

Supplementary Information for Comparative Analysis of Mammalian Y Chromosomes Illuminates Ancestral Structure and Lineage-Specific Evolution

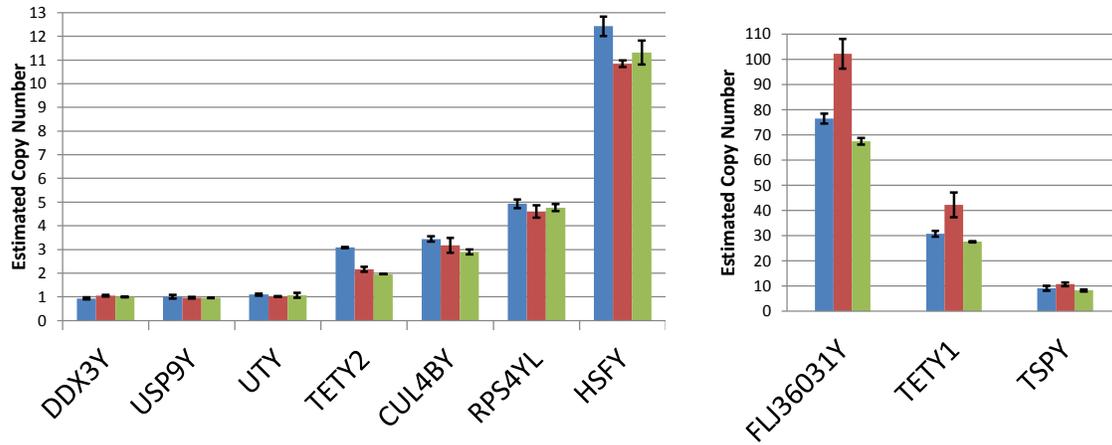
Gang Li, Brian W. Davis, Terje Raudsepp, Alison J. Pearks Wilkerson, Victor C. Mason, Malcolm Ferguson-Smith, Patricia C. O'Brien, Paul Waters and William J. Murphy.

Figures 1-12	Pgs. 2-16
Tables 1-3 *	Pgs. 17-19
Supplemental Methods	Pgs. 20-24

*Supplementary Table 4 is attached as a separate Excel file

Figure 1. Copy number analysis of felid Y chromosome genes using qPCR and FISH

A. Quantitative PCR analysis in three domestic cats



B. FISH mapping of domestic cat multicopy cDNA probes on snow leopard (*Panthera uncia*) Y chromosomes



Figure 2. Summary of copy number and expression information for domestic cat and dog Y chromosome genes.

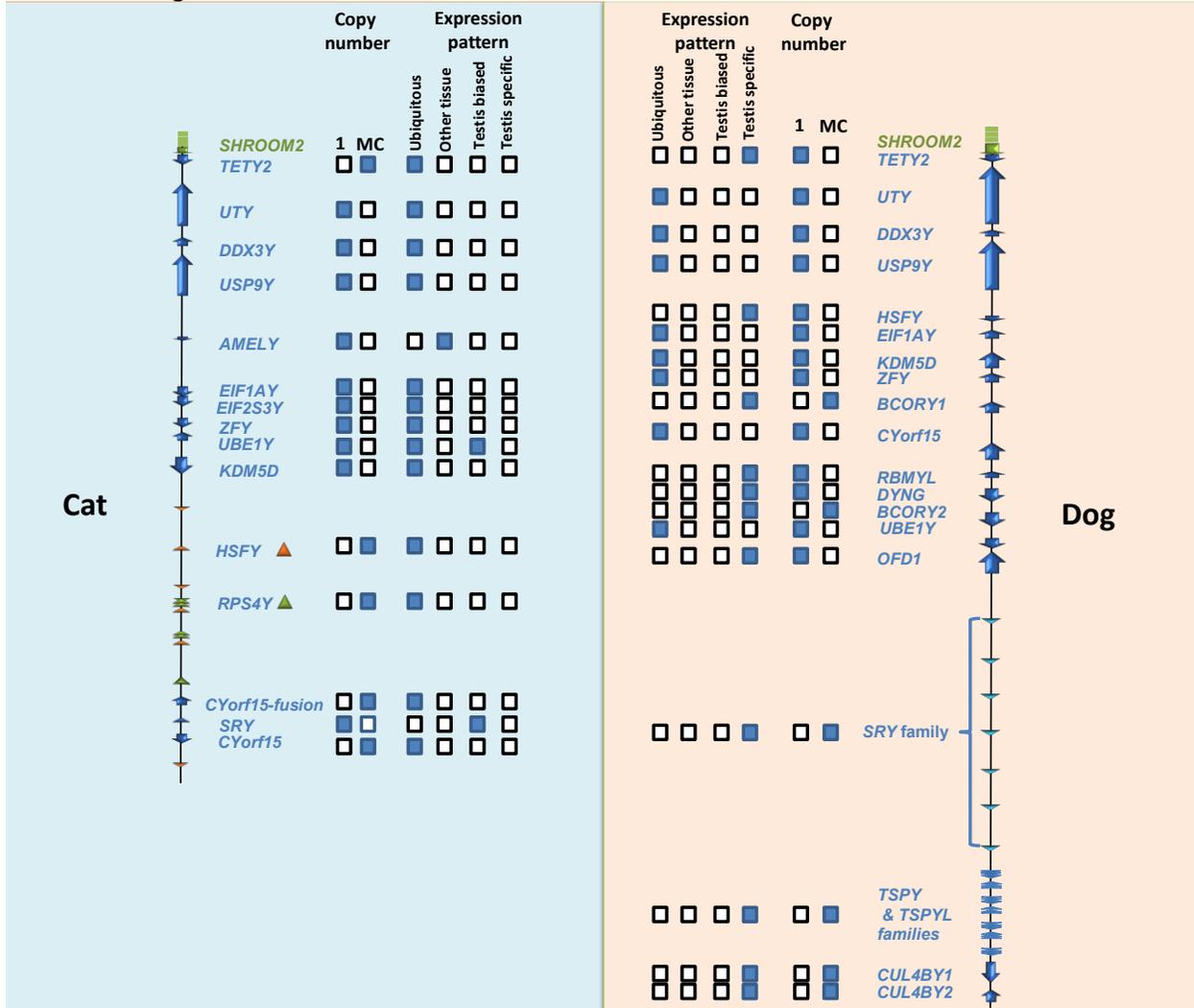
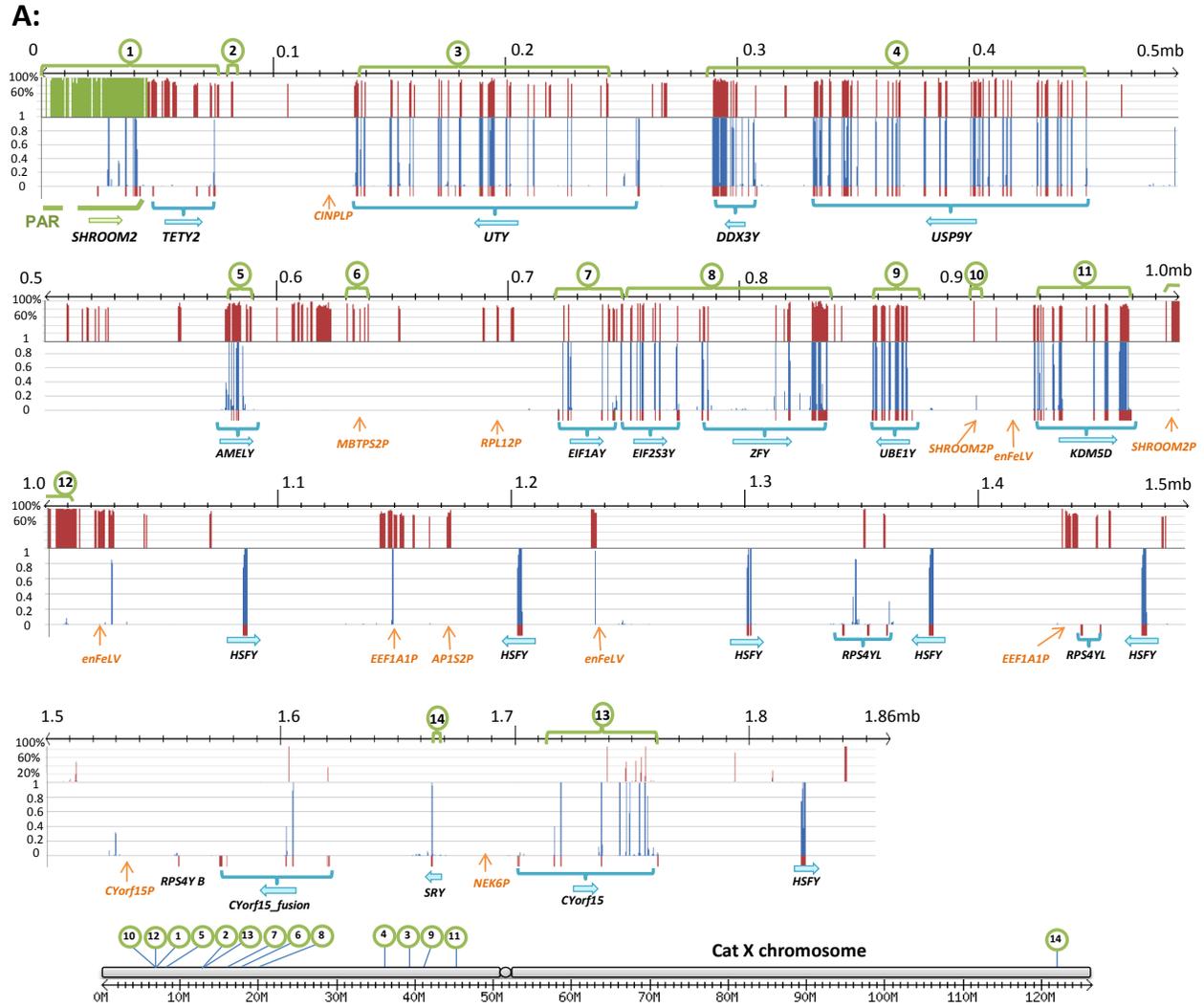


Figure 3: Comparison of sequence conservation among 5 Y mammalian chromosomes (human, chimpanzee, rhesus, dog, and cat) using the sequenced cat (A) and dog (B, next page) Y chromosome sequence as a reference. Multispecies sequence conservation scores are shown by blue columns (bottom panel), and pairwise sequence identity between X and Y chromosome sequence is shown by red columns (top panel). Green circles with numbers indicate X-Y orthologous regions, the order of which is shown on a schematic of the X chromosome at the bottom. Gene and transcript information is listed below each panel. Functional genes are shown in black, while pseudogenes are shown in orange. Red bars at the bottom of each panel indicate known gene exons.



B:

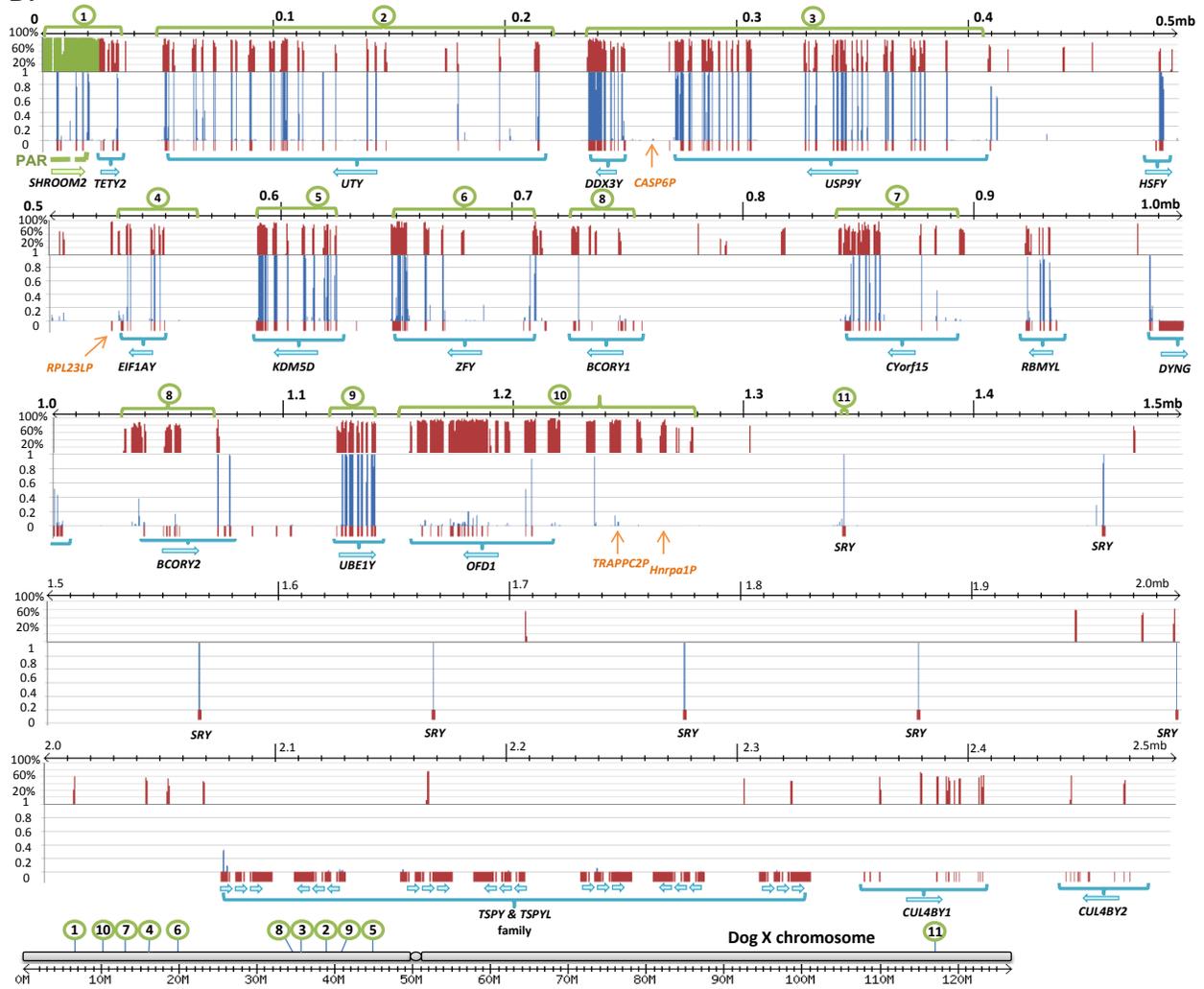


Figure 4. Triangular self dot plots of self-self DNA sequence identities within the dog and cat MSY sequences. The horizontal panel below each figure displays the average % self-self identity, based on a sliding window analysis (window size=50 nucleotides, distance between each window step=5 nucleotides). Sequences with 100 percent identity present a dot in the triangular dot plot. Horizontal and vertical lines represent direct and inverted repeats. The MSY schematic with scale is shown below the triangular plot with colored background representing the different gene regions. Gene bodies are represented within the colored blocks as directional arrows, drawn to scale.

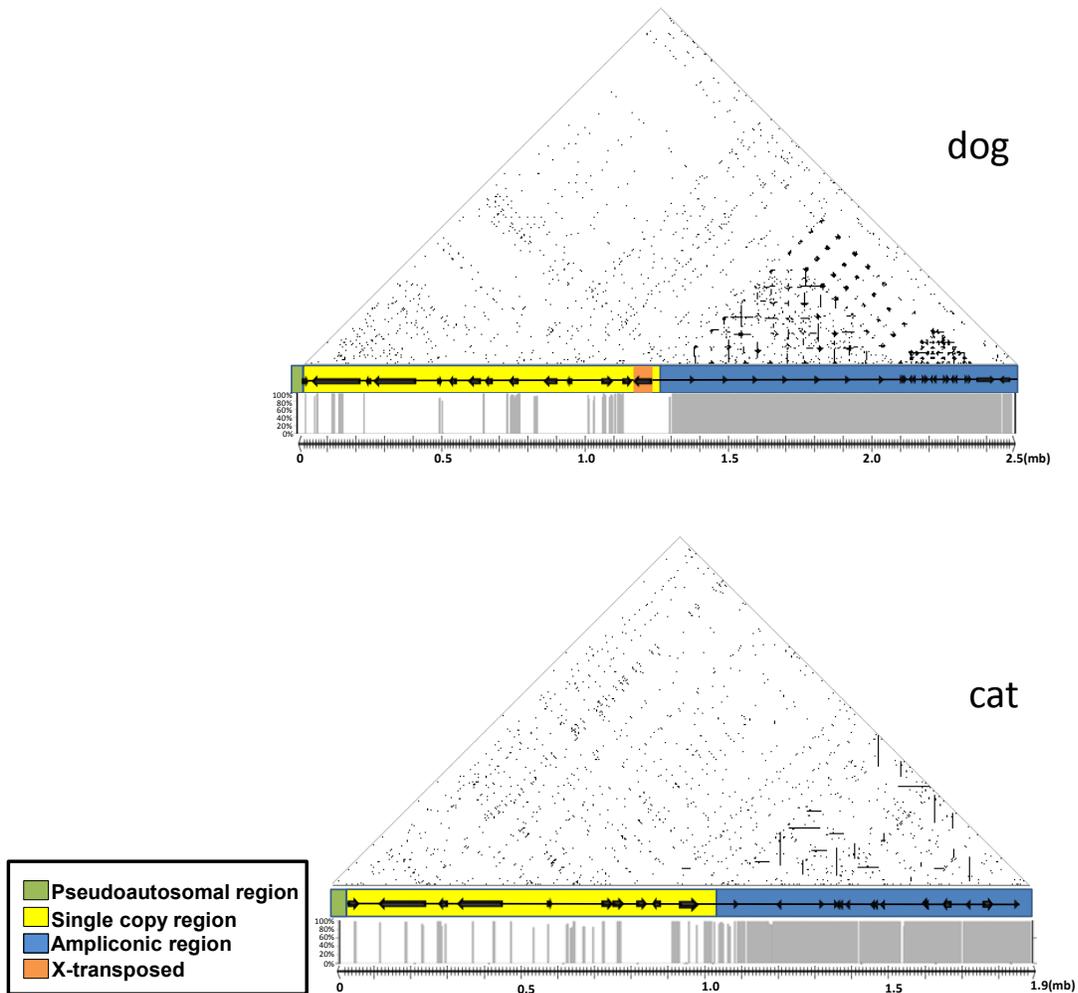


Figure 5. A) Gene structure and expression profile (RNASeq) of the novel coding gene, *DYNG*, found on the domestic dog Y chromosome. **B)** Self-self blast result of the transcript identifies the internal repeats within the open reading frame.

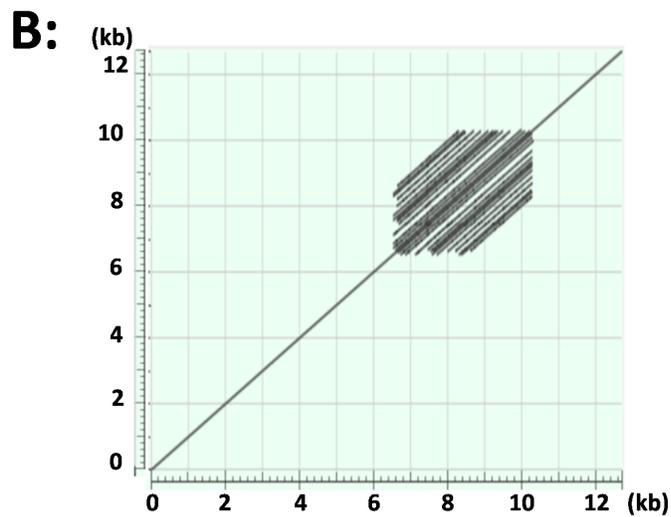
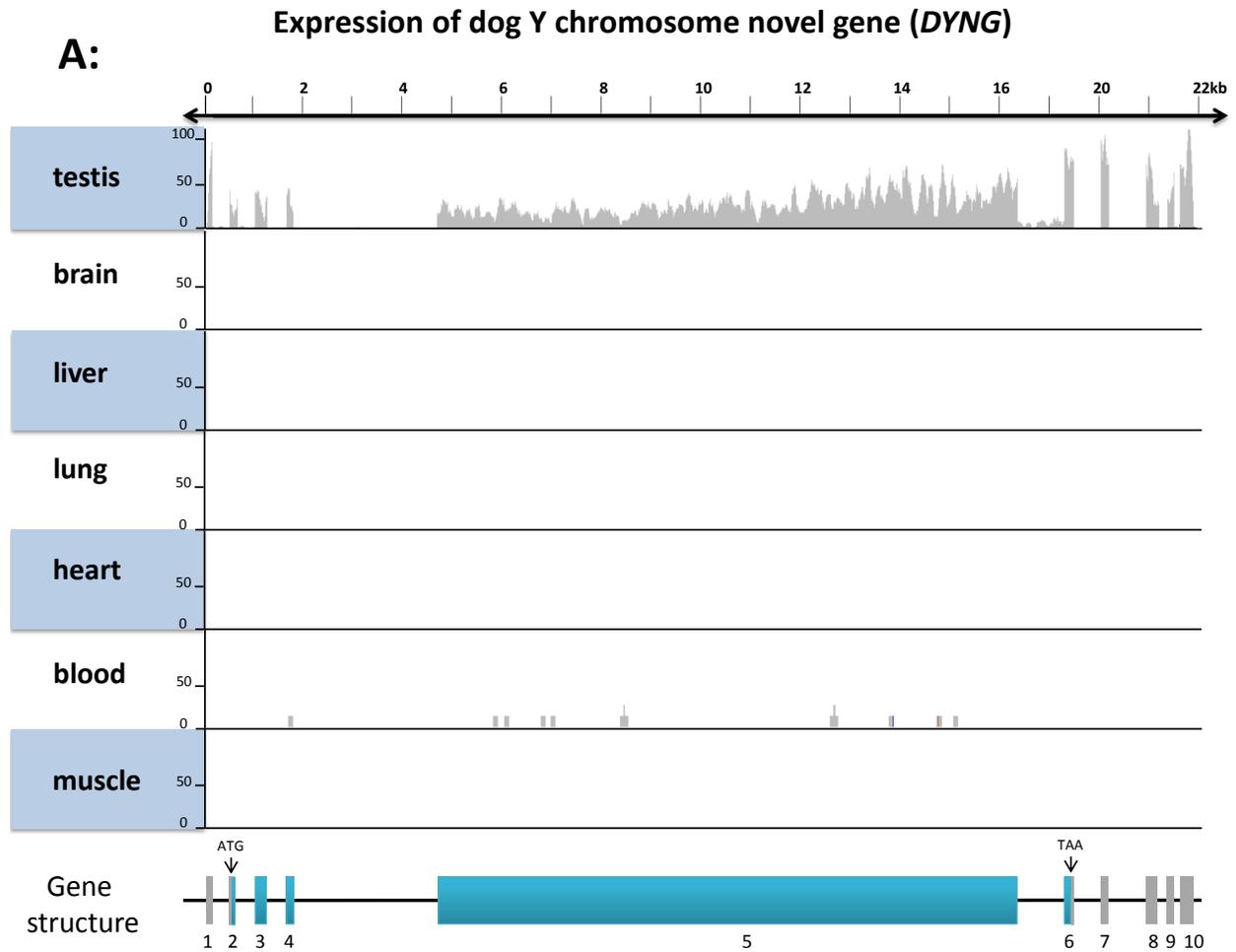


Figure 6. Evolution of the *CYorf15* gene family in select mammals. **A)** Structural comparison of human *CYorf15* gene family, with **B)** dog *CYorf15* gene, cat *CYorf15* fusion gene and the X chromosome gene Gamma taxilin, (*TXLNG*). Dashed arrows and boxes indicate regions of sequence homology between the different gene family members, based on pairwise sequence comparisons.

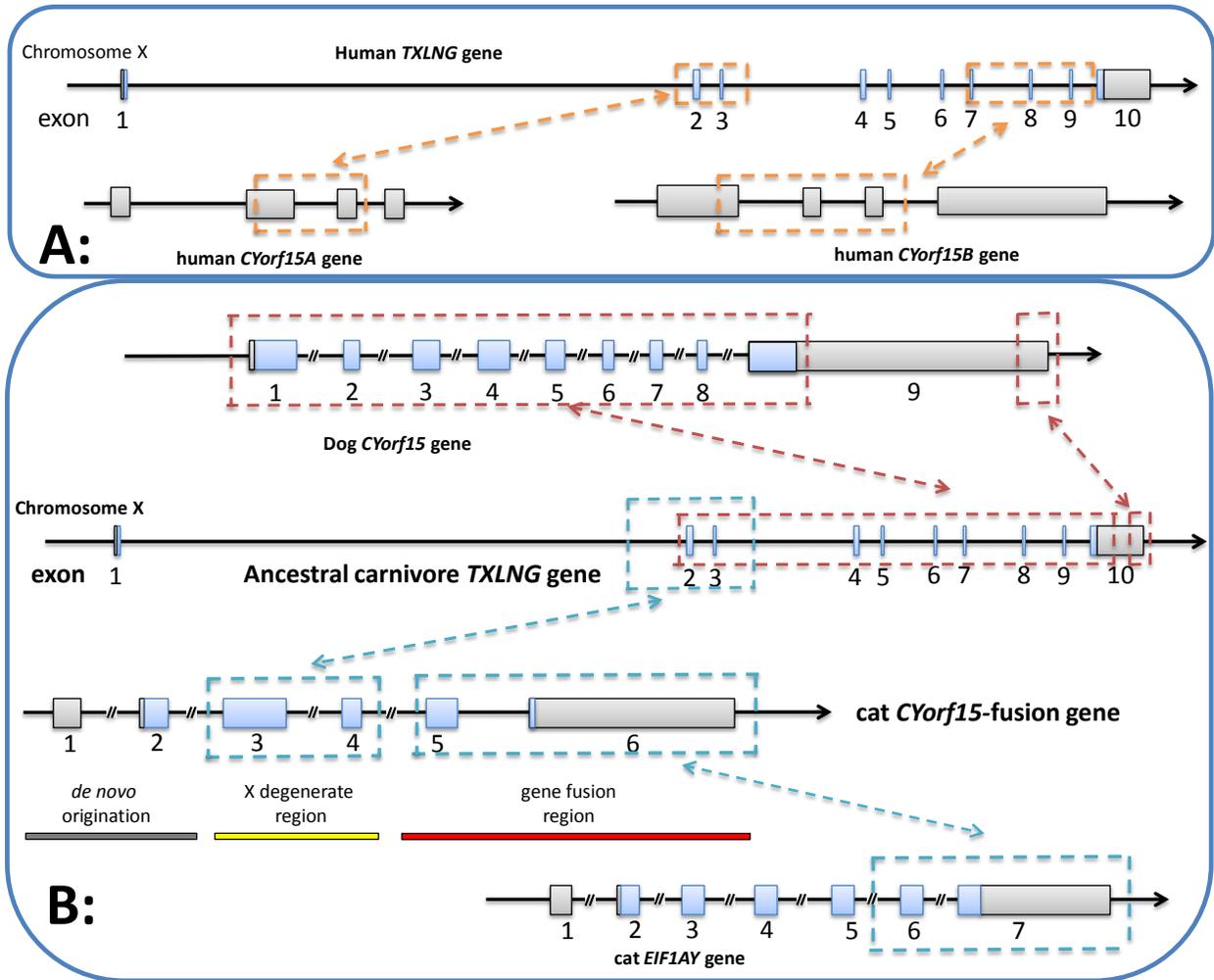


Figure 7. *OFD1* gene tree estimated using maximum likelihood method under a GTR+G+I model. Numbers near each node on the tree are bootstrap values (500 pseudoreplicates: only values higher than 80 have been shown). Taxon names are denoted by (species common name)/(ENSEMBL transcript name or NCBI accession number)/(available chromosome information). The chimpanzee Y chromosome *OFD1* pseudogene is denoted by the start and end locations within the published Y chromosome sequence.

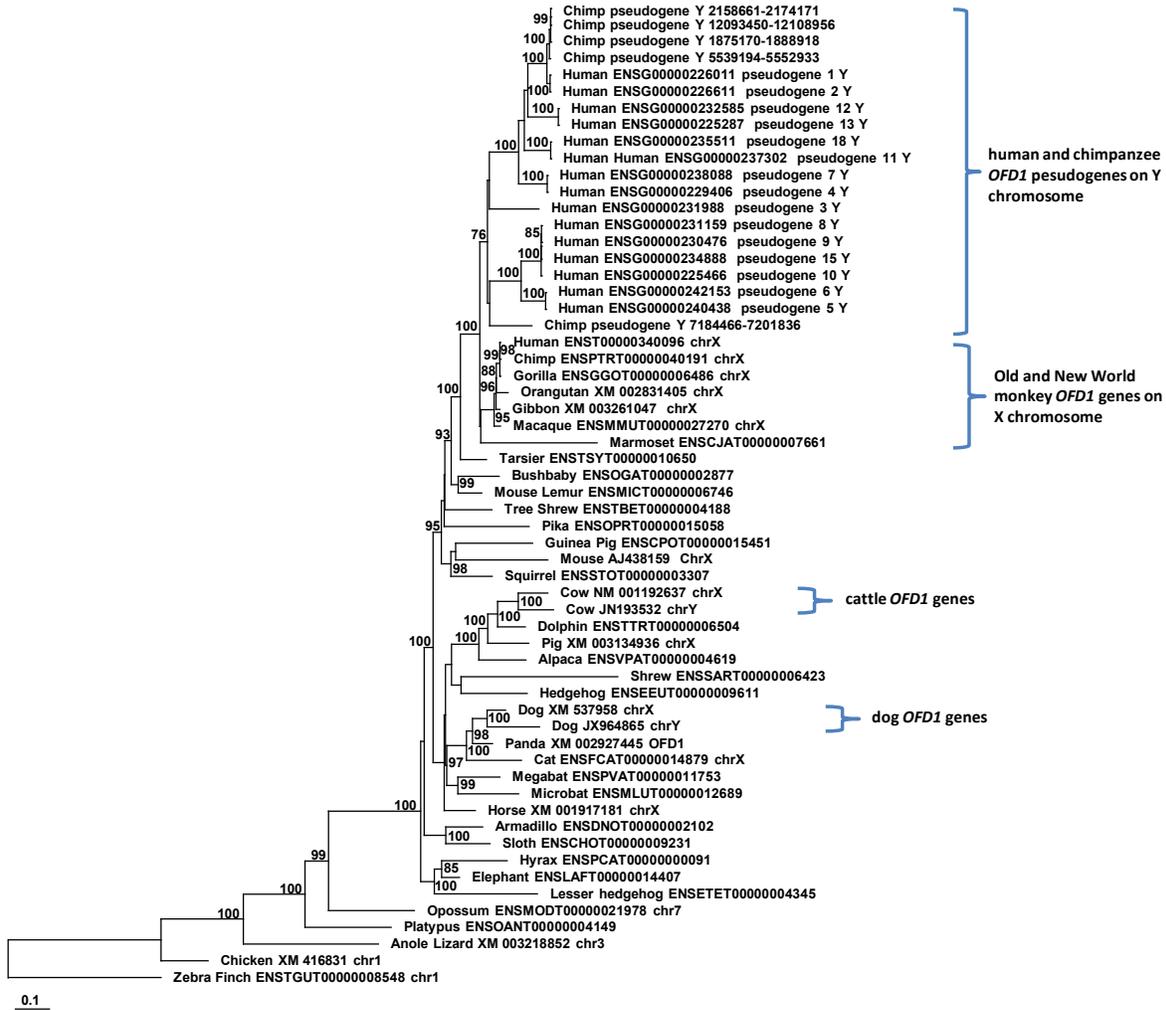


Figure 8. Results of positive selection test for the branch leading to domestic dog *OFDIY* gene. **A)** *OFDI* gene tree, with branch lengths based on non-synonymous substitution rates. **B)** Bayesian posterior probabilities for protein sites under positive natural selection (branch-site models comparison). Colored bars underneath indicate different structural domains. **C)** Results of likelihood ratio tests for two pairs of model comparisons.

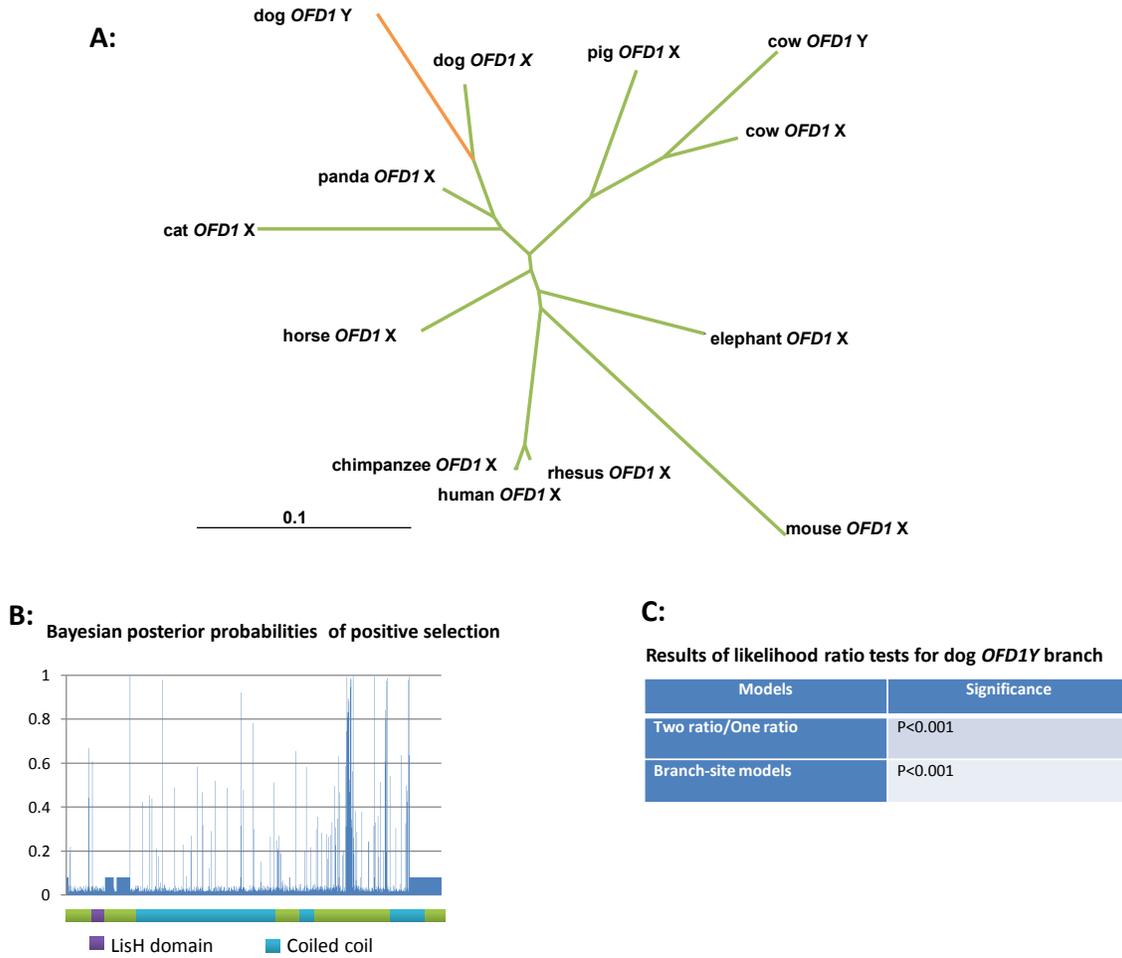


Figure 9. Location of the dog (top) and cat (bottom) pseudoautosomal boundaries (dashed vertical line), inferred from X-Y sequence similarity plots. Gene annotations are shown at the bottom of each figure.

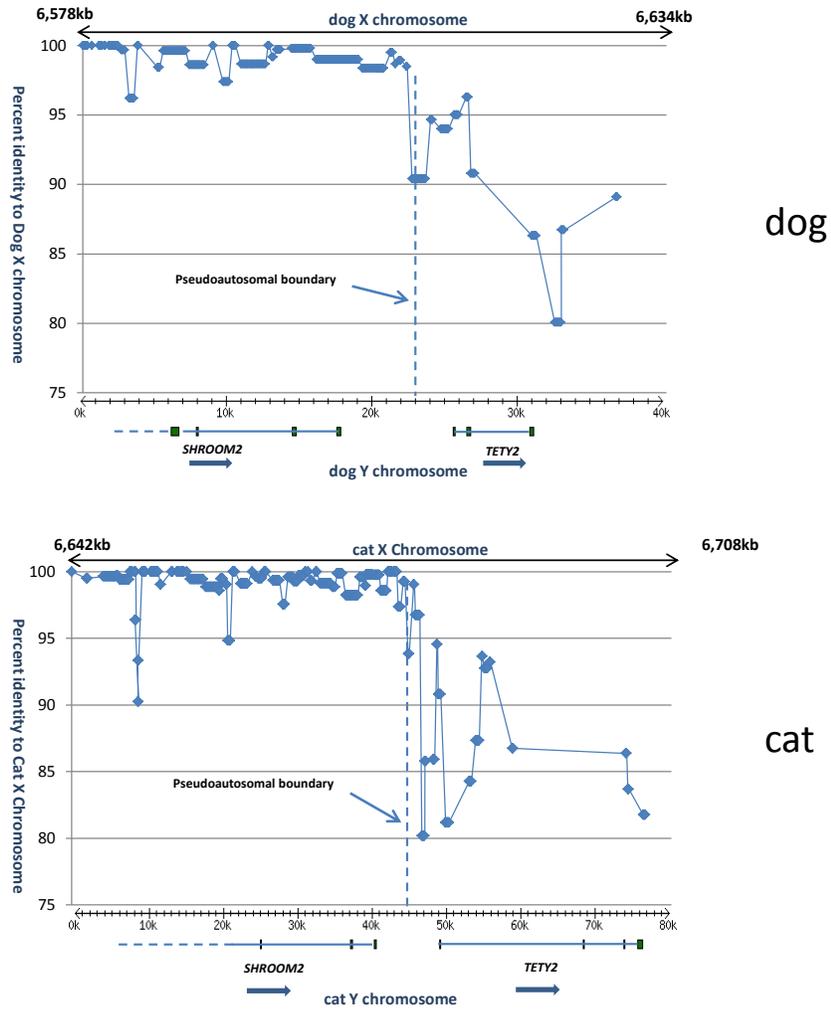


Figure 10. **A)** Maximum likelihood tree of *TETX2* orthologues sampled from different eutherian X chromosomes, estimated under a GTR+G+I model. **B)** Schematic showing the origination and divergence of the Y-linked *TETY2* gene on dog and cat Y chromosomes from the X chromosome progenitor gene, resulting from several mechanisms, including PAB movement and exon loss and *de novo* exon gain.

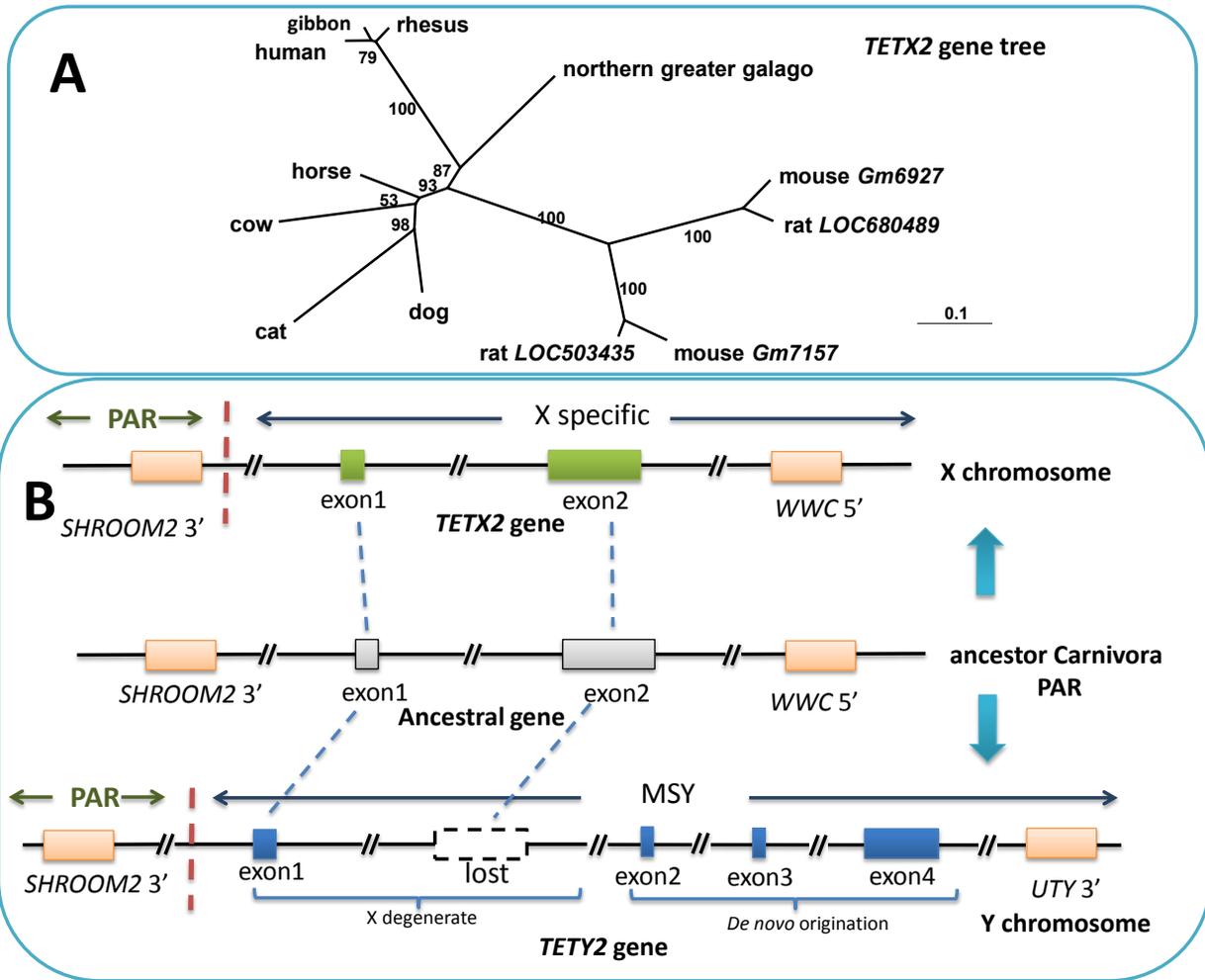


Figure 11. X-Y dot plots show regions of sequence conservation (>97%). The blue box highlights additional remnants of *SHROOM2* X-Y similarity, well inside the MSY and distant from the PAR, reflecting ancient chromosomal inversions that likely occurred on the ancestral felid branch.

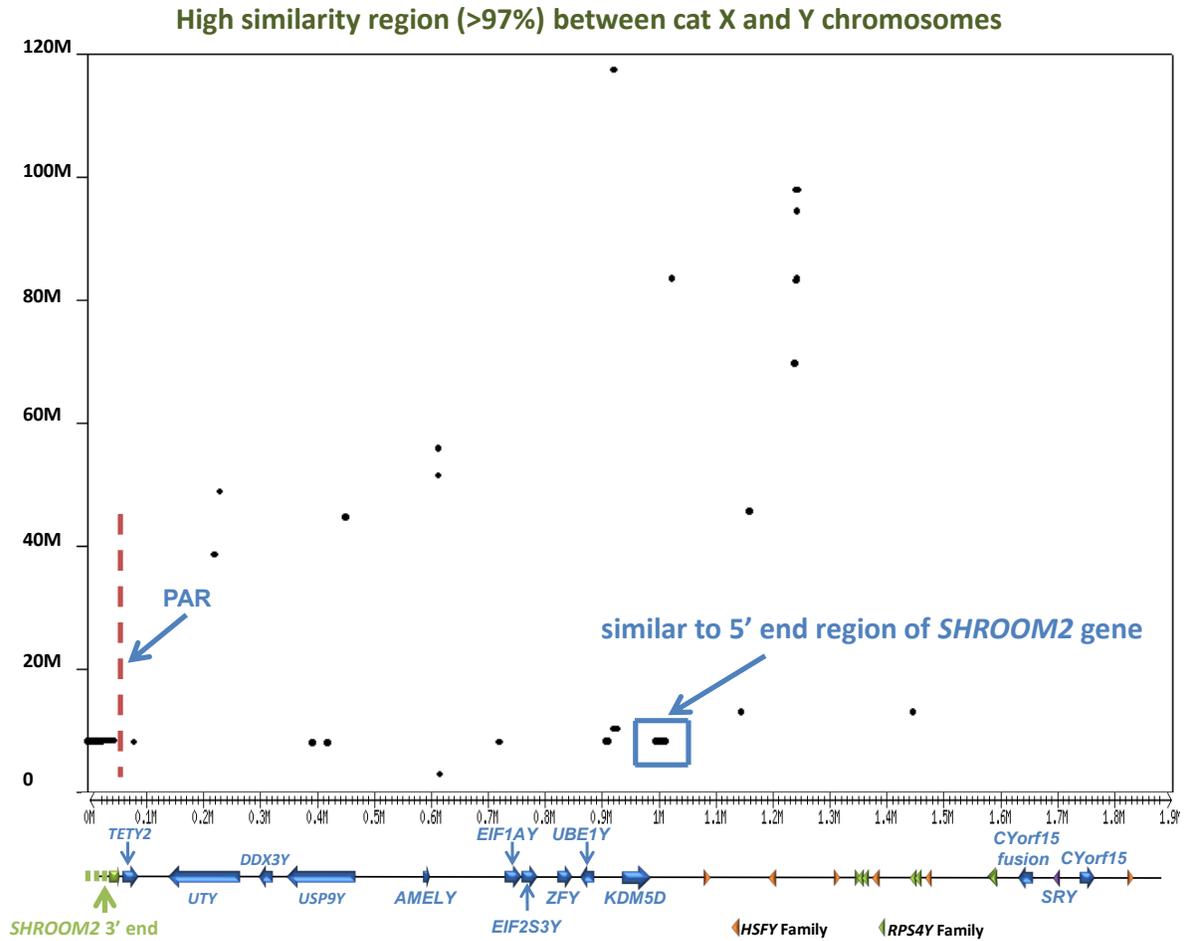


Figure 12. Pairwise sequence comparisons of four Y chromosomes (dog versus cat, rhesus, and human), using a BLAST-based method on repeatmasked sequences. Each dot represents a 50 nucleotide window with at least 50 percent identity between two species. Green boxes indicate the *UTY/DDX3Y/USP9Y* conserved gene region. The schematic on the x and y axes of each figure indicates the size and location of each MSY gene (green bars).

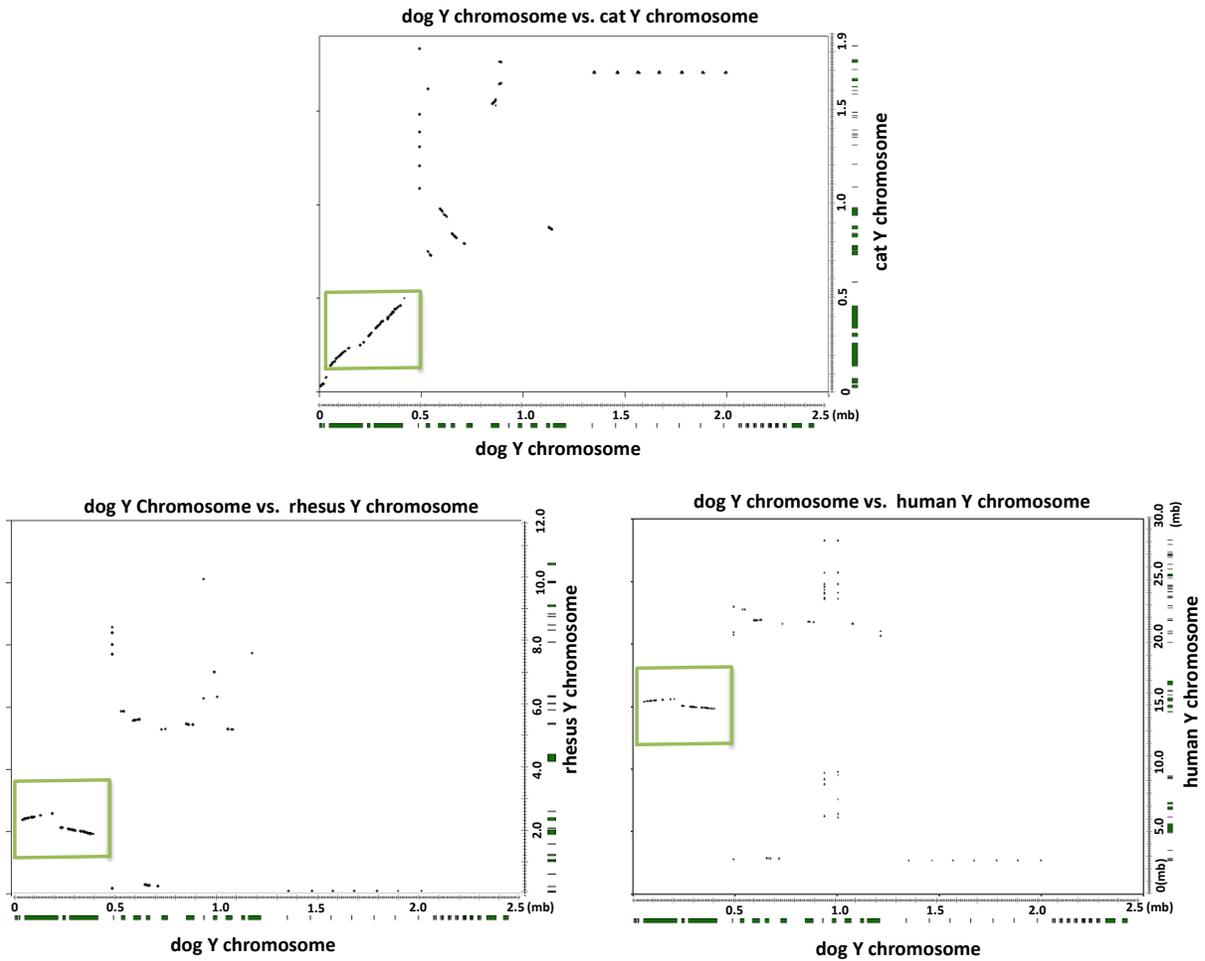


Figure 13. Plots showing rates of MSY gene loss during the evolution of mammals and continuing in the carnivore lineage, assuming an exponential model of decay with baseline (Hughes et al. 2012). Plots are divided into two groups of X-degenerate gene classes: Stratum I (Murphy et al. 2006; Hughes et al. 2012), and stratum II/III (following Pearks Wilkerson et al. 2008). Estimated ancestral gene numbers (Hughes et al. 2012), and active genes inferred for different eutherian ancestors (with divergence times drawn from Meredith et al. 2011) are shown as blue diamonds. Age of the ancestral Stratum I genes (a) is based upon the inferred ~218 Mya prototherian/therian divergence time (Meredith et al. 2011), whereas the age of the origin of Stratum II/III genes (a) is derived from Pearks Wilkerson et al. 2008. Other nodes are as follows: (b) primate/rodent/carnivore ancestor, (c) primate/rodent ancestor, (d) carnivore ancestor, (e) catarrhine primate ancestor.

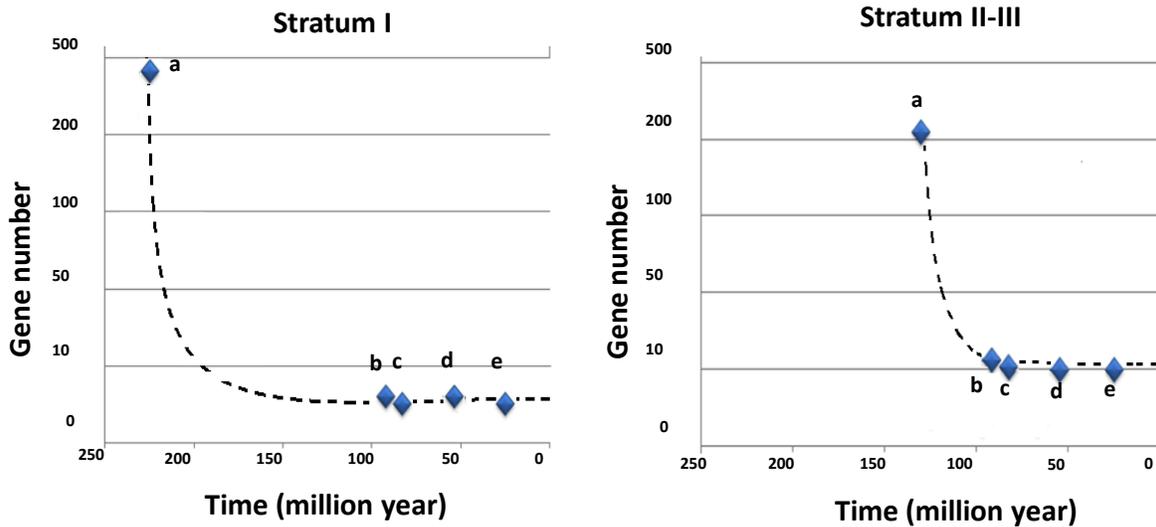
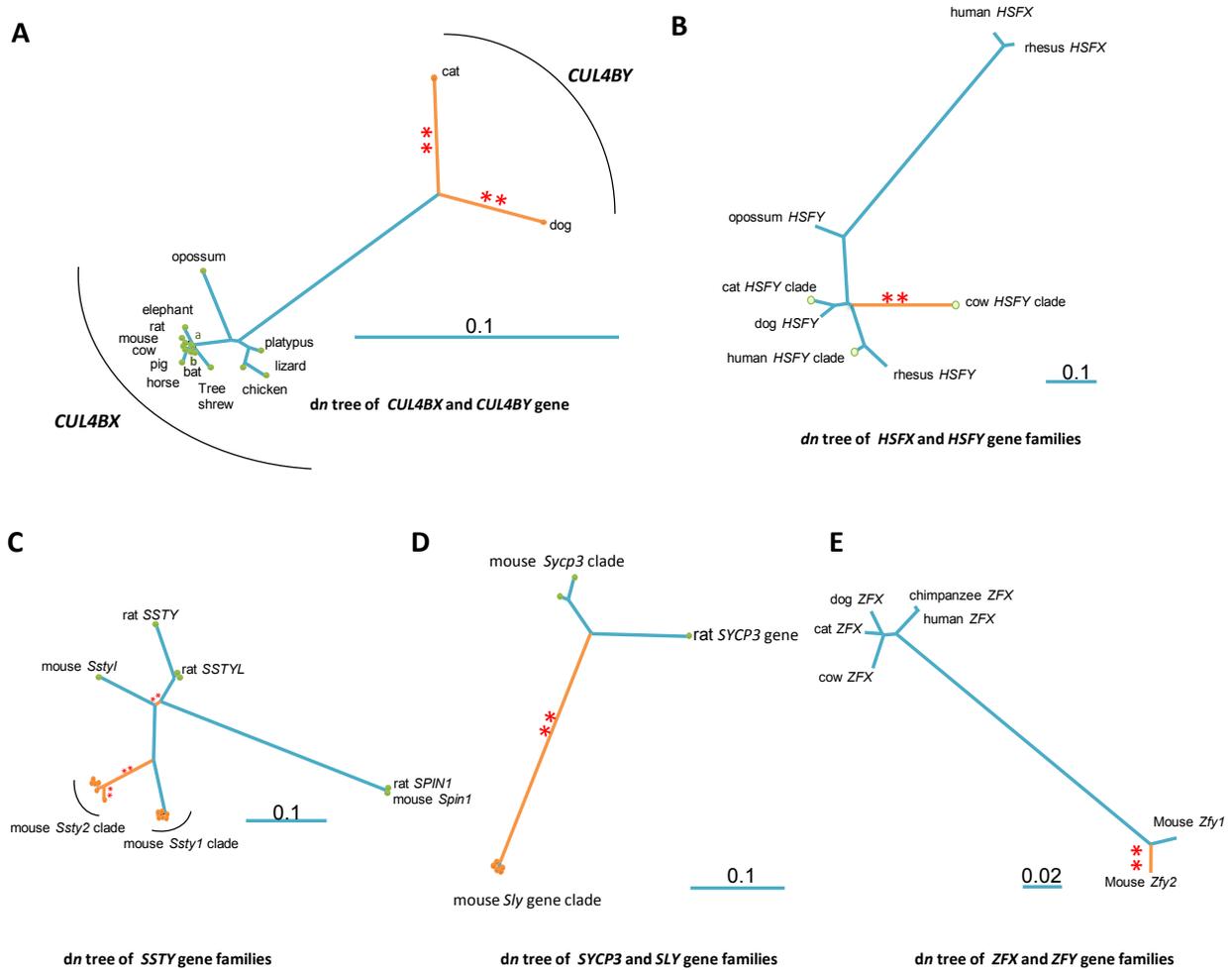


Figure 14. Results of positive selection tests for several multicopy MSY gene families, performed using the branch-site model. Asterisks indicate significant ($p < 0.01$) positive selection detected with likelihood ratio tests.



Supplementary Table 1. Structural domains of the domestic cat and dog Y chromosomes, and relative sequence coverage. Estimated sizes are based on cytogenetic and physical mapping data (Murphy et al. 2006, Pearks-Wilkerson et al. 2008).

region	PAR		Y single copy		Multicopy/ampliconic		NOR	
	estimated size	sequenced	estimated size	sequenced	estimated size	sequenced	estimated size	sequenced
cat	6.6 Mb Xp/Yp	23 Kb ^a Yp	1-2 Mb Yp	0.7 Mb Yp	~35 Mb Yp+Yq	1 Mb Yp	NA*	NA
dog	6.6 Mb Xp/Yq	45 Kb ^a Yq	1-2 Mb Yq	1.1 Mb Yq	1-2 Mb Yq	1.5 Mb Yq	8-10 Mb Yp	0 bp Yp

^alength of corresponding Y-PAR.

^bNA=not applicable

Supplementary Table 2. Comparison of synonymous substitution rates between canine X-Y paralogs.

Orthologous genes		Synonymous substitution rate	average sequence divergence within intron region
X copy	Y copy		
<i>OFD1X</i>	<i>OFD1Y</i>	0.155	0.149
<i>TRAPPC2</i>	<i>TRAPPC2P</i>	NA	0.146
<i>DDX3X</i>	<i>DDX3Y</i>	0.527	not alignable
<i>UTX</i>	<i>UTY</i>	0.334	not alignable
<i>KDM5C</i>	<i>KDM5D</i>	0.751	not alignable
<i>UBA1</i>	<i>UBE1Y</i>	0.560	not alignable
<i>USP9X</i>	<i>USP9Y</i>	0.561	not alignable
<i>ZFX</i>	<i>ZFY</i>	0.266	not alignable
<i>TXLNG</i>	<i>CYorf15</i>	0.498	not alignable
<i>BCOR</i>	<i>BCORY1/2</i>	0.774/0.687	not alignable
<i>EIF1AX</i>	<i>EIF1AY</i>	0.532	not alignable

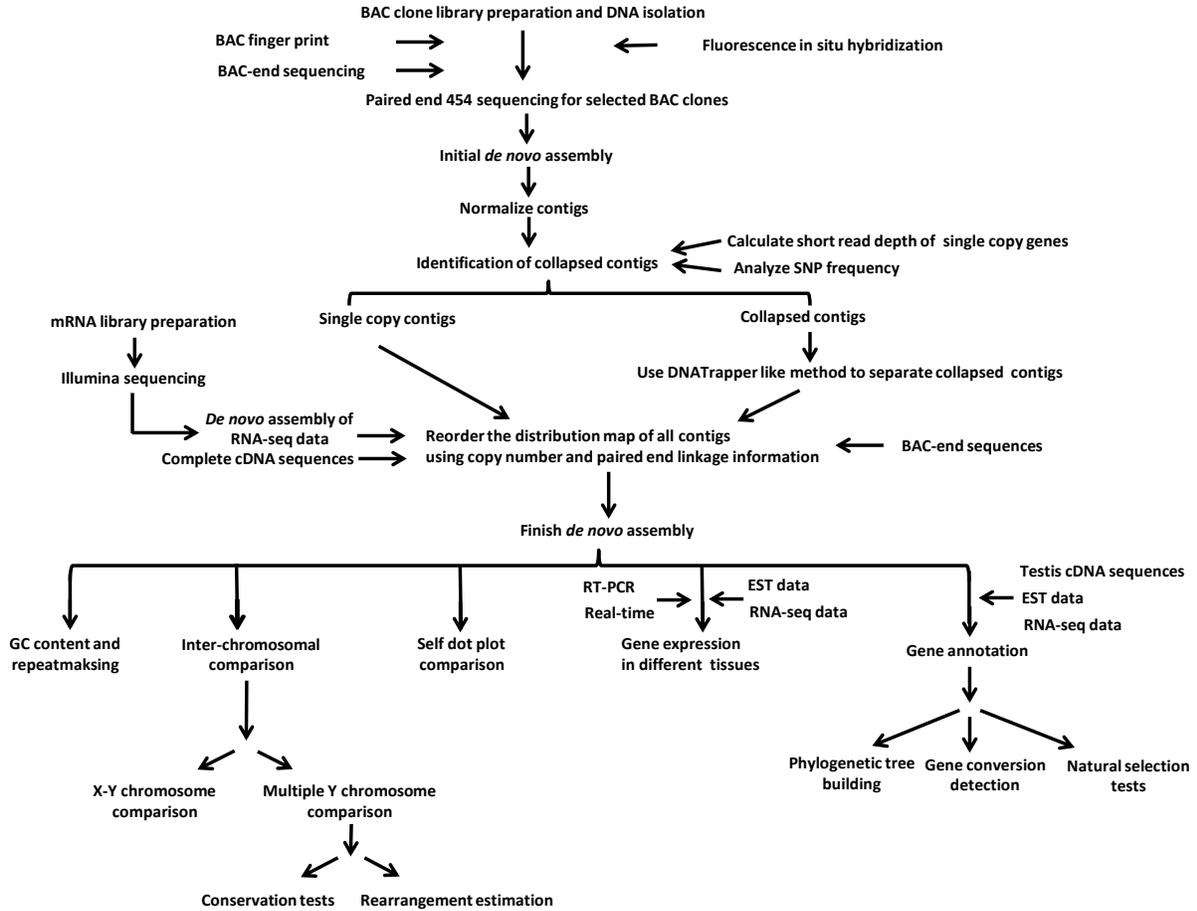
Supplementary Table 3. Sequences involved in selection tests.

Gene	species name	Accession number	Gene	species name	Accession number
<i>AMELY</i>	pig	AB091792	<i>UTY</i>	chimpanzee	ENSPTRT00000042001
<i>AMELY</i>	horse	AB032194	<i>UTY</i>	mouse	NM_009484
<i>AMELY</i>	cow	M63500	<i>ZFY</i>	mouse	AK076618
<i>AMELY</i>	human	ENST00000215479	<i>ZFY</i>	mouse	M24401
<i>AMELY</i>	chimpanzee	ENSPTRT00000041965	<i>ZFY</i>	cow	AF032867
<i>DDX3Y</i>	human	ENST00000360160	<i>ZFY</i>	human	ENST00000383052
<i>DDX3Y</i>	chimpanzee	ENSPTRT00000041996	<i>ZFY</i>	chimpanzee	ENSPTRT00000041960
<i>DDX3Y</i>	mouse	ENSMUST00000091190	<i>HSFY</i>	human	AF332227
<i>EIF1AY</i>	human	ENST00000361365	<i>HSFY</i>	human	AK058182
<i>EIF1AY</i>	chimpanzee	ENSPTRT00000042030	<i>HSFY</i>	rhesus	FJ527015
<i>EIF1AY</i>	rhesus	FJ527014	<i>HSFY</i>	opossum	GQ253469
<i>EIF1AY</i>	cow	FJ195366	<i>HSFY</i>	cow	NM_001077006
<i>EIF1AY</i>	pig	AY609402	<i>HSFY</i>	cow	XM_001250229
<i>EIF2S3Y</i>	rat	FJ775732	<i>HSFY</i>	cow	XM_003585244
<i>EIF2S3Y</i>	mouse	AJ006584	<i>TSPY</i>	human	ENST00000451548
<i>KDM5D</i>	human	ENST00000317961	<i>TSPY</i>	human	ENST00000457222
<i>KDM5D</i>	chimpanzee	ENSPTRT00000042025	<i>TSPY</i>	human	ENST00000428845
<i>KDM5D</i>	mouse	NM_011419	<i>TSPY</i>	human	ENST00000426950
<i>SRY</i>	human	AM884750	<i>TSPY</i>	human	ENST00000426950
<i>SRY</i>	chimpanzee	DQ977342	<i>TSPY</i>	human	ENST00000457222
<i>SRY</i>	rhesus	FJ527022	<i>TSPY</i>	human	ENST00000287721
<i>SRY</i>	rabbit	AY785433	<i>TSPY</i>	chimpanzee	ENSPTRT00000055849
<i>SRY</i>	mouse	U03645	<i>TSPY</i>	chimpanzee	ENSPTRT00000041983
<i>SRY</i>	rat	AF274872	<i>TSPY</i>	chimpanzee	ENSPTRT00000055827
<i>SRY</i>	cow	AB039748	<i>TSPY</i>	chimpanzee	ENSPTRT00000055807
<i>SRY</i>	pig	U49860	<i>TSPY</i>	pig	ENSSSCT00000022749
<i>SRY</i>	horse	AB004572	<i>TSPY</i>	cow	ENSBTAT00000049474
<i>SRY</i>	panda	AB292070	<i>CDY</i>	human	ENST00000426790
<i>UBE1Y</i>	mouse	AF150963	<i>CDY</i>	human	ENST00000306882
<i>UBE1Y</i>	rat	EF690356	<i>CDY</i>	human	ENST00000544303
<i>UBE1Y</i>	opossum	GQ253467	<i>CDY</i>	human	ENST00000306609
<i>USP9Y</i>	human	ENST00000338981	<i>CDY</i>	chimpanzee	ENSPTRT00000076216
<i>USP9Y</i>	chimpanzee	ENSPTRT00000041991	<i>CDY</i>	chimpanzee	ENSPTRT00000073911
<i>USP9Y</i>	mouse	NM_148943	<i>CDY</i>	chimpanzee	ENSPTRT00000076243
<i>USP9Y</i>	cow	NM_001145509	<i>CDY</i>	chimpanzee	ENSPTRT00000072418
<i>UTY</i>	human	ENST00000331397	<i>CDY</i>	chimpanzee	ENSPTRT00000072527

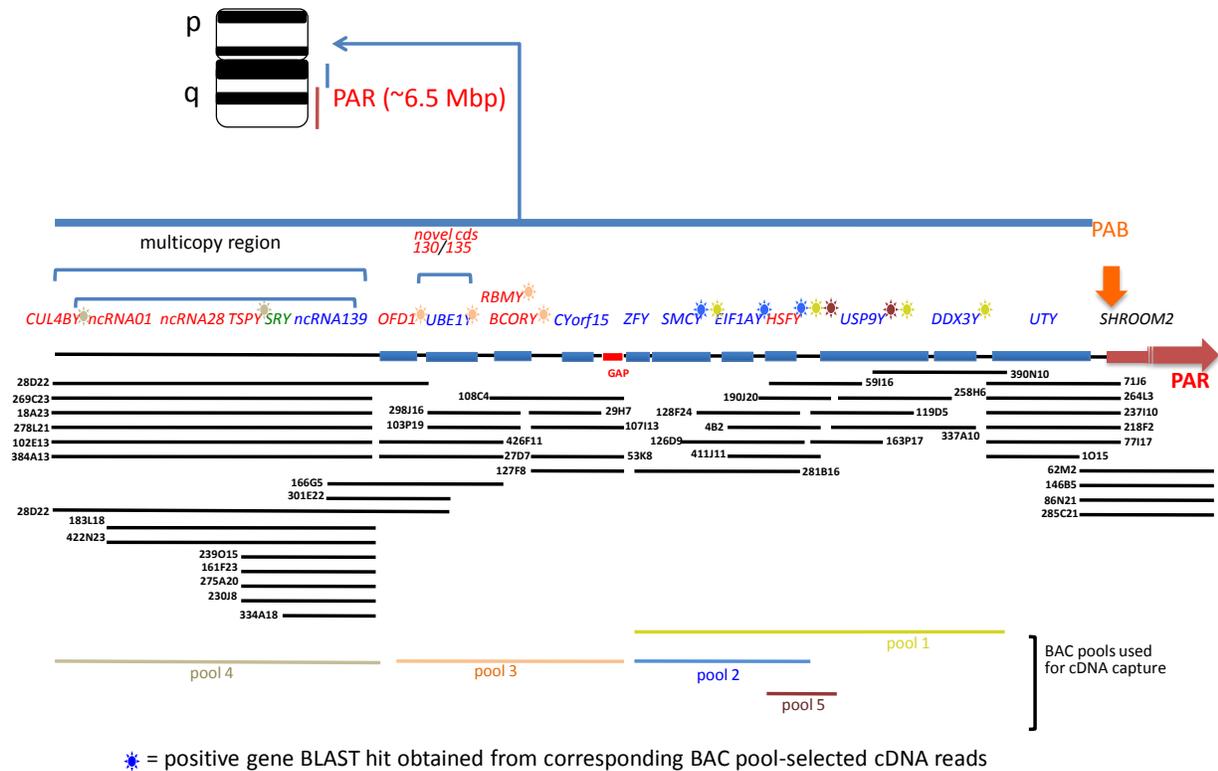
Gene	species name	Accession number	Gene	species name	Accession number
<i>CDY</i>	rhesus	FJ527011	<i>SSTYIL</i>	rat	XM_222167
<i>DAZ</i>	chimpanzee	ENSPTRT00000062513	<i>SSTYIL</i>	rat	XM_001078732
<i>DAZ</i>	chimpanzee	ENSPTRT00000062535	<i>SstyL</i>	mouse	XM_144574
<i>DAZ</i>	chimpanzee	ENSPTRT00000062545	<i>SstyI</i>	mouse	X05260
<i>DAZ</i>	chimpanzee	ENSPTRT00000055821	<i>SstyI</i>	mouse	Mm#S38822574 (EST)
<i>DAZ</i>	human	ENST00000382440	<i>SstyI</i>	mouse	Mm#S38822973
<i>DAZ</i>	human	ENST00000405239	<i>SstyI</i>	mouse	Mm#S38822990
<i>DAZ</i>	human	ENST00000382365	<i>SstyI</i>	mouse	Mm#S8039565
<i>DAZ</i>	human	ENST00000440066	<i>SstyI</i>	mouse	Mm#S7116079
<i>DAZ</i>	rhesus	FJ648738	<i>SstyI</i>	mouse	Mm#S38821906
<i>Cyorf15A</i>	human	ENST00000407724	<i>SstyI</i>	mouse	Mm#S32784017
<i>Cyorf15A</i>	chimpanzee	NM_001009080	<i>SstyI</i>	mouse	Mm#S26680069
<i>Cyorf15A</i>	rhesus	FJ527012	<i>SstyI</i>	mouse	Mm#S26702624
<i>Cyorf15A</i>	gorilla	FJ532256	<i>SstyI</i>	mouse	Mm#S40811957
<i>Cyorf15B</i>	human	ENST00000382832	<i>Ssty2</i>	mouse	Mm#S7227131
<i>Cyorf15B</i>	rhesus	FJ648737	<i>Ssty2</i>	mouse	Mm#S52502287
<i>Cyorf15B</i>	gorilla	FJ532257	<i>Ssty2</i>	mouse	Mm#S7227130
<i>CUL4BX</i>	microbat	ENSMLUT00000013243	<i>Ssty2</i>	mouse	Mm#S38822045
<i>CUL4BX</i>	dog	ENSCAFT00000029435	<i>Ssty2</i>	mouse	Mm#S19760504
<i>CUL4BX</i>	platypus	ENSOANT00000004088	<i>Ssty2</i>	mouse	Mm#S14960481
<i>CUL4BX</i>	pig	ENSSSCT00000013782	<i>Sly</i>	mouse	BC049626
<i>CUL4BX</i>	lizard	ENSACAT00000010801	<i>Sly</i>	mouse	Mm#S38821277 (EST)
<i>CUL4BX</i>	chicken	ENSGALT00000013942	<i>Sly</i>	mouse	Mm#S38821362
<i>CUL4BX</i>	marmoset	ENSCJAT00000005037	<i>Sly</i>	mouse	Mm#S38822154
<i>CUL4BX</i>	human	ENST00000404115	<i>Sly</i>	mouse	Mm#S38822342
<i>CUL4BX</i>	tree shrew	ENSTBET00000007462	<i>Sly</i>	mouse	Mm#S60235098
<i>CUL4BX</i>	mouse	NM_001110142	<i>Sly</i>	mouse	Mm#S38822478
<i>CUL4BX</i>	cat	ENSFCAT00000010442	<i>Sly</i>	mouse	Mm#S21508819
<i>CUL4BX</i>	elephant	ENSLAFT00000034221	<i>Sly</i>	mouse	Mm#S14960473
<i>CUL4BX</i>	cow	ENSBTAT00000024713	<i>Sycp3</i>	mouse	ENSMUST00000033458
<i>CUL4BX</i>	horse	ENSECAT00000025736	<i>Sycp3</i>	mouse	ENSMUST00000067763
<i>CUL4BX</i>	rat	ENSRNOT00000065023	<i>Sycp3</i>	rat	ENSRNOT00000056700
<i>CUL4BX</i>	rabbit	ENSOCUT00000015646	<i>FLJ36031Y</i>	cat (clone a260)	DQ329513
<i>CUL4BX</i>	opossum	XM_001372861	<i>FLJ36031Y</i>	cat (b321)	DQ329514
<i>Spin1</i>	mouse	AK168544	<i>FLJ36031Y</i>	cat (c254)	DQ329515
<i>SPIN1</i>	rat	BC097359	<i>FLJ36031Y</i>	cat (e268)	DQ329517
<i>SSTY</i>	rat	AABR06008269	<i>FLJ36031Y</i>	cat (d340)	DQ329516

Supplementary Methods

Pipeline of data generation and analyses.



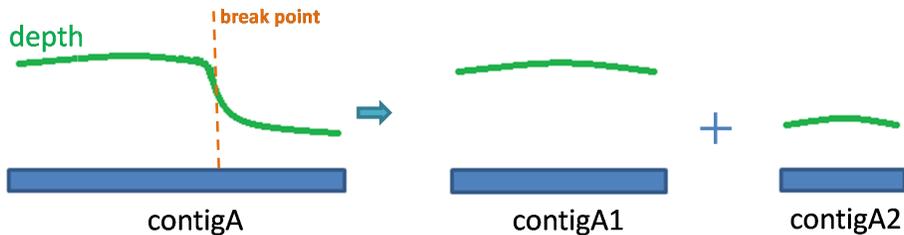
Domestic dog Y chromosome BAC clone map.



De novo assembly of cat and dog Y chromosome sequences

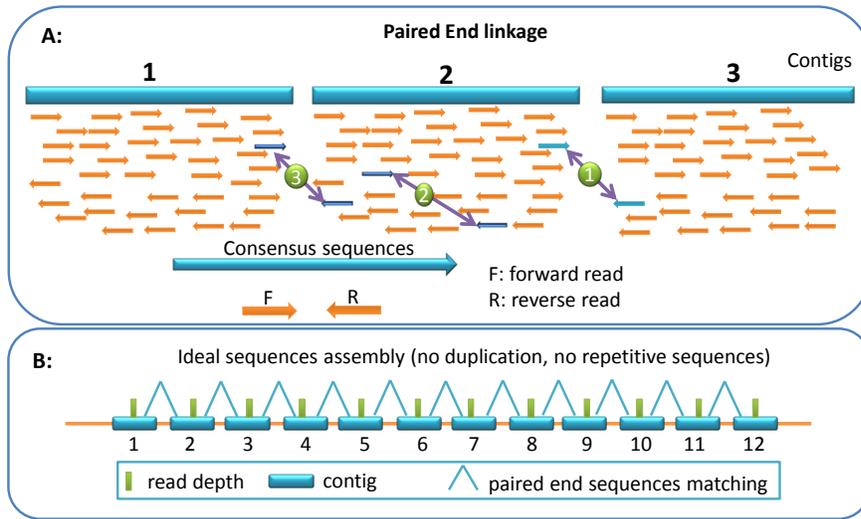
The initial sequence assembly was subjected to custom manual annotation, as follows:

- We normalized the initial *de novo* contigs containing significantly uneven depth of coverage by identifying regions with significant shifts in coverage, and cutting the contigs into new sub-contigs with even read depth coverage (see below).

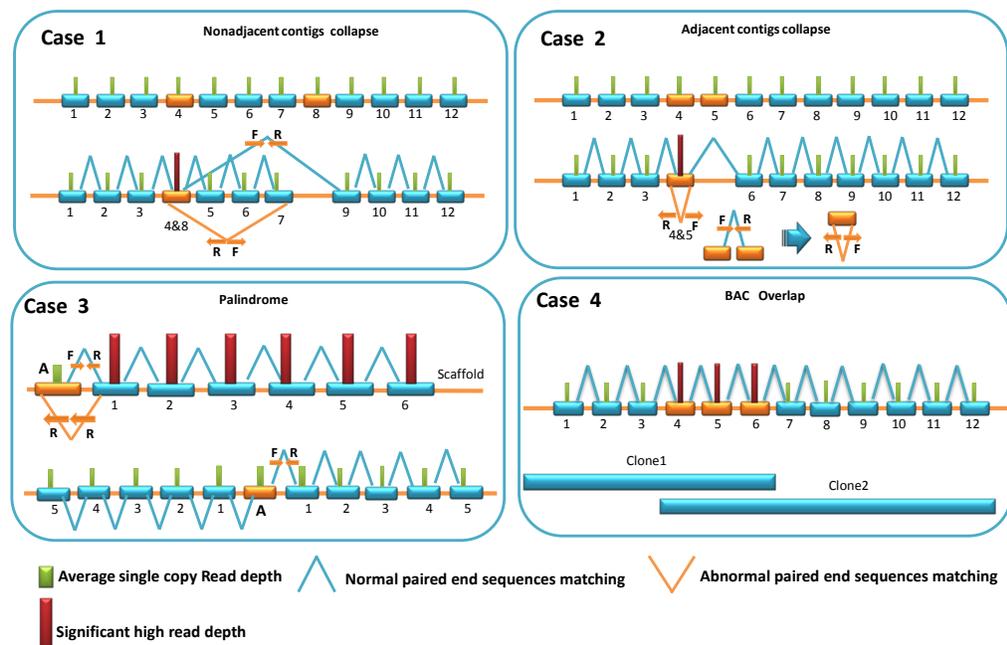


- Paired end sequences are favorable for sequence assembly by utilizing upstream and

downstream read recognition (see Fig. A below). An ideal *de novo* sequence assembly based on NGS data will exhibit relatively even read-depth within each contig and a contiguous distribution of matching paired-end reads (Fig. B below).



Using this information, we identified collapsed contigs (caused by gene/sequence duplication) in the initial *de novo* assembly. Several examples of different assembly artifacts are shown below. For example, in a case where two regions with highly similar, paralogous sequences collapse into one contig with greater than average read coverage, this results in two distinct and abnormal patterns of paired end read distribution (Case 1 and Case 2 below).

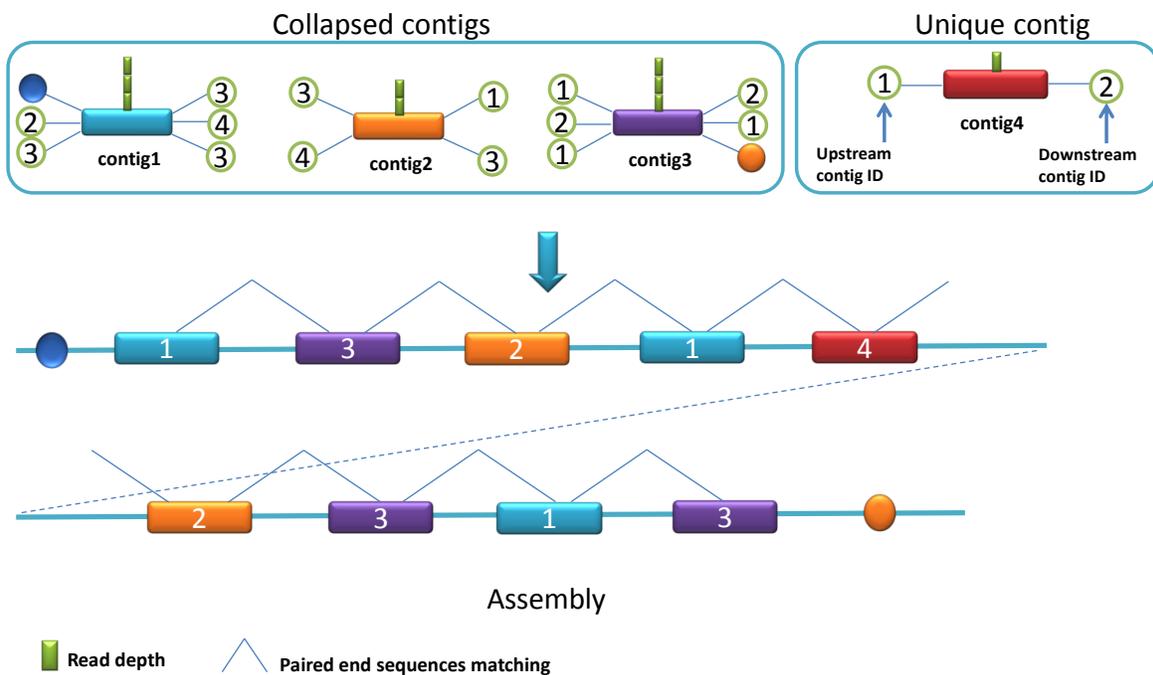


Case 3 indicates a palindrome structure of highly similar sequences with read depth and paired end reads matching. Regions of BAC clone overlap also can cause significantly higher read

depth (Case 4), but paired end sequences patterns that are distinct from collapsed contigs caused by gene duplication.

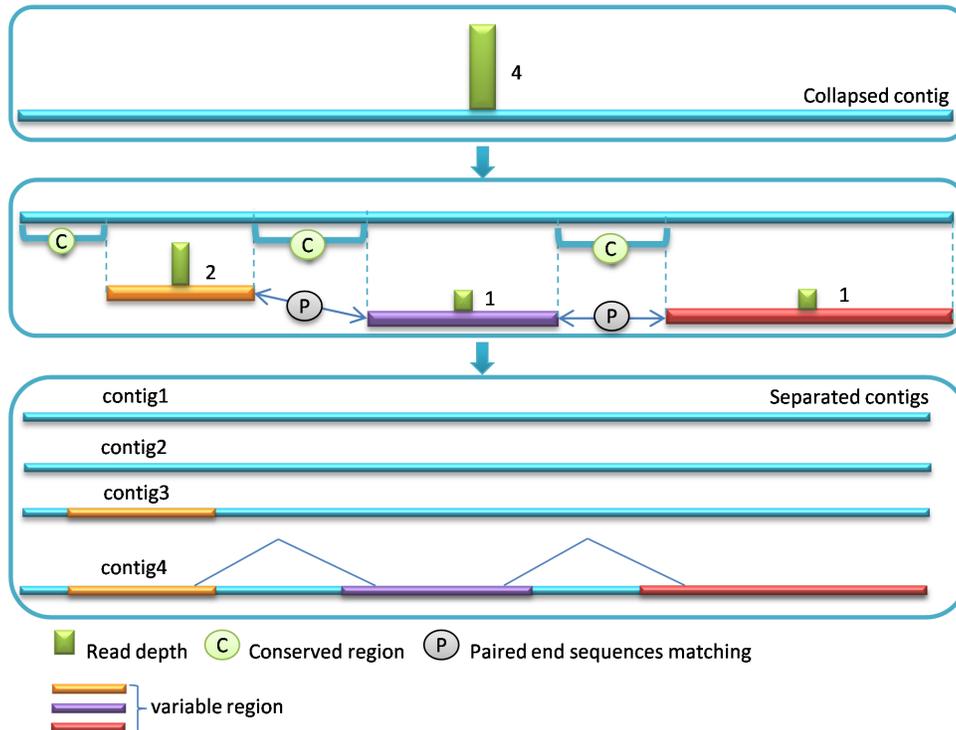
Using this information we reanalyzed the initial *de novo* sequence assembly. First we evaluated the copy number of each contig that was identified as potential collapsed contigs containing sequences from highly similar, duplicated Y chromosome regions. In this step, two criteria were used. We estimated the average expected depth of sequence coverage in contigs which contain known single copy coding genes (e.g. *UTY/DDX3Y/USP9Y* gene coding region). All normalized contigs were then compared to this value to provide an initial estimate of the potential collapsed copy number. Second, because the sequenced MSY region is haploid, only sequencing error or collapsed sequences could account for polymorphisms found in the initial *de novo* assembled contigs. We analyzed the SNP frequency of polymorphic sites found in contigs exhibiting a significant excess of coverage depth to adjust our initial copy number estimate.

c). Using paired-end sequence information and the estimated copy number of collapsed contigs, we calculated the maximum parsimony distribution map of contigs to rebuild scaffolds (below).



d). We used a method similar to that applied in the software DNPtrapper(Arner et al. 2006) to separate collapsed contigs into unique contigs. Collapsed contigs include highly conserved regions and polymorphic region that result from collapse of nearly identical sequences distinguished by rare sequence variants. Variable region sequences are separated from the original contig and then sequence coverage is recalculated. Using paired end information the new

separated contigs are rebuilt.



The *de novo* assembly is completed by repositioning the formerly collapsed contigs into unique positions in the existing assembly. Copy number information inferred from read depth data was validated with quantitative PCR.

References

- Arner E, Tammi MT, Tran AN, Kindlund E, Andersson B. 2006. DNPTrapper: an assembly editing tool for finishing and analysis of complex repeat regions. *BMC Bioinformatics* 7:155.
- Hughes JF, Skaletsky H, Brown LG, Pyntikova T, Graves T, Fulton RS, Dugan S, Ding Y, Buhay CJ, Kremitzki C et al. 2012. Strict evolutionary conservation followed rapid gene loss on human and rhesus Y chromosomes. *Nature* 483(7387): 82-86.
- Meredith RW, Janecka JE, Gatesy J, Ryder OA, Fisher CA, Teeling EC, Goodbla A, Eizirik E, Simao TLL, Stadler T et al. 2011. Impacts of the Cretaceous Terrestrial Revolution and KPg Extinction on Mammal Diversification. *Science* 334(6055): 521-524.
- Murphy WJ, Wilkerson AJP, Raudsepp T, Agarwala R, Schaffer AA, Stanyon R, Chowdhary BP. 2006. Novel gene acquisition on carnivore Y chromosomes. *Plos Genet* 2(3): 353-363.
- Pearks Wilkerson AJ, Raudsepp T, Graves T, Albracht D, Warren W, Chowdhary BP, Skow LC, Murphy WJ. 2008. Gene discovery and comparative analysis of X-degenerate genes from the domestic cat Y chromosome. *Genomics* 92(5): 329-338.