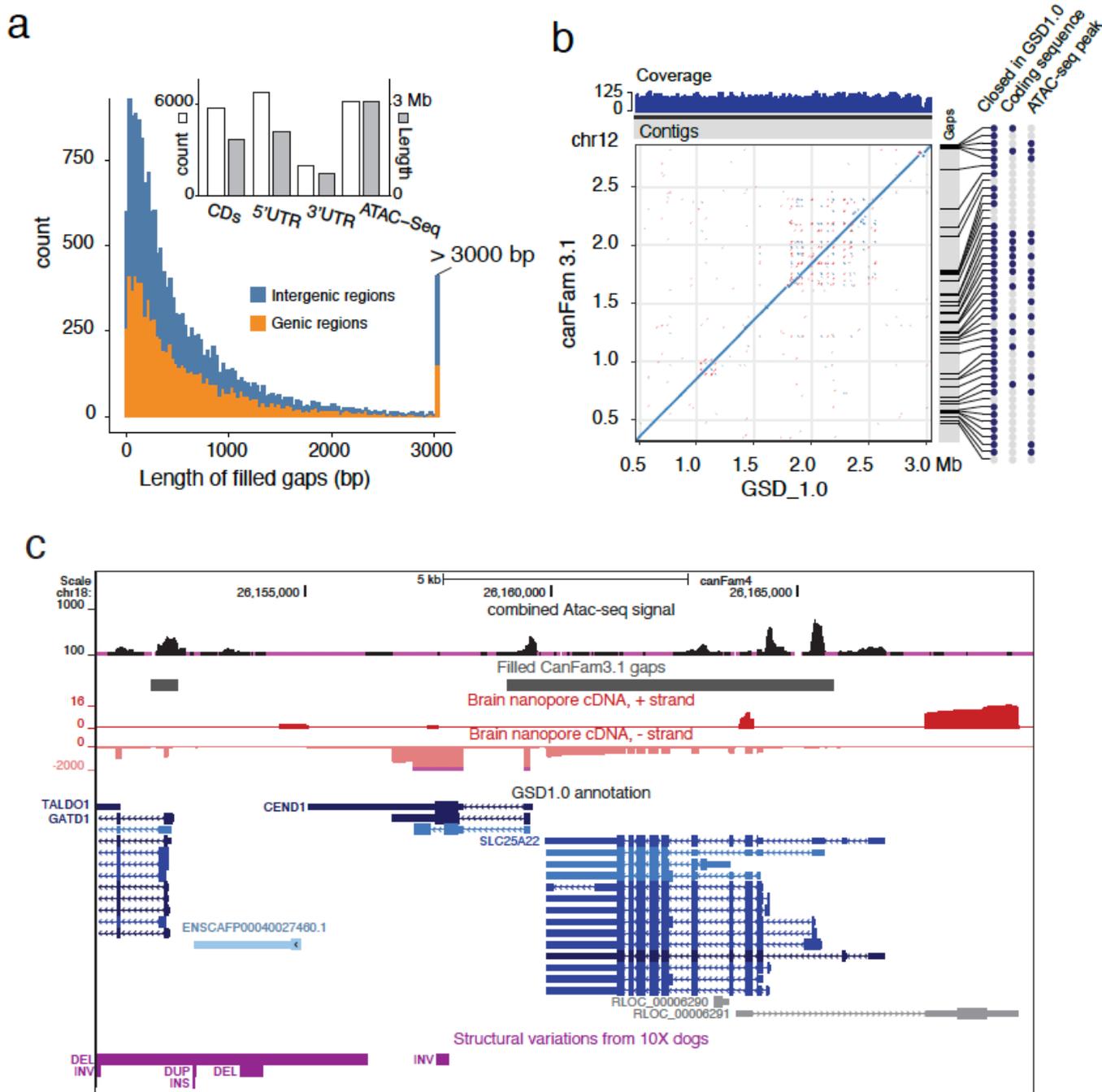


Figure 2

Gap closure and functionality. a) Size distribution and overlap with exons and promoters for filled gaps. b) Sequence comparison of DLA on chr 12 between CanFam3.1 and GSD_1.0. c) Representative GSD_1.0 annotation from the UCSC track hub highlighting available data and an example of a gene hidden in CanFam3.1.

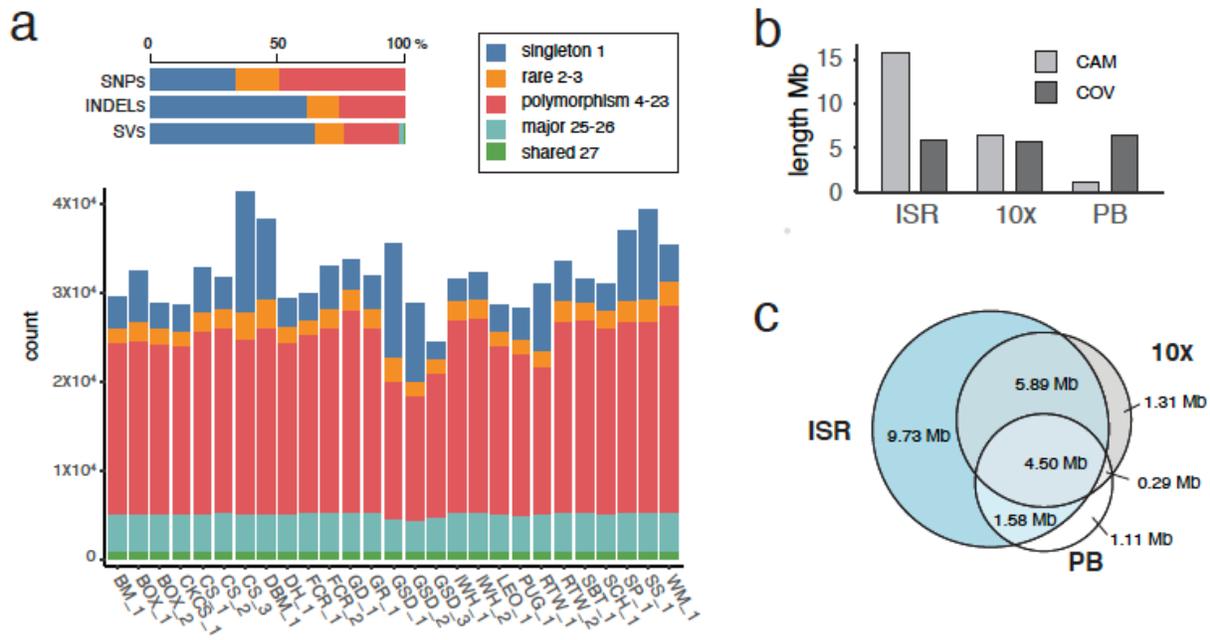


Figure 3

Genome variation and dark regions. a) SNPs, indels and structural variations shared among Mischka and the 27 10x sequenced dogs. b) The total length of dark regions detected from Illumina short reads (ISR), 10x and PacBio sequencing. c) Intersection of dark regions from different datasets.

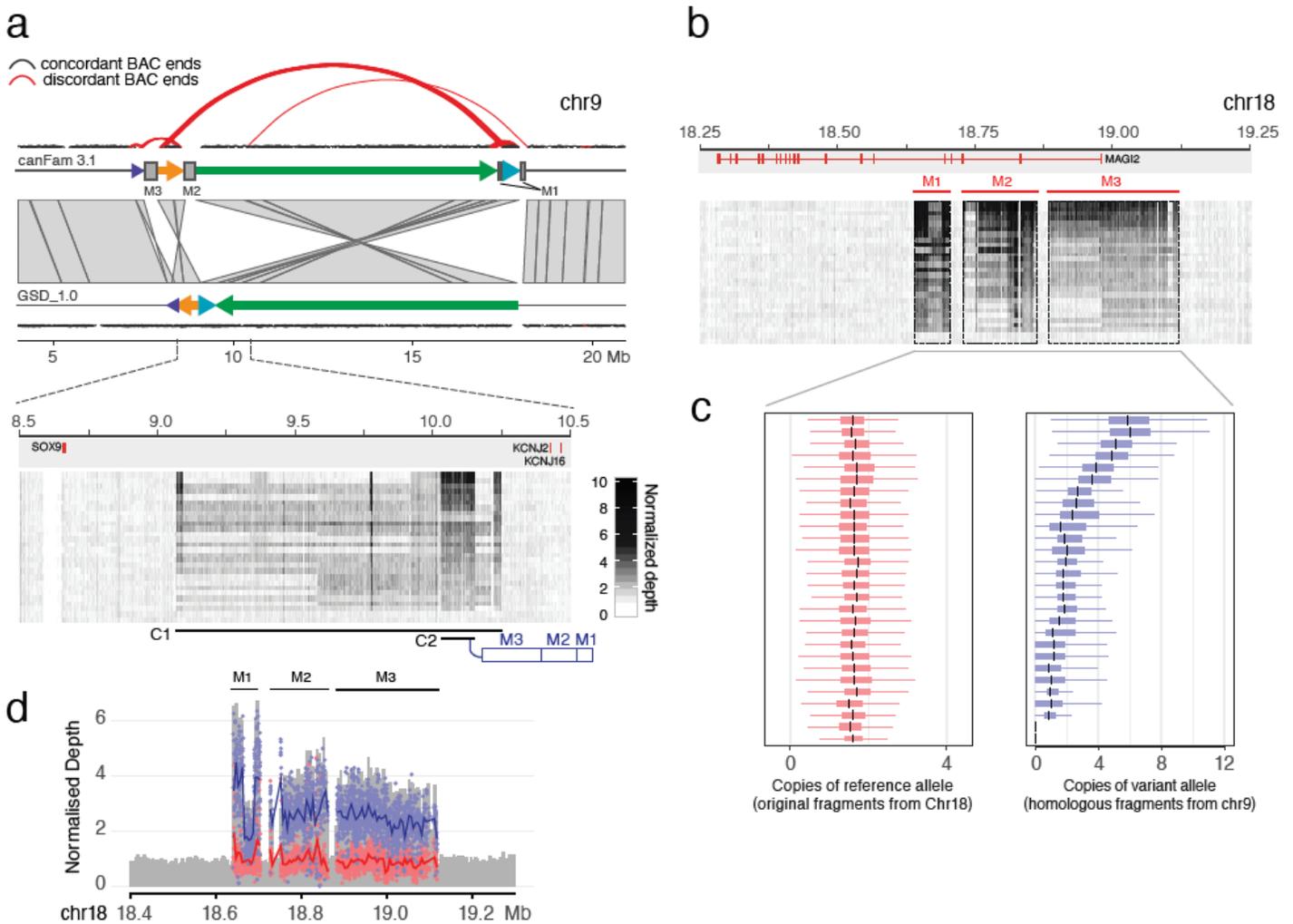


Figure 4

Correction of an inverted region in chromosome 9. a) Four fragments from the region were rearranged in GSD_1.0. The corrected order was confirmed using CanFam3.1 BAC clone (CH82) the end sequences of 49 discordant end pairs (red curves; > 500 kb or not in a forward-reverse direction) were found at the edge of rearranged fragments in CanFam3.1, whereas these were properly mapped in GSD_1.0. From this region, three homologous chr 18 fragments spanning MAGI2 (M1, M2 and M3) were present on chr 9 of CanFam3.1, but missing in the GSD_1.0. We proposed that those homologous fragments should be located together with C2 (chr9:10.03-10.16 Mb) in a large duplicated C1 (chr9:9.07-10.25 Mb) region. b) Reads from both original and homologous M1, M2 and M3 fragments were mapped to chr 18 of GSD_1.0. c) Mischka and all 10x dogs have only two original chr 18 copies M1, M2 and M3, but carry between 0-6 copies of the chr 9 homologous fragments. d) The example plot of normalised depth illustrates how the copy number of the reference alleles and variant alleles were measured to distinguish the original (red) and homologous (blue) of M1, M2 and M3.

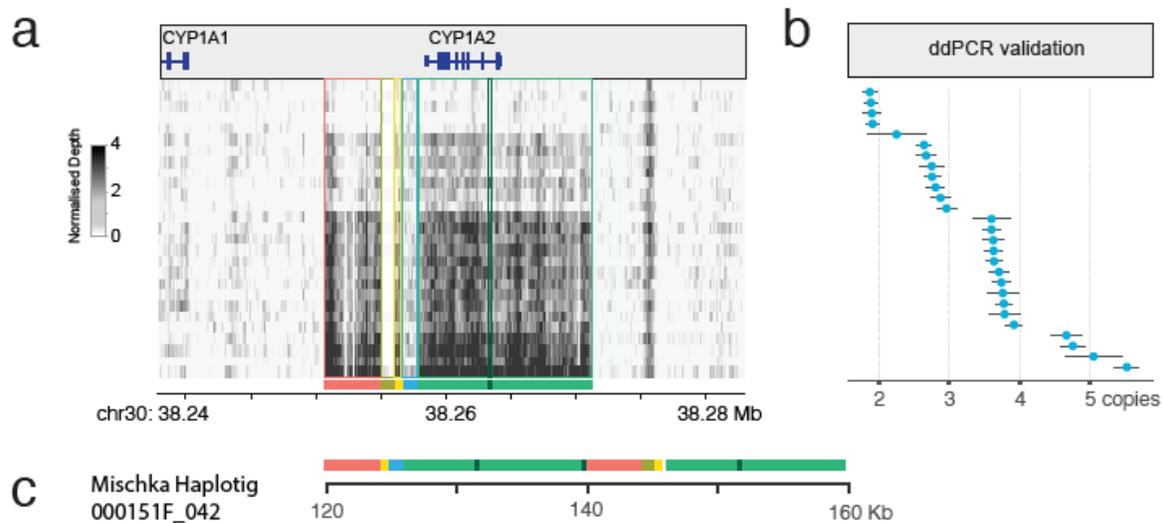


Figure 5

Copy number expansion encompassing CYP1A2. a) A duplication identified on chr 30 consists of six segments and contains the CYP1A2 gene. b) The duplication was validated in the 10x sequenced individuals using ddPCR. c) The individual pieces from the reference are plotted as they appear in the alternative haplotig sequence (000151F_042) for Mischka (CNV = 3). Sequence was extracted from the FALCON assembly. Full length CYP1A2 sits within copies of the green fragment.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [Supplementaryinformation.pdf](#)
- [rs.pdf](#)