

An Evolutionary Driver of Interspersed Segmental Duplications in Primates

Stuart Cantsilieris^{1,2}, Susan M. Sunkin³, Matthew E. Johnson⁴, Fabio Anaclerio⁵, John Huddlestone^{6,7}, Carl Baker¹, Max L. Dougherty¹, Jason G. Underwood⁸, Arvis Sulovari¹, PingHsun Hsieh¹, Yafei Mao¹, Claudia Rita Catacchio⁵, Maika Malig^{1,9,10}, AnneMarie E. Welch^{1,11}, Melanie Sorensen¹, Katherine M. Munson¹, Weihong Jiang¹², Santhosh Girirajan¹³, Mario Ventura⁵, Bruce T. Lamb¹⁴, Ronald A. Conlon¹², and Evan E. Eichler^{1,15*}

¹Department of Genome Sciences, University of Washington School of Medicine, Seattle, WA 98195, USA.

²Centre for Eye Research Australia, Department of Surgery (Ophthalmology), University of Melbourne, Royal Victorian Eye and Ear Hospital, East Melbourne, VIC 3002, Australia.

³Allen Institute for Brain Science, Seattle, WA, USA.

⁴Center for Spatial and Functional Genomics, Children's Hospital of Philadelphia, Philadelphia, PA 19104, USA.

⁵Department of Biology-Genetics, University of Bari, Bari, Italy.

⁶Vaccine and Infectious Disease Division, Fred Hutchinson Cancer Research Center, Seattle, WA 98109, USA.

⁷Molecular and Cellular Biology Program, University of Washington, Seattle, WA 98195, USA.

⁸Pacific Biosciences (PacBio) of California, Incorporated, Menlo Park, CA 94025, USA.

⁹Department of Molecular and Cellular Biology, University of California, Davis, CA 95616, USA.

¹⁰Integrative Genetics and Genomics Graduate Group, University of California, Davis, CA 95616, USA.

¹¹Brain and Mitochondrial Research, Murdoch Children's Research Institute, Royal Children's Hospital, Melbourne, VIC, Australia.

¹²Case Transgenic and Targeting Facility, Department of Genetics and Genome Sciences, School of Medicine, Case Western Reserve University, Cleveland, OH 44106, USA.

¹³Department of Biochemistry and Molecular Biology, Department of Anthropology, Pennsylvania State University, University Park, PA 16802, USA.

¹⁴Stark Neurosciences Research Institute, Indiana University School of Medicine, Indianapolis, IN 46202, USA.

¹⁵Howard Hughes Medical Institute, University of Washington, Seattle, WA 98195, USA.

#Bold represent current affiliations.

*Corresponding author: Evan E. Eichler, Ph.D.

University of Washington School of Medicine

Howard Hughes Medical Institute

3720 15th Ave NE, S413C

Box 355065

Seattle, WA 98195-5065

Phone: (206) 543-9526

E-mail: eee@gs.washington.edu

Methods

Transgenic Mouse Lines

We established two independent lines for each transgenic. A15.26, A15.3, O14.20, O14.23, H15.1, and H15.2 corresponded to three human BACs (RP11-344H15, RP11-1381A15, and RP11-1236O14) and lines I13.43 and I13.49 corresponded to a single BAC from the baboon (RP41-285I13). RP11-1381A15 (205 kbp), which was used to generate transgenic mouse line A15, represented the transcriptional splice form of *NPIP* (*NPIPA7*), which contained all eight exons and produces a 1,070 bp transcript in humans. The BAC also contains the NODAL modulator 3 gene (*NOMO3*) downstream of *NPIP*. RP11-1236O14 (218 kbp), which was used to generate transgenic mouse line O14, represents the transcriptional splice form *NPIPA1* and contains all eight exons but had a 54 bp larger spliceform of exon 4. RP11-1236O14 also contains a full-length gene transcript of NODAL modulator 1 (*NOMO1*) located upstream of *NPIPA1*. RP11-344H15 (185 kbp), which was used to generate transgenic mouse line H15, expresses a variation of *NPIP* containing two copies of exon 2 in the full-length transcript. We identified three other full-length genes; the first two, RNA polymerase I transcription factor (*RRN3*) and N-terminal Asnamidase (*NTANI*), reside downstream of the *NPIP* transcript. The third gene, polycystic kidney disease 1 (*PKD1*), is upstream of the *NPIP* transcript. Finally, transgenic mice were created by using the ancestral single-copy *NPIP* gene from baboon RP41-285I13 (184 kbp) and this line was referred to as I13. The BAC also contained two homologous genes (*NTANI* and *RRN3*) downstream of the single-copy *NPIP* gene. We investigated expression patterns across two independent lines and in all transgenics (3 x Human: A15, O14, H15; 1 x Baboon: I13). RT-PCR experiments were performed from cDNA sourced from a panel of seven tissues that included brain, heart, kidney, liver, lung, muscle, and testis (data not shown). In each independent line the expression results were near identical. For all human BAC integrants we observed a pattern of ubiquitous expression, while both independent lines generated from the baboon showed expression specific to the testis.

In situ Hybridization (ISH)

For the ISH experiments, we proceeded with four human BAC transgenic lines representing two independent integrations (A15.26 and A15.3) of the same BAC, RP11-344H15, and two independent integrations (I13.43 and I13.49) of the baboon BAC, RP41-285I13. A control probe from the mouse locus, *Drd1a*, was used as a positive control for the ISH experiments and was hybridized to tissue sections from each of mouse transgenic lines. *NPIP* probes specific to the corresponding sequence in both baboon and human were designed based on cDNA sequence available for representative loci. In the case of baboon, we specifically designed a 659 bp probe derived from macaque cDNA (XM_001109190), which is 99% identical. In the case of human, we selected two probes based on available human cDNA (Additional file 2: Table S9). A homologous *NPIP* gene and/or the gene family does not appear to exist in the mouse. Sequence analysis of the corresponding homologous locus shows that it is highly divergent and lacks any significant sequence homology with the transcripts identified amongst primates. The three independent paralogs of the human *NPIP* versions gave near identical ISH results consistent with the observations from the primary tissue (Additional file 1: Fig. S17). We do note, however, there was some evidence of differences in levels of expression although no difference in the overall pattern with respect to the primary tissue was observed.

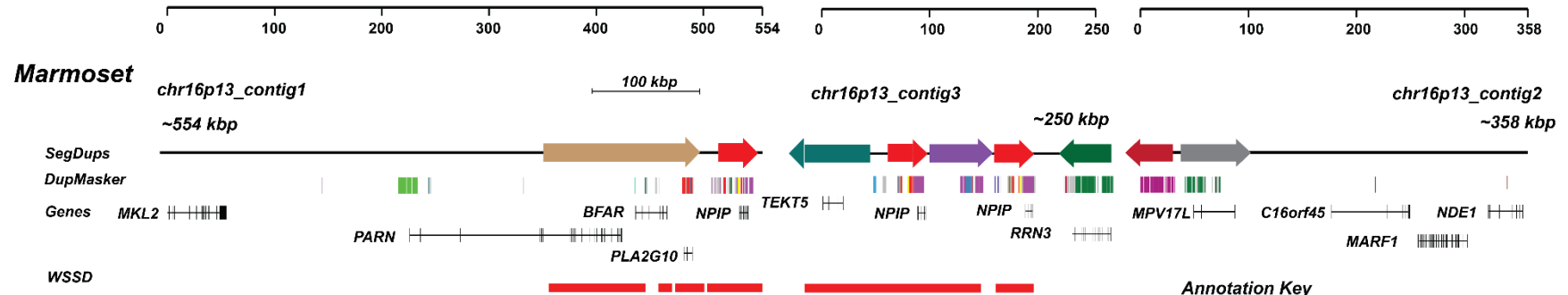
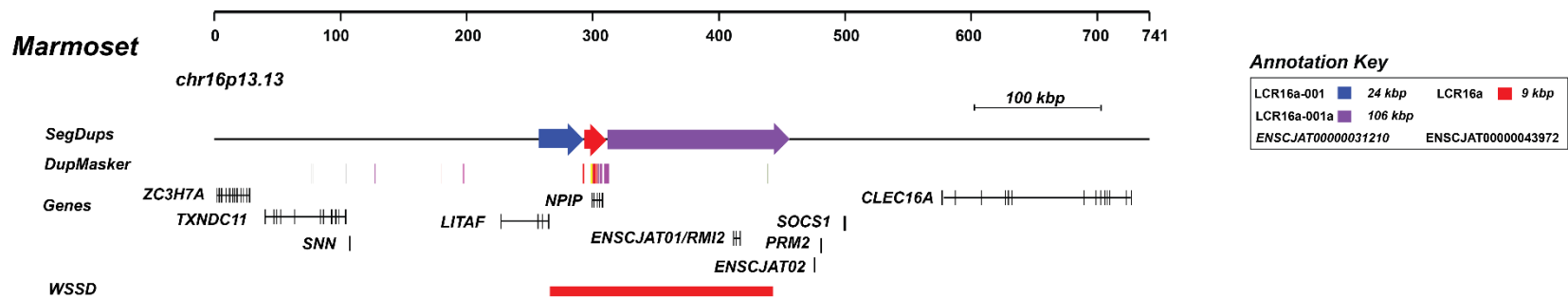
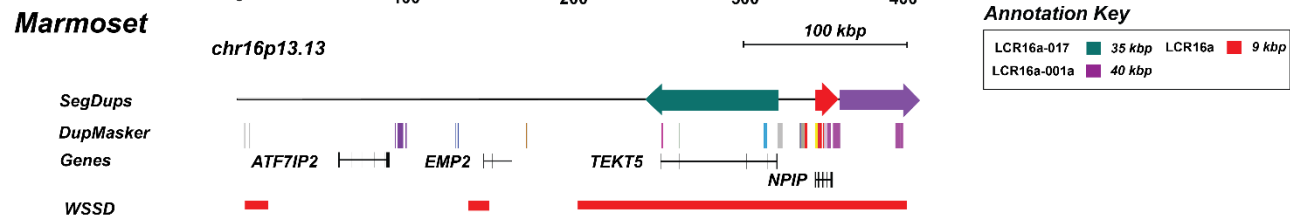
Maximum likelihood analysis using PAML

We performed a maximum likelihood analysis using the branch models of PAML [19] and the entire gene model with the exclusion of exons 1 and 8, which are highly variable among the *NPIP* copies from dog (n=2), macaque (n=1), baboon (n=1), orangutan (n=5), gorilla (n=5), chimpanzee (n=5), and human (n=7) (Additional file 1: Fig. S11-S12 and S18-S19). Note that we excluded the marmoset paralogs in this analysis due to the dramatic restructuring of the New World gene model. The input phylogeny is inferred using these 26 *NPIP* paralogs and the maximum likelihood-based method in IQ-TREE [48]. We applied the codon substitution models (codeml) implemented in PAML to test positive selection and constructed a null model under neutrality ($w=1$ on all branches; H0) for the entire inferred phylogeny (Additional file 1: Fig. S11). The significance of each test was evaluated using the likelihood ratio chi square test.

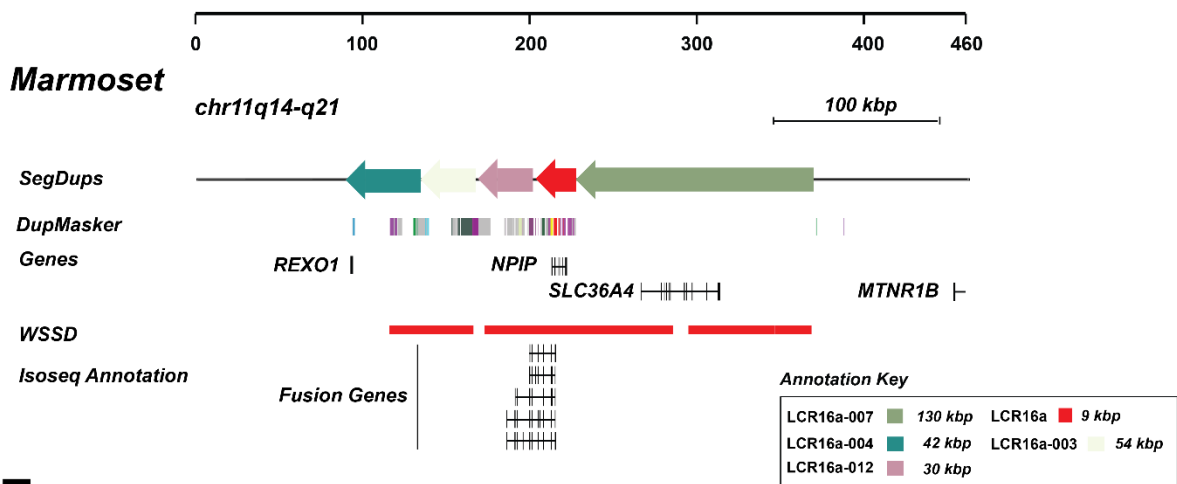
Selection analysis in canine lineages: To search for evidence for selection in the canine lineage, we set $w=1$ for all branches except the canine lineages (CAN1 and CAN2; Additional file 1: Fig. S11) and estimate w_0 (H1).

Selection analysis in primate lineages: We first built a model (H2) to infer selection on all ape lineages (a single free parameter $w:=w_1=w_2$, and set $w_0=1$; Additional file 1: Fig. S11), then considered two additional models: selection on the African ape sequences (H3: a free parameter w_2 , and set $w_0=w_1=1$) or only within the Old World monkey (OWM) and orangutan (H4: a free parameter w_1 , and set $w_0=w_2=1$). We also tested whether there was evidence for selection on either subfamily A (H5_A; Additional file 1: Fig. S18) or the subfamily B (H5_B; Additional file 1: Fig. S18) in the African ape lineages.

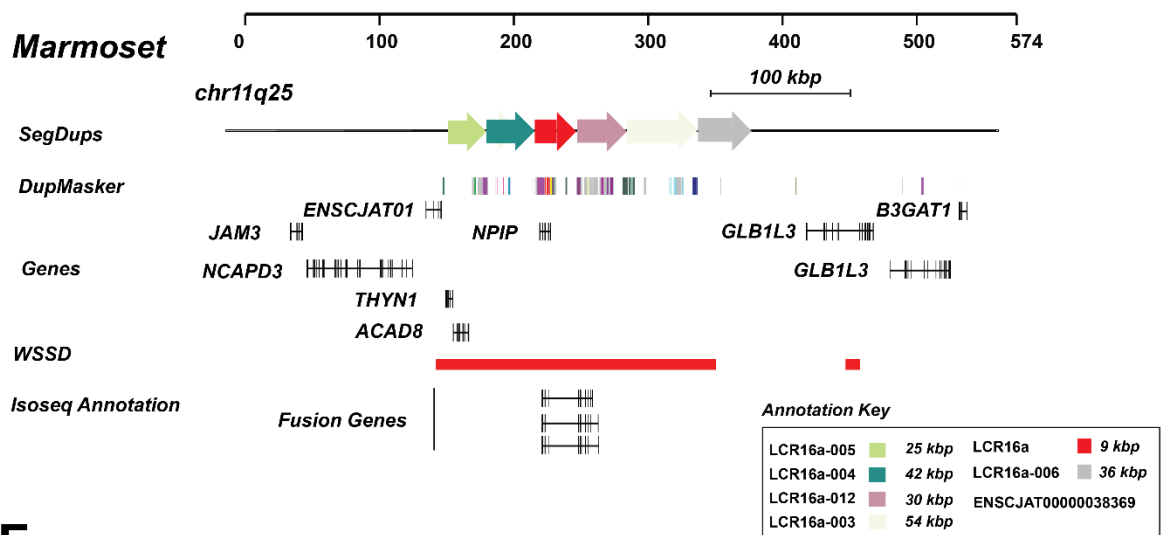
Positive selection sites model: To identify sites that are likely the targets of positive selection, we applied the branch-site test of positive selection implemented in PAML (model=2, NSsites=2, fix_omega=0) to the same 26 sequences from exons 2 to 7 with the same configuration in the models H3 (Additional file 1: Fig S12) and H5_B (Additional file 1: Fig. S19).

A**B****C**

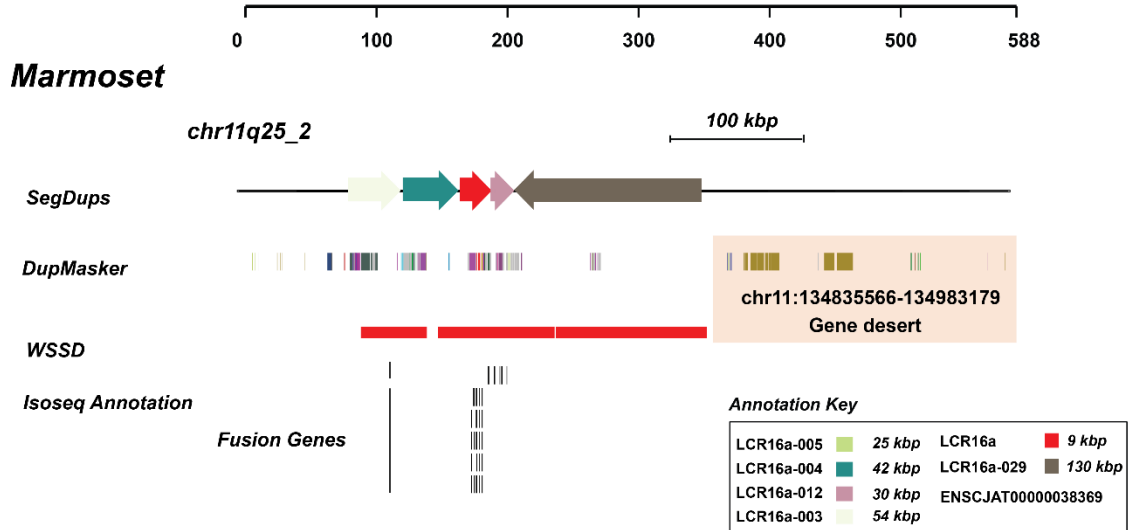
D



E



F



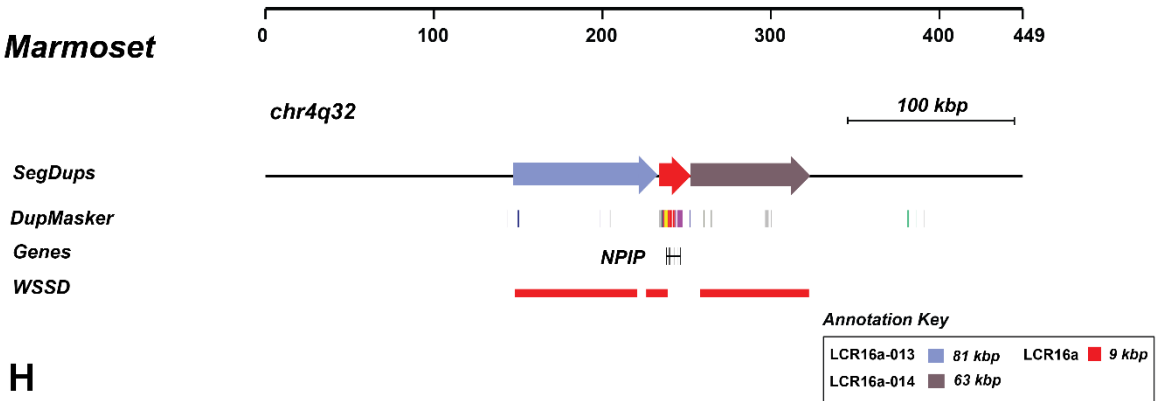
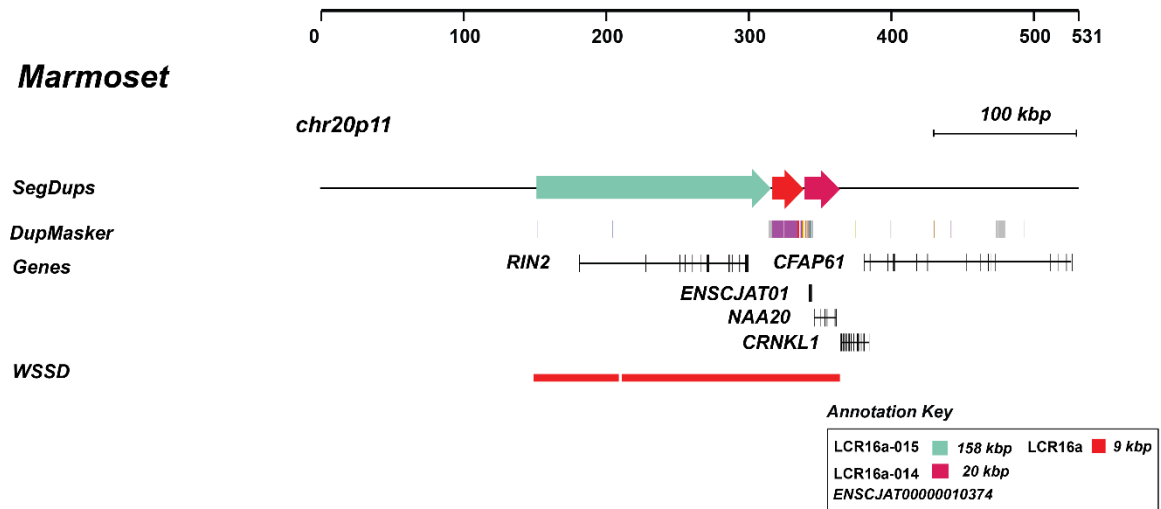
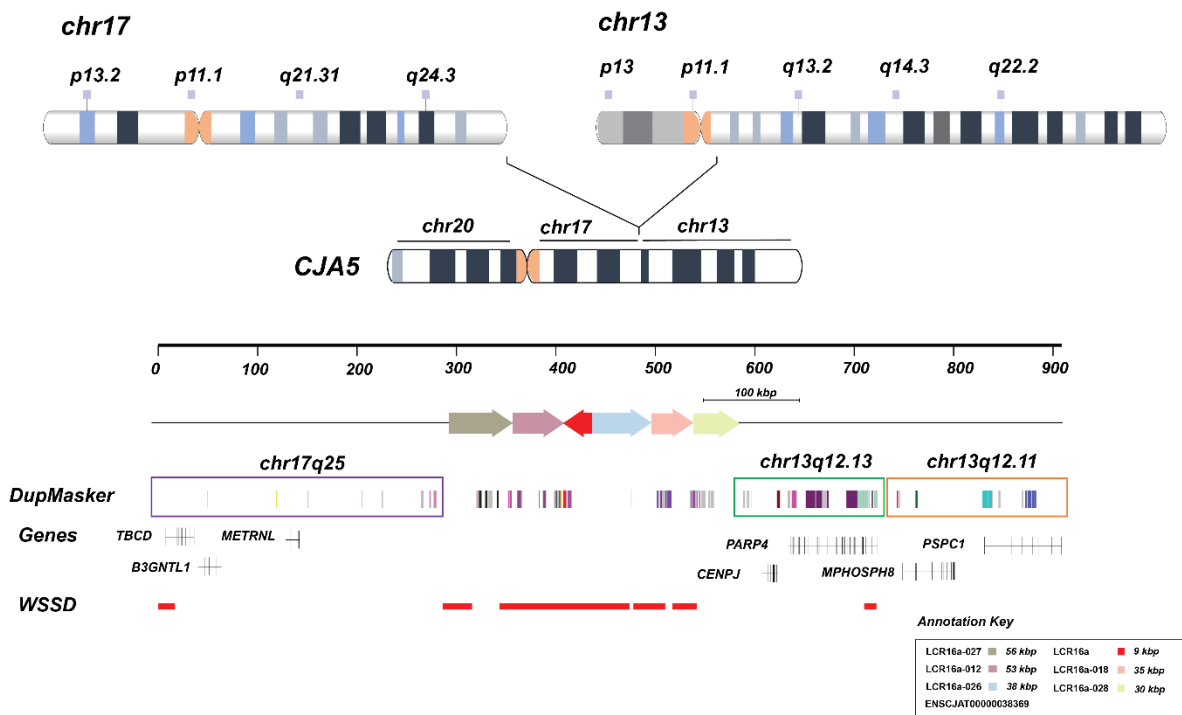
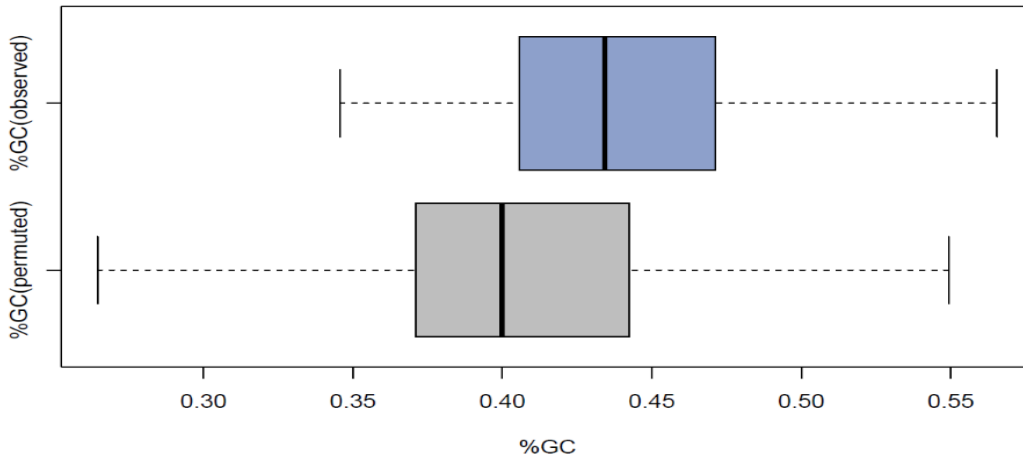
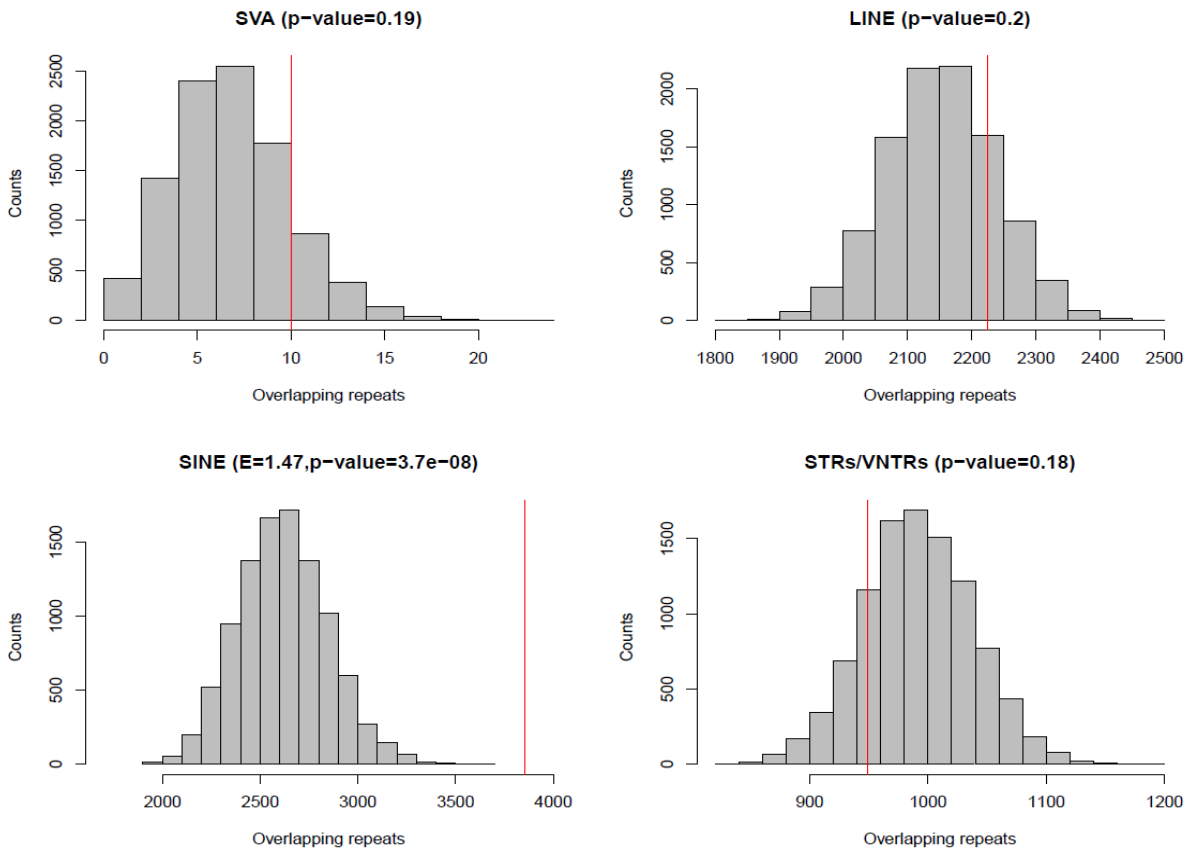
G**Marmoset****H****Marmoset****I**

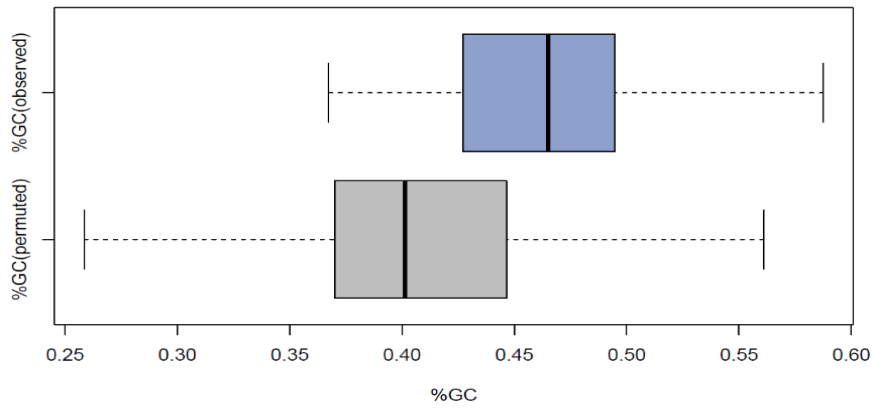
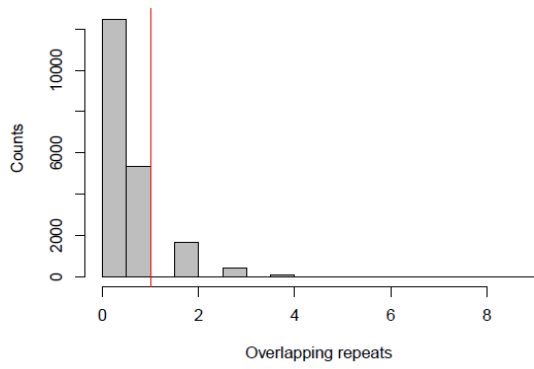
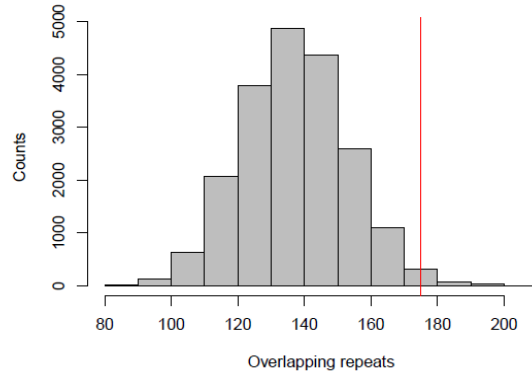
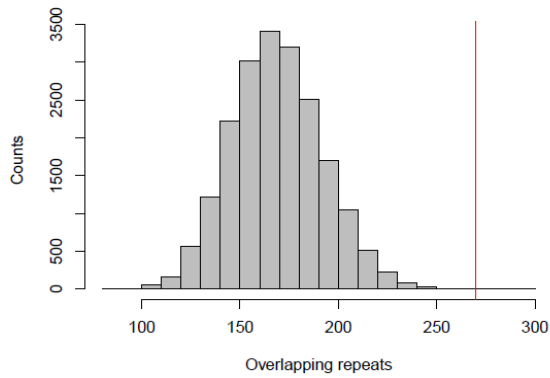
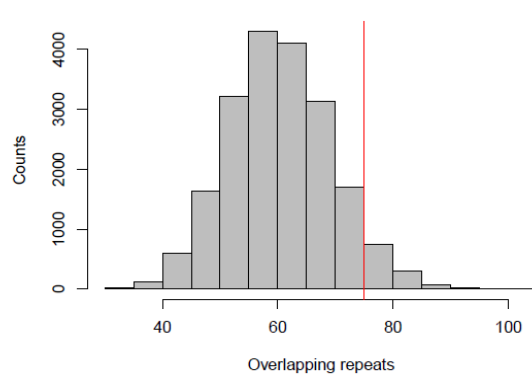
Figure S1: Eleven LCR16a marmoset insertions anchored to the GRCh38 reference genome. The schematics represents the organization nine LCR16a marmoset insertions based on sequence and assembly of large-insert clones. Supercontigs (396–920 kbp) were created to span duplicated sequences and are anchored in unique sequence (>20 kbp). Segmental duplication (SD) organization is depicted using coloured arrows; regions of increased read depth (WSSD) identify duplicated regions (red bars). Gene models are predicted using GMAP based on Iso-Seq and Ensembl transcript data. Additional annotations include DupMasker. The chromosomal map location of the insertions (GRCh38 coordinates) map to **A)** 16p13.12, **B)** 16p13.13, **C)** 16p13.13, **D)** 11q14, **E)** 11q25, **F)** 11q25, **G)** 4q32, **H)** 20p11, and **I)** 17q25_13q14.

A GC content (KS p-value = $1.6e-05$)



B



C**GC content (KS p-value = 2e-09)****D****SVA (p-value=0.27)****LINE (E=1.27, p-value=0.009)****SINE (E=1.6, p-value=7.9e-06)****STRs/VNTRs (p-value=0.06)**

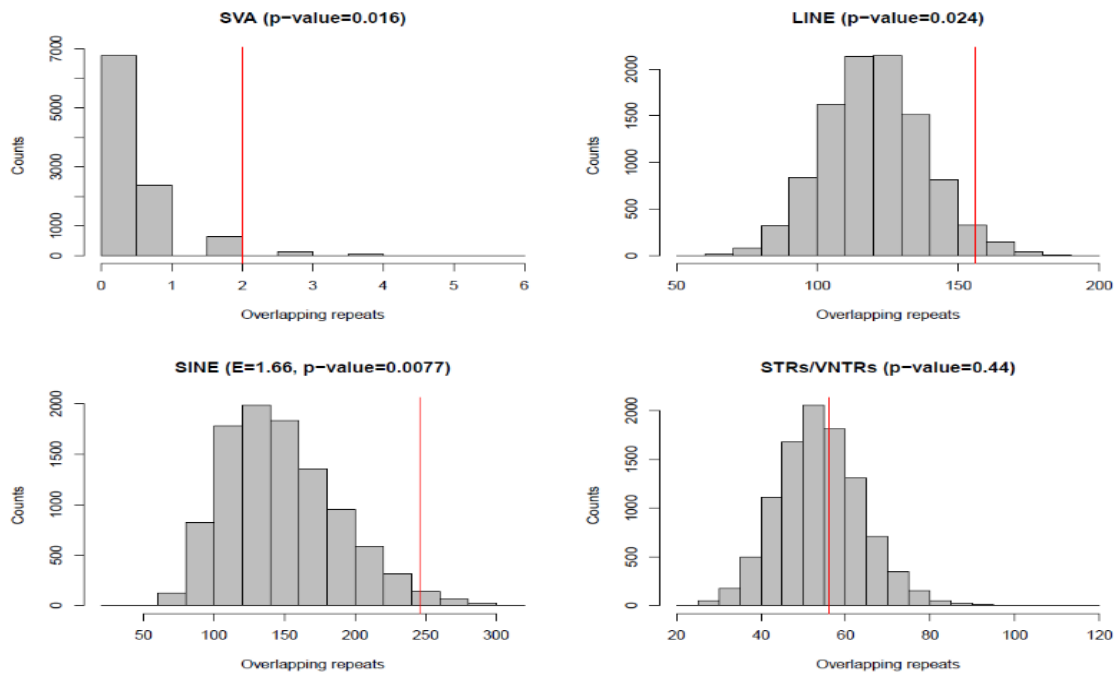
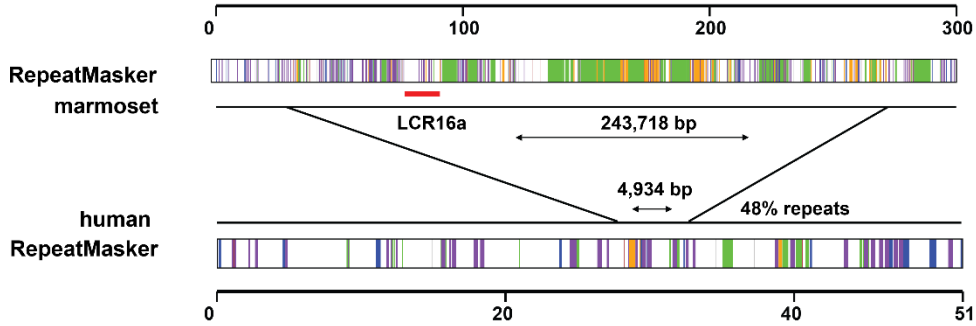
E

Figure S2: Sequence properties and enrichment analysis of donor and acceptor regions in association with LCR16a. **A)** GC content is elevated among donor duplication sites (n=63) compared to the null distribution (10,000 permutations). **B)** Analysis of common repeat content in donor duplication sites shows a relative enrichment of SINE elements compared to the null distribution. No evidence of SVA, LINE, or VNTR enrichment is identified at these sites. **C)** GC content is elevated at acceptor duplication sites (n=27) compared to the null distribution (10,000 permutations). **D)** Acceptor sites show a relative enrichment of SINE and LINE elements compared to the null distribution. **E)** LCR16a integration sites (n=13) show a relative enrichment of SINE elements compared to the null distribution, consistent with sequence resolved breakpoints among primates.

B

chr11:134089171-134426493



centromeric breakpoint

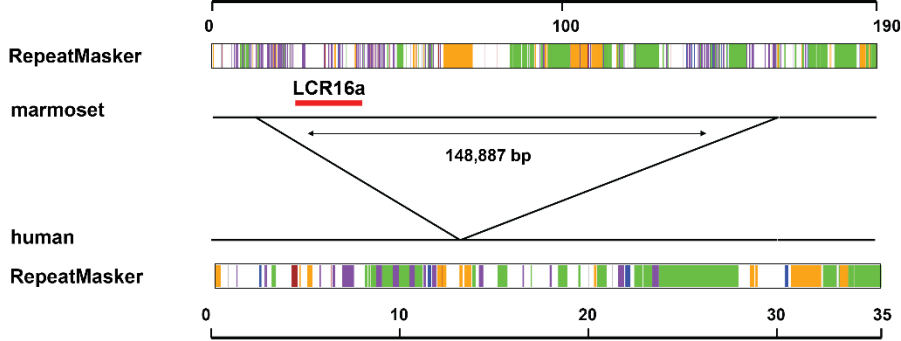
telomeric breakpoint

marmoset	AAGACCTCCCTACTAGTCAAGCTGCTTCCTGATCTTCGAGTCTGAGAGGCTTCTATTTT	-----GCACAGAGC--ACCTGGGAGCTGGGGCTCTTCT-----	
human	AAGACCTTCTACTAGTCAAGCTGCTTCCTGCAATTCGAGCCTGAGAGGCTTCTATTTT	CCCCAGGAGTTGGAGACCACTGGGCAACAGCTGTGAGACCCCATCTCTAGAAAAATAC	AluJb
	*****	* * * *	
marmoset	GACTCAAAGCAAAACAATGCAATGTTAATGTCTTCCTATGTTGTTCCCTCATAAGAT	-GGCATTTCAGTGAGCTATGATTAGGCAACAGAGTGGGACCCCATCTCTAAAAAAGA	
human	GACTCAAAGCAAAACAATGCTATGTTAATGTCTTCCTATGTTGTTCCCTCATAAGAT	AAAAATTAGCCTGGCATGATGGTGTGACCTCTGCTCCCAAGCGCTGGGATTACAGG	AluJb
	*****	* * * *	
marmoset	TTTGATGAAATGATAGTATTAATAAGTGCATTTCTCTCACACCTGTAATCCAGCACTTT	TTTTGAAAAAACATAAAATCAATTAAAGTAGACTTTTAAAGCCTCTAATATTGGA	AluSx
human	TTTGGTGAATCATATTAATAAGTGCATTTCTCAGTATCAGTCTCAGTTATTAGAGATTT	TGTGAGCCACCAGCCAGCCGAAAAACAGAGAAATTTTAAAGCCTCTAATATTGGA	
	*****	* * * *	
marmoset	GGGAGCCGAGGCAGGTGGATCACGAGGTCAGGAGATCAAGACCAGCCTGGCCAAACATGG	TATAGTTTCCATGTTATGTTAGTTTATTTATAGACCTGATGGTTAGAAGCATTTCAC	AluSx
human	CAGTTCGCCAGTACCCCTGGCAAGGATTTGACGAACGTAGCAAGAAGTCGTGAACCTTTG	TATAGTTTCCATGTTCCCATAGTTTATTTATAGACCTGATGGTTAGAAGCATTTCAC	
	* * * *	*****	
marmoset	TGAAGC-----TCGATTTCTACAAAAAATACAAAAATTAGCCGGGTCTGGTGG	TATCCCTTTTGCTAGTCAGCTGGTTATACCTTTAATTTGAAATGCAATTTGTACTTTTCC	AluSx
human	AGCTAATTTATATCATTTCTCATTTCTCTTTAATAAATAAAA-----GCAAGAAAACC	TTCCGCTTTGCTAGCCAGCTGGTTATACCTTTAATTTAAGTCAATTTGCTGCTTTTCC	
	* * *	* * * *	
marmoset	TGGGCATCTGTAATCC	TAAAG	AluSx
human	AGAGCAA-----	TGAAG	
	* * *	* * *	



C

chr11:92975146-93322890

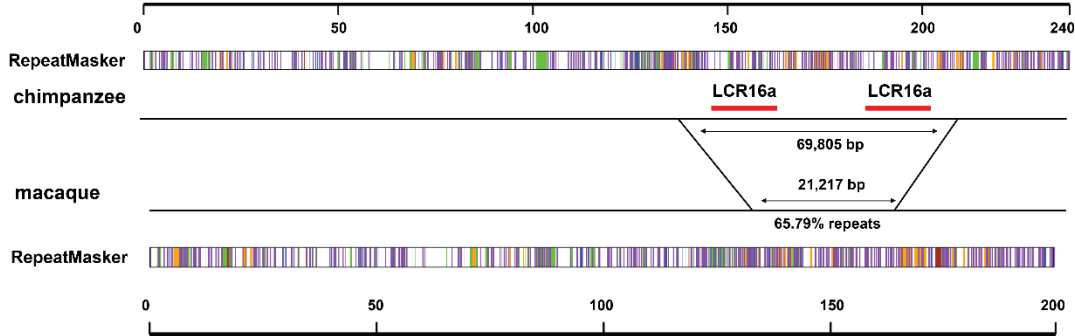


telomeric breakpoint		centromeric breakpoint	
marmoset	GAGAGCAAAATGTGGCAGAGACTAGTGTTCAGCCAGAAGATGCAGACAGGGCCCCAG	marmoset	-----AGAGATCACTCAATGCATTCGGCCCTGGTGACAGAGTGAGACTGTCTCAA
human	--GAGCAAAATGTGTAGAGACTAGGTTTTCAGCTGGCAGATGCAGCCAGGGGACCAAG	human	AGCAAAATGTGTAGAGACTAGGTTTTCAGCTGG-----CAGATGCAGGGGACCAAA
	*****		*****
marmoset	CCAGACAAGGGTTAATCCAGGATCTGACAGAGCCAGAGAGGGTCTGAGTGACACTAAC	marmoset	AAAAAAAAAAAAAGATCAAAATTAATGCAAAATCCCAATGGAGCCAGAATA-----
human	CCAGCCAAGGGTCAATCCAGGATCTGACAGAGCCAGAGAGGGTCTGAGTGACACTAGC	human	GCCAGCCAAGGGTCAATCCAGGATCTGAC-AGAGCCAGCAGGGTCTGAGTGACACTA
	*****		*****
marmoset	ACAGACTTTACAGCACTAGAGAGGATGA-GGGCACTTCTCGAAGCCAGAGACAGAA	marmoset	---AAAATTTAAATCGAGGTGGGTTCAGGAATTCGATTTTCCTTTGGCTTGGCTCC
human	ACAGACTTTACAGCACTAGAGAGGATGA-GGGCACTTCTCGAAGCCAGAGACAGAA	human	GCACAGACTTTACAGCAGATAGAGGACAA-GAGCCAC-----TTCTCTGGAGAC
	*****		*****
marmoset	AGAACACATTAGAGACCCCACTCTCCCTGTGCTCATGTCTGAAATGACACTTC	marmoset	AAATGCTCTCAGAGCAGTGAATGATCCCACTTTGAGATAAAATCTGATCATTTGG
human	AGGACACATTAGAGACCCCACTCTCCCTGTGCTCATGTCTGAAATGACACTTC	human	CAGAGACAGAGAGGACACTTAGAGAGACCCCACTCTCCCTGTGCTCATTTGCT
	*****		*****
marmoset	CTCCCTGCAATGTCTACTTTCCCTCAT-AAACCAAGATGCTAGTTTGGAAATCAAACA	marmoset	TCACCAATGGCTTTTCAGATTCATTTAAGTTTTATAAGAGTGCATTAATAATTTT
human	CTCCCTGCAACCCCTACTATCCCCCATAAACCCAGATGCTAGTTTCGGAACCAAAACA	human	TGGAATGAC-----ACCTCTCCCTGCAACCCCTACTATCCCCCATAAACCCAGAT
	*****		*****
	LTR-GYP		LTR-GYP
marmoset	GAATTCACGCTGTAAATCCAGCACTTTGGGAGCCAGACATGTGATCGGAGGTCAGGA	marmoset	ATAGGCCAGTCCCTCAAAACAGAGGATGGGAGGCCAAGCTTAAAGCTGAACATTTGC
human	GAATTCAGGGAGGGCAACTGAAAGCTGAAC-----ATTTCGTGGGATCTGAG	human	GCTAGTTTCGGAACCAAAACAGAGGAGGGGAGGCCAAGCTTAAAGCTGAACATTTGC
	*****		*****
	LTR-GYP		LTR-GYP
marmoset	GTTCAAGCCAGACTACCAACATAGTAAACCCCACTCTACTAAAAATCAAAAAATTAG	marmoset	TGGGGAATTCGAACTTGCCTGATGAGAATGTATAAAATTTGGATCAAAATGGAGTTT
human	CATT--GCCGATTGAGAAATGTATAAAATATTG-----ATTGAATGG-----AGTTTA	human	TGGGGAATTCGAGCATTGCCTGATTGAGAAATGTATAAAATTTGGATGAAATGGAGTTT
	*****		*****
	LTR-GYP		LTR-GYP
marmoset	CTGGCATGTGGCAGAGCCCTAAGTCCAGTACTCTGGAGGCTGAGCAGGGAATTC	marmoset	AAAGACATTGGATATTAGTAAATTTCCCACTAGCAGCTGAGGGTTTGTGAGGCAAGAG
human	ATGACATTGG--ACATTAATTAATTTCCCACTAGCAGTTGAGGGTTCAAGAGGCAAGA	human	AATGACATTGGACATTAGTAAATTTCCCACTAGCAGTTGAGGGTTTCAAGAGGCAAGAG
	*****		*****
	LTR-GYP		LTR-GYP
marmoset	ACTTGAACCCAGAGGTTCAAGTTGCAGTGAAGTAAATTTTCGCATGCATCCA	marmoset	GATAA---CAATATTAAGATTTTCCTTTTTCCTGTCTTAAATAGGGTTTAAATATG
human	GTATAACAACAAATATTAACATTTCTTTTCTGTGTTCTTAACGAGAGTTTAAATGTT	human	TATAACAACAAATATTAAGATTTTCTTTTTCCTGTCTTAAATAGGGTTTAAATATG
	*****		*****
	LTR-GYP		LTR-GYP
marmoset	GCTTGGCAGAGAGTAAAGACTCTTCGGTCCCTCTCTC--TCTCTCTCTCAGCATAT	marmoset	CCACGTACACTGTATTATCCACTCTCAGGCTGCTATAAAGAACTACAGAGACTGG
human	GCCACATACACACTGTATTATCCACTCTCAGGCTGCTATAAAGAACTATCTGAGACTGG	human	CCACATACACTGTATTATCCACTCTCAGGCTGCTATAAAGAACTATCTGAGACTGG
	*****		*****
	LTR/ERVL		LTR/ERVL
marmoset	A-----T--ATAATA	marmoset	GAAATTTAGAAAAAAGAG
human	GAAATTTAGAAAAAAGATTTA	human	GAATTTAGAAAAA-----
	*****		*****
	LTR/ERVL		LTR/ERVL

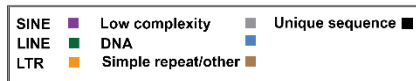


G

chr16:14662753-15595746



telomeric breakpoint		centromeric breakpoint	
chimpanzee	-----TTCTAAGTCTCATTTTCTCTTCAGTTCOCATTAAGTAAATCCATTTGAGCT	chimpanzee	-----GTGGCTGAGACAGGAGAATTGCTTGAACCTGGG-AGGAATAGGTTGTAGTG
macaque	AAGCCTCTAAGTCTCATTTTCTCTTCAGTTCOCATTAAGTAAATCCATTTGAGCT	macaque	TGAGCCACCCTGGCTGGG-CTTAGAGTCTCTTTTAAAGCTTAATGCTGGTCAATGGTGTCT
chimpanzee	GCTCCACAAAGAACACCATTGGATATTAGAGGATATTGCTCTAGAGATTGCTACCAATC	chimpanzee	AGCCAAATAATATAAAGTATTAGAGTCTTAACAGAGAAAGTTCCACATGATCACCTT
macaque	GTTCCACAAAGAACACCAGGATATTAGAGGATATTGCTCTACAGATTGCTACCAATC	macaque	TGATTCTAGAGGGAAGAA-----GGTATAAATGAGGCACTTTG-----ACACC
chimpanzee	CCTGGCTCTCTGGGTTCCTAGGAGACTGACTTGTGCAGAGGTGCCAGTAARITCTCA	chimpanzee	TTAGCTTTGAATGAATGCA---GAGCAATTTGCACAGTTACTTGTAAATAAATAATGAG
macaque	CCTGGCTCTCTGGGTTCCTAGGAGACTGACTTGTGCAGAGGTGCCAGTAARITCTCA	macaque	TCCCTCTCCCAACAGGCTGAGCTGGTTTTTTCAGTTTCTTTGGAAAGTCTTGGCC-
chimpanzee	GTTTCATCCCTTTGGCTTTGCCAGATTAGCAAATAATTTTTTCTTACTTTTTTTAA	chimpanzee	ATCATTCACAAATTTATCAGTCTTTTTTGTGTACAAAACATAATAAGTATTCTCTT
macaque	GTTTCATCCCTTTGGCTTTGCCAGATTAGCAAACGATTTTTTATTATGTTTTTAA	macaque	AACAGGA-----AGGCTATTAGTTC-GAGTGG--GTGGCTAGAAATTTTATT
chimpanzee	TTTTTAATTTTTTGTTTTTTTTTTCCAGAGTGGAGTTTGCCTCTTGTGGCCAGGCT	chimpanzee	TTGCTGTATTGTGACTGGTTTGGCTGGAAGG-----GAATTA-----GGCACT
macaque	TTTTTTGTTTTT-----TATTTTTTTTCCAGACTGAGTTTGCCTCTTGTGGCCAGGTT	macaque	TGGCTTAACTCAAACATACACGCAAAAAGTAAGGAGGGAATTTGCTACTGTCTAGT
chimpanzee	GGAGTCTCTCTGCTTGGACTCCAGAGTAG-----AluSp	chimpanzee	CGAGAGTTTGTGTGTTAAAGTTTTCTGGCCAGGCAGGTTGCTCATGCCGTGAATCC
macaque	GGAGTCTAATGGCATGATCTCAGCTCAGTCACTGCAACTCTGCTCCCGGGTTCAAGCGATT	macaque	TAAACTGGTTAATGCAGAAAGGAGTCTTGGCCGGCGCGTGGCTCAAGCTGTAAATCC
chimpanzee	-----CTGGACTACAGGTTGCACCACCTGCCAGTAAT	chimpanzee	CAGCACATTCGGAGGCTGAGGCGAGGGTATCCTCTGAGGTCAGGAGTCTAGACCAGCT
macaque	TCTGCTCAGCCTCTGCATAGCTGGGATACAGGCGCTGCCACCAGCCAGCTAAT	macaque	CAGCACTTTGGAGGCGGAGATGGCGGATC--AGAGGTCAGGAGTCTGATCATCCT
chimpanzee	TAAACAATTTTTTTTTTTTTTTAGATGAGAGTGTATCAGTCCGTTCTTGGATTG	chimpanzee	GACCAACATGGAGAAACCCCATCTCTACTAAAATAACAAAATAGCTGGCGTGGTGGT
macaque	TT-----TGTATTTTAGTAGAGATGGGTTTACCATGT-----TGGCCAGG	macaque	GGCTAACACGGTGAACCCCGTCTCTACTAGAAATAACAAAATAGCCGCGAGGTTGG
chimpanzee	CTATAAAGAACTCCCTGAGACTGGGTAATTTATAAAGAAAGAGGTTAATTGACTACA	chimpanzee	ACGTGCTGTATCCAGCTACTCGGGAGGCTGAGGCGAGGAAATGCTTGAACCTGGGA
macaque	CTGGTCTCAACTCTGT-----AGTCAAGTGAATGCTGCTCAGC	macaque	CAGCGCTGTATCCAGCTACTCGGGAGGCTGAGGCGGAGAAATGGCTGAACCCGGGA
chimpanzee	GTTCCACAGGCTCTCCAGSAGACARCTTGGGAGGCT-----CAGSAACTTTGA	chimpanzee	GGCGGTGTTGGTGCAGCTGAGATGTGCCATTGCATCCAGCCTGGGCATAAGAGTG
macaque	CTCCAA-----GTCCAGGATTAACAGGCTGAGCCACTGTGCTGCTAATTGTTTAA	macaque	GGCGGACTTGGCAGTGCAGTCCGCTGAGTCCGCTGCTGCTCCTCCAGCTGGCTCAGAGCA
chimpanzee	ATCAT-----GGCGAAGGCAAGGGGAAGCAGCACATCTTACATAGCTGGAGCAGGAG	chimpanzee	AAACTGTCTCAAATAAATAAATAA--- 600 AluSp
macaque	AATTTTTTATCTGTTTTGATCAGTGCAGAAATTAATTTTTT----- 601 AluYh1	macaque	GACT-CCGCTCAAATAAATAAATAA--- 601 AluYh1



H

chr2:100489259-100624778

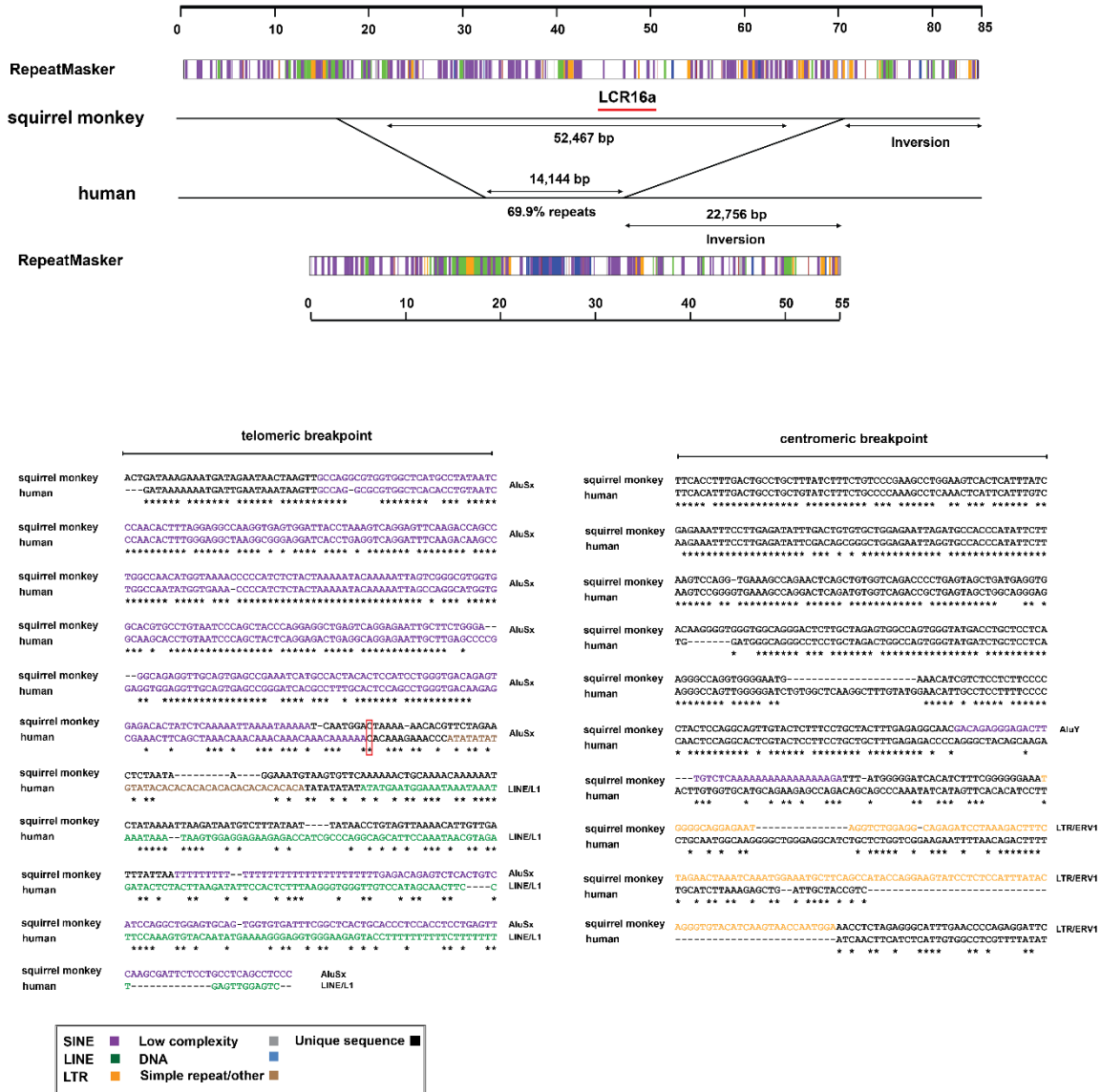
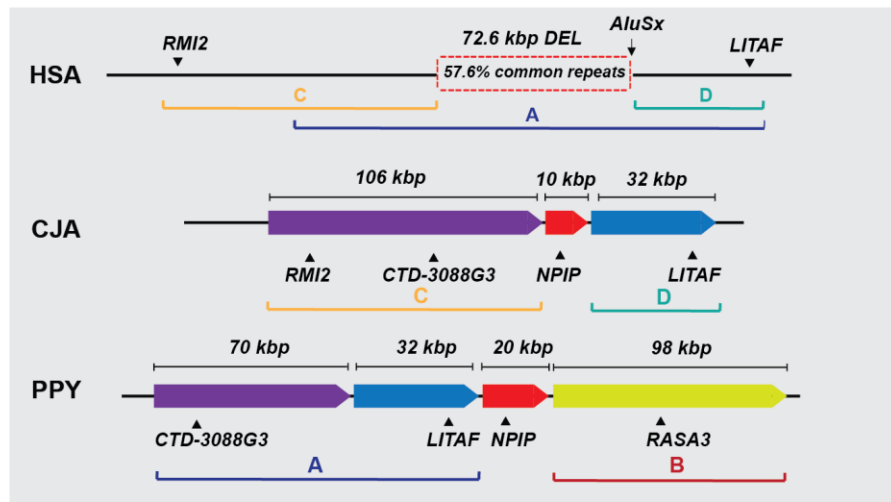
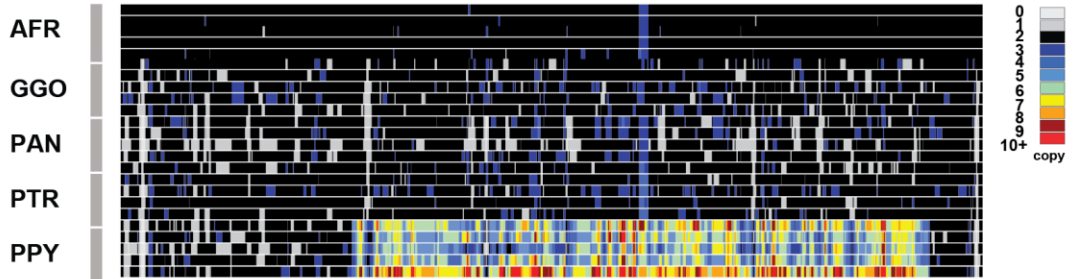
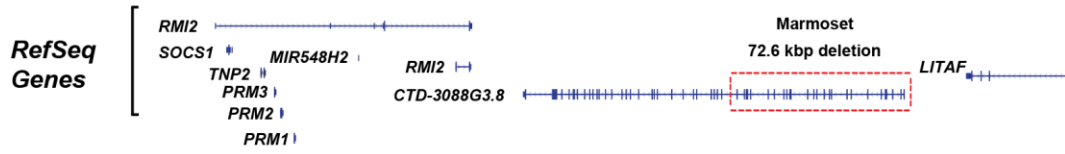
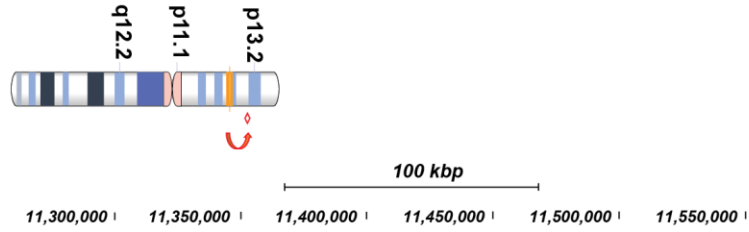


Figure S3: Breakpoint resolution from eight LCR16a insertions. The schematic depicts eight SD insertions corresponding to six marmoset (A-F), one chimpanzee (G), and one squirrel monkey (H). The chromosomal map location of the insertions (GRCH38 coordinates) map to A) 16p13.13, B) 11q25, C) 11q14, D) 17q25_13q14, E) 4q32, F) 20p11, G) 16p13, and H) 2q11. All comparisons were made against GRCh38, except in panel (G) whereby the chimpanzee assembly was compared to a custom sequence contig derived from the macaque CH250 BAC library. Starting from the top, tracks correspond to RepeatMasker annotation, size of the duplication insertion, LCR16a location, and corresponding sequence deleted from the insertion site. Pairwise sequence alignments taken upstream and downstream of the telomeric and centromeric breakpoints are depicted, with common repeats annotated at the breakpoints (red box). 64% of LCR16a-associated insertions map to an AluS element and in all but one case (C) where we find coordinated deletion of repeat-rich sequence (average 67.56%) ranging from 3.4 to 72.6 kbp in length.

A

chr16
recurrent site



A	chr16:11.46 -11.56 Mbp	D	chr16:11.52 -11.55 Mbp
B	chr13:11.40 -11.41 Mbp	-	chr16:11.45-11.52 Mbp
C	chr16:11.31 -11.43 Mbp		

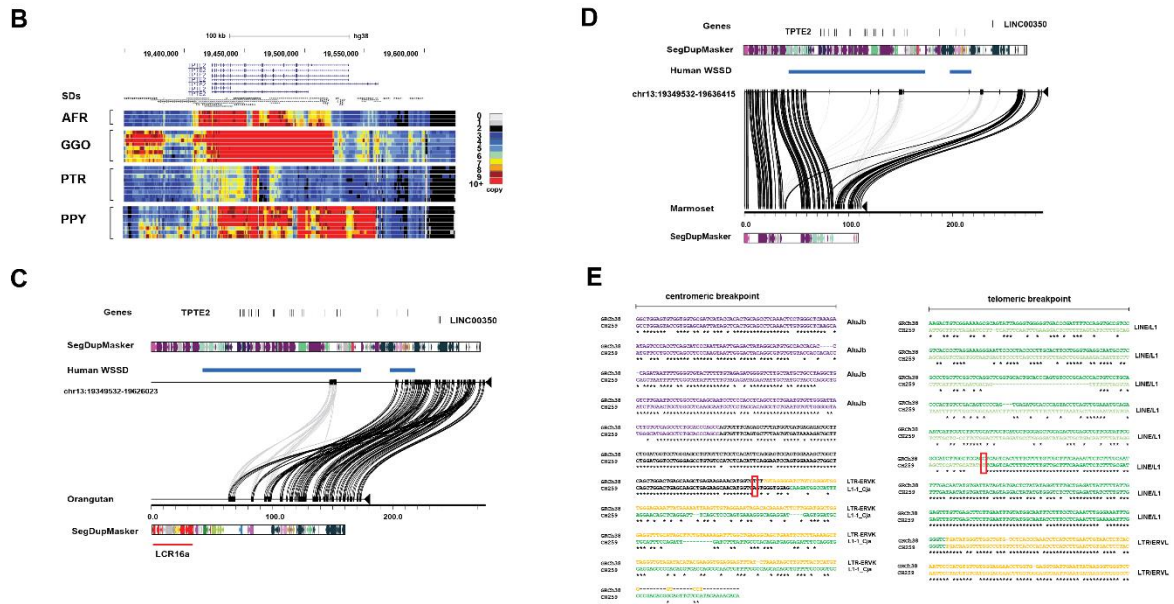


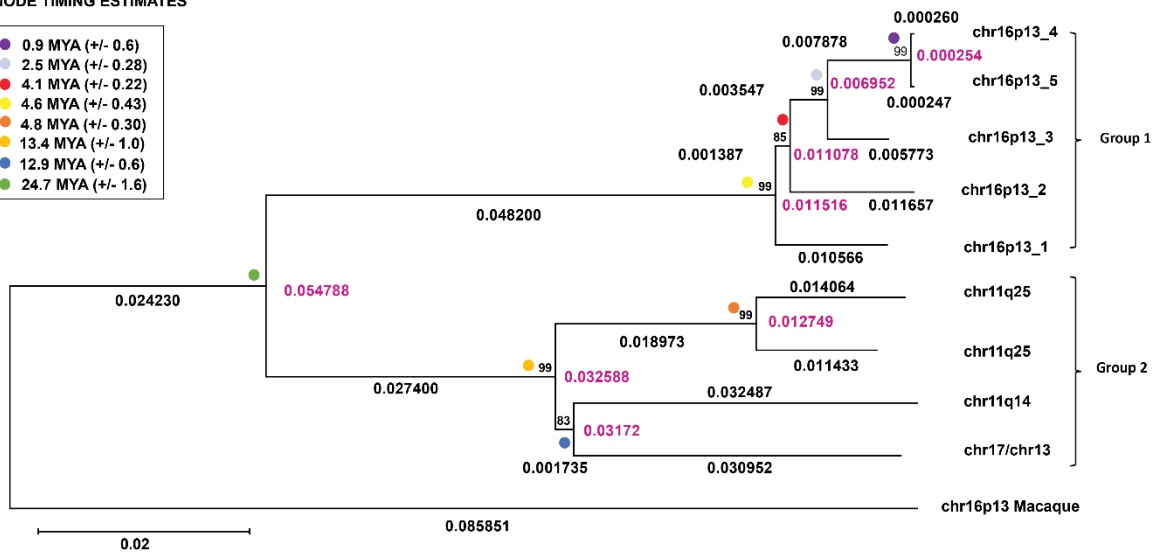
Figure S4: Evidence of recurrent LCR16a duplication during primate evolution. **A)** A UCSC Genome Browser snapshot of a recurrent LCR16a duplication block on chr16p13.13. Annotation tracks include RefSeq gene annotation and copy number (CN) heatmaps (CN index shown) produced from Illumina read-depth profiles from modern humans and nonhuman primates (NHPs). All orangutans carry a large >220 kbp duplication block, which includes the carboxy-terminus of *LITAF* and a long noncoding RNA LOC101927131. At the bottom, a schematic of the locus in human, marmoset, and orangutan is provided based on comparative analysis of large-insert clones. SDs were identified within the contigs of marmoset and orangutan shown as coloured arrows. The marmoset contains a 176.73 kbp duplication block, which includes duplications of *RM12*, *ENST00000598234.5*, and *LITAF*. The orangutan has a larger >220 kbp duplication block that includes *ENST00000598234.5*, *LITAF*, and a gene from 13q24, *RASA3*. All of these duplications in marmoset and orangutan are in association with LCR16a. **B)** A UCSC Genome Browser snapshot of a recurrent LCR16a duplication block on chr13q12.1. CN heatmaps show that this region has been subject to multiple rounds of rearrangement during ape evolution. The CN heatmaps include representations of AFR=African Humans, GGO=gorilla, PTR=chimpanzee, PPY=orangutan. **C)** A Miropeats comparison of the human and orangutan contigs shows the pairwise differences between the orthologous regions. Annotations include whole-genome shotgun sequence detection (WSSD) in human indicating duplicated regions identified by sequence read depth, DupMasker, and exons of genes. Miropeats identifies a bifurcated alignment that includes >80 kbp of orthologous sequence and a >100 kbp of duplication block absent from the human assembly adjacent to LCR16a. **D)** A Miropeats comparison of human and marmoset chr13q12.1 region shows a large ~208 kbp duplication block missing from the marmoset assembly that includes the ~80 kbp duplicate gene *TPTE2*. LCR16a is located ~210 kbp downstream of this site. **E)** Pairwise sequence alignment flanking the breakpoint of the ~80 kbp *TPTE2* deletion in marmoset shows an AluJb element at the centromeric breakpoint and a LINE/L1 element at the telomeric breakpoint.

A

MARMOSET

NODE TIMING ESTIMATES

- 0.9 MYA (+/- 0.6)
- 2.5 MYA (+/- 0.28)
- 4.1 MYA (+/- 0.22)
- 4.6 MYA (+/- 0.43)
- 4.8 MYA (+/- 0.30)
- 13.4 MYA (+/- 1.0)
- 12.9 MYA (+/- 0.6)
- 24.7 MYA (+/- 1.6)

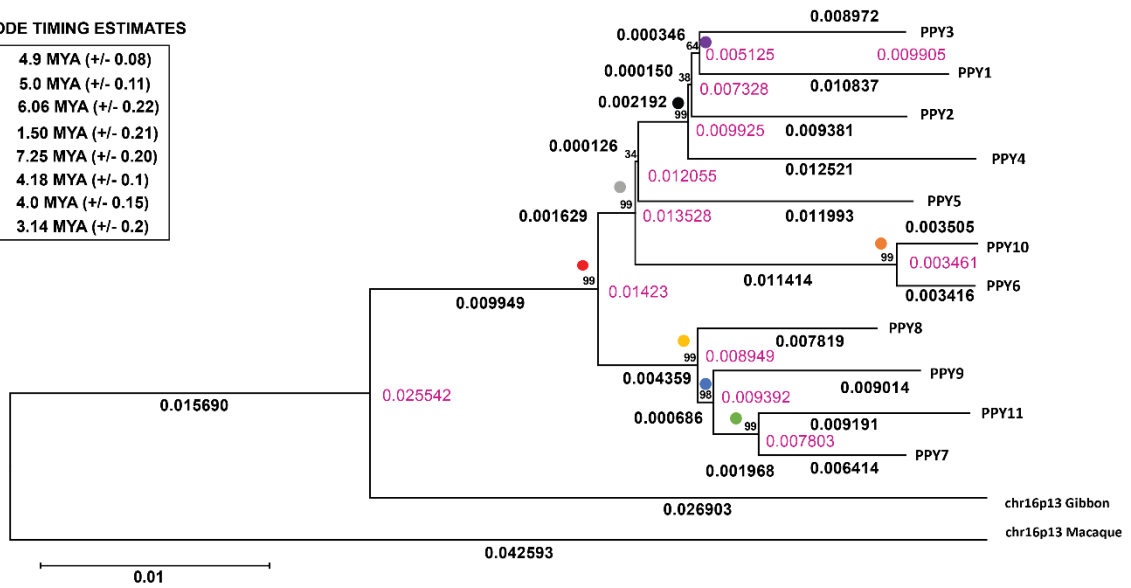


B

ORANGUTAN

NODE TIMING ESTIMATES

- 4.9 MYA (+/- 0.08)
- 5.0 MYA (+/- 0.11)
- 6.06 MYA (+/- 0.22)
- 1.50 MYA (+/- 0.21)
- 7.25 MYA (+/- 0.20)
- 4.18 MYA (+/- 0.1)
- 4.0 MYA (+/- 0.15)
- 3.14 MYA (+/- 0.2)

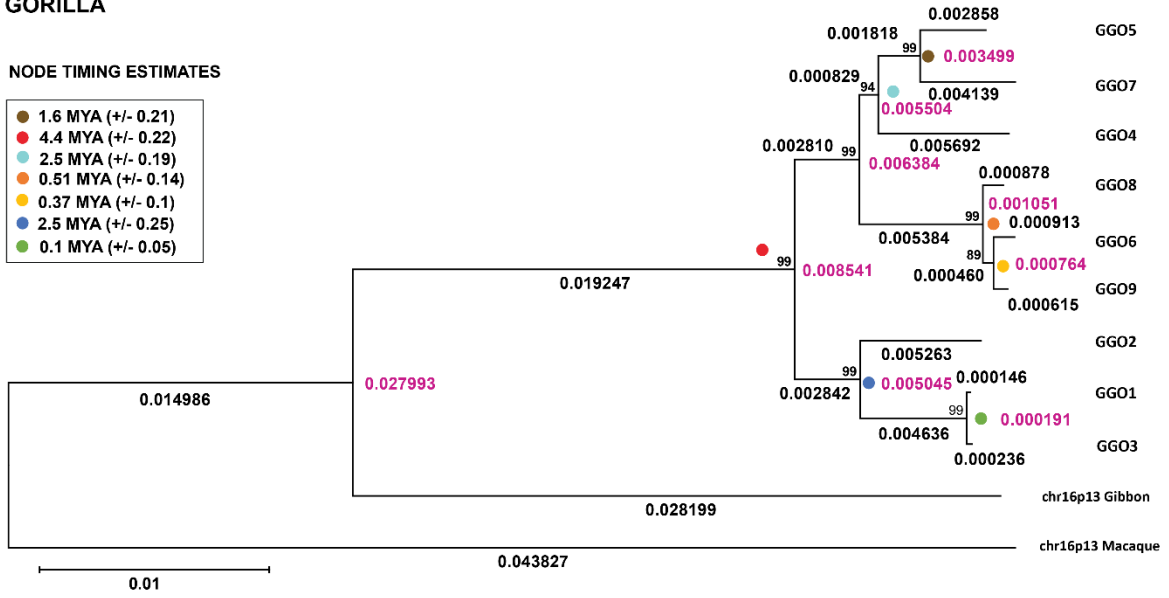


C

GORILLA

NODE TIMING ESTIMATES

- 1.6 MYA (+/- 0.21)
- 4.4 MYA (+/- 0.22)
- 2.5 MYA (+/- 0.19)
- 0.51 MYA (+/- 0.14)
- 0.37 MYA (+/- 0.1)
- 2.5 MYA (+/- 0.25)
- 0.1 MYA (+/- 0.05)

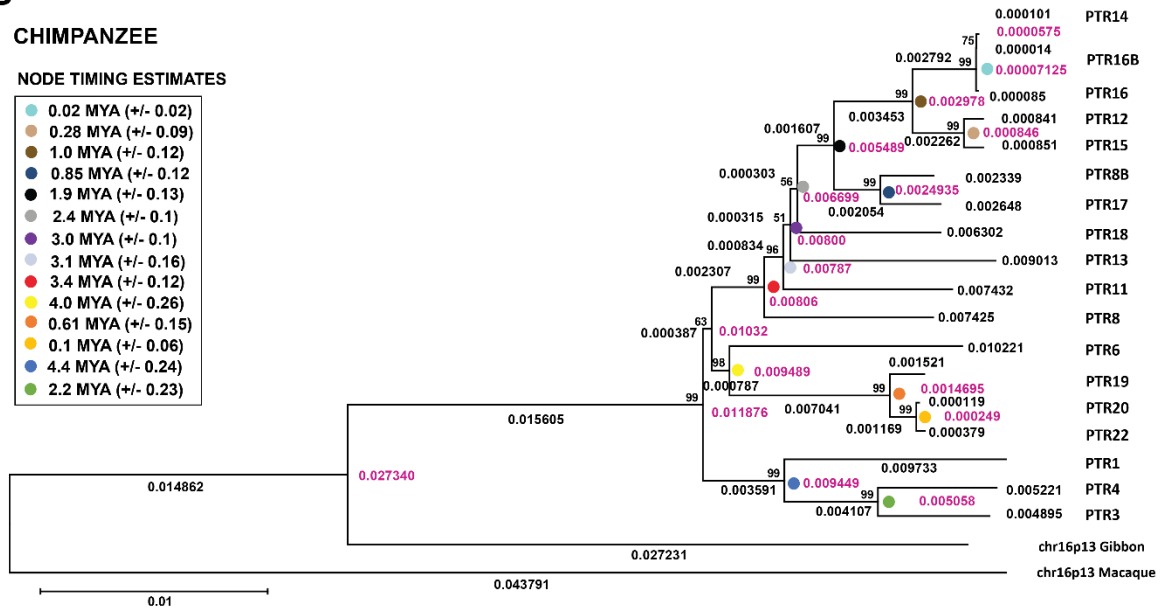


D

CHIMPANZEE

NODE TIMING ESTIMATES

- 0.02 MYA (+/- 0.02)
- 0.28 MYA (+/- 0.09)
- 1.0 MYA (+/- 0.12)
- 0.85 MYA (+/- 0.12)
- 1.9 MYA (+/- 0.13)
- 2.4 MYA (+/- 0.1)
- 3.0 MYA (+/- 0.1)
- 3.1 MYA (+/- 0.16)
- 3.4 MYA (+/- 0.12)
- 4.0 MYA (+/- 0.26)
- 0.61 MYA (+/- 0.15)
- 0.1 MYA (+/- 0.06)
- 4.4 MYA (+/- 0.24)
- 2.2 MYA (+/- 0.23)



E

HUMAN

NODE TIMING ESTIMATES

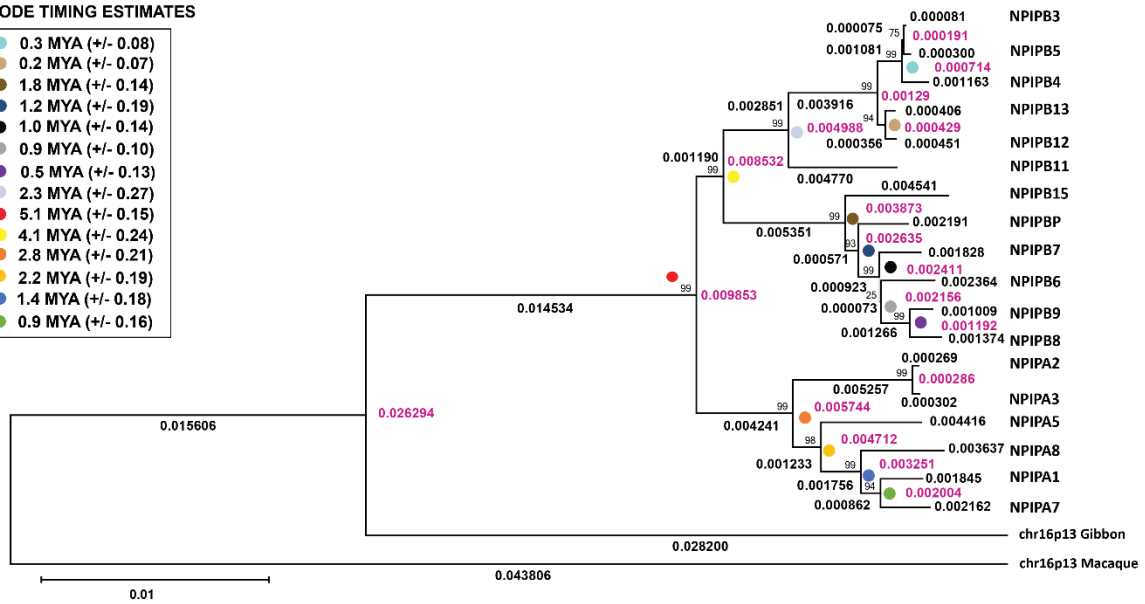
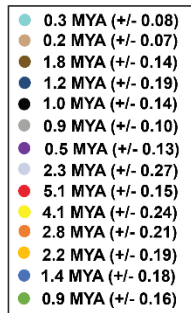
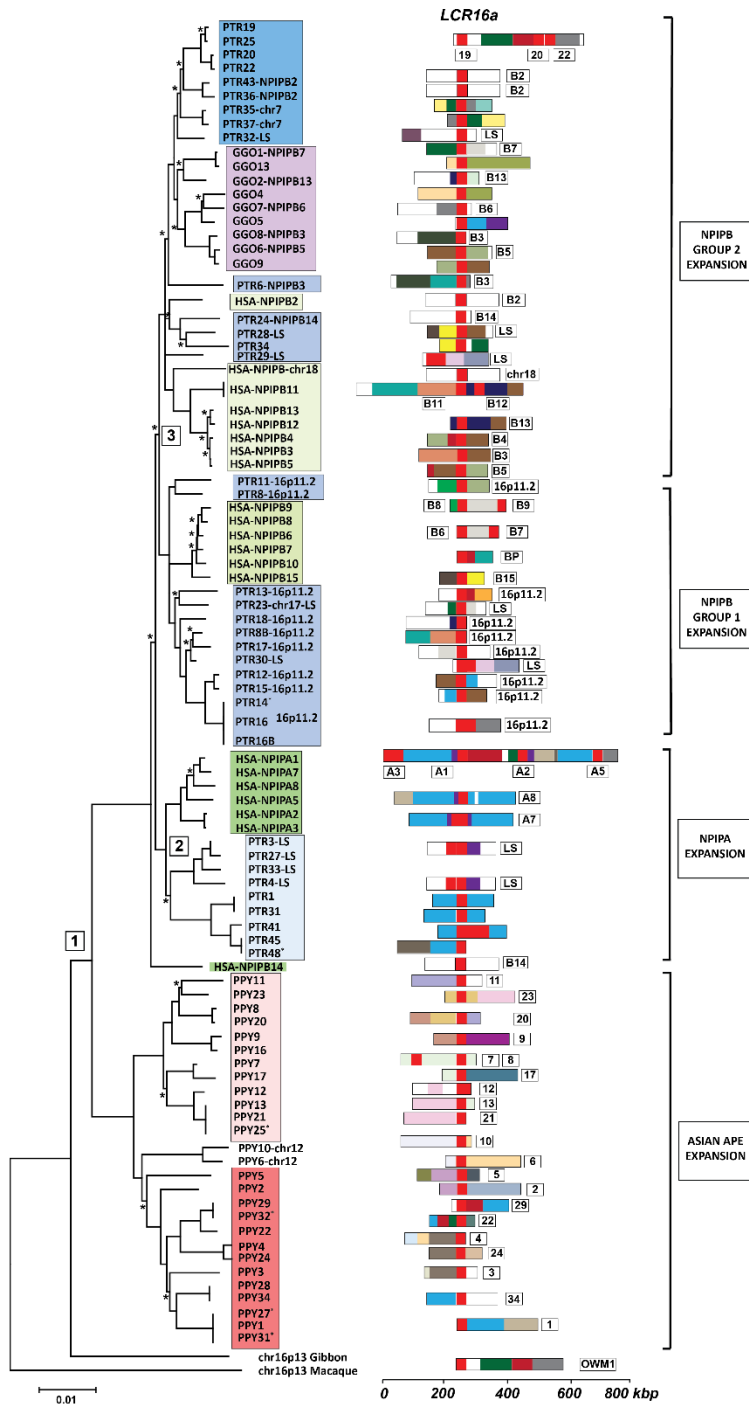


Figure S5: Evolutionary analysis and timing estimates of LCR16a copies in primates. An unrooted neighbour-joining tree was constructed using the MEGA5 complete deletion option based on ~9 kbp of aligned sequence representing LCR16a paralogs identified in marmoset (A), orangutan (B), gorilla (C), chimpanzee (D) and human (E), including the ancestral sequence in macaque as an outgroup. The evolutionary distances were computed using the Kimura 2-parameter method and timing estimates were performed using a divergence time of 25 mya for macaque and 35 mya for macaque/marmoset. Coloured dots represent the timing estimates of each node identified in the phylogeny.



Annotation Key

PPY	GGO	PTR	HSA
chr13	MIPEPF3	NPIPA	NPIPA-exp
chr16	ZDHHC20	NPIPB-exp2	NPIPB-exp1
RASA3	PARP4	INTS4P2	RRN3
NOMO	MIPEP	PDXDC1	PDXDC1
LITAF	ATP8A2	MPV17L	NOMO
MYH11	LINC01048	POPK1	PKD1
ABCC1, FOPNL	chr16:1472810-14816195	MPV17L	MPV17L
RM2	Unique Sequence	EIF3C	EIF3C
RRN3		SMG1	SMG1
PDXDC1		CLEC18	CLEC18
UGRC2, PDZD9		PDPDR	PDPDR
USP31		BANP	BANP
COG7		PKD1	EIF3C
SCGB2A2		chr7:53527933-83567877	SULT1A
SLC01B1		RRN3	SMG1
TPTE2		PDXDC1	OTOAP1
		SNX28P1	SNX28P1
		Unique Sequence	BOLA2
			Unique Sequence

1 Speciation node between asian and african apes
2 NPIPA expansion in african apes
3 NPIPB expansion in african apes

Figure S6: Evolutionary analysis of LCR16a copies in great apes. An unrooted neighbour-joining tree was constructed using the MEGA5 complete deletion option based on ~6 kbp of aligned sequence representing 86 LCR16a paralogs identified in the great ape lineage (human = HSA, chimpanzee = PTR, gorilla = GGO, PPY = orangutan). The ancestral sequence in macaque is used as an outgroup and evolutionary distances were computed using the Kimura 2-parameter method. Mosaic duplication block architecture is represented by coloured blocks and individual duplicons are referred to based on their gene content. HSA LCR16a copies are identified based on RefSeq nomenclature. Key nodes defining independent LCR16a clades are defined by boxed numbers. Independent LCR16a expansion in among Asian apes (box 1) includes two major clades corresponding to chromosome 13 (pink) and chromosome 16 (red) LCR16a duplications. Among African apes, three additional clades correspond to expansions of *NPIPA* (box 2) and two separate expansions of *NPIPB* (box 3). The LCR16a ‘core element’ is shared amongst all duplication blocks (red).

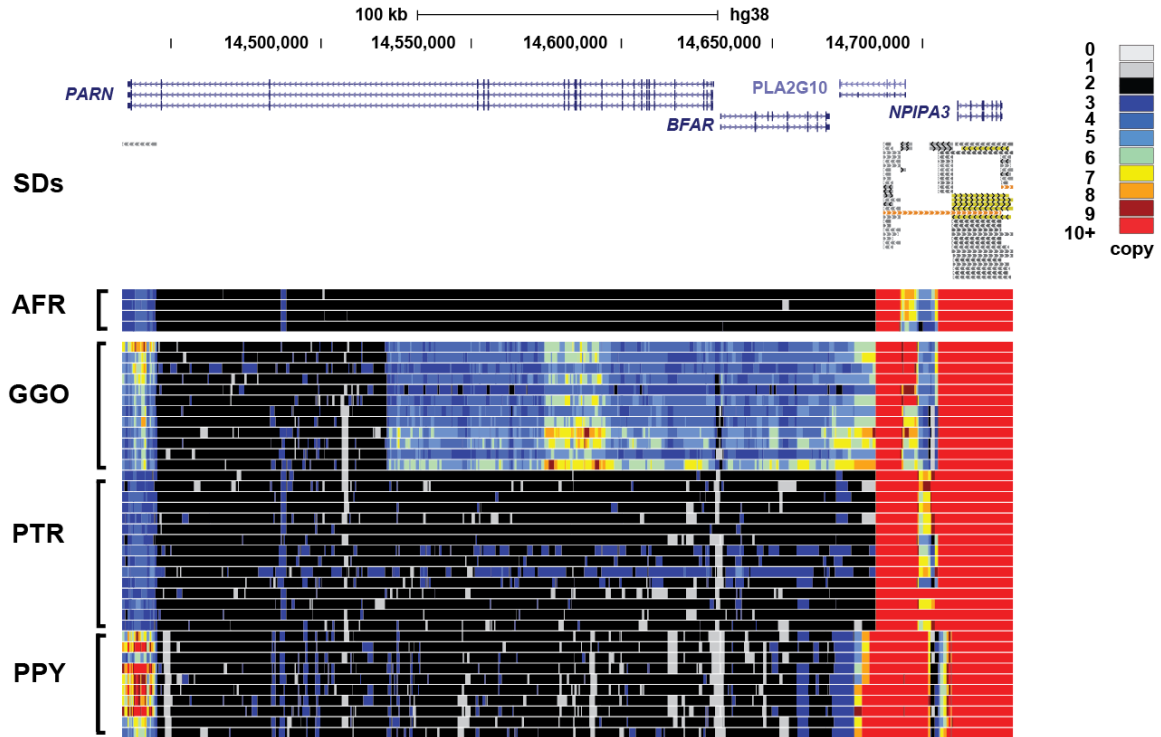


Figure S7: Lineage-specific duplicate genes. Pictured are UCSC Genome Browser snapshots (human build GRCh38) of three genes *PARN*, *BFAR*, and *PLA2G10*. Copy number (CN) heatmaps, with CN index shown, produced from Illumina read-depth predictions representative of modern humans from the Human Genome Diversity Project (HGDP) cohort, and nonhuman primates (NHPs). A large 130 kbp lineage-specific duplication of *PARN*, *BFAR*, and *PLA2G10* is identified in the gorilla, which is flanked by a copy of LCR16a.

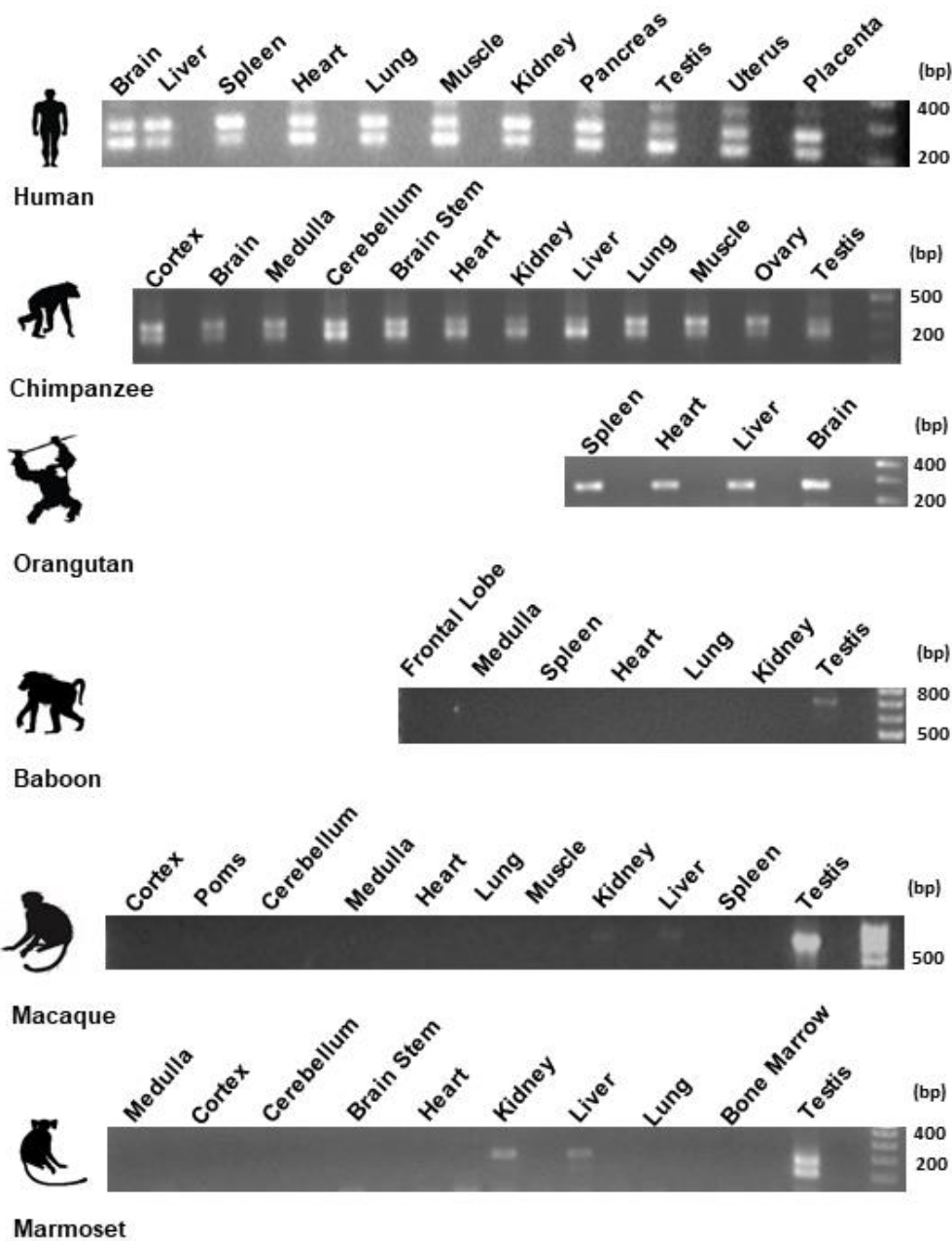


Figure S8: *NPIP* expression in a diversity panel of tissues/subtissues originating from human and NHP primary source material. RT-PCR reactions were performed using primers (Additional file 2: Table S7) designed to the canonical *NPIP* transcript described previously [13]. cDNA was prepared from mRNA generated from tissue source material (Methods) originating from chimpanzee, orangutan, baboon, macaque, and marmoset. Reactions were visualized on a 2% agarose gel, with primers designed to the ubiquitously expressed gene *UBE1* used as a control. A pattern of ubiquitous expression is observed in all great apes, while tissue-specific expression, largely limited to the testis, is observed in NWM and OWM.

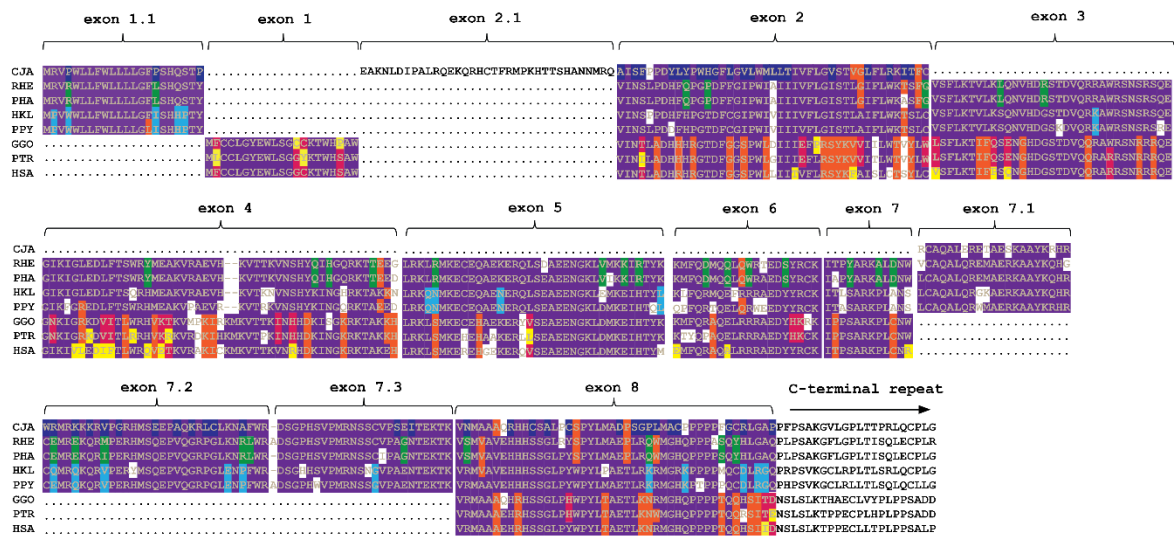


Figure S9: Gene amino acid structure of *NPIP* isoforms throughout primate evolution. A multiple sequence alignment compares *NPIP* isoforms across eight primate lineages. Individual exons are translated based on the putative open reading frames (ORFs) generated from gene models predicted using PacBio Iso-Seq and RefSeq gene annotations. Purple shading shows conserved amino acids. Marked changes at the N- and C-termini are shown during primate evolution.

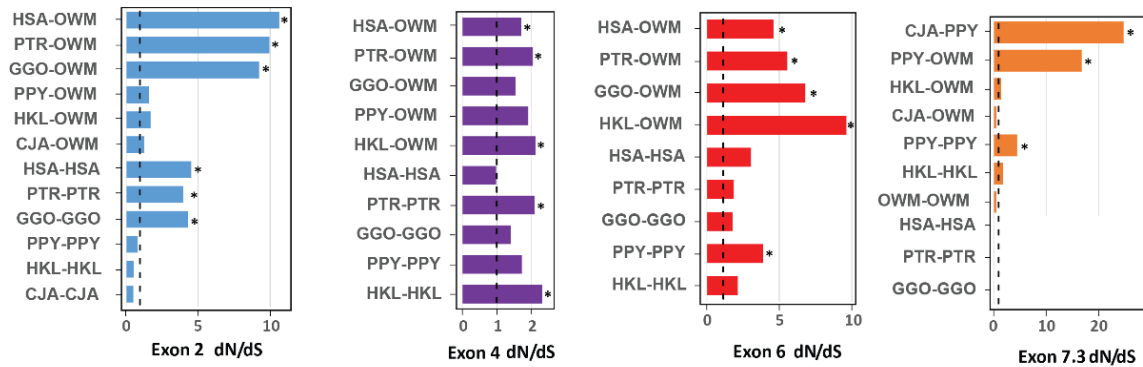
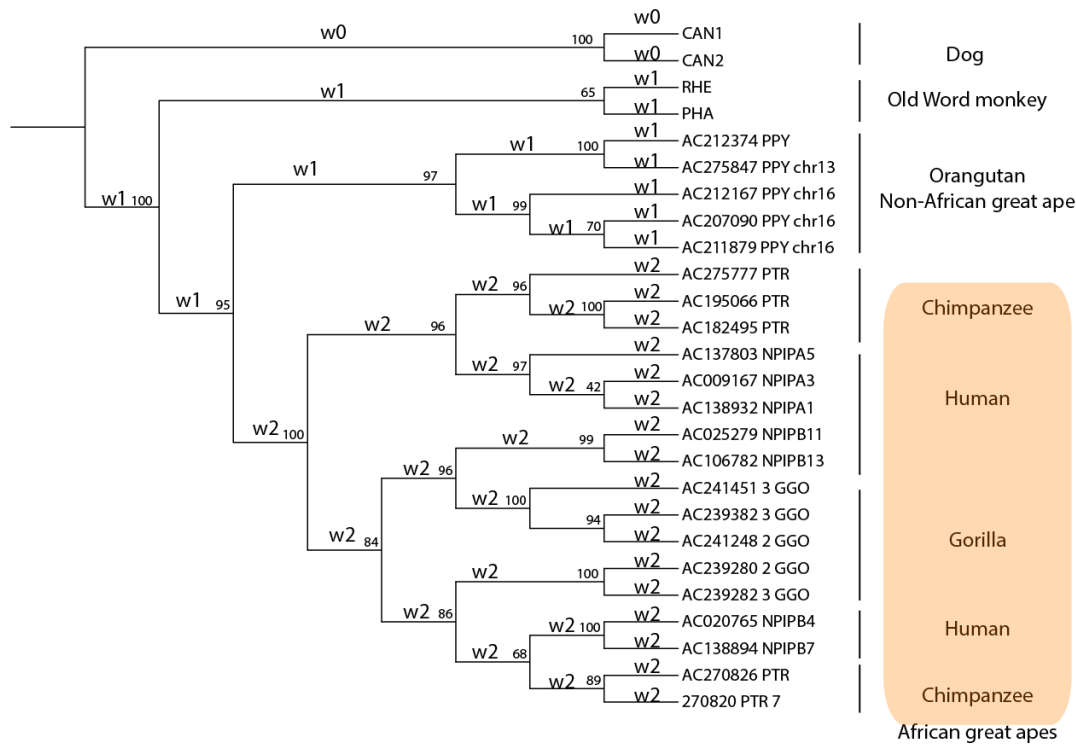


Figure S10: Selection analysis across four *NPIP*-coding exons based on an excess of nonsynonymous amino acid replacements. Graphs represent the ratio of dN/dS from individual exons between species. For each exon (exons 1-7.3 based on analysis of different primate lineages), we compared the number of nonsynonymous substitutions per site (dN) to the number of synonymous substitutions per site (dS) (Additional file 2: Table S5). Within and between species, comparisons were made (HSA=human, PTR=chimpanzee, GGO=gorilla, PPY=orangutan, HKL=gibbon, OWM=Old World monkey-baboon and macaque). Among OWM species a single copy of *NPIP* exists and thus all comparisons are orthologous. dN/dS ratios significantly greater than 1.0 are taken as evidence of positive selection (Z test transformation) (asterisks). Exons 2 and 6 show the most extreme levels of positive selection. Specifically, in exon 2 we find selection restricted to the African ape lineage (blue graph); in comparison, positive selection is observed in exon 6 (red graph) and exon 4 (blue graph) among all great apes. Note, selection is restricted to the orangutan lineage in exon 7.3 (orange graph) indicating species-specific differences for *NPIP* exons.



H0: $w_0=w_1=w_2=1$; (neutral evolution on all lineages?)
 $\ln L = -3718.735169$ (#parameters=51)

H1: $w_1=w_2=1$; **w_0 free**; (Selection on the dog lineages?)
 $\ln L = -3718.000394$ (#parameters=52); $w_0=0.7828$; H0 vs H1: $p_value = 0.2254$

H2: $w_0=1$; **$w_1=w_2$ free**; (Selection on the primate lineages?)
 $\ln L = -3709.921393$ (#parameters=52); $w_2=1.7961$; H0 vs H2: $p_value = 2.6867e-5$

H3: $w_0=w_1=1$; **w_2 free**; (Selection on the African great ape lineages?)
 $\ln L = -3710.773721$ (#parameters=52); $w_2=1.96412$; H0 vs H3: $p_value = 6.5975e-5$

H4: $w_0=w_2=1$; **w_1 free**; (Selection on the orangutan and Old World monkey lineages?)
 $\ln L = -3717.525219$ (#parameters=52); $w_1=1.46481$; H0 vs H4: $p_value = 0.1198$

Figure S11: Evidence for positive selection in the African ape lineages using branch model analysis (PAML). The gene tree was reconstructed using a maximum likelihood method (IQ-TREE), from sequences representing dog ($n=2$), macaque ($n=1$), baboon ($n=1$), orangutan ($n=5$), gorilla ($n=5$), chimpanzee ($n=5$), and human ($n=7$). The w 's in black and red represent fixed and free parameters, respectively, in the PAML analysis. Numbers below branches are the percentage of bootstrap supports for the inferred gene tree. P values are computed using the likelihood ratio test.

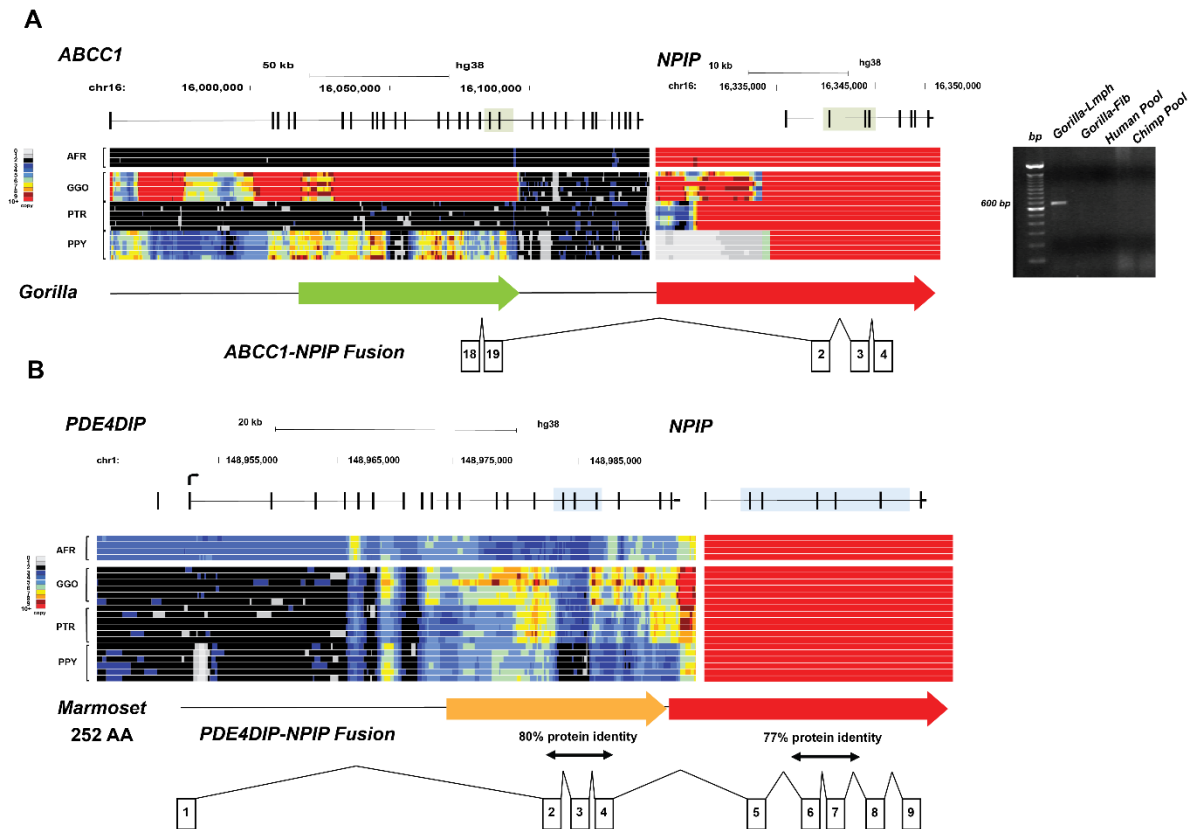


Figure S13: *NPIP* gene fusion transcripts detected using PacBio Iso-Seq and RT-PCR. **A)** An *ABCC1-NPIP* fusion transcript detected in gorilla. A schematic of an LCR16a duplication block on chr16p13.11 (UCSC Genome Browser snapshot). Annotation tracks include RefSeq gene annotation and copy number (CN) heatmaps (CN index shown) produced from Illumina read-depth profiles from modern humans and NHPs. Lineage-specific duplications are observed in *ABCC1* for gorilla and orangutan. BAC sequence analysis shows that the *ABCC1* duplication sits adjacent to a copy of LCR16a in gorilla (green and red arrows). An *ABCC1-NPIP* fusion transcript (exons 18 and 19 of *ABCC1* and exons 2-4 of *NPIP*) is detected from sequenced RT-PCR products originating from cDNA generated from gorilla lymphoblast source material. The corresponding product is missing from pooled cDNA (human and chimpanzee source tissue). **B)** A *PDE4DIP-NPIP* fusion transcript is detected from PacBio Iso-Seq experiments. A 252 amino acid putative ORF is predicted. The transcript originates from a lineage-specific duplication of *PDE4DIP* and LCR16a (orange and red arrows) in marmoset, creating a fusion transcript that includes exons 15-17 of *PDE4DIP* and exons 3-6 of *NPIP-S*.

	Ex2	Ex3	Ex4	Ex5
Dog 2	V HLLPGHFFPLGSDGVRLLWTLTIIISLGLICAFLLIYLMWRTILA	LQDFVEDPELRLVLRSKKCDLE	ESKKVLEDKCNALSSVKTMEKAELEKQLKKKVDTLVEMVEFYEQKRM	RKLEKTHYELVAKKNQFSTKEEFFKVATEEIDRY
	V HLLPGHFFPLGSDGV ILWTLT++ SLGICAFLLIYLMWRTI A	LQD FEDPELRLVLRSKK DLE	+SKKVLEDKCNALSSVKT KEAELEKQLKKK DTL EFFEQ	KL KTHYELVAKKNQ S E KVATEEID Y
Dog 1	V HLLPGHFFPLGSDGVEILWTLTVVTSLSGLICAFLLIYLMWRTIFA	LQDEFEDPELRLVLRSKKYDLE	DSKKVLEDKCNALSSVKTKEAELEKQLKKKMDTLAEFFEQTKVAAE	EKLEKTHYELVAKKNQLSATEKNLKVATEEIDKY
	V + LP HF G D I W + ++ LGI I+LW+T F	D D + R RS E	K LED + ++ AE+ ++ K+++ + + Q K E	KL+ E K+ QLS E+n K+ ++I Y
Macaque	VINSLPDHFPGDFGDFGIPW-IAIIIVFLGISTLGIELWRTSFG	-HRSRTDVQRRAWRSNSRSQE	GIKIGLEDLFTSWRYMEAKVRAEVHKVTTKVNSHYQIHGQRKTTEE	GKLRMKECEQAEKERQLSDAEENGLVMKKIRTY

	Ex6	Ex7	Ex8	Ex9
Dog 2	RQVQEMQEQLOESELTFRRHQ	IAIHEKNAQDNW	VKAQIWEWETIAQOSREKAYLKH	RLGILEGRMLPERHRRQVLMGRPEMQN
	QQVQEMQEQLOESELTFRR Q	IA+HEKNAQDNW	VKA IWE EIAQOSREKAYLKH	RLG LE LPERHRRQ I PGRPE QN
Dog 1	RQVQEMQEQLOESELTFRRHQ	IADVHEKNAQDNW	VKARIWERETIAQOSREKAYLKH	RLGLE-ESLPERHRRQELIPGRPEIQN
	K+ Q+MQ+ +E ++R +	I + + A DNW	V A+ +RE+A+ R+ AY +H	+ E + +PERH QE + GRP ++N
Macaque	KKMFQDMQQLQWRTEDSYRCK	ITFYARKALDNW	VCAQALQREMAE--RKAAKQKH	GCEMREKQRMPEHMSQEFVQGRGLKN

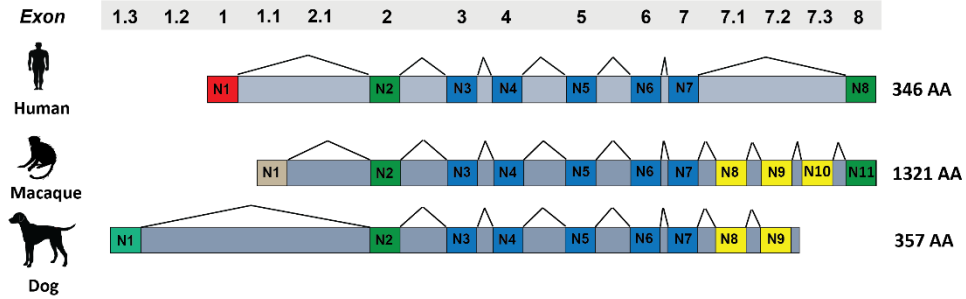


Figure S14: NPIP protein alignment and exon comparison between dog, macaque, and human. The *NPIP* ancestral gene structure is largely conserved between orthologs. Translation of putative protein-coding exons between macaque and dog *NPIP* copies identifies 8/11 exons shared between the two gene models (note three predicted canine-specific exons not shown).

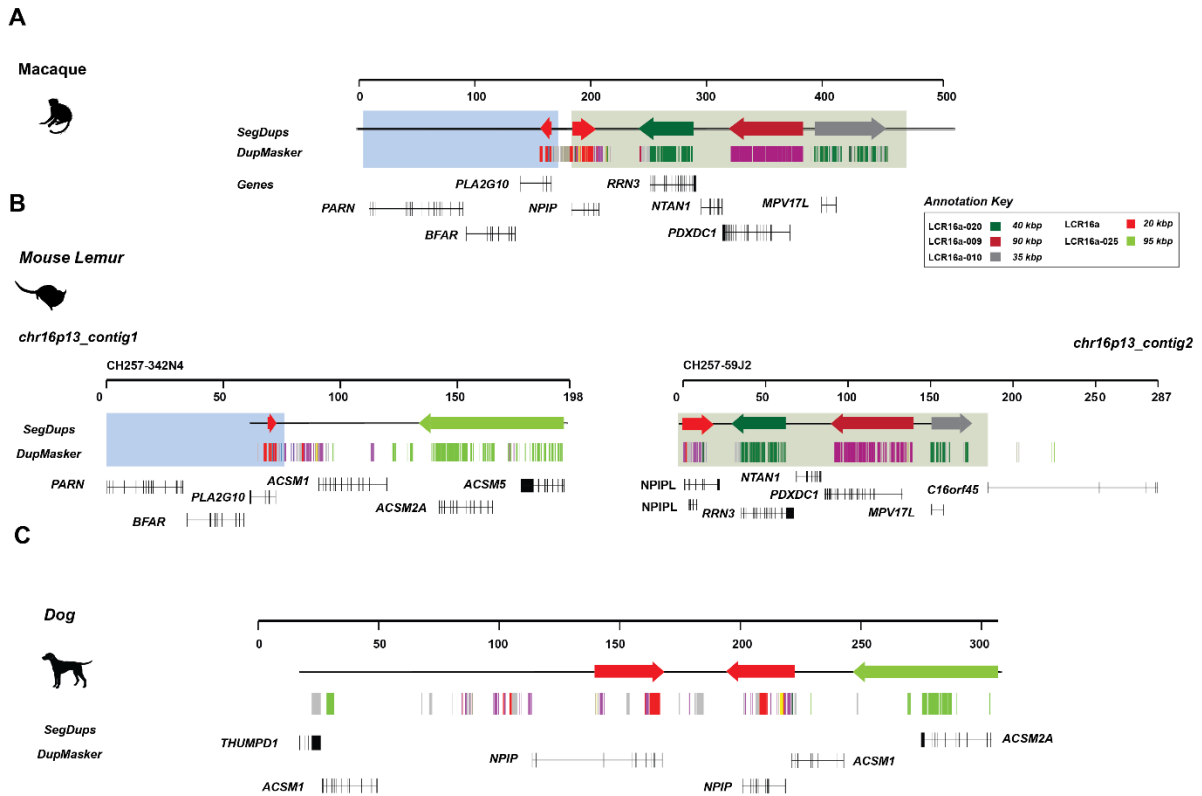


Figure S15: Cross-species comparison of the predicted ancestral *NPIP* organization between macaque, mouse lemur, and dog. Annotations include SDs (coloured arrows), gene models (black arrows), and DupMasker [50]. **A)** Large-insert clone-based assembly of the macaque locus demonstrates three ancestral duplicons, plus LCR16a (red arrow) representing a single copy of *NPIP*. **B)** Contigs mapping to the ancestral locus in the mouse lemur show that the organization is largely conserved (blue and cream shading) with the exception of a lineage-specific duplication, which includes *ACSM2A* and *ACSM5* (green arrow). **C)** Analysis of the canFAM3.0 reference assembly (chr6:24589985-24896899) identifies two unannotated copies of *NPIP* organized in an inverted orientation. Similar to the mouse lemur, an *ACSM2A* duplication is identified adjacent to *NPIP*.

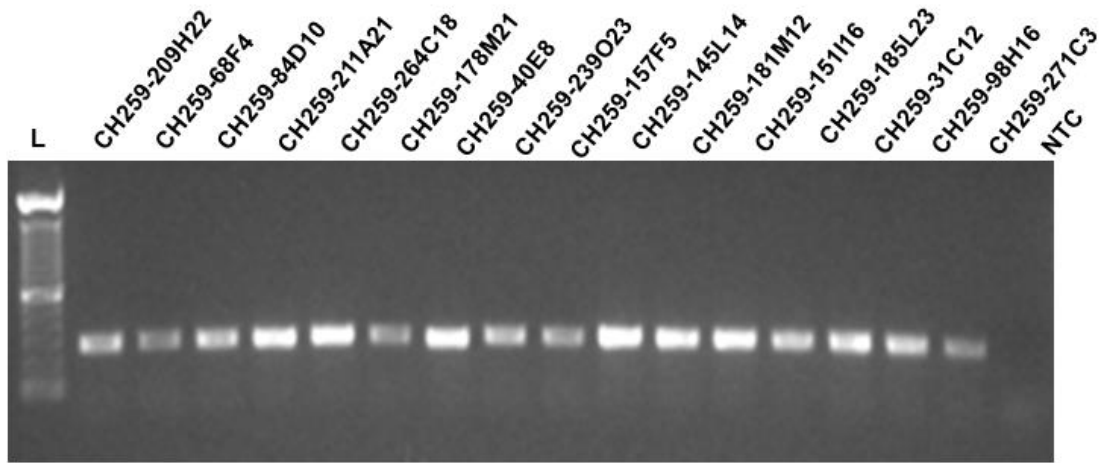


Figure S16: PCR-based testing of LCR16a-positive BAC clones derived from the CH259 large-insert clone library. DNA extracted from 16 LCR16a-positive BAC clones by hybridization were PCR confirmed for the presence of the LCR16a core duplicon. A ~300 bp product was visualized on a 2% agarose gel.

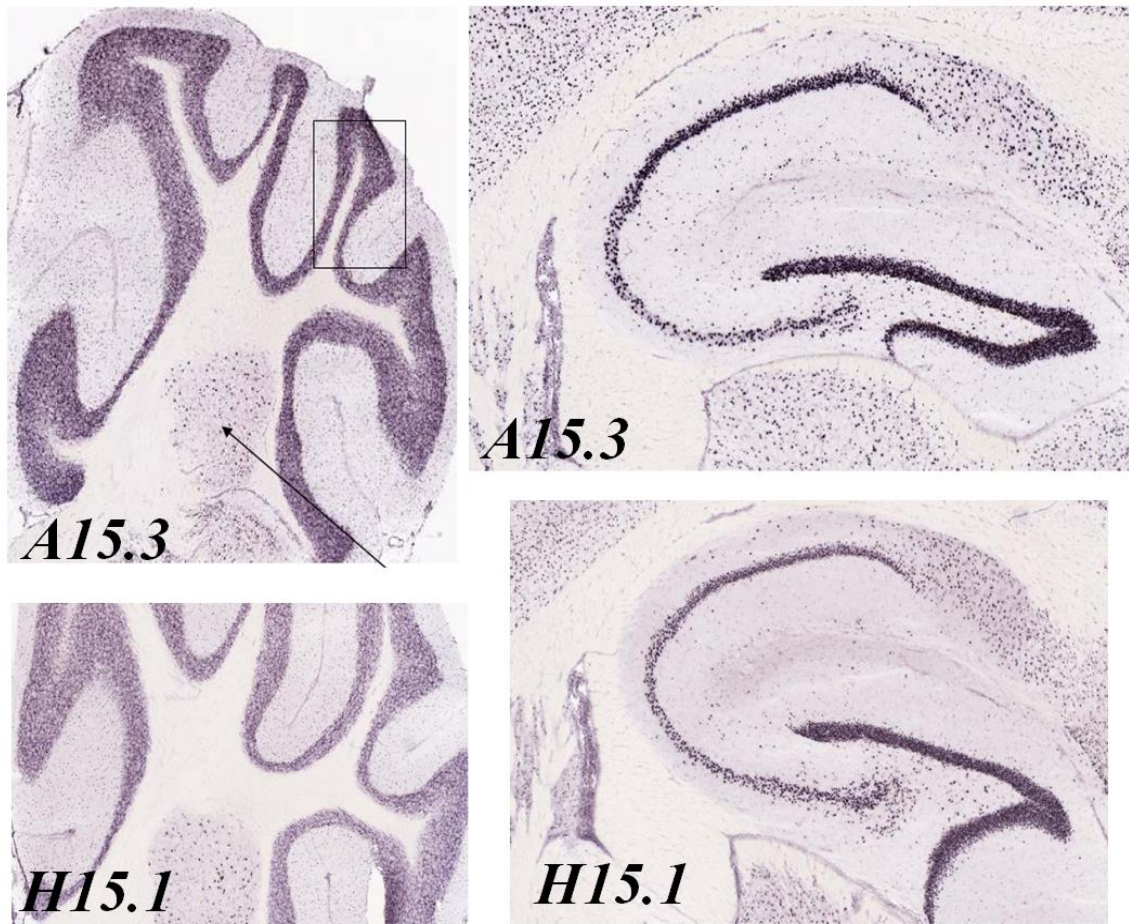
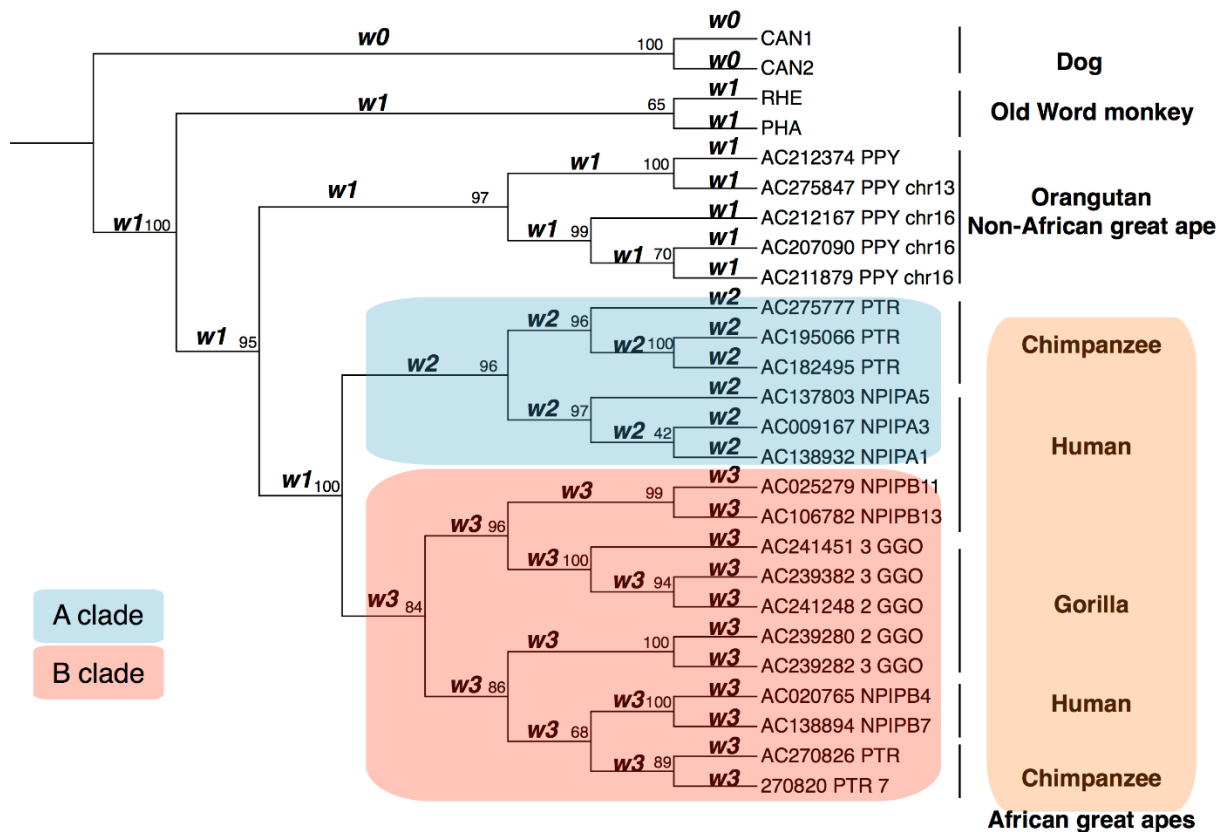


Figure S17: ISH expression analysis for BAC transgenic mice. Patterns of neuronal expression are consistent between BAC transgenic lines A15.3 and H15.1 in the hippocampus (left) and dentate gyrus (right). A15.3 and H15.1 correspond to different *NPIP* copies (RP11-344H15 and RP11-1236O14).



H0: $w_0=w_1=w_2=w_3=1$; (neutral evolution on all lineages?)
 $\ln L = -3718.735169$ (#parameters=51)

H5_A clade free: $w_0=w_1=w_3=1$; w_2 free; (Selection on the A clade of African great apes?)
 $\ln L = -3717.877165$ (#parameters=52); $w_2=1.57116$; H0 vs H5_A clade: $p_value = 0.1902$

H5_B clade free: $w_0=w_1=w_2=1$; w_3 free; (Selection on the B clade of African great apes?)
 $\ln L = -3712.902393$ (#parameters=52); $w_3=2.16776$; H0 vs H5_B clade: $p_value = 0.0006$

Figure S18: Evidence for positive selection in the *NPIP*B subtype using branch model analysis (PAML). The gene tree was reconstructed using a maximum likelihood method (IQ-TREE). The w 's in black and red represent fixed and free parameters, respectively, in the PAML analysis. Numbers below branches are the percentage of bootstrap supports for the inferred gene tree. P values are computed using the likelihood ratio test.

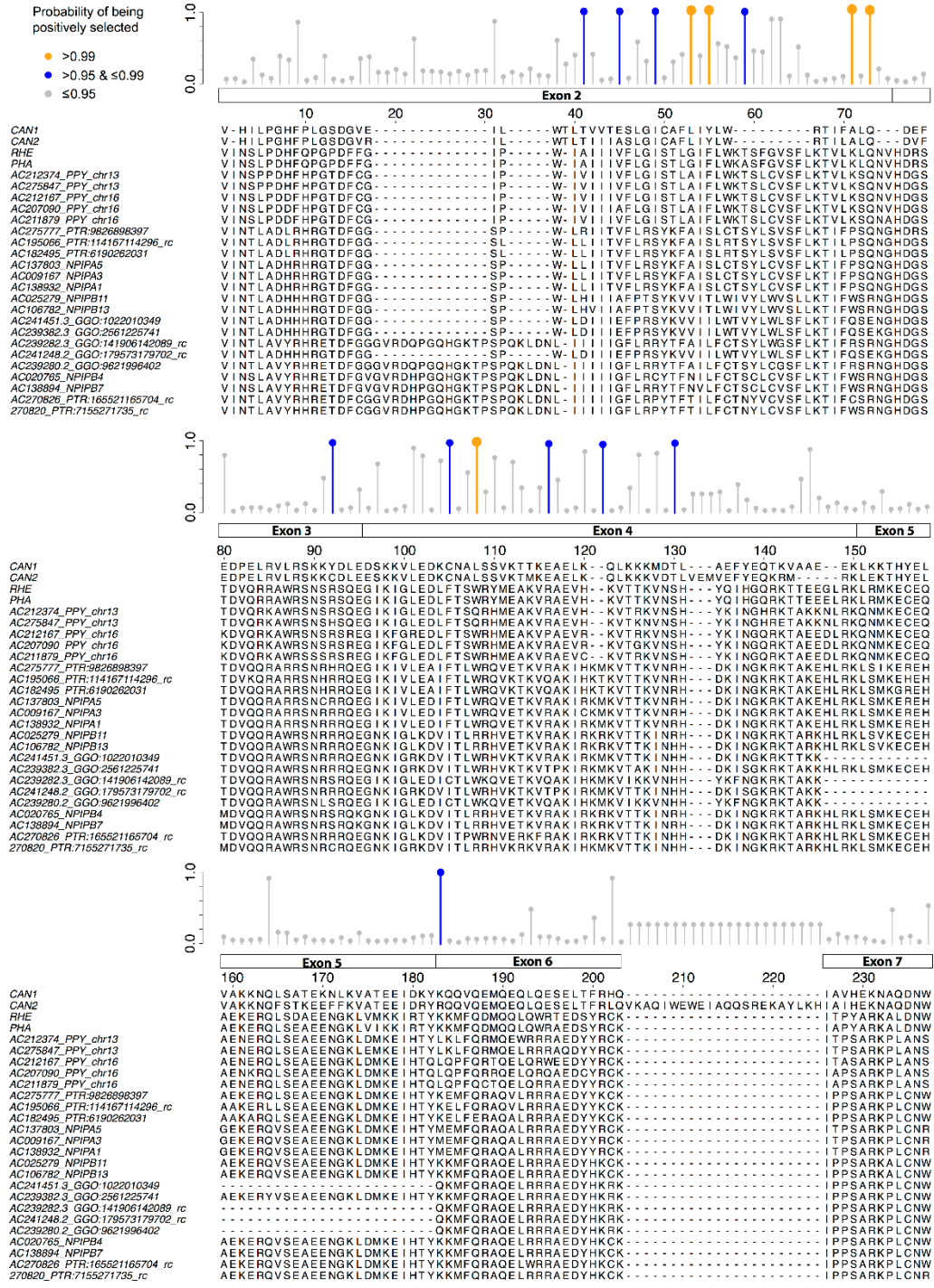


Figure S19: Evidence for positive selection relating to the *NPIPb* subtype within the African ape lineage. The branch-site test of positive selection (PAML v14.9) identifies 15 positively selected sites in exons 2, 3, 4, and 6, using 26 *NPIP* sequences (dog (n=2), macaque (n=1), baboon (n=1), orangutan (n=5), gorilla (n=5), chimpanzee (n=5), and human (n=7)).