

40 **ABSTRACT**

41 To better understand the pattern of primate genome structural variation, we sequenced and
42 assembled using multiple long-read sequencing technologies the genomes of eight nonhuman
43 primate species, including New World monkeys (owl monkey and marmoset), Old World
44 monkey (macaque), Asian apes (orangutan and gibbon), and African ape lineages (gorilla,
45 bonobo, and chimpanzee). Compared to the human genome, we identified 1,338,997 lineage-
46 specific fixed structural variants (SVs) disrupting 1,561 protein-coding genes and 136,932
47 regulatory elements, including the most complete set of human-specific fixed differences.
48 Across 50 million years of primate evolution, we estimate that 819.47 Mbp or ~27% of the
49 genome has been affected by SVs based on analysis of these primate lineages. We identify
50 1,607 structurally divergent regions (SDRs) wherein recurrent structural variation contributes
51 to creating SV hotspots where genes are recurrently lost (*CARDs*, *ABCD7*, *OLAH*) and new
52 lineage-specific genes are generated (e.g., *CKAP2*, *NEK5*) and have become targets of rapid
53 chromosomal diversification and positive selection (e.g., *RGPDs*). High-fidelity long-read
54 sequencing has made these dynamic regions of the genome accessible for sequence-level
55 analyses within and between primate species for the first time.

56

57 INTRODUCTION

58 An early and still unmet grand challenge of the Human Genome Project has been to
59 reconstruct the evolutionary history of every base pair of the human reference sequence¹⁻⁵.
60 To do so requires both a diverse sampling of nonhuman primate (NHP) genomes but also a
61 more complete assembly of those genomes so that all forms of variation can be assessed
62 without bias introduced from a superior quality reference⁶⁻¹³. Early attempts to sequence
63 closely related ape species focused primarily on characterizing simpler forms of variation
64 (e.g., single-nucleotide variants, (SNVs)) from portions of the genome that could be readily
65 aligned to human^{7-10,13}. As long-read sequence assemblies began to emerge, our ability to
66 catalog larger forms of structural variation significantly improved resulting in a series of
67 more contiguous NHP genomes. These new references, L
68 assemblies where allelic variation was collapsed and the most complex forms of gene-rich
69 structural variants (SVs) were still not resolved, including recently duplicated sequence¹⁴⁻¹⁹.
70 Advances in long-read sequencing technology over the last three years now allow for most of
71 these regions to be accurately sequenced and assembled to a degree where both paralogous
72 and allelic variation can be readily distinguished²⁰⁻²³. Numerous studies focused on the
73 human lineage have shown that such regions are incubators for the emergence of new genes,
74 adaptive evolution while also contributing to disease, and disease susceptibility²⁴⁻²⁶.

75
76 To better characterize SVs and these complex genic SV regions, we generated genome
77 assemblies of eight NHP genomes using two long-read sequencing platforms. Our plan was
78 twofold: First, we wanted to broaden the phylogenetic diversity by sequencing additional
79 NHP genomes using the same sequencing platform (in this case continuous long-read
80 sequencing or PacBio CLR) that had been initially applied to the other ape references to
81 minimize sequencing technology biases. This included sequence and assembly of primate
82 genomes representing gibbon (*Nomascus leucogenys*), marmoset (*Callithrix jacchus*), and
83 one owl monkey (*Aotus nancymaae*) (Table 1). Second, we wanted to leverage the higher
84 accuracy and assembly contiguity of HiFi (high-fidelity) sequencing data by sequence and
85 assembly of all NHP genomes where haplotypic differences could be distinguished. These
86 served as a means to validate all fixed structural variation events as well as provide complete
87 haplotype-resolved access to any particular regions of interest without the need to construct
88 and annotate these different NHP genomes for yet a third time.

89
90

91 **RESULTS**

92 ***Genome assembly of NHP genomes***

93 Building on our previous analysis of African great ape genomes^{14,17,19}, we first sequenced
94 and assembled three additional female NHP genomes using CLR sequencing, namely, white-
95 cheeked gibbon (*Nomascus leucogenys*), the common marmoset (*Callithrix jacchus*), and
96 owl monkey (*Aotus nancymaae*). Each genome was sequenced to high depth (>56-fold
97 coverage), assembled, and error corrected as described previously^{14,16,17,19} (Supplementary
98 Figure 1 and Supplementary Table 1). We generated highly contiguous (contig N50=9.9 to
99 25 Mbp) squashed assemblies of ~2.84-2.9 Gbp with an overall sequence accuracy of
100 >99.98% (Table 1 and Supplementary Table 1). Next, to further reduce sequencing error and
101 increase our ability to investigate more complex regions, we sequenced the same eight NHP
102 samples using PacBio HiFi sequencing^{17,27} (Table 1; Supplementary Figure 2 and
103 Supplementary Table 1). We used hifiasm to produce haplotype-resolved genomes that were
104 substantially smaller among monkeys (5.84 to 6.23 Gbp, diploid) when compared to
105 nonhuman apes²¹ (6.12 to 6.98 Gbp). These HiFi assemblies are estimated to be more
106 accurate (QV=42 to 58 or 99.9937% to 99.9998% accuracy) and significantly more
107 contiguous (contig N50=19 to 104 Mbp) when compared to the CLR draft genome
108 assemblies (Table 1 and Supplementary Figure 3).

109

110 ***NHP sequence divergence and incomplete lineage sorting (ILS)***

111 As a baseline for sequence divergence among the lineages, we mapped the HiFi sequence
112 data from each NHP back to human and computed single-nucleotide divergence (Methods).
113 The mean autosomal sequence divergence ranged from 1.3% to 9.83%, consistent with the
114 expected phylogeny, and was predictably higher than that of the X chromosome (0.99% to
115 8.24%; Figure 1a and 1b, Supplementary Table 2). We note that these estimates are also
116 slightly higher than earlier reports likely because a great fraction of repetitive DNA is being
117 included among NHPs^{8,19}. For example, among the apes ~92% of the human genome is
118 aligned in contrast to the New World monkey lineages where 64% and 59.7% of the
119 sequence from owl monkey and marmoset are unambiguously aligned (Supplementary Table
120 3). An assembly-based comparison yields similar results but involves a smaller fraction of
121 the genome due to extensive and more complex forms of structural variation (Supplementary
122 Figure 4 and Supplementary Table 3).

123

124 We used these data to generate a time-calibrated phylogeny for the nine primate species,
125 including human (Figure 1a and 1b; Supplementary Tables 4-6). We constructed more than
126 one million complete multiple sequence alignments (MSAs) at a resolution of 500 bp (518.9
127 Mbp of aligned sequence). While the majority of trees (52.7%) are consistent with the
128 generally accepted phylogeny, the fraction of alternate topologies is, once again, greater than
129 previous estimates^{9,13,17,28} (Figure 1c, Supplementary Table 4). Most of the difference can be
130 attributed to potential ILS during African ape or great ape speciation as gene tree
131 concordance factors show the lowest values in these two nodes (gene tree concordance=64.3
132 and 62, respectively)²⁹. Lineage-specific branch lengths are generally balanced with one
133 notable exception: the owl monkey branch length is significantly shorter and divergence to
134 human significantly lower when compared to marmoset (Figure 1a). An analysis of 16,244
135 gene trees using human as an outgroup to both owl monkey and marmoset shows that the owl
136 monkey evolves significantly slower ($p=0$ autosome, $p=6.85 \cdot 10^{-185}$ for the X chromosome)
137 (Supplementary Figure 5). Excluding potential sites of ILS, we estimated split times of the
138 species and find that mean split times of the apes better match the lower bounds of previous
139 estimates³⁰⁻³⁶ (Supplementary Table 7).

140

141 *Primate lineage-specific versus shared SVs*

142 We applied a three-pronged approach to discover and validate SVs (> 50 bp) mapping to the
143 euchromatic portion of the primate lineages^{37,38}. Using read-based and assembly-based
144 callers (pbsv, Sniffles and PAV), we first compared the eight NHP genomes against the
145 human reference genome, including three additional human genomes (CHM13, HG00733
146 and NA19240) to mitigate the effect of human polymorphism and missing variants in a
147 particular reference (Supplementary Table 8). In total, we identified 2.23 million putative
148 insertions and 1.89 million deletions in these nine lineages. Using both HiFi sequence data
149 and genome assemblies, we validated 1.85 million insertions and 1.63 million deletions
150 (mean validation rate: 86.79% and 89.37%, respectively) (Supplementary Table 9). We note
151 that genome-based HiFi and CLR SV calling are highly congruent (>95%) although HiFi
152 tended to recover larger insertions (Supplementary Figure 6). Finally, we generated Oxford
153 Nanopore Technologies (ONT) data from the same primate DNA samples and manually
154 inspected a subset (900 SV events) for confirmation using this orthogonal sequencing
155 platform estimating a false positive rate and a false negative rate of ~2.6% and 11.4%,
156 respectively (Supplementary Table 10).

157 To distinguish fixed from polymorphic events, we further genotyped (Methods) the validated
158 SVs against Illumina whole-genome sequence (WGS) data from a panel of 120 genomes (30
159 humans and 90 NHPs, Supplementary Table 11)³⁹⁻⁴³. We projected the 1,338,997 fixed
160 events (441,453 deletions and 897,544 insertions) onto the primate phylogeny (Figure 2a;
161 Supplementary Tables 12 and 13) classifying events as shared or lineage-specific¹⁷
162 (Methods). The number of SV events correlates strongly with evolutionary genetic distances
163 separating species (Figure 2b) with characteristic insertion peaks at ~6 kbp and 300 bp full-
164 length *LI* and *Alu* mobile element insertions (Supplementary Figure 7 and Supplementary
165 Table 14). Remarkably, we estimate that 27.2% of the genome (819.47 Mbp) has been
166 subjected to structural variation across these nine lineages with fixed insertions
167 outnumbering deletions approximately two to one (the total length of shared and lineage-
168 specific insertions is ~524.8 Mbp versus ~294.68 Mbp of deletions) (Figure 2a). The excess
169 of insertions is greatest for the ancestral ape and African great ape lineages (~2- to 3-fold)
170 (Figure 2a and Supplementary Table 13) and this twofold excess is still observed when
171 calibrating for the number of fixed SNV differences^{44,45} (Figure 2b; Supplementary Figures 8
172 and 9).

173

174 A small fraction of fixed primate SVs affect genes (~18.78 Mbp of deletions and ~1.31 Mbp
175 insertions). Using human gene annotation as a guide, we annotated the fixed SVs against the
176 human gene models (GRCh38, RefSeq) and the regulatory element database (ENCODE V3)
177 with Variant Effect Predictor (VEP)^{46,47}. These fixed SVs intersect 6,067 genes, including
178 1,561 protein-coding genes, and 136,932 regulatory elements. The latter includes 2,389
179 promoter-like (PLS) and 16,455 proximal enhancer-like signatures (pELS) potentially
180 disrupted by 16,671 fixed SVs (Supplementary Table 15). We estimate that 244 genes and
181 1,759 regulatory elements are novel and several are likely to confer functional effect
182 (Supplementary Figures 10 and 11). Such is the case for the 3,741 bp *LIPA5* insertion shared
183 in apes mapping to the last exon of the neuronal-function gene, *astrotactin 2 (ASTN2)*, which
184 encodes a glycoprotein that guides neuronal migration during the development of the central
185 nervous system^{48,49}. The insertion creates a novel transcript isoform resulting in a new exon
186 in human (NM_1884735) and this innovation is accompanied by a 1 base-pair deletion in this
exon, which in gibbon, orangutan, and gorilla is incapable of read

191 chondroitin sulfate attachment domain (Supplementary Figure 13). In gibbons, we identify a
192 large ~42.7 kbp deletion of the neurogenesis-associated gene, trace-amine associated
193 receptor 2, (*TAAR2*) along with seven of its enhancers (Figure 2d and Supplementary Figure
194 14). Loss of this brain-expressed gene in knockout mice has been shown to result in higher
195 levels of dopamine and lower levels of norepinephrine in the striatum and hippocampus
196 respectively⁵¹. A complete list of these gene and gene-regulatory fixed SVs is provided along
197 with additional discussion (e.g., *AR*, *SPATA1*, *ELN*, and *MAGEB16*) (Supplementary Tables
198 16 and 17, Supplementary Figures 15-18, and Supplementary Discussion).

199

200 We also reassessed human-specific changes and the effect of potential reference biases in
201 discovery. Importantly, 7,169 human-specific SVs have been reclassified, in part, because of
202 the inclusion of more outgroup species in addition to the use of more accurate sequence
203 aligner (minimap2 vs. blasr) that improves alignment within repetitive regions such as
204 subtelomeres^{52,53} (Supplementary Figures 19 and 20). Nevertheless, we identified 13
205 additional genes and 252 additional regulatory elements as potentially disrupted compared to
206 our previous report¹⁹ (Supplementary Figures 21 and 22). This includes, for example, a 90-
207 base pair deletion within the third exon of N-acetyltransferase 16 (*NAT16*) resulting in 30
208 amino acid loss in human lineage with respect to all other NHPs. The event was confirmed in
209 all humans by genotyping and by full-length transcript sequencing (Figure 2e and
210 Supplementary Figure 23). *NAT16* is highly expressed in the brain and pituitary and is
211 L -acetylhistidine synthesis, but its biological function remains unknown.

212

213 To assess the effect of using a human reference genome to classify such events, we repeated
214 ape-specific SV analyses using an assembled African human genome and the orangutan ape,
215 instead as the reference genomes to base the comparison. As expected, the analyses
216 reclassified approximately 34 gene-disruption events and led to a reduction of SVs most
217 notably with respect to insertions (Supplementary Figure 24). For example, using orangutan
218 as a reference reduces the number of lineage-specific insertions in orangutan (56,389 vs.
219 77,933), chimpanzee (2,020 vs. 4,471), bonobo (3,108 vs. 5,886), and human (13,446 vs.
220 16,696) lineage-specific insertions (Supplementary Figures 25 and 26, Supplementary Table
221 18). The intersect of these two sets provides the most conservative set of lineage-specific
222 changes on each branch. Consistent with the previous analyses, we find that the number of
223 insertions is ~2-3 times than that of deletions.

224

225 ***Structurally Divergent Regions (SDRs)***

226 In addition to increased accuracy and haplotype resolution, another major advantage of HiFi-
227 based assemblies is their 4- to 6-fold increase in sequence contiguity (Table 1). During our
228 comparison of monkey and ape chromosomes, we identified much larger, structurally
229 divergent regions (SDRs) that had been missed or incompletely assayed by our standard SV
230 analyses (Supplementary Figures 27 and 28). These regions were often gene-rich but had
eluded complete characterization due to the limitations of standard SV

259 thiol-dependent ubiquitinyl hydrolase activity ($p=1.9 \cdot 10^{-24}$), antimicrobial activity ($p=2.2 \cdot 10^{-5}$), innate immune response ($p=5 \cdot 10^{-5}$), neurotransmitter receptor activity ($p=2.5 \cdot 10^{-4}$), etc.
260
261 (Supplementary Table 22). Notably, most of these enrichments are associated with core
262 duplicons including *DEFBs*, *NPIPs*, *RGPDs*, *CYPs*, *NBPFs*, *GOLGAs*, *UGTs*, *RHDs*, and
263 *USPs*^{60,61} (Supplementary Table 23).

264

265 A few examples of these hotspot regions are illustrative. We confirmed, for example, that the
266 *CARD18* (caspase recruitment domain family member 18) was lost in the ancestral *Pan*
267 lineage by ~60 kbp deletion event⁷. We identified, however, a larger and independent
268 deletion of ~190 kbp in the gibbon lineage that completely removes the entire gene cluster
269 *CARD16* (pLI=0.04), *CARD17* (pLI=0), and *CARD18* (pLI=0.05). A third independent
270 deletion of ~150 kbp removed yet another member, *CARD17*, in the owl monkey suggesting
271 that this entire gene family has been under relaxed selection during primate evolution (Figure
272 3b and Supplementary Figure 34). Other hotspots are more complex, such as the *OLAH-*
273 *ACBD7* region showing evidence of both gain and loss of genes (Figure 3c). In gorilla,
274 *OLAH* (pLI=0) is deleted by a ~32 kbp deletion (Supplementary Figure 35) whereas in
275 macaque the locus has been the target of ~190 kbp duplication that truncates *OLAH* in that
276 lineage but also creates a new copy of *ACBD7*, which is actively transcribed as a fusion gene
277 (Figure 3c). In *Pan*, the same region has been the target of a ~250 kbp SD that originated
278 from chromosome 12 and produces a *Pan*-specific transcript with an open-reading frame
279 (ORF) of 97 amino acids whose promoter region is hypomethylated (Figure 3d,
280 Supplementary Figure 36). This large insertion of an SD in the *Pan* lineage also had the
281 benefit of removing one of two directly orientated duplications flanking *MEIG*
282 (meiosis/spermiogenesis associated 1), theoretically eliminating recurrent
283 microdeletion/microduplication of *MEIG1* in the *Pan* lineage (Figure 3d and 3e). A 28 kbp
284 genomic duplication region has been depleted in orangutans, but this has not resulted in any
285 alteration of gene content (Supplementary Figure 37). *MEIG1* (pLI=0.05) is a
286 spermiogenesis-related gene and *MEIG1* deficiency severely disrupts mouse spermatogenesis
287 and is potentially associated in human infertility⁶²⁻⁶⁴.

288

289 In order to test the potential for SDRs to serve as cradles for gene innovation, we repeated
290 our SDR analysis in a more distantly related primate. Using our graph-based approach, we
291 compared human and marmoset and identified 697 SDRs (~38.45 Mbp) that could not be

292 orthologously aligned to the complete human reference genome. Next, we manually
293 clustered them into 270 distinct SDR events since these two genomes are too divergent
294 (Supplementary Table 24). For the purpose of gene discovery, we also generated ~5.13
295 million full-length cDNA transcripts from 10 distinct primary tissues from the common
296 marmoset (Table 1 and Supplementary Table 25). We identified five regions that showed
297 evidence of novel or structurally divergent transcripts that lacked orthologous counterparts in
298 the human genome (Supplementary Figure 38 and Supplementary Table 24). Of particular
299 interest was a gene-rich region of human chromosome 13 that had been subject to a series of
300 inversions and duplications increasing by ~350 kbp in size and adding nine putative
marmoset-specific genes (Fig. 4a). We found 13 (c) genes expressed in 121 (90.75%) of 121 (100%)
marmoset-specific genes (Fig. 4a). We found 13 (c) genes expressed in 121 (90.75%) of 121 (100%)

326 5a and 5b, Supplementary Figure 43). For example, none of the gibbon or orangutan
327 duplicate copies map syntenically to each other or other African great apes thus, although
328 orangutan has multiple *RGPDs*, all originated independently and none have orthologs among
329 the other apes and group as distinct clade within the tree (Figure 5b and Supplementary
330 Figure 43). We identify only one paralogous gene, *hRGPD2*, that is syntenic and orthologous
331 among the African great apes. Within the five different ape lineages, we estimate ~20
332 independent mutation events (total length: ~1.2 Mbp) representing one of the most extreme
333 examples of homoplasy (Figure 5a and Supplementary Figure 42).

334
335 Most of the *RGPD* interspersed SDs were accompanied by both local restructuring of the
336 duplication blocks as well as larger scale structural rearrangements of the chromosome 2
337 flanking sequence especially in association with large-scale inversions in different NHP
338 lineages (Figure 5c and Supplementary Figure 44). Haplotype-resolved sequence assemblies
339 allowed the origin and spread of lineage-specific copies to be distinguished phylogenetically
340 (Figure 5b). Human *RGPD3* and *RGPD4* are not phylogenetically, for example, orthologs of
341 chimpanzee *RGPD3* and *RGPD4* even though they appear syntenic (Figure 5b and
342 Supplementary Figure 43) suggesting potential gene conversion. In addition, the emergence
343 of many *RGPDs* in apes appears to have been driven by recurrent large-scale inversions,
344 duplicative transpositions, and deletions within a ~7 Mbp genomic region over the last 15
345 million years of evolution creating unique configurations and distinct copies in each ape
346 lineage (Supplementary Figure 44).

347
348 *RGPD1* is a human-specific paralog predicted to have arisen ~570 thousand years ago (kya)
349 within the *Homo* lineage at ~0.57 mya (Figure 5b). This specific copy has several amino acid
350 replacements at the protein N-terminus with respect to all other human *RGPDs* this change
351 is predicted to alter the protein structure between *hRGPD1* and its antecedent *hRGPD2*⁶⁶
352 (Figure 5d). In this regard, it is interesting that the *hRGPD1* genomic region shows a dearth
353 of genetic diversity based on the analysis of Human Pangenome Reference Consortium
354 (HPRC) haplotype-resolved assemblies (π value = 4.65×10^{-5} , $p < 0.05$, TajimaD = -1.98)
355 (Figure 5e and Supplementary Figure 45) consistent with the region potentially being
356 subjected to a selective sweep specifically and recently in the human lineage.

357
358 In comparison to human, most of the copies mapping to bonobo and chimpanzee
359 chromosome 2 represent independent expansions from ancestral *RANBP2* that also gave rise

360 to human *RGPD5*, *RGPD6*, and *RGPD8* (Supplementary Figure 43). Of note, *RGPD6* is a
361 human-specific gene copy that arose via segmental duplication or gene conversion from
362 human *RGPD5* most recently (~5.2 kya, 95% CI [0.002,16.08]) (Figure 5b). The interval
363 between these human-specific copies, which includes *NPHPI*, is subjected to both inversion
364 toggling and microdeletion associated with Joubert syndrome and juvenile nephronophthisis
365 as a result of nonallelic homologous recombination (NAHR) between inverted and directly
366 orientated duplications⁶⁷⁻⁶⁹, respectively (Figure 5f and Supplementary Figure 46). We
367 examined 94 human phased haplotypes from the HPRC and Human Genome Structural
368 Variation Consortium^{38,69-71} and identified 11 distinct structural configurations four
369 predisposing to microdeletion (Figure 5g; Supplementary Figures 46-50 and Supplementary
370 Table 27). We also identified as single pathogenic allele deleting *NPHPI* (HG00733) and
371 confirmed maternal transmission (Supplementary Figures 51-53). A maximum likelihood
372 phylogenetic analysis identified the most closely related (non-deleted) haplotype and
373 breakpoint analysis confirms that the deleted allele arose from one of the haplotypes
374 predisposing to microdeletion (Supplementary Figure 51). Given the recent evolutionary
375 restructuring of this region of chromosome 2, it follows that this predisposition to
376 microdeletion is specific to the human lineage.

377

378 **DISCUSSION**

379 Using three long-read sequencing platforms across multiple primate genera, we present a
380 comprehensive analysis of SVs within euchromatic DNA of the primate order^{15,19}. The use of
381 HiFi data and inclusion of additional NHP species as well as genotyping in population
382 samples significantly improves earlier surveys of fixed SV events³⁹⁻⁴³ and extends the
383 analysis deeper within the primate phylogeny. Among the great apes for example, we
384 identify 13 genes and 1,759 regulatory elements not previously reported¹⁹ (Supplementary
385 Figures 21 and 22). The addition of other primate genomes identified lineage-specific SDR
386 events in the gibbon (n=680), macaque (n=219), and marmoset (n=697) lineages
387 (Supplementary Figure 30). Similarly, while we identify all 16 previously identified ape-
388 specific genic SVs; 13/16 are no longer classified as (great) ape-specific SVs (Supplementary
389 Table 28) due to the inclusion of other NHP lineages¹⁵. Finally, the use of a highly
390 contiguous orangutan genome as an alternate reference, helped reduce earlier human genome
391 reference biases by refining and polarizing the set of fixed SVs that occurred specifically
392 since humans diverged from the other ape lineages (Supplementary Table 18). Among the
393 6,067 genes (both coding and noncoding) and 136,932 regulatory DNA associated with fixed

394 SVs, we find a significant enrichment in transcription regulation ($p=1.1 \cdot 10^{-9}$), sensory
395 transduction ($p=6.3 \cdot 10^{-3}$), cell division ($p=2.3 \cdot 10^{-2}$), and vocal learning ($3.4 \cdot 10^{-3}$)
396 (Supplementary Table 29). These data serve as a rich resource for the characterization of
397 gene expression differences and candidate mutations for adaptation among NHPs.

398

399 The overall topology of the primate phylogenetic tree is consistent with previous
400 expectations with the proportion of ILS generally increasing as more of the repetitive content
401 is accessed by long-read sequencing technology¹⁷ (Figure 1). Our comparison of two New
402 World monkeys lineages, however, reveals significant acceleration of the marmoset SNV
403 branch length when compared to that of the owl monkey (branch length: 0.024 vs. 0.017).
404 This finding is also consistent with the shorter blocks of synteny in the marmoset lineage
405 when compared to the human genome (only 102 regions >500 kbp compared to 169 regions
406 >500 kbp in the owl monkey) and the significant increase in the number of recent SDs (165.7
407 Mbp in marmoset vs. 125.7 Mbp in owl monkey) (Supplementary Table 30). The slower
408 evolution of the owl monkey lineage compared to marmoset may simply be a consequence of
409 differences in reproductive longevity as has been proposed⁴⁰ or changes in the generation
410 time of the two lineages during evolution. The three major clades of New World monkeys,
411 however, are thought to have diverged over a short time frame (19-24 mya)^{35,36,72,73} (Figure
412 1a). Studying multi-generational pedigrees, Thomas and colleagues showed a 32.5%
413 reduction in the rate of *de novo* mutation in owl monkey when compared to that of apes with
414 an overall mutation rate of $0.81 \cdot 10^{-8}$ per site per generation⁴⁰. Our results suggest that this
415 reduced mutation rate may be longstanding property of the *Aotinae* with the net consequence
416 that the owl monkey genome is less derived when compared to marmoset. These findings
417 have some practical considerations regarding the use of these different New World monkeys
418 as models for human disease⁷⁴⁻⁷⁶.

419

420 The greater accuracy afforded by HiFi sequencing allowed more complex regions of genetic
421 variation to be assembled contiguously across the primates (e.g., MHC). We developed a
422 graph-based approach to systematically identify 1,604 SDRs among apes and macaque
423 (Figure 3) of which a third ($n=557$) showed evidence of recurrent structural variation and
424 were enriched for SDs. We hypothesize that these hotspots of recurrent structural variation
425 and their associated 631 genes (mean pLI=0.133) demarcate either regions of the ape genome
426 no longer under selection (e.g., *CARD18*, *OLAH*, etc.) or regions where rapid structural

427 diversification has facilitated the emergence of new genes showing signatures of positive
428 selection (e.g., *RGPD*, *NPIP*, *NPF*)⁷⁷⁻⁷⁹ (Figure 5) and/or important for adaptive
429 specializations in different primate lineages^{24,80,81}. Ironically, the innovations often come at a
430 cost with respect to fitness as the SDRs are associated with human disease susceptibility
431 regions (e.g., 1q22.3, 2q13, 16p11.2, 10p13), such as the human-specific duplication of
432 *RGPD5* and Joubert syndrome deletion alleles (Figure 5).

433
434 Our analysis also suggests that SDRs are common in the primate genome though with few
435 exceptions these regions have not been considered as part of previous large-scale sequencing
436 efforts because of 1) difficulties in their assembly and 2) challenges they pose in alignment
437 even among closely related species when fully resolved. We identified, for example, SDRs in
438 marmoset compared to owl monkey giving rise to marmoset-specific duplicate genes (Figure
439 4). Using our resource of ~5.13 million full-length transcripts, we show that these duplicate
440 genes are expressed in the brain, maintain an ORF, and emerged specifically since marmoset
441 diverged from other owl monkey ~20 mya (Supplementary Figure 54). The ancestral genes
442 have critical functions: *NEK5*, for example, is member of NimA family of serine/threonine
443 protein kinases involved in cell differentiation while *CKAP2* (cytoskeleton associated protein
444 2) is involved in cell division^{82,83}. These findings caution against simply using human gene
445 models to annotate NHP genomes or to assess NHP gene expression differences from single-
446 cell RNA sequencing experiments. Understanding the gene innovations in such previously
447 inaccessible complex regions of primate genomes will be critical to realizing the full
448 potential of these species as models of human genetic disease⁷⁴⁻⁷⁶.

449

450 **Materials and Methods**

451 We sequenced and assembled eight NHP reference genomes using long-read PacBio HiFi
452 and ONT sequencing chemistry and the hifiasm genome assembler²¹. All samples, with one
453 exception, were female and correspond to the same samples used in previous studies as
454 references, namely; Central chimpanzee (Clint)⁷, bonobo (Mhudiblu)¹⁷, Western gorilla
455 (Kamilah)¹³, Sumatran orangutan (Susie)⁸, Northern white-cheeked gibbon (Asia)¹⁰, rhesus
456 macaque (AG07107)¹⁶, common marmoset (CJ1700), and owl monkey (86718) (Table 1).
457 We used pbsv, Sniffles, and PAV to characterize SVs and merged SVs using the SVPOP
458 pipeline^{37,38}. The merged calls were validated with HiFi sequencing data and assembly of
459 select regions; ONT sequence data from the same specimens were used to calculate the false
460 positive rate and validate assembly of select regions in our data set. The validated SVs were

461 genotyped by Paragraph using Illumina WGS data from 120 population samples^{16,39-43,84}.
462 VEP was used to annotate the functional disruption of SVs⁴⁶. In addition to SVs (<20 kbp)
463 identified by the three callers, we used a graph-based aligner (Mashmap) to identify large
464 structural changes across apes and Old World monkey⁵⁵, defined here as SDRs. SDR
465 validation was based on haplotype-resolved assemblies and ONT data. The ONT data also
466 were used to call methylation by Guppy⁸⁵. We also generated full-length Iso-Seq data
467 specifically from 10 diverse marmoset tissues and from a gibbon immortalized lymphoblast
468 line. In the case of the marmoset, full-length RNA was prepared from 10 distinct tissues
469 obtained upon necropsy from a different specimen (*Callithrix jacchus*). Genomic divergence
470 analyses were based on HiFi sequencing data and genomes, respectively. Syntenic regions
471 across New World monkey to apes and MSAs were constructed with minimap2 and
472 mafft^{53,86}. The phylogenetic analyses were performed using TREEasy, IQTREE, and
473 BEAST2⁸⁷⁻⁸⁹.

474 **Acknowledgments**

475 We thank T. Brown for manuscript proofreading and editing. This article is subject to
476 L L L usly granted a
477 nonexclusive CC BY 4.0 license to the public and a sublicensable license to HHMI in their
478 research articles. Pursuant to those licenses, the author-accepted manuscript of this article
479 can be made freely available under a CC BY 4.0 license immediately upon publication.

480

481 **Funding**

482 This work was supported, in part, by National Institutes of Health (NIH) grants HG002385,
483 HG010169, and HG009081 to E.E.E.; GM147352 to G.A.L.; R01HG010485,
484 U41HG010972 and U01HG010961 to B.P.; R01-AI-137011 and DP1-DA-046108 to S.L.S.;
485 by Shanghai Pujiang Program (22PJ1407300) and Shanghai Jiao Tong University 2030
486 Program (WH510363001-7) to Y.M.; by National Natural Science Foundation of China
487 grants 82001372 to X.Y.; L.C. is supported by the P51 OD011092 (to the Oregon National
488 Primate Research Center); E.E.E. is an investigator of the Howard Hughes Medical Institute.

489

490 **Author contributions**

491 Y.M. and E.E.E. conceived the project; Y.M., W.T.H., K.M.M., K.H., A.P.L., P.A.A., A.R.,
492 D.S.G., G.A.L., P.C.D., and E.E.E. generated sequencing data, assembled genomes, analyzed
493 the data, and performed quality control analyses; X.Y., R.R., V.L.B., W.T.F., G.K.W., G.F.,

494 S.L.S., and W.C.W. contributed the marmoset and owl monkey samples; L.C. contributed the
495 bonobo and gibbon samples; Y.M. performed the SNV divergence and ILS analyses; Y.M.,
496 W.T.H., P.A.A., S.Z., G.A.L., H.J., and E.E.E. performed SV analyses; Y.M. performed
497 SDR analyses; M.H., and B.P. generated gene model annotations; Y.M., D.P., and E.E.E.
498 performed *NPHP1* haplotype analyses; Y.M., X.W., and Q.L. performed the protein structure
prediction analyses. Y.M. and E.E.E. drafted the manuscript.

- 535 12 Juan, D., Santpere, G., Kelley, J. L., Cornejo, O. E. & Marques-Bonet, T. Current
536 advances in primate genomics: novel approaches for understanding evolution and
537 disease. *Nature Reviews Genetics*, 1-18 (2023).
- 538 13 Scally, A. *et al.* Insights into hominid evolution from the gorilla genome sequence.
539 *Nature* 483, 169-175 (2012).
- 540 14 Gordon, D. *et al.* Long-read sequence assembly of the gorilla genome. *Science* 352,
541 aae0344 (2016).
- 542 15 He, Y. *et al.* Long-read assembly of the Chinese rhesus macaque genome and
543 identification of ape-specific structural variants. *Nature communications* 10, 4233
544 (2019).
- 545 16 Warren, W. C. *et al.* Sequence diversity analyses of an improved rhesus macaque
546 genome enhance its biomedical utility. *Science* 370, eabc6617 (2020).
- 547 17 Mao, Y. *et al.* A high-quality bonobo genome refines the analysis of hominid
548 evolution. *Nature* 594, 77-81 (2021).
- 549 18 Yang, C. *et al.* Evolutionary and biomedical insights from a marmoset diploid
550 genome assembly. *Nature* 594, 227-233 (2021).
- 551 19 Kronenberg, Z. N. *et al.* High-resolution comparative analysis of great ape genomes.
552 *Science* 360, eaar6343 (2018).
- 553 20 Logsdon, G. A., Vollger, M. R. & Eichler, E. E. Long-read human genome
554 sequencing and its applications. *Nature Reviews Genetics* 21, 597-614 (2020).
- 555 21 Cheng, H., Concepcion, G. T., Feng, X., Zhang, H. & Li, H. Haplotype-resolved de
556 novo assembly using phased assembly graphs with hifiasm. *Nature methods* 18, 170-
557 175 (2021).
- 558 22 Rautiainen, M. *et al.* Telomere-to-telomere assembly of diploid chromosomes with
559 Verkko. *Nature Biotechnology*, 1-9 (2023).
- 560 23 Mao, Y. & Zhang, G. A complete, telomere-to-telomere human genome sequence
561 presents new opportunities for evolutionary genomics. *Nature methods* 19, 635-638
562 (2022).
- 563 24 Dennis, M. Y. *et al.* Evolution of human-specific neural SRGAP2 genes by
564 incomplete segmental duplication. *Cell* 149, 912-922 (2012).
- 565 25 Fiddes, I. T. *et al.* Human-specific NOTCH2NL genes affect notch signaling and
566 cortical neurogenesis. *Cell* 173, 1356-1369. e1322 (2018).
- 567 26 Kawanishi, K. *et al.* Human species-specific loss of CMP-N-acetylneuraminic acid
568 hydroxylase enhances atherosclerosis via intrinsic and extrinsic mechanisms.
569 *Proceedings of the National Academy of Sciences* 116, 16036-16045 (2019).
- 570 27 Logsdon, G. A. *et al.* The structure, function and evolution of a complete human
571 chromosome 8. *Nature* 593, 101-107 (2021).
- 572 28 Mailund, T., Munch, K. & Schierup, M. H. Lineage sorting in apes. *Annual review of*
573 *genetics* 48, 519-535 (2014).
- 574 29 Minh, B. Q., Hahn, M. W. & Lanfear, R. New methods to calculate concordance
575 factors for phylogenomic datasets. *Molecular biology and evolution* 37, 2727-2733
576 (2020).
- 577 30 Steiper, M. E. & Young, N. M. Primate molecular divergence dates. *Mol Phylogenet*
578 *Evol* 41, 384-394, (2006).
- 579 31 Wilkinson, R. D. *et al.* Dating primate divergences through an integrated analysis of
580 palaeontological and molecular data. *Systematic biology* 60, 16-31, (2011).
- 581 32 Prado-Martinez, J. *et al.* Great ape genetic diversity and population history. *Nature*
582 499, 471-475, (2013).
- 583 33 Pozzi, L. *et al.* Primate phylogenetic relationships and divergence dates inferred from
584 complete mitochondrial genomes. *Mol Phylogenet Evol* 75, 165-183, (2014).

- 585 34 de Manuel, M. *et al.* Chimpanzee genomic diversity reveals ancient admixture with
586 bonobos. *Science* 354, 477-481, (2016).
- 587 35 Vanderpool, D. *et al.* Primate phylogenomics uncovers multiple rapid radiations and
588 ancient interspecific introgression. *PLoS biology* 18, e3000954, (2020).
- 589 36 Álvarez-Carretero, S. *et al.* A species-level timeline of mammal evolution integrating
590 phylogenomic data. *Nature* 602, 263-267, (2022).
- 591 37 Sedlazeck, F. J. *et al.* Accurate detection of complex structural variations using
592 single-molecule sequencing. *Nature methods* 15, 461-468 (2018).
- 593 38 Ebert, P. *et al.* Haplotype-resolved diverse human genomes and integrated analysis of
594 structural variation. *Science* 372, eabf7117 (2021).
- 595 39 Prado-Martinez, J. *et al.* Great ape genetic diversity and population history. *Nature*
596 499, 471-475 (2013).
- 597 40 Thomas, G. W. *et al.* Reproductive longevity predicts mutation rates in primates.
598 *Current Biology* 28, 3193-3197. e3195 (2018).
- 599 41 Rogers, J. *et al.* The comparative genomics and complex population history of *Papio*
600 baboons. *Science Advances* 5, eaau6947 (2019).
- 601 42 Okhovat, M. *et al.* Co-option of the lineage-specific LAVA retrotransposon in the
602 gibbon genome. *Proceedings of the National Academy of Sciences* 117, 19328-19338
603 (2020).
- 604 43 Byrska-Bishop, M. *et al.* High-coverage whole-genome sequencing of the expanded
605 1000 Genomes Project cohort including 602 trios. *Cell* 185, 3426-3440. e3419
606 (2022).
- 607 44 Marques-Bonet, T. *et al.* A burst of segmental duplications in the genome of the
608 African great ape ancestor. *Nature* 457, 877-881 (2009).
- 609 45 Sudmant, P. H. *et al.* Evolution and diversity of copy number variation in the great
610 ape lineage. *Genome research* 23, 1373-1382 (2013).
- 611 46 McLaren, W. *et al.* The ensembl variant effect predictor. *Genome biology* 17, 1-14
612 (2016).
- 613 47 Moore, J. E. *et al.* Expanded encyclopaedias of DNA elements in the human and
614 mouse genomes. *Nature* 583, 699-710 (2020).
- 615 48 Behesti, H. *et al.* ASTN2 modulates synaptic strength by trafficking and degradation
616 of surface proteins. *Proceedings of the National Academy of Sciences* 115, E9717-
617 E9726 (2018).
- 618 49 Bauleo, A. *et al.* Rare copy number variants in ASTN2 gene in patients with
619 neurodevelopmental disorders. *Psychiatric Genetics* 31, 239-245 (2021).
- 620 50 *et al.* Heterozygous aggrecan variants are associated with
621 short stature and brachydactyly: description of 16 probands and a review of the
622 literature. *Clinical endocrinology* 88, 820-829 (2018).
- 623 51 Efimova, E. V. *et al.* Trace amine-associated receptor 2 is expressed in the limbic
624 brain areas and is involved in dopamine regulation and adult neurogenesis. *Frontiers*
625 *in Behavioral Neuroscience* 16 (2022).
- 626 52 Chaisson, M. J. & Tesler, G. Mapping single molecule sequencing reads using basic
627 local alignment with successive refinement (BLASR): application and theory. *BMC*
628 *bioinformatics* 13, 1-18 (2012).
- 629 53 Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* 34,
630 3094-3100 (2018).
- 631 54 Porubsky, D. *et al.* Gaps and complex structurally variant loci in phased genome
632 assemblies. *bioRxiv*, (2022).

- 633 55 Jain, C., Koren, S., Dilthey, A., Phillippy, A. M. & Aluru, S. A fast adaptive
634 algorithm for computing whole-genome homology maps. *Bioinformatics* 34, i748-
635 i756 (2018).
- 636 56 Yang, X. *et al.* A refined characterization of large-scale genomic differences in the
637 first complete human genome. *bioRxiv*, (2022).
- 638 57 Sekar, A. *et al.* Schizophrenia risk from complex variation of complement component
639 4. *Nature* 530, 177-183 (2016).
- 640 58 Cantsilieris, S. *et al.* Recurrent structural variation, clustered sites of selection, and
641 disease risk for the complement factor H (CFH) gene family. *Proceedings of the*
642 *National Academy of Sciences* 115, E4433-E4442 (2018).
- 643 59 Thamadilok, S. *et al.* Human and nonhuman primate lineage-specific footprints in the
644 salivary proteome. *Molecular biology and evolution* 37, 395-405 (2020).
- 645 60 Vollger, M. R. *et al.* Segmental duplications and their variation in a complete human
646 genome. *Science* 376, eabj6965 (2022).
- 647 61 Jiang, Z. *et al.* Ancestral reconstruction of segmental duplications reveals punctuated
648 cores of human genome evolution. *Nature genetics* 39, 1361-1368 (2007).
- 649 62 Khan, N. *et al.* Crystal structure of human PACRG in complex with MEIG1 reveals
650 roles in axoneme formation and tubulin binding. *Structure* 29, 572-586. e576 (2021).
- 651 63 Zhang, Z. *et al.* MEIG1 is essential for spermiogenesis in mice. *Proceedings of the*
652 *National Academy of Sciences* 106, 17055-17060 (2009).
- 653 64 Du, R. *et al.* Efficient typing of copy number variations in a segmental duplication-
654 mediated rearrangement hotspot using multiplex competitive amplification. *Journal*
655 *of human genetics* 57, 545-551 (2012).
- 656 65 Ciccarelli, F. D. *et al.* Complex genomic rearrangements lead to novel primate gene
657 function. *Genome research* 15, 343-351 (2005).
- 658 66 Jumper, J. *et al.* Highly accurate protein structure prediction with AlphaFold. *Nature*
659 596, 583-589 (2021).
- 660 67 Parisi, M. A. *et al.* The NPHP1 gene deletion associated with juvenile
661 nephronophthisis is present in a subset of individuals with Joubert syndrome.
662 *American journal of human genetics* 75, 82-91, (2004).
- 663 68 Gana, S., Serpieri, V. & Valente, E. M. Genotype-phenotype correlates in Joubert
664 syndrome: A review. *Am J Med Genet C Semin Med Genet* 190, 72-88, (2022).
- 665 69 Porubsky, D. *et al.* Recurrent inversion polymorphisms in humans associate with
666 genetic instability and genomic disorders. *Cell* 185, 1986-2005. e1926 (2022).
- 667 70 Liao, W.-W. *et al.* A draft human pangenome reference. *bioRxiv*, 2022.2007.
668 2009.499321 (2022).
- 669 71 Wang, T. *et al.* The Human Pangenome Project: a global resource to map genomic
670 diversity. *Nature* 604, 437-446 (2022).
- 671 72 Schneider, H. The current status of the New World monkey phylogeny. *Anais da*
672 *Academia Brasileira de Ciências* 72, 165-172 (2000).
- 673 73 Perelman, P. *et al.* A molecular phylogeny of living primates. *PLoS genetics* 7,
674 e1001342 (2011).
- 675 74 Baer, J. F., Weller, R. E. & Kakoma, I. *Aotus: the owl monkey*. (Academic Press,
676 2012).
- 677 75 Okano, H., Hikishima, K., Iriki, A. & Sasaki, E. in *Seminars in fetal and neonatal*
678 *medicine*. 336-340 (Seminars in fetal and neonatal medicine).
- 679 76 Grillner, S. *et al.* Worldwide initiatives to advance brain research. *Nature*
680 *neuroscience* 19, 1118-1122 (2016).
- 681 77 Nuttle, X. *et al.* Emergence of a Homo sapiens-specific gene family and chromosome
682 16p11.2 CNV susceptibility. *Nature* 536, 205-209 (2016).

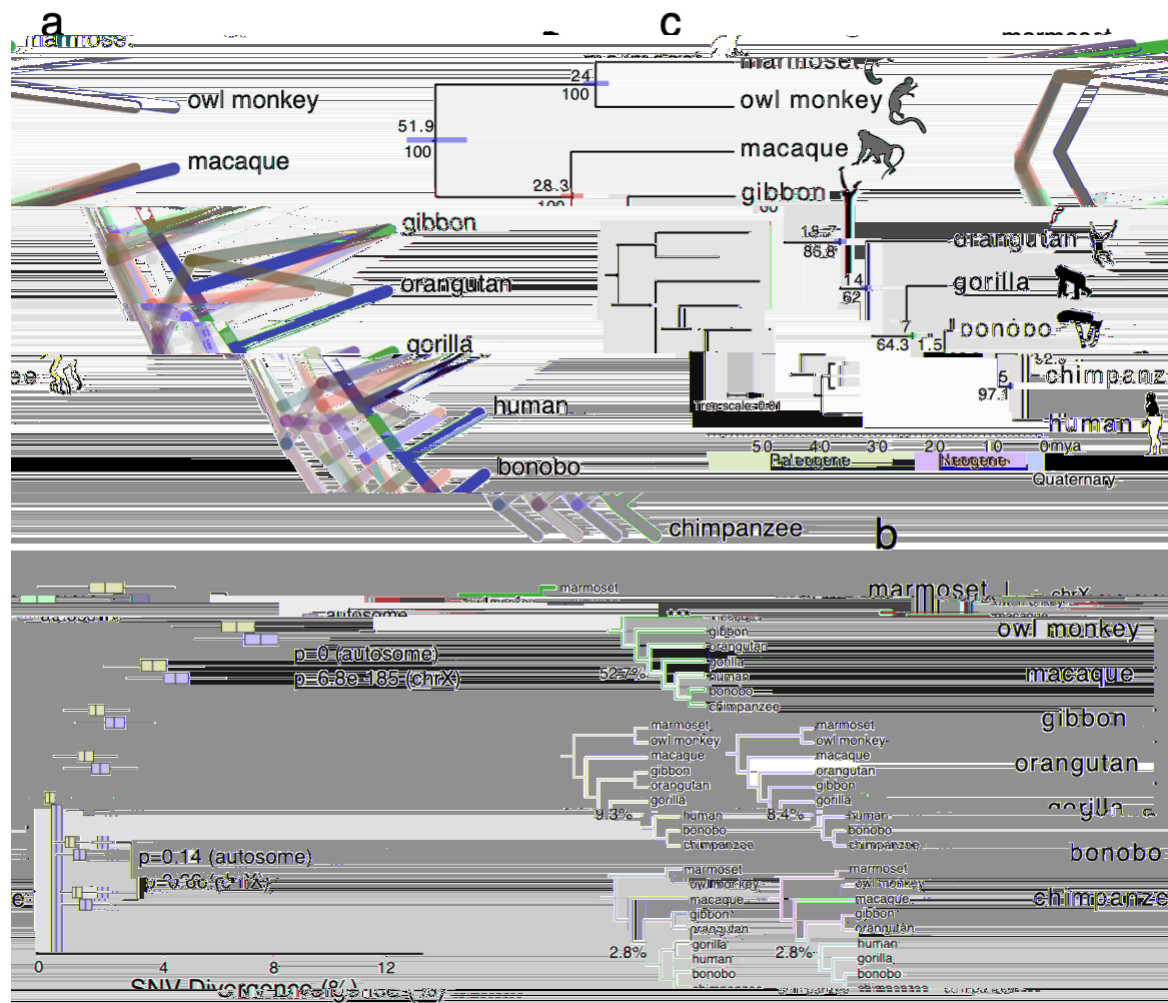
- 683 78 Hsieh, P. *et al.* Evidence for opposing selective forces operating on human-specific
684 duplicated TCAF genes in neanderthals and humans. *Nature communications* 12,
685 5118 (2021).
- 686 79 Hsieh, P. *et al.* Adaptive archaic introgression of copy number variants and the
687 discovery of previously unknown human genes. *Science* 366, eaax2083 (2019).
- 688 80 Ju, X.-C. *et al.* The hominoid-specific gene TBC1D3 promotes generation of basal
689 neural progenitors and induces cortical folding in mice. *Elife* 5, e18197 (2016).
- 690 81 Dennis, M. Y. *et al.* The evolution and population diversity of human-specific
691 segmental duplications. *Nature ecology & evolution* 1, 0069 (2017).
- 692 82 Prosser, S. L., Sahota, N. K., Pelletier, L., Morrison, C. G. & Fry, A. M. Nek5
693 promotes centrosome integrity in interphase and loss of centrosome cohesion in
694 mitosis. *Journal of Cell Biology* 209, 339-348 (2015).
- 695 83 McAlear, T. S. & Bechstedt, S. The mitotic spindle protein CKAP2 potently
696 increases formation and stability of microtubules. *Elife* 11, e72202 (2022).
- 697 84 Chen, S. *et al.* Paragraph: a graph-based structural variant genotyper for short-read
698 sequence data. *Genome biology* 20, 1-13 (2019).
- 699 85 Liu, Y. *et al.* DNA methylation-calling tools for Oxford Nanopore sequencing: a
700 survey and human epigenome-wide evaluation. *Genome biology* 22, 1-33 (2021).
- 701 86 Katoh, K. & Standley, D. M. MAFFT multiple sequence alignment software version
702 7: improvements in performance and usability. *Molecular biology and evolution* 30,
703 772-780 (2013).
- 704 87 Bouckaert, R. *et al.* BEAST 2: a software platform for Bayesian evolutionary
705 analysis. *PLoS computational biology* 10, e1003537 (2014).
- 706 88 Nguyen, L.-T., Schmidt, H. A., Von Haeseler, A. & Minh, B. Q. IQ-TREE: a fast and
707 effective stochastic algorithm for estimating maximum-likelihood phylogenies.
708 *Molecular biology and evolution* 32, 268-274 (2015).
- 709 89 Mao, Y., Hou, S., Shi, J. & Economo, E. P. TREEasy: An automated workflow to
710 infer gene trees, species trees, and phylogenetic networks from multilocus data. *Mol*
711 *Ecol Resour* 20, (2020).
- 712

713 **Table 1. Primate genome sequence and assembly**
714

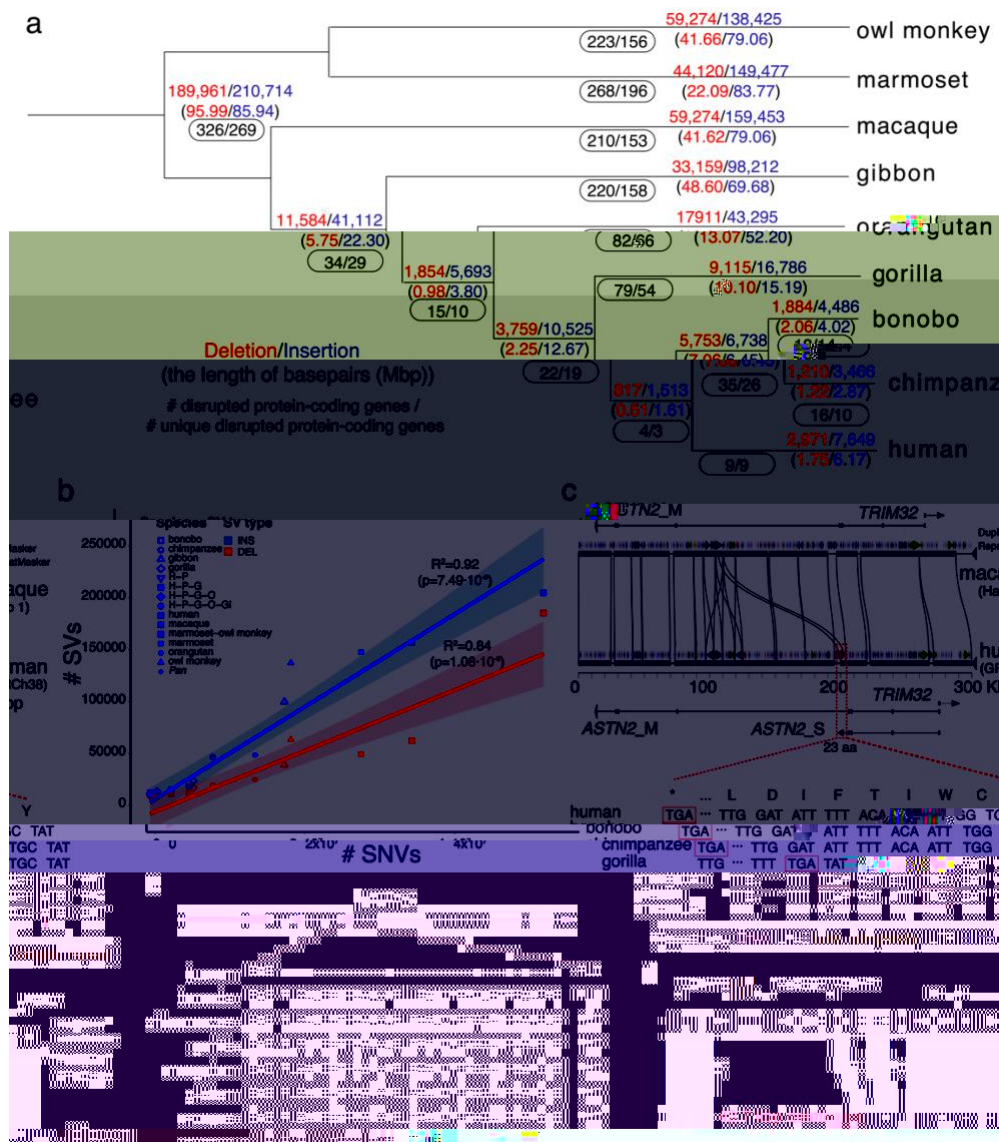
Common name	Scientific name	Individual ID	Sex	CLR raw data and assembly			HiFi raw data and assembly			Iso-Seq (Gbp)	ONT (Gbp)
				reads (coverage)	assembly (contig N50, Mbp)	QV	reads (coverage)	assembly (contig N50, Mbp)	QV hap1/hap2		
Chimpanzee	<i>Pan troglodytes</i> (common chimpanzee)	Clint_PTR	M	117	12.27	39.19	37	66.89 /49.98*	45/44	1.94	294 (178*)
Bonobo	<i>Pan paniscus</i> (pygmy chimpanzee)	Mhudiblu_PP A	F	74	15.06	39.25	39	50.45 /36.22*	47/47	1.38	124*
Gorilla	<i>Gorilla gorilla gorilla</i> (western lowland gorilla)	Kamilah_GGO	F	84.3	9.52	38.72	31	38.19 /37.87*	46/46	1.84	264*
Orangutan	<i>Pongo abelii</i> (Sumatran orangutan)	Susie_PAB	F	94.9	11.07	34.83	43	62.38/ 58.39*	42/42	1.09	272 (126*)
Gibbon	<i>Nomascus leucogenys</i> (northern white-cheeked gibbon)	Asia_NLE	F	92.5*	12.78*	38.65	31*	44.67 /34.99*	43/43	15.25*	97*
Macaque	<i>Macaca mulatta</i> (Rhesus monkey)	AG07107_M MU	F	66	46.61	36.18	29	18.81 /19.01*	51/52	104.58	329 (231*)
Marmoset	<i>Callithrix jacchus</i> (white-tufted-ear marmoset)	CJ1700_CJA	F	66*	25.23*	42.95*	39*	103.97 /87.06*	58/58	18.43*	NA
Owl monkey	<i>Aotus nancymaae</i>	86718_ANA	F	56.3*	9.85*	37.4*	31*	55.92 /44.99*	57/57	NA	91*

715 * New data in this study

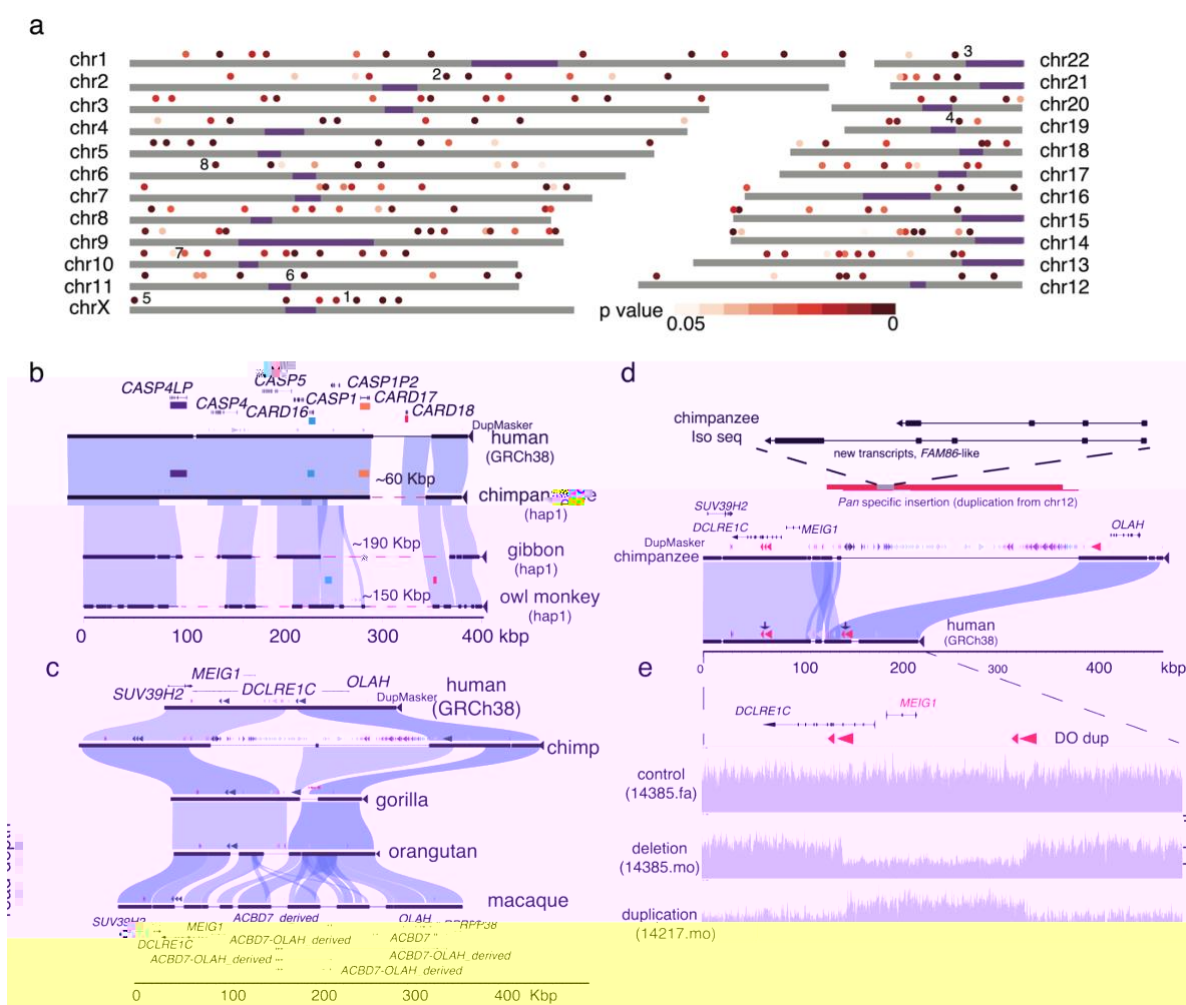
716 **Figures**



717
 718 **Figure 1. Primate phylogeny and SNV divergence between NHPs and humans.** (a) A
 719 primate time-calibrated phylogeny was constructed from a multiple sequence alignment
 720 (MSA) of 81.63 Mbp of autosomal sequence from nine genomes. The estimated species
 721 divergence time (above node) with 95% confidence interval (CI, horizontal blue bar) was
 722 calculated using BEAST2. All nodes have 100% posterior possibility support, and the gene
 723 tree concordance factor (gCF) is indicated (below node). The inset (gray) depicts a maximum
 724 likelihood phylogram generated using IQ-TREE2, which reveals a significantly shorter
 725 branch length in owl monkey, with respect to marmoset. (b) SNV divergence calculated by
 726 mapping HiFi sequence reads to human GRC38 separately for autosomes and the X
 727 chromosome (excluding pseudoautosomal regions). Approximately 85% of the genome was
 728 aligned for Old World monkey and apes and ~60% for New World monkey. The owl
 729 monkey shows significantly less divergence compared to human than the marmoset
 730 (Wilcoxon rank sum test). An analysis using 20 kbp nonoverlapping segments from the
 731 assembly gives almost identical results (Supplementary Figure 4). (c) The percent of trees
 732 showing an alternate tree topology are indicated (percentages are drawn from a total of
 733 302,575 gene trees): 159,546 (52.7%) support the primate topology depicted in panel a.
 734

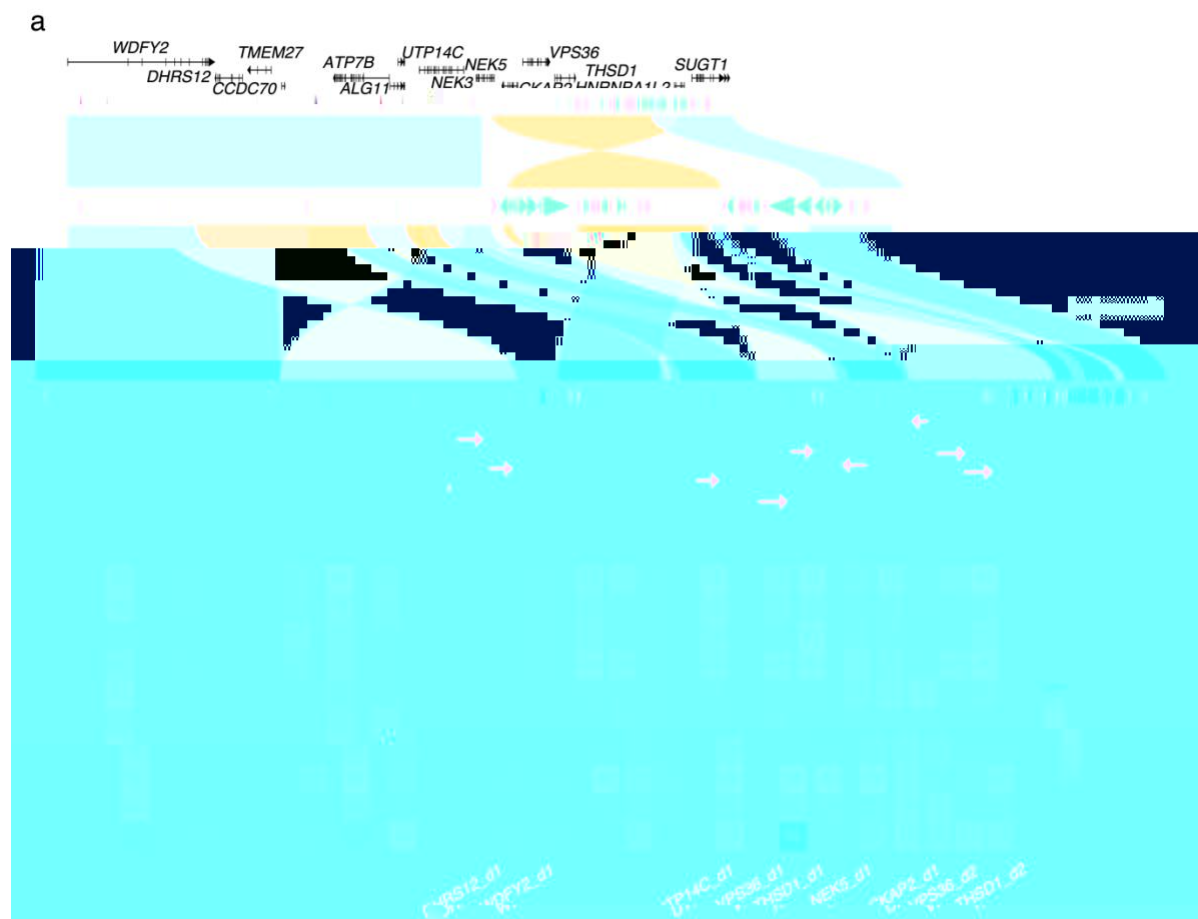


735
 736 **Figure 2. Primate genome structural variation.** (a) The number of fixed structural variants
 737 (SVs) including deletions (red) and insertions (blue) are shown for each branch of the
 738 primate tree (number of events above the line and number of Mbp below). The number of
 739 protein-coding genes based on human RefSeq models are also indicated (black
 740 oval) with the total number of events (first number) and the subset specific to each lineage
 741 (second number). (b) The number of fixed SVs correlates with the accumulation of SNVs in
 742 each lineage (comparison to GRCh38) for both deletions (red) and insertions (blue). (c) An
 743 ape-specific fixed *LI* insertion (shown with a red dashed line box) in the human genome but
 744 not in the macaque genome (Miropeats alignment) serves as an exapted exon of the short
 745 isoform of astrotactin 2, *ASTN2*, in human. The coding sequences of the exon are shown in
 746 the bottom panel. The red triangles represent 1 bp insertion resulting in a frameshift in
 747 gorilla, orangutan, and gibbon. The red box represents the stop codon. (d) A 42.7 kbp
 748 lineage-specific deletion in the gibbon genome (red dashed line) deletes *TAAR2* and seven
 749 enhancers (shown in orange) compared to the human (GRCh38) (Miropeats comparison).
 750 (e) A 90 bp deletion (30 amino acids) human-specific deletion of *NAT16* (NM_001369694)
 751 removes 30 amino acids in humans compared to all other NHPs.



752
753 **Figure 3. Structurally divergent regions (SDRs) of the primate genome.** (a) A schematic
754 of human chromosomes (T2T-CHM13) depicts SDR hotspots where recurrent
755 rearrangements occur in excess. Heat map indicates significance based on simulation model
756 (dark (p=0) to light red (p=0.05)). Centromeres are depicted in purple. Enumerated regions
757 identify specific gene families or regions of biomedical interest (1: *UPRT*, 2: *RGPDs*, 3:
758 *USP41*, 4: *ZNFs*, 5. *IL3RA_2*, 6: *CARDs*, 7: *OLAHA*, and 8: *MHC*). (b) Recurrent deletion of
759 the caspase recruitment domain (*CARD*) gene family. SafFire plot
760 (<https://github.com/mrvollger/SafFire>) shows a ~58 kbp deletion of *CARD18* (orange) in the
761 *Pan* lineage, multiple deletions (~190 kbp total) in gibbon of *CARD16* (blue), *CARD17* (red)
762 and *CARD18*, and multiple deletions ~150 kbp, including *CARD17* (red), in marmoset.
763 (c) SafFire plot of SDR mapping to genes *OLAHA*, *MEIG1*, and *ABCD7* in human shows a
764 large ~250 kbp insertion of segmental duplications (SDs; colored arrowheads) in chimpanzee
765 within the intergenic region between *MEIG1* and *OLAHA*. *OLAHA* is deleted in gorilla by an
766 independent lineage-specific deletion (~30 kbp). Multiple independent insertion events in
767 macaque add ~190 kbp of sequence, including a duplication of *OLAHA* in macaque. Full-
768 length transcript sequencing of macaque using Iso-Seq supports the formation of five novel
769 transcripts, including four *OLAHA-ABCD* fusion events and a derived *ABCD7* (macaque gene
770 models below). (d) The chimpanzee-specific 250 kbp SD from chromosome 12 creates a
771 novel multi-exonic gene model supported by Iso-Seq transcript sequencing in chimpanzee
772 (upper panel) with an unmethylated promoter (Supplementary Figure 36). The insertion

773 simultaneously deletes one of two directly orientated (DO) SDs in chimpanzee. (e) In
774 humans, the DO repeats associate with the breakpoints of recurrent deletions and
775 duplications of the spermiogenesis gene *MEIG1*. Two females carrying a deletion and a
776 duplication (as measured by sequence read depth) are depicted from a population sample of
777 19,584 genomes (CCDG, <https://ccdg.rutgers.edu/>). The carrier frequencies for microdeletion
778 and microduplication in control samples are 0.026% and 0.189%, respectively.
779



780
781 **Figure 4. Marmoset-specific genes in a SDR.** (a) SafFire plot comparing the organization
782 of a gene-rich region of ~1.1 Mbp in human (middle), owl monkey (top), and marmoset
783 (bottom) genomes. Human and marmoset differ mainly by a large 250 kbp inversion (orange)
784 associated with the addition of 150 kbp of SD at the boundary of the inversion in humans
785 (colored arrowheads). The corresponding region in marmoset has expanded by ~400 kbp due
786 to inversion and marmoset-specific SDs creating marmoset-specific paralogs (red arrows) of
787 *CCDC70*, *TMEM27*, *DHR12*, *UTP14C*, *THSD1*, *VPS36*, *NEK5* and *CKAP2*. (b) Iso-Seq
788 full-length non-chimeric transcript sequencing from 10 marmoset primary tissues confirms
789 transcription of 8/10 of the paralogous copies and the maintenance of an open-reading frame
790 in at least six of these marmoset-specific gene candidates.
791

