# The germline mutational process in rhesus macaque and its implications for phylogenetic dating

Lucie A. Bergeron [1*], Søren Besenbacher [2], Jaco Bakker [3], Jiao Zheng [4,5], Panyi Li [4], George Pacheco [6], Mikkel-Holger S. Sinding [7,8], Maria Kamilari [1], M. Thomas P. Gilbert [6,9], Mikkel H. Schierup [10] and Guojie Zhang [1,4,11,12*]

[1] Section for Ecology and Evolution, Department of Biology, University of Copenhagen, Copenhagen, Denmark

[2] Department of Molecular Medicine, Aarhus University, Aarhus, Denmark

[3] Animal Science Department, Biomedical Primate Research Centre, Rijswijk, Netherlands

[4] BGI-Shenzhen, Shenzhen 518083, Guangdong, China

[5] BGI Education Center, University of Chinese Academy of Sciences, Shenzhen 518083, Guangdong, China

[6] Section for Evolutionary Genomics, The GLOBE Institute, University of Copenhagen, Copenhagen, Denmark

[7] Trinity College Dublin, Dublin, Ireland

[8] Greenland Institute of Natural Resources, Nuuk, Greenland

[9] Department of Natural History, NTNU University Museum, Norwegian University of Science and Technology (NTNU), NO-7491 Trondheim, Norway

[10] Bioinformatics Research Centre, Aarhus University, Aarhus, Denmark

[11] State Key Laboratory of Genetic Resources and Evolution, Kunming Institute of Zoology, Chinese Academy of Sciences, Kunming 650223, China

[12] Center for Excellence in Animal Evolution and Genetics, Chinese Academy of Sciences, Kunming 650223, China

[*] Corresponding author

E-mail: guojie.zhang@bio.ku.dk or lucie.a.bergeron@gmail.com

1

## Abstract

Understanding the rate and pattern of germline mutations is of fundamental importance for understanding evolutionary processes. Here we analyzed 19 parent-offspring trios of rhesus macaques (*Macaca mulatta*) at high sequencing coverage of ca. 76X per individual, and estimated an average rate of $0.77 \times 10^{-8}$ *de novo* mutations per site per generation (95 % CI: $0.69 \times 10^{-8}$ - $0.85 \times 10^{-8}$). By phasing 50 % of the mutations to parental origins, we found that the mutation rate is positively correlated with the paternal age. The paternal lineage contributed an average of 81 % of the *de novo* mutations, with a trend of an increasing male contribution for older fathers. About 3.5 % of *de novo* mutations were shared between siblings, with no parental bias, suggesting that they arose from early development (postzygotic) stages. Finally, the divergence times between closely related primates calculated based on the yearly mutation rate of rhesus macaque generally reconcile with divergence estimated with molecular clock methods, except for the Cercopithecidae/Hominoidea molecular divergence dated at 52 Mya using our new estimate of the yearly mutation rate.

## Introduction

Germline mutations are the source of heritable disease and evolutionary adaptation. Thus, having precise estimates of germline mutation rates is of fundamental importance for many fields in biology, including searching for *de novo* disease mutations (Acuna-Hidalgo et al. 2016; Oliveira et al. 2018), inferring demographic events (Lapierre et al. 2017; Zeng et al. 2018), and accurate dating of species divergence times (Teeling et al. 2005; Ho and Larson 2006; Pulquério and Nichols 2007). Over the past ten years, new sequencing techniques have allowed deep sequencing of individuals from the same pedigree, enabling direct estimation of the *de novo* mutation rate for each generation, and precise estimation of the individual parental contributions to germline mutations across the whole genome. Most such studies have been conducted on humans, using large pedigrees with up to 3000 trios (Jónsson et al. 2017; Halldorsson et al. 2019), leading to a consensus estimate of ~1.25 x $10^{-8}$ *de novo* mutation per site per generation,

2

60    with an average parental age of ~ 29 years, leading to a yearly rate of 0.43 x $10^{-9}$ *de novo*

61    mutation per site per year and most variation between trios explained by the age of the parents

62    (Awadalla et al. 2010; Roach et al. 2010; Kong et al. 2012; Neale et al. 2012; Wang and Zhu

63    2014; Besenbacher et al. 2015; Rahbari et al. 2016; Jónsson et al. 2017; Maretty et al. 2017).

64    The observed increases in the mutation rate with paternal age in humans and other primates

65    (Venn et al. 2014; Jónsson et al. 2017; Thomas et al. 2018) has generally been attributed to

66    errors during replication (Li et al. 1996; Crow 2000). In mammalian spermatogenesis, primordial

67    germ cells go through meiotic divisions, to produce stem cells by the time of puberty. After this

68    time, stem cell divisions occur continuously throughout the male lifetime. Thus, human

69    spermatogonial stem cells have undergone 100 to 150 mitoses in a 20 years old male, and ~ 610

70    mitoses in a 40 years old male (Acuna-Hidalgo et al. 2016), leading to an additional 1.51 *de novo*

71    mutations per year increase in the father's age (Jónsson et al. 2017). Female age also seems to

72    affect the mutation rate in humans, with 0.37 mutations added per year (Jónsson et al. 2017).

73    This maternal effect cannot be attributed to replication errors, as different from spermatogenesis,

74    female oocytogenesis occurs during embryogenesis process and is already finished before birth

75    (Byskov 1986). Moreover, there seems to be a bias towards males in contribution to *de novo*

76    mutations, as the paternal to maternal contribution is 4:1 in human and chimpanzee (Venn et al.

77    2014; Jónsson et al. 2017). One recent study proposed that damage-induced mutations might be a

78    potential explanation for the observation of both the maternal age effect and the male-bias also

79    present in parents reproducing right after puberty when replication mutations should not have

80    accumulated yet in the male germline (Gao et al. 2019). Parent-offspring analyses can also be

81    used to distinguish mutations that are caused by gametogenesis from mutations that emerge in

82    postzygotic stages (Acuna-Hidalgo et al. 2015; Scally 2016). While germline mutations in

83    humans are relatively well studied, it remains unknown how much variability exists among

84    primates on the contribution of replication errors to *de novo* mutations, the parental effects, and

85    the developmental stages at which these mutations are established (postzygotic or

86    gametogenesis).

87    Up until now, the germline mutation rate has only been estimated using pedigrees in few non-

88    human primate species, including chimpanzee (*Pan troglodytes*) (Venn et al. 2014; Tatsumoto et

89    al. 2017; Besenbacher et al. 2019), gorilla (*Gorilla gorilla*) (Besenbacher et al. 2019), orangutan

3

90  (*Pongo abelii*) (Besenbacher et al. 2019), African green monkey (*Chlorocebus sabaeus*) (Pfeifer

91  2017), owl monkey (*Aotus nancymaae*) (Thomas et al. 2018) and recently rhesus macaque

92  (*Macaca mulatta*) (Wang et al. 2020). The mutation rate of baboon (*Papio anubis*) (Wu et al.

93  2019) and grey mouse lemur (*Microcebus murinus*) (Campbell et al. 2019) have also been

94  estimated in preprinted studies. To precisely call *de novo* mutations in the offspring, collecting

95  and comparing the genomic information of the pedigrees is a first essential step for detecting

96  mutations only present in offspring but not in either parent. Next, the *de novo* mutations need to

97  be separated from sequencing errors or somatic mutations, which cause false-positive calls.

98  Because mutations are rare events, detecting *de novo* mutations that occur within a single

99  generation requires high sequencing coverage in order to cover a majority of genomic regions

100  and identify the false-positives. Furthermore, the algorithms used to estimate the mutation rate

101  should take false-negative calls into account. However, a considerable range of sequencing depth

102  (ranging from 18X (Pfeifer 2017) to 120X (Tatsumoto et al. 2017)) has been applied in many

103  studies for estimation of mutation rate. Different filtering methods have been introduced to

104  reduce false-positives and false-negatives but the lack of standardized methodology makes it

105  difficult to assess whether differences in mutation rate estimates are caused by technical or

106  biological variability. In addition, most studies on non-human primates used small pedigrees

107  with less than ten trios, which made it difficult to detect any statistically significant patterns over

108  *de novo* mutation spectra.

109  Studying non-human primates could help us understanding whether the mutation rate is affected

110  by life-history traits such as mating strategies or the age of reproduction. The variation in

111  mutation rate among primates will also be useful for re-calibrating the speciation times across

112  lineages. The sister group of Hominoidea is Cercopithecidae, including the important biomedical

113  model species, rhesus macaque (*Macaca mulatta*), which share 93 % of its genome with humans

114  (Gibbs et al. 2007). This species has a generation time estimate of ~ 11 years (Xue et al. 2016),

115  and their sexual maturity is much earlier than in humans with females reaching maturity around

116  three years old, while males mature around the age of 4 years (Rawlins and Kessler 1986). While

117  female macaques generally start reproducing right after maturation, males rarely reproduce in the

118  wild until they reach their adult body size, at approximately eight years old (Bercovitch et al.

119  2003). They are also a promiscuous species, and do not form pair bonds, but reproduce with

4

120    multiple individuals. These life-history traits, along with their status as the closest related

121    outgroup species of the hominoid group, make the rhesus macaque an interesting species for

122    investigating the differences and common features in mutation rate processes across primates.

123    In this study, we, produced high depth sequencing data for 33 rhesus macaque individuals (76X

124    per individual) representing 19 trios. This particular dataset consists of a large number of trios,

125    each with high coverage sequencing, and allowed us to test different filter criteria and choose the

126    most appropriate ones to estimate the species mutation rate with high confidence. With a large

127    number of *de novo* mutations phased to their parents of origins, we can statistically assess the

128    parental contribution and the effect of the parental age. We characterize the type of mutations

129    and their location on the genome to detect clusters and shared mutations between siblings.

130    Finally, we use our new estimate to infer the effective population size and date their divergence

131    time from closely related primate species.

132

133    **Results**

134

135    **Estimation of mutation rate for 19 trios of rhesus macaques**

136    To produce an accurate estimate for the germline mutation rate of rhesus macaques, we

137    generated high coverage (76 X per individual after mapping, min 64 X, max 86 X) genome

138    sequencing data for 19 trios of two unrelated families (Fig 1). The first family consisted of two

139    reproductive males and four reproductive females, and the second family had one reproductive

140    male and seven reproductive females. In the first family, the pedigree extended over a third

141    generation in two cases. The promiscuous mating system of rhesus macaques allowed us to

142    follow the mutation rates in various ages of reproduction, and compare numerous full siblings

143    and half-siblings.

144    We developed a pipeline for single nucleotide polymorphisms (SNP) calling with multiple

145    quality control steps involving the filtering of reads and sites (see Methods). For each trio, we

146    considered candidate sites as *de novo* mutations when i) both parents were homozygotes for the

147    reference allele, while the offspring was heterozygous with 30 % to 70 % of its reads supporting

148    the alternative allele, and ii) the three individuals passed the depth and genotype quality filters (see

5

149 Methods). These filters were calibrated to ensure a low rate of false-positives among the

150 candidate *de novo* mutations.

151 We obtained an unfiltered set of 12,785,386 average candidate autosomal SNPs per trio (se =

152 26,196), of which a total of 177,227 were potential Mendelian violations (average of 9,328 per

153 trio; se= 106). Of these, 744 SNPs passed the filters as *de novo* mutations, ranging from 25 to 59

154 for each trio and an average of 39 *de novo* mutations per trio (se = 2) (see S1 Table). We

155 manually curated all mutations using IGV on bam files and found that 663 mutations

156 convincingly displayed as true positives. This leaves a maximum of 10.89 % (81 sites) that could

157 be false-positives due to the absence of the variant in the offspring or presence of the variant in

158 the parents (see S1 Fig and the 81 curated mutations in supplementary). Most of those sites were

159 in dinucleotide repeat regions or short tandem repeats (56 sites), while others were in non-

160 repetitive regions of the genome (25 sites). The manual curation may have missed the

161 realignment executed during variant calling. Thus, in the absence of objective filters, we decided

162 to keep these regions in the estimate of mutation rate but corrected the number of mutations for

163 each trio with a false-positive rate (see equation 1 in Methods section).

164 To confirm the authenticity of the *de novo* mutations, we performed PCR experiments for all

165 candidate *de novo* mutations from one trio before manual correction. We designed primers to a

166 set of 39 *de novo* candidates among which 3 *de novo* mutations assigned as spurious from the

167 manual inspection. Of these, 24 sites were successfully amplified and sequenced for all three

168 individuals i.e mother, father, and offspring, including 1 of the spurious sites. Among those

169 sequenced sites, 23 were correct, only one was wrong (S2 Fig). This invalidated candidate was

170 the spurious candidate removed by manual curation, therefore supporting our manual curation

171 method. The PCR validation results suggested a lower false-positive rate of 4.2 % before manual

172 curation. As the PCR validation was done only on 24 candidates we decided to keep a strict

173 false-positive rate of 10.89 % found by manual curation.

174 We then estimated the mutation rate, per site per generation, as the number of mutations

175 observed, and corrected for false-positive calls, divided by the number of callable sites. The

176 number if callable sites for each trio ranged from 2,334,764,487 to 2,359,040,186, covering on

177 average 88 % of the autosomal sites of the rhesus macaque genome. A site was defined as

178  callable when both parents were homozygotes for the reference allele, and all individuals passed

179  the depth and genotype quality filters at that site. As callability is determined using the base-pair

180  resolution vcf file, containing every single site of the genome, all filters used during calling were

181  taken into account during the estimation of callability, except for the site filters and the allelic

182  balance filter. We then corrected for false-negative rates, calculated as the number of "good"

183  sites that could be filtered away by both the site filters and allelic balance filters - estimated at

184  4.02 % (see equation 1 in Methods section). Another method to estimate the false-negative rate

185  is to simulate mutations on the bam files and evaluate the detection rate after passing through all

186  filters. On 552 randomly simulated mutations among the 19 offsprings, 545 were detected as *de*

187  *novo* mutations, resulting in a false-negative rate of 1.27 %. The 7 remaining mutations were

188  filtered away by the allelic balance filter only, which can be explained by the reads filtering in

189  the variant calling step. This result might be underestimated due to the methodological limitation

190  of simulating *de novo* mutations, yet, it ensures that a false-negative rate of 4.02 % is not out of

191  range. Thus, the final estimated average mutation rate of the rhesus macaques was $0.77 \times 10^{-8}$ *de*

192  *novo* mutations per site per generation (95 % CI $0.69 \times 10^{-8}$ - $0.85 \times 10^{-8}$). We removed the 81

193  sites that, based on manual curation, could represent false-positive calls from the following

194  analyses (see the 663 *de novo* mutations in S2 Table).

195

**Parental contribution and age impact to the *de novo* mutation rate**

197  We observed a positive correlation between the paternal age and the mutation rate in the

198  offspring (adjusted $R^2 = 0.23$; $P = 0.021$; regression: $\mu = 1.022 \times 10^{-9} + 5.393 \times 10^{-10} \times$

199  $age_{paternal}$; $P = 0.021$; Fig 2A). We also detected a slight positive correlation with the maternal

200  age, though not significant (adjusted $R^2 = 0.09$; $P = 0.111$; regression: $\mu = 6.200 \times 10^{-9} + 1.818 \times$

201  $10^{-10} \times age_{maternal}$; $P = 0.111$; Fig 2B). A multiple regression of the mutation rate on paternal and

202  maternal age resulted in this formula: $\mu_{Rhesus} = 1.355 \times 10^{-9} + 7.936 \times 10^{-11} \times age_{maternal} + 4.588 \times$

203  $10^{-10} \times age_{paternal}$ (P = 0.06), where $\mu_{Rhesus}$ is the mutation rate for the species.

204  We were able to phase 337 mutations to their parent of origin, which accounted for more than

205  half of the total number of *de novo* mutations (663). There is a significant male bias in the

7

206    contribution of *de novo* mutations, with an average of 80.6 % paternal *de novo* mutations (95 %

207    CI 76.6 % - 84.6 %; T = 22.62, DF = 36, P < $2.2 \times 10^{-16}$; Fig 2C). Moreover, with more than half

208    of the *de novo* mutations phased to their parent of origin, we were able to disentangle the effect

209    of the age of each parent on mutation rate independently (Fig 2D). By assuming that the ratio of

210    mutations phased to a particular parent was the same in the phased mutations than in the

211    unphased ones, we could predict the total number of mutations given by each parent. For

212    instance, if an offspring had 40 *de novo* mutations and only half were phased, with 80 % given

213    from its father, we would apply this ratio to the total number of mutations in this offspring,

214    ending up with 32 *de novo* mutations from its father and eight from its mother. This analysis

215    suggested a stronger male age effect to the number of mutations (adjusted $R^2$ = 0.41, P = 0.002),

216    and a similar, non significant maternal age effect (adjusted $R^2$ = -0.01, P = 0.38). The two

217    regression lines meet around the age of sexual maturity (3 years for females and 4 years for

218    males), which is consistent with a similar accumulation of *de novo* mutations during the

219    developmental process from birth to sexual maturity in both sexes, but the variances on the

220    regression line slopes are large (see Fig 2C and S3 Fig for the same analysis with a Poisson

221    regression). Using these two linear regressions, we can predict the number of *de novo* mutations

222    in the offspring based on the age of each parent at the time of reproduction: *nb of mutations $_{Rhesus}$*

223    $= 4.6497 + 0.3042 \times age_{maternal} + 4.8399 + 1.8364 \times age_{paternal}$, where *nb of mutations $_{Rhesus}$* is the

224    number of *de novo* mutations for the given trio. The expected mutation rates calculated using the

225    two different regression models show similar correlations with the observed mutation rate ($R^2$ =

226    0.54, P = 0.016 for the first regression and $R^2$ = 0.54, P = 0.016 for the upscaled one, see S4 Fig).

227    However, on the first regression on the mutation rate, the maternal age effect may be confounded

228    by the paternal age, as maternal and paternal age are correlated in our dataset, yet, non-

229    significantly ($R^2$ = 0.38, P = 0.106, see S5 Fig). The upscaled regression unravels the effect of

230    the parental age independently from each other. This regression can also be used to infer the

231    contribution of each parent at different reproductive age. For instance, if both parents reproduce

232    at 5 years old, based on the upscaled regression, the father is estimated to give ~ 14 *de novo*

233    mutations (95 % CI:6 - 22) and the mother ~ 6 *de novo* mutations (95 % CI:3 - 10),

234    corresponding to a contribution ratio from father to mother of 2.3:1 at 5 years old. If they

235    reproduce at 15 years old, this ratio would be 3.6:1 with males giving ~ 32 *de novo* mutations

8

236   (95 % CI: 29 – 36) and females ~ 9 *de novo* mutations (95 % CI: 4 – 14). It seems that the male

237   bias increases with the parental age, yet, our model was based on too few data points in early

238   male reproductive ages to reach a firm conclusion. For the two extended trios for which a second

239   generation is available, we looked at the proportion of *de novo* mutations in the first offspring

240   that were passed on to the third generation - the third generation inherited a heterozygote

241   genotype with the alternative allele being the *de novo* mutation. In one case, 66 % of the *de novo*

242   mutations in the female (Heineken) were passed to her daughter (Hoegaarde), while in another

243   case, 40 % of the *de novo* mutations in the female (Amber) were passed to her son (Magenta).

244   These deviations from the expected 50 % inheritance rate are not statistically significant

245   (Binomial test; $P_{Hoegaarde}$ = 0.14 and $P_{Magenta}$ = 0.27).

246

### Characterizations of *de novo* mutations

248   We characterized the type of *de novo* mutations and found that transition from a strong base to

249   weak base ($G > A$ and $C > T$) were most common (332/663), with 43 % of those mutations

250   located in CpG sites (Fig 3A). In total, 23.2 % (154/663) of the *de novo* mutations were located in

251   CpG sites. This is slightly higher than what has been found in humans, for which 19 % of the *de*

252   *novo* mutations are in CpG sites (Besenbacher et al. 2015), but not significantly (human: $X^2 =$

253   2.774, df = 1, P = 0.096). Moreover, 32.1 % (144/448) of the transition mutations ($A > G$ and $C$

254   $> T$) were in CpG sites, higher than what has been found in chimpanzee, with 24 % of the

255   transition *de novo* mutations in CpG sites (Venn et al. 2014). The transition to transversion ratio

256   (ti/tv) was 2.08, which is similar to the ratio observed in other species (human: ti/tv ~ 2.16 (Yuen

257   et al. 2016); human ti/tv ~ 2.2 (Wang and Zhu 2014); chimpanzee: ti/tv ~ 1.98 (Tatsumoto et al.

258   2017). The 663 *de novo* mutations showed some clustering in the genome (Fig 3B and S6 Fig).

259   Across all trios, we observed 11 clusters, defined as windows of 20,000 bp where more than one

260   mutation occurred in any individual, involving 23 mutations. Four clusters were made of

261   mutations from a single individual, accounting for eight mutations (Fig 3B). Overall, 3.47 % of

262   the *de novo* mutations were located in clusters, and 1.21 % were mutations within the same

263   individual located in a cluster, which is significantly lower than the 3.1 % reported in humans

264   (Besenbacher et al. 2016) ($X^2$ = 7.35, DF = 1, P = 0.007; S7 Fig, S3 Table). We observed 23

9

265 mutations occurring recurrently in more than one related individual (Table 1), which accounted

266 for 3.5 % of the total number of *de novo* mutations (23/663) and 1.5 % of sites (10/650 unique

267 sites). Four *de novo* mutations (2 sites) were shared between half-siblings on the maternal side,

268 and 19 (8 sites) were shared between half-siblings on the paternal side. However, there was no

269 significant difference between the proportion of mutations shared between pairs of individuals

270 related on the maternal side (9 pairs, 0.70 % shared), and pairs related on their paternal side (53

271 pairs, 0.80 % shared; Fisher's exact test P = 1). In 6 sites, the phasing to the parent of origin

272 confirmed that the mutation was coming from the common parent for at least one individual

273 (Tab. 1). Moreover, the phasing was never inconsistent by attributing a shared *de novo* mutation

274 to the other parent than the parent in common. However, 5 shared sites did appear as mosaic in

275 the common parent, with a maximum of 5 % of the reads of the father supporting the alternative

276 allele (4 out of 80 reads).. Nine of the *de novo* mutations (1.4 % of the total *de novo* mutations)

277 were located in coding sequences (CDS regions), which is close to the overall proportion of

278 coding sequences region (1.2%) in the whole macaque genome. Eigth out of those eight

279 mutations were non-synonymous.

280 **Table 1 – Six mutations shared between related individuals**.

| Chrom | Position | Ref | Alt | Sibling 1 | Phasing [a] | Sibling 2 | Phasing | Sibling 3 | Phasing | Sibling 4 | Phasing | Common parent | Name parent |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| chr2 | 101979137 | G | A | Khan | P | Delta | U | | | | | father | Smack |
| chr6 | 132663101 | A | T | Amber | U | Babet | M | | | | | mother | Mayke |
| chr7 | 60635102 | G | T | Sir | U | Honoria | U | | | | | father | Smack |
| chr7 | 116648579 | G | A | Amber | M | Babet | U | | | | | mother | Mayke |
| chr9 | 32544257 | C | T | Hoegaarde | U | Djembe | U | Babet | U | Magenta | U | father | Poseidon |
| chr10 | 65163492 | G | A | Khan | P | Delta | P | | | | | father | Smack |
| chr15 | 35463257 | C | T | Leffe | U | Babet | U | Magenta | U | | | father | Poseidon |
| chr17 | 88174686 | C | A | Hoegaarde | P | Babet | P | | | | | father | Poseidon |
| chr19 | 7047030 | C | T | Leffe | U | Djembe | U | | | | | father | Poseidon |
| chr19 | 15861061 | C | T | Bavaria | P | Lithium | U | | | | | father | Noot |

281 a: P: paternal; M: maternal; U: unphased

282

### Molecular dating with trio-based mutation rate

284 Based on our inferred mutation rate and the genetic diversity of Indian rhesus macaques ($\pi =$

285 0.00247) estimated using whole genomic sequencing data from more than 120 unrelated wild

286 individuals (Xue et al. 2016), we calculated the effective population size ($N_e$) of rhesus

287 macaques to be 79,874. This is similar to the $N_e = 80,000$ estimated previously using $\mu = 0.59 \times$

288 $10^{-8}$ from hippocampal transcriptome and H3K4me3-marked DNA regions from 14 individuals

10

289    (Yuan et al. 2012), yet higher than $N_e = 61,800$ estimated using $\mu = 1 \times 10^{-8}$ with 120 individual

290    full genome data (Xue et al. 2016). Assuming a generation time of 11 years and an average

291    reproduction age of 10 years for females and 12 years for males, the yearly mutation rate of

292    rhesus macaques was calculated based on our regression model of the number of mutations given

293    by males and females independently, and the average callability (see equation 2 in the Methods

294    section). As captive animals usually reproduce later than in the wild, which could impact the

295    average mutation rate per generation, we used the regression instead of the mutation rate per

296    generation to correct for this possible bias. The yearly mutation rate of rhesus macaques in our

297    calculation was $0.62 \times 10^{-9}$ per site per year, almost 1.5 times that of humans (Jónsson et al.

298    2017).

299    Given a precise evolutionary mutation rate is essential for accurate calibration of molecular

300    divergence events between species, we used the mutation rate we inferred for rhesus macaques to

301    re-date the phylogeny of closely-related primate species with full genome alignment available

302    (Moorjani et al. 2016) (Fig 4A). The molecular divergence time ($T_D$) is the time since an

303    ancestral lineage started to split into two descendant lineages, and can be inferred from the

304    genetic divergence between the two descendant lineages and the mutation rate. The speciation

305    time ($T_S$) is a younger event that implies no more gene flow between lineages  (Steiper and

306    Young 2008). On the backward direction, the alleles of two descendant lineages are randomly

307    sampled from their parents until going back to the most recent common ancestor (Rosenberg and

308    Nordborg 2002). This stochastic event, known as the coalescent, depends on the population

309    sizes, being slower in a large population (Kingman, 1982). Thus, from the divergence time, the

310    speciation time can be inferred given the rate of coalescence (see equation 3 in the Method

311    section). We also compared our results to those of previous dating attempts based on molecular

312    phylogenetic trees calibrated with fossils records (Fig 4B). We found that the two methods concur

313    for the most recent events. Specifically, we estimated that the *Macaca mulatta* and *Macaca*

314    *fascicularis* genomes had already diverged around 3.90 million years ago (Mya) (95 % CI: 3.46

315    – 4.46), which is slightly older than previous estimates using the molecular clock calibrated with

316    fossils, as the molecular divergence of the two species has been estimated at 3.44 Mya with

317    mitochondrial data (Pozzi et al. 2014) and 3.53 Mya from nuclear data (Perelman et al. 2011).

11

318    We estimated a speciation event between the two species 2.14 Mya after the coalescent time,

319    also consistent with previous findings of a most common recent ancestor to the two populations

320    of the rhesus macaque, the Chinese and the Indian population, around 1.94 Mya based on

321    coalescent simulations (Hernandez et al. 2007). For the next node, the molecular clock seems to

322    differ between mitochondrial and nuclear data, as the divergence time for the Papionini group

323    into the *Papio* and *Macaca* genera has been estimated to 8.13 Mya using nuclear data (Perelman

324    et al. 2011), and 12.17 Mya with mitochondrial data (Pozzi et al. 2014). We estimated a

325    divergence time between these two genera of 13.17 Mya (95 % CI: 11.70 – 15.07). For earlier

326    divergence events, our estimated divergence times are more ancient than previous reports. For

327    instance, we estimated that the Cercopithecini and Papionini diverged 19.86 Mya (95 % CI:

328    17.64 – 22.71), while other studies had calculated 11.55 Mya using nuclear data (Perelman et al.

329    2011), and 14.09 Mya using mitochondrial data (Pozzi et al. 2014). Finally, the divergence

330    between Cercopithecidae and Hominoidea has been reported between 25 and 30 Mya (Stewart

331    and Disotell 1998; Moorjani et al. 2016), with an estimation of 31.6 Mya using the nuclear

332    molecular clock (Perelman et al. 2011) and 32.12 Mya using the mitochondrial one (Pozzi et al.

333    2014). Our dating of the divergence time between the Cercopithecidae and Hominoidea of 52.31

334    Mya (95 % CI: 46.47 – 59.83) is substantially older than previous estimates. However, the

335    estimated speciation time inferred based on the ancestral population size, suggested a speciation

336    of the Catarrhini group into two lineages 44.50 Mya (Fig 4B).

337

## Discussion

339

340    Despite many efforts to accurately estimate direct *de novo* mutation rates, it is still a

341    challenging task due to the rare occurrence of *de novo* mutations, and the small sample size that

342    is often available. Sequencing coverage is known to be a significant factor in affecting false-

343    positive (FP), and false-negative (FN) calls when detecting *de novo* mutation (Acuna-Hidalgo

344    et al. 2016; Tatsumoto et al. 2017). A minimal sequencing coverage at 15X was recommended

345    for SNPs calling (Song et al. 2016). However, such coverage cannot provide sufficient power

346    to reduce FPs because the lower depth threshold cannot preclude Mendelian violations due to

12

347   sequencing errors. Moreover, a larger portion of the genome would be removed in the

348   denominator at low depth in order to reduce the FN. While most studies on direct estimation of

349   mutation rate use 35-40X coverage (Jónsson et al. 2017; Thomas et al. 2018; Besenbacher et al.

350   2019), their methods to reduce FP and FN differ. Some studies use the deviation from 50 % of

351   the *de novo* mutation pass to the next generation to infer the false-positive rate (Jónsson et al.

352   2017; Thomas et al. 2018). Others use probabilistic methods to access the callability

353   (Besenbacher et al. 2019), or simulation of known mutation to control the pipeline quality

354   (Pfeifer 2017). Differences in methods likely impact the calculated rate. Here, we produced

355   sequences at 76X coverage, which allows us to apply conservative filtering processes, while

356   still obtaining high coverage (88 %) of the autosomal genome region when inferring *de novo*

357   mutations. To our knowledge, only one other study has used very high coverage (120X per

358   individuals), on a single trio of chimpanzees (Tatsumoto et al. 2017). Such high coverage

359   allowed us to achieve a false-positive rate below 10.89 % and within the regions we deemed

360   callable, we calculated a low false-negative rate of 4.02 %.

361   Our estimated rate is higher than the $0.58 \times 10^{-8}$ *de novo* mutations per site per generation

362   estimated in a preprint report (Wang et al. 2020). The difference should be mainly attributed to

363   the fact that they sequenced the offspring of younger parents (average parental age of 7.1 years

364   for females and 7.8 years for males compared to 8.4 years for females and 12.4 years for males

365   in this study). Using our regression from the phased mutation, we estimated a mutation rate of

366   $0.51 \times 10^{-8}$ per site per generation, when males reproduce at 7.8 years and females reproduce at

367   7.1 years old. Moreover, using their regression based on the age of puberty and the increase of

368   paternal mutation per year, Wang and collaborators estimated a per generation rate of $0.71 \times$

369   $10^{-8}$ mutations when males reproduce at 11 years, and a yearly rate of $0.65 \times 10^{-9}$ mutations

370   per site per year, which is approx 5 % higher than our estimate of $0.62 \times 10^{-9}$ (2020). This

371   difference may be due to any combination of stochasticity, differences in *de novo* mutation rate

372   pipelines (callability estimate, false-negative rate, and false-positive rate estimate) and different

373   models for converting pedigree estimates to yearly rates. Our combination of high coverage

374   data and a large number of trios allowed us to gain high confidence estimates of the germline

375   mutation rate of rhesus macaques at around $0.77 \times 10^{-8}$ *de novo* mutation per site per

376   generation, ranging from $0.49 \times 10^{-8}$ to $1.16 \times 10^{-8}$. This is similar to the mutation rate

13

377  estimated for other non-Hominidae primates; $0.81 \times 10^{-8}$ for the owl monkey (*Aotus*

378  *nancymaae*) (Thomas et al. 2018) and $0.94 \times 10^{-8}$ for the African green monkey (*Chlorocebus*

379  *sabaeus*) (Pfeifer 2017), while all Hominidae seem to have a mutation rate that is higher than 1

380  $\times 10^{-8}$ *de novo* mutation per site per generation (Jónsson et al. 2017; Besenbacher et al. 2019).

381  However, if we count for the *de novo* mutation per site per year, the rate of rhesus macaque

382  $(0.62 \times 10^{-9})$ is almost 1.5-fold the human one of $0.43 \times 10^{-9}$ mutation per sites per year

383  (Jónsson et al. 2017).

384  One of the main factors affecting the mutation rate within the species is the paternal age at the

385  time of reproduction, which was attributed to the accumulation of replication-driven mutations

386  during spermatogenesis (Drost and Lee 1995; Li et al. 1996; Crow 2000), and has been

387  observed in many other primates (Venn et al. 2014; Jónsson et al. 2017; Maretty et al. 2017;

388  Thomas et al. 2018; Besenbacher et al. 2019). In rhesus macaques, the rate at which germline

389  mutation increases with paternal age seems faster than in humans; we inferred 1.84 mutations

390  more per year for the rhesus macaque father (95% CI 0.77 – 2.90  for an average callable

391  genome of 2.35 Mb), compared to 1.51 in humans (95% CI 1.45–1.57 for an average callable

392  genome of 2.72 Mb) (Jónsson et al. 2017). For females, there is less difference, with 0.30 more

393  mutations per year for the mother in rhesus macaque (95% CI  -0.41 – 1.02), and 0.37 more per

394  year in human mothers (95% CI 0.32–0.43) (Jónsson et al. 2017). In rhesus macaques, males

395  produce a larger number of sperm cells per unit of time ($23 \times 10^{6}$ sperm cells per gram of testis

396  per day (Amann et al. 1976)) than humans ($4.4 \times 10^{6}$ sperm cells per gram of testis per day

397  (Amann and Howards 1980)). This could imply a higher number of cell division per unit of

398  time in rhesus macaques and thus more replication error during spermatogenesis. This is also

399  consistent with the generation time effect which stipulates that an increase in generation time

400  would decrease the number of cell division per unit of time as well as the yearly mutation rate

401  assuming that most mutations arise from replication errors (Wu and Lit 1985; Goodman et al.

402  1993; Ohta 1993; Li et al. 1996; Ségurel et al. 2014; Scally 2016). Indeed, humans have a

403  generation time of 29 years, while it is 11 years for rhesus macaques. Another explanation for a

404  higher increase of mutation rate with paternal age could be differences in the replication

405  machinery itself. Due to higher sperm competition in rhesus macaque, the replication might be

406  under selective pressure for fast production at the expense of replication fidelity, leading to less

14

407    DNA repair mechanisms. As in other primates, we found a male bias in the contribution of *de*

408    *novo* mutations, as the paternal to maternal ratio is 4.2:1.This ratio is higher than the 2.7:1 ratio

409    observed in mice (Lindsay et al. 2019) and slightly higher than the 4:1 ratio observed in

410    humans (Goldmann et al. 2016; Jónsson et al. 2018; Lindsay et al. 2019). Similarly to the wild,

411    the males of our dataset reproduced from 10 years old, which did not allow us to examine if the

412    contribution bias was also present just after maturation. Moreover, the promiscuous behavior of

413    rhesus macaque leads to father reproducing with younger females. Using our model to compare

414    the contribution of each parent reproducing at the similar age, it seems that the male bias

415    increases with the parental age, with a lower difference in contribution at the time of sexual

416    maturation (2.3:1 for parents of 5 years old) and an increase in male to female contribution

417    with older parents (3.6:1 for parents of 15 years old). This result differs from humans, where

418    the male bias seems constant over time (Gao et al. 2019), but more time points in macaque

419    would be needed to interpret the contribution over time. In rhesus macaques, the ratio of

420    paternal to maternal contribution to the shared mutations between related individuals is 1:1,

421    similarly to what has been shown in mice (Lindsay et al. 2019), highlighting that those

422    mutations probably occur during primordial germ cell divisions in postzygotic stages. Our

423    study shows many shared patterns in the *de novo* mutations among non-Hominid primates.

424    More estimation of mammals could help understanding if these features are conserved across a

425    broad phylogenetic scale. Moreover, further work would be needed to understand if some

426    gamete production stages are more mutagenic in some species than others.

427    An accurate estimation of the mutation rate is essential for the precise dating of species divergence

428    events. We used the rhesus macaque mutation rate to estimate its divergence time with related

429    species for which whole-genome alignments are already available and their molecular divergence

430    times have been investigated before with other methods (Moorjani et al. 2016). The results of our

431    direct dating method, based on molecular distances between species and *de novo* mutation rate,

432    matched those of traditional molecular clock approaches for speciation events within 10 to 15

433    million years. However, it often produced earlier divergence times for more ancient nodes than the

434    molecular clock method. This incongruence might be attributed to the fossils that were used for

435    calibration with the clock method, which has many limitations (Heads 2005; Pulquério and

436    Nichols 2007; Steiper and Young 2008). A fossil used for calibrating a node is usually selected

15

437    to represent the oldest known specimen of a lineage. Still, it cannot be known if real even older

438    specimens existed (Heads 2005). Thus, a fossil is usually assumed to be younger than the real

439    divergence time of the species (Benton et al. 2015). Moreover, despite the error associated with

440    the dating of a fossil itself, determining its position on a tree can be challenging and have

441    effects on the inferred ages across the whole tree (Pulquério and Nichols 2007; Steiper and

442    Young 2008). For instance, the Catarrhini node, marking the divergence between the

443    Cercopithecidae and the Hominoidea, is often calibrated in primate phylogenies (Heads 2005).

444    This node has been calibrated to approx. 25 Mya using the oldest known Cercopithecidae fossil

445    (*Victoriapithecus*), and the oldest known Hominoidea fossil (*Proconsul*), both around 22 My

446    old (Goodman et al. 1998). However, if the oldest Catarrhini fossil (*Aegyptopithecus*) of 33 to

447    34 My age is used, this node could also be calibrated to 35 Mya (Stewart and Disotell 1998).

448    Finally, instead of being an ancestral specimen of the Catarrhini, *Aegyptopithecus* has been

449    suggested as a sister taxon to Catarrhini, which would lead to an even older calibration time for

450    this node (Stewart and Disotell 1998).

451    On the other hand, the direct mutation rate estimation could have produced overestimated

452    divergence times for the Catarrhini node age compared to previous estimates (Perelman et al.

453    2011; Pozzi et al. 2014), because the mutation rate and generation time might change cross-

454    species and over time. It is possible that the Catarrhini ancestor would have had a faster yearly

455    mutation rate, and/or a shorter generation time than the recent macaques. Since fossil

456    calibration could underestimate real divergence times, molecular-based methods could

457    overestimate it, especially by assuming a unique mutation rate to an entire clade (Steiper and

458    Young 2008).

459    To obtain more confidence in the estimation of divergence time, it would be necessary to have an

460    accurate estimation of the mutation rate for various species. The estimates available today for

461    primates vary from $0.81 \times 10^{-8}$ per site per generation for the Owl monkey (*Aotus nancymaae*)

462    to $1.66 \times 10^{-8}$ per site per generation for Orangutan (*Pongo abelii*). However, the different

463    methods and sequencing depth make it difficult to compare between species and attribute

464    differences to biological causes or methodological ones. Therefore, more standardized methods

465    in further studies would be needed to allow for cross-species comparison.

466

16

## Methods

### Samples

Whole blood samples (2 mL) in EDTA (Ethylenediaminetetraacetic acid) were collected from 53 Indian rhesus macaques (Macaca mulatta) during routine health checks at the Biomedical Primate Research Centre (BPRC, Rijswijk, Netherlands). Individuals originated from two groups, with one or two reproductive males per group. After ensuring the relatedness with a test based on individual genotypes (Manichaikul et al. 2010), we ended up with 19 trios formed by 33 individuals and two extended trios (for which a second generation was available). In our dataset males reproduced from 10 years old to 14.5 years old ($\male$ reproductive range: 4.5 years), and females from 3.5 years old to 15.7 years old ($\female$ reproductive range: 12.2 years). Genomic DNA was extracted using DNeasy Blood and Tissue Kit (Qiagen, Valencia-CA, USA) following the manufacturer's instructions. BGIseq libraries were built in China National GeneBank (CNGB), Shenzhen, China. The average insert size of the samples was 230 base pairs. Whole-genome pair-ended sequencing was performed on BGISEQ500 platform, with a read length of 2x100 bp. The average coverage of the raw sequences before trimming was 81X per sample (SE = 1.35). Whole-genome sequences have been deposited in NCBI (National Center for Biotechnology Information) with BioProject number PRJNA588178 and SRA submission SUB6522592.

### Reads mapping, SNPs calling, and filtering pipeline

Adaptors, low-quality reads, and N-reads were removed with SOAPnuke filter (Chen et al. 2017). Trimmed reads were mapped to the reference genome of rhesus macaque Mmul 8.0.1 using BWA-MEM version 0.7.15 with the estimated insert size option. Only reads mapping uniquely were kept and duplicates were removed using Picard MarkDuplicates. The average coverage after mapping was 76X per individuals (SE = 1.16). Variants were called using GATK 4.0.7.0 (Poplin et al. 2018); calling variants for each individual with HaplotypeCaller in BP-RESOLUTION mode; all gVCF files per sample were combined into a single one per trio using CombineGVCFs per autosomal chromosomes; finally joint genotyping was applied with

17

496    GenotypeGVCF. Because *de novo* mutations are rare events, variant quality score recalibration

497    (VQSR) is not a suitable tool to filter the sites as *de novo* mutations are more likely to be filtered

498    out as low-quality variants. Instead we used a site filtering with the following parameters: $QD <$

499    $2.0, FS > 20.0, MQ < 40.0, MQRankSum < -2.0, MQRankSum > 4.0,$ ReadPosRankSum < -

500    3.0, ReadPosRankSum > 3.0 and SOR > 3.0. These filters were chosen by first, running the

501    pipeline with the site filters recommended by GATK (QD < 2.0; FS > 60.0; MQ < 40.0;

502    MQRankSum < -12.5; ReadPosRankSum < -8.0 ; SOR > 3.0), then, doing a manual

503    curation of the candidates *de novo* mutations on the Integrative Genome Viewer (IGV).

504    Finally, we identified the common parameters within the apparent false-positive calls and

505    decided to adjust the site filter to remove as many false-positives without losing much true

506    positive calls (see the pipeline S8 Fig).

507

508    **Detection of *de novo* mutations**

509    The combination of high coverage (76X) and stringent filters reduced false-positive - calling a *de*

510    *novo* mutation while it is not there. Thus, for each trio, we applied the following filters:

511    (a) Mendelian violations were selected using GATK SelectVariant and refined to only keep
512    sites where both parents were homozygote reference (HomRef), and their offspring was
513    heterozygote (Het).

514    (b) In the case of a *de novo* mutation, the number of alternative alleles seen in the offspring
515    should account for ~ 50 % of the reads. Our allelic balance filter allowed the alternative
516    allele to be present in 30 % to 70 % of the total number of reads (applying the same 30%
517    cutoff as in other studies (Kong et al. 2012; Besenbacher et al. 2015; Francioli et al.
518    2015; S9 Fig).

519    (c) The depth of the three individuals was filtered to be between $0.5 \times m_{depth}$ and $2 \times m_{depth}$, with
520    $m_{depth}$ being the average depth of the trio. Most of the Mendelian violations are due to
521    sequencing errors in regions of low sequencing depth; therefore, we applied a stricter
522    threshold on the minimum depth to avoid the peak of Mendelian violations around 20X
523    (S10 Fig).

18

524

525   (d) Finally, after analyzing each trio with different genotype quality GQ cutoff (from 10 to

526          90), we set up a filter on the genotype quality of 60 to ensure the genotypes of the

527          HomRef parents and the Het offspring (S11 Fig).

528   From 242,922,329 autosomal SNPs (average of 12,785,386 per trio), 2,251,363 were potential

529   Mendelian violations found by GATK (average of 118,493 per trio), 177,227 were filtered

530   Mendelian violations with parents HomRef and offspring Het (average of 9,328 per trio) (a),

531   78,339 passed the allelic balance filter (average of 4,123 per trio) (b), 13,251 passed the depth

532   filter (average of 697 per trio) (c) and 744 the genotype quality filter (average of 39 per trio)

533   (d) (see S4 Table for details on each individual). We also remove sites where a *de novo*

534   mutation was shared among non-related individuals (1 site shared between 4 unrelated

535   individuals). This allowed us to detect the number of *de novo* mutations observed per trio

536   called m. We manually checked the reads mapping quality for all *de novo* mutations sites in the

537   Integrative Genome Viewer (IGV). And we found possible false-positive calls in 10.89 % of

538   the sites for which the variant was absent from the offspring or also present in a parent (see S1

539   Fig). We kept those sites for the estimation of the mutation rate, and corrected for false-positive

540   ($\beta = 0.1089$), but removed them for downstream pattern analysis. We experimentally validated

541   the *de novo* candidates from the trio Noot (father), Platina (mother), and Lithium (offspring).

542   Primers were designed for 39 candidates (S5 Table). PCR amplification and Sanger sequencing

543   were conducted on each individual (protocol in Supplementary materials). On 24 sites the PCR

544   amplification and sequencing returned high-quality results for all three individuals. A candidate

545   was considered validated when both parents showed homozygosity for the reference allele and

546   the offspring showed heterozygosity (S2 Fig). All sequences generated for the PCR validation

547   have been deposited in Genbank with accession numbers MT426016 - MT426087 (S4 Table).

548   **Estimation of the mutation rate per site per generation**

549   From the number of *de novo* mutations to an estimate of the mutation rate per site per

550   generation, it is necessary to also correct for false-negatives - not calling a true *de novo* mutation

551   as such. To do so, we estimated two parameters: the false-negative rate and the number of

19

552  callable sites, *C*, ie. the number of sites in the genome where we would be able to call a *de*

553  *novo* mutation if it was there. We used the BP_RESOLUTION option in GATK to call variants

554  for each position and thus get the exact genotype quality for each site in each individual - also

555  sites that are not polymorphic. So unlike other studies, we do not have to rely on sequencing

556  depth as a proxy for genotype quality at those sites. Instead, we can apply the same genotype

557  quality threshold to the non-polymorphic sites as we do for *de novo* mutation candidate sites.

558  This should lead to a more accurate estimate of the number of callable sites. For each trio, *C* is

559  the sum of all sites where: both parents are HomRef, and the three individuals passed the depth

560  filter (b) and the genotype quality filter (d). To correct for our last filter, the allelic balance (c),

561  we estimated the false-negative rate *α,* defined as the proportion of true heterozygotes sites

562  (one parent HomRef, the other parent HomAlt and their offspring Het) outside the allelic

563  balance threshold (S9 Fig). We also implemented in this parameter the false-negative rate of

564  the site filters following a normal distribution (FS, MQRankSum, and ReadPosRankSum). For

565  all trios combined, the rate of false-negatives caused by the allele balance filter and the site

566  filters was 0.0402. To validate this false-negative rate estimation we also used a simulation

567  method, used in other studies (Keightley et al. 2015; Pfeifer 2017). With BAMSurgeon (Ewing

568  et al. 2015), 552 mutations were simulated across the 19 trios at random callable sites. The

569  false-negative rate was calculated as 1 – (number of detected mutations/number of simulated

570  mutations), after running the pipeline from variant calling. The mutation rate per sites per

571  generation can then be estimated per trio with the following equation:

572
$$\mu = \frac{m \times (1 - \beta)}{(1 - \alpha) \times 2 \times C} \qquad (1)$$

575

576  **Sex bias, ages, and relatedness**

577  *De novo* mutations were phased to their parental origin using the read-backed phasing method

578  described in Maretty et al. 2017 (script available on GitHub:

579  https://github.com/besenbacher/POOHA). The method uses read-pairs that contain both a *de*

580  *novo* mutation and another heterozygous variant, the latter of which was used to determine the

20

581     parental origin of the mutation if it is present in both offspring and one of the parents. The

582     phasing allowed us to identify any parental bias in the contribution of the *de novo* mutations.

583     Pearson's correlation test was performed between the mutation rate and ages of each parent, as

584     well as a linear regression model for father and mother independently. A multiple linear regression

585     model was performed to predict the mutation rate from both parental ages as predictor variables.

586     The phased mutations were used to dissociate the effect of the parental age from one another.

587     Because the total number of SNPs phased to the mother or the father may differ, we divided the

588     phased *de novo* mutations found in a parent by the total SNPs phased to this parent. Only a

589     subset of the *de novo* mutations in an offspring was phased. Thus, we applied the paternal to

590     maternal ratio to the total number of mutations in a trio, referred to as 'upscaled' number of

591     mutations, to predict the number of total mutations given by each parent at different ages. The

592     two extended trios, analyzed as independent trios, also allowed us to determine if ~ 50 % of the

593     *de novo* mutations observed in the first trio were passed on to the next generation.

594

595     **Characterization of *de novo* mutations**

596     From all the *de novo* mutations found, the type of mutations and their frequencies were

597     estimated. For the mutations from a C to any base we determined if they were followed by a G to

598     detect the CpG sites (similarly if G mutations were preceded by a C. We defined a cluster as a

599     window of 20,000 bp to qualify how many mutations were clustered together; over all

600     individuals, looking at related individuals, and within individuals. We simulated 663 mutations

601     following a uniform distribution to compare with our dataset. We investigated the mutations that

602     are shared between related individuals. Finally, we looked at the location of mutations in the

603     coding region using the annotation of the reference genome.

604

605     **Molecular dating using the new mutation rate**

606     We calculated the effective population size using Watterson's estimator $\theta = 4N_e\mu$ (Watterson

607     1975). We estimated $\theta$ with the nucleotide diversity $\pi = 0.00247$ according to a recent population

608     study (Xue et al. 2016). Thus, we calculated the effective population size as $N_e = \frac{\pi}{4\mu}$ with $\mu$ the

609     mutation rate per site per generation estimated in our study. To calculate divergence time, we

21

610 converted the mutation rate to a yearly rate based on the regression model of the number of

611 mutations given by each parent regarding their ages and the average callability $C =$

612 2,351,302,179. Given the maturation time and the high mortality due to predation, we assumed

613 an average age of reproduction in the wild at 10 years old for females and 12 years old for males

614 and a generation time of 11 years, also reported in another study (Xue et al. 2016). Thus, the

615 yearly mutation rate was:

616 $$\mu = \frac{4.6497 + 0.3042 \times agematernal + 4.8399 + 1.8364 \times agepaternal \times (1 - \beta)}{(1 - \alpha) \times 2 \times C} \tag{2}$$

617 The divergence time between species was then calculated using $T_{divergence} = \frac{d}{2\mu}$ with $d$ the

618 genetic distance between species which were calculated from the whole-genome comparison

619 (Moorjani et al. 2016) and μ the yearly mutation rate of rhesus macaques. We also used the

620 confidence interval at 95% of our mutation rate regression to compute the confidence interval on

621 divergence time. Based on the coalescent theory (Kingman, 1982), the time to coalescence is

622 2NeG with G the generation time and Ne the ancestral effective population size, assumed

623 constant over time, as shown in a previous study (Xue et al. 2016). Thus, we dated the speciation

624 event as previously done by Besenbacher et al. 2019 with:

625 $$T_{speciation} = T_{divergence} - 2 \times N_{e\ ancestor} \times G \tag{3}$$

626

## Acknowledgments

628 We would like to thank GenomeDK at Aarhus University for providing computational resources

629 and supports to this study. We also thank Josefin Stiller for helpful comments on the manuscript.

630

## Data Availability Statement

632 Whole-genome sequences have been deposited in NCBI (National Center for Biotechnology

633 Information) with BioProject number PRJNA588178 and SRA submission SUB6522592. All

634 sequences generated for the PCR validation have been deposited in Genbank with accession

635 numbers MT426016 - MT426087.

636

## References

638    Acuna-Hidalgo R, Bo T, Kwint MP, Van De Vorst M, Pinelli M, Veltman JA, Hoischen A, Vissers
639        LELM, Gilissen C. 2015. Post-zygotic Point Mutations Are an Underrecognized Source of de Novo
640        Genomic Variation. Am. J. Hum. Genet. 97:67–74.

641    Acuna-Hidalgo R, Veltman JA, Hoischen A. 2016. New insights into the generation and role of de novo
642        mutations in health and disease. Genome Biol. 17.

643    Amann RP, Howards SS. 1980. Daily spermatozoal production and epididymal spermatozoal reserves of
644        the human male. J. Urol. 124:211–215.

645    Amann RP, Johnson L, Thompson DL, Pickett BW. 1976. Daily Spermatozoal Production, Epididymal
646        Spermatozoal Reserves and Transit Time of Spermatozoa Through the Epididymis of the Rhesus
647        Monkey. Biol. Reprod. 15:586–592.

648    Awadalla P, Gauthier J, Myers RA, Casals F, Hamdan FF, Griffing AR, Côté M, Henrion E, Spiegelman
649        D, Tarabeux J, et al. 2010. Direct measure of the de novo mutation rate in autism and schizophrenia
650        cohorts. Am. J. Hum. Genet. 87:316–324.

651    Benton MJ, Donoghue PCJ, Asher RJ, Friedman M, Near TJ, Vinther J. 2015. Constraints on the
652        timescale of animal evolutionary history. Palaeontol. Electron. 18:1–106.

653    Bercovitch FB, Widdig A, Trefilov A, Kessler MJ, Berard JD, Schmidtke J, Nürnberg P, Krawczak M.
654        2003. A longitudinal study of age-specific reproductive output and body condition among male
655        rhesus macaques, Macaca mulatta. Naturwissenschaften 90:309–312.

656    Besenbacher S, Hvilsom C, Marques-Bonet T, Mailund T, Schierup MH. 2019. Direct estimation of
657        mutations in great apes reconciles phylogenetic dating. Nat. Ecol. Evol. [Internet] 3:286–292.
658        Available from: http://www.nature.com/articles/s41559-018-0778-x

659    Besenbacher S, Liu S, Izarzugaza JMG, Grove J, Belling K, Bork-Jensen J, Huang S, Als TD, Li S,
660        Yadav R, et al. 2015. Novel variation and de novo mutation rates in population-wide de novo
661        assembled Danish trios. Nat. Commun. [Internet] 6:5969. Available from:
662        http://www.nature.com/doifinder/10.1038/ncomms6969

663    Besenbacher S, Sulem P, Helgason A, Helgason H, Kristjansson H, Jonasdottir A, Jonasdottir A,
664        Magnusson OT, Thorsteinsdottir U, Masson G, et al. 2016. Multi-nucleotide de novo Mutations in
665        Humans. Petrov DA, editor. PLOS Genet. [Internet] 12:e1006315. Available from:
666        http://dx.plos.org/10.1371/journal.pgen.1006315

667    Byskov AG. 1986. Differential of mammalian embryonic gonad. Physiol. Rev. 66:71–117.

668    Campbell CR, Tiley GP, Poelstra JW, Hunnicutt KE, Larsen PA, dos Reis M, Yoder AD. 2019. Pedigree-
669        based measurement of the de novo mutation rate in the gray mouse lemur reveals a high mutation
670        rate, few mutations in CpG sites, and a weak sex bias. bioRxiv [Internet]. Available from:
671        http://dx.doi.org/10.1101/724880

672    Chen Y, Chen Y, Shi C, Huang Z, Zhang Y, Li S, Li Y, Ye J, Yu C, Li Z, et al. 2017. SOAPnuke: A
673        MapReduce acceleration-supported software for integrated quality control and preprocessing of
674        high-throughput sequencing data. Gigascience 7:1–6.

675    Crow JF. 2000. The origins, patterns and implications of human spontaneous mutation. Nat. Rev. Genet.

23

676      1:40–47.

677  Drost JB, Lee WR. 1995. Biological basis of germline mutation: Comparisons of spontaneous germline
678      mutation rates among drosophila, mouse, and human. Environ. Mol. Mutagen. [Internet] 25:48–64.
679      Available from: http://doi.wiley.com/10.1002/em.2850250609

680  Ewing AD, Houlahan KE, Hu Y, Ellrott K, Caloian C, Yamaguchi TN, Bare JC, P'ng C, Waggott D,
681      Sabelnykova VY, Kellen MR, et al. 2015. Combining tumor genome simulation with crowdsourcing
682      to benchmark somatic single-nucleotide-variant detection. Nature methods. 12(7), 623-630.

683  Francioli LC, Polak PP, Koren A, Menelaou A, Chun S, Renkens I. 2015. Genome-wide patterns and
684      properties of de novo mutations in humans. Nat. Genet. [Internet] 47:822. Available from:
685      https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4485564/pdf/nihms-679155.pdf

686  Gao Z, Moorjani P, Sasani TA, Pedersen BS, Quinlan AR, Jorde LB, Amster G, Przeworski M. 2019.
687      Overlooked roles of DNA damage and maternal age in generating human germline mutations. Proc.
688      Natl. Acad. Sci. U. S. A. [Internet] 116:9491–9500. Available from:
689      http://www.ncbi.nlm.nih.gov/pubmed/31019089

690  Gibbs RA, Rogers J, Katze MG, Bumgarner R, Weinstock GM, Mardis ER, Remington KA, Strausberg
691      RL, Venter JC, Wilson RK, et al. 2007. Evolutionary and biomedical insights from the rhesus
692      macaque genome. Science (80-. ). [Internet] 316:222–234. Available from:
693      http://science.sciencemag.org/

694  Goldmann JM, Wong WSW, Pinelli M, Farrah T, Bodian D, Stittrich AB, Glusman G, Vissers LELM,
695      Hoischen A, Roach JC, et al. 2016. Parent-of-origin-specific signatures of de novo mutations. Nat.
696      Genet.

697  Goodman M, Porter CA, Czelusniak J, Page SL, Schneider H, Shoshani J, Gunnell G, Groves CP. 1998.
698      Toward a Phylogenetic Classification of Primates Based on DNA Evidence Complemented by
699      Fossil Evidence. Mol. Phylogenet. Evol.

700  Goodman MF, Creighton S, Bloom LB, Petruska J, Kunkel TA. 1993. Biochemical Basis of DNA
701      Replication Fidelity. Crit. Rev. Biochem. Mol. Biol. [Internet] 28:83–126. Available from:
702      https://www.tandfonline.com/action/journalInformation?journalCode=ibmg20

703  Halldorsson B V., Palsson G, Stefansson OA, Jonsson H, Hardarson MT, Eggertsson HP, Gunnarsson B,
704      Oddsson A, Halldorsson GH, Zink F, et al. 2019. Characterizing mutagenic effects of recombination
705      through a sequence-level genetic map. Science (80-. ). 363.

706  Heads M. 2005. Dating nodes on molecular phylogenies: A critique of molecular biogeography.
707      Cladistics 21:62–78.

708  Hernandez RD, Hubisz MJ, Wheeler DA, Smith DG, Ferguson B, Rogers J, ... & Muzny D. 2007.
709      Demographic histories and patterns of linkage disequilibrium in Chinese and Indian rhesus
710      macaques. Science, 316(5822), 240-243.

711  Ho SYW, Larson G. 2006. Molecular clocks: When timesare a-changin'. Trends Genet. 22:79–83.

712  Jónsson H, Sulem P, Arnadottir GA, Pálsson G, Eggertsson HP, Kristmundsdottir S, Zink F, Kehr B,
713      Hjorleifsson KE, Jensson BÖ, et al. 2018. Multiple transmissions of de novo mutations in families.

24

714    Nat. Genet. [Internet] 50:1674. Available from: http://www.nature.com/articles/s41588-018-0259-9

715    Jónsson H, Sulem P, Kehr B, Kristmundsdottir S, Zink F, Hjartarson E, Hardarson MT, Hjorleifsson KE,
716        Eggertsson HP, Gudjonsson SA, et al. 2017. Parental influence on human germline de novo
717        mutations in 1,548 trios from Iceland. Nature [Internet] 549:519–522. Available from:
718        http://www.nature.com/doifinder/10.1038/nature24018

719    Keightley PD, Pinharanda A, Ness RW, Simpson F, Dasmahapatra KK, Mallet J, … & Jiggins CD. 2015.
720        Estimation of the spontaneous mutation rate in Heliconius melpomene. Molecular biology and
721        evolution, 32(1), 239-243.

722    Kingman JFC. 1982. The coalescent. Stochastic Processes and Their Applications, 13(3), 235–248.
723        https://doi.org/10.1016/0304-4149(82)90011-4

724    Kong A, Frigge ML, Masson G, Besenbacher S, Sulem P, Magnusson G, Gudjonsson SA, Sigurdsson A,
725        Jonasdottir A, Jonasdottir A, et al. 2012. Rate of de novo mutations and the importance of father's
726        age to disease risk. Nature [Internet] 488:471. Available from:
727        https://www.nature.com/articles/nature11396.pdf

728    Lapierre M, Lambert A, Achaz G. 2017. Accuracy of demographic inferences from the site frequency
729        spectrum: The case of the yoruba population. Genetics 206:139–449.

730    Li WH, Ellsworth DL, Krushkal J, Chang BHJ, Hewett-Emmett D. 1996. Rates of nucleotide substitution
731        in primates and rodents and the generation-time effect hypothesis. Mol. Phylogenet. Evol. 5:182–
732        187.

733    Lindsay SJ, Rahbari R, Kaplanis J, Keane T, Hurles ME. 2019. Similarities and differences in patterns of
734        germline mutation between mice and humans. Nat. Commun. 10:1–12.

735    Manichaikul A, Mychaleckyj JC, Rich SS, Daly K, Sale M, Chen WM. 2010. Robust relationship
736        inference in genome-wide association studies. Bioinformatics 26:2867–2873.

737    Maretty L, Jensen JM, Petersen B, Sibbesen JA, Liu S, Villesen P, Skov L, Belling K, Theil Have C,
738        Izarzugaza JMG, et al. 2017. Sequencing and de novo assembly of 150 genomes from Denmark as a
739        population reference. Nature [Internet] 548:87–91. Available from:
740        http://www.nature.com/doifinder/10.1038/nature23264

741    Moorjani P, Amorim CEG, Arndt PF, Przeworski M. 2016. Variation in the molecular clock of primates.
742        Proc. Natl. Acad. Sci. U. S. A. [Internet] 113:10607–10612. Available from:
743        http://www.ncbi.nlm.nih.gov/pubmed/27601674

744    Neale BM, Devlin B, Boone BE, Levy SE, Lihm J, Buxbaum JD, Wu Y, Lewis L, Han Y, Boerwinkle E,
745        et al. 2012. Patterns and rates of exonic de novo mutations in autism spectrum disorders. Nature
746        485:242–246.

747    Ohta T. 1993. An examination of the generation-time effect on molecular evolution. Proc. Natl. Acad.
748        Sci. USA 90:10676–10680.

749    Oliveira S, Cooper DN, Azevedo L. 2018. De Novo Mutations in Human Inherited Disease. In: eLS. John
750        Wiley & Sons, Ltd. p. 1–7.

25

751  Perelman P, Johnson WE, Roos C, Seuánez HN, Horvath JE, Moreira MAM, Kessing B, Pontius J,
752      Roelke M, Rumpler Y, et al. 2011. A Molecular Phylogeny of Living Primates.Brosius J, editor.
753      PLoS Genet. [Internet] 7:e1001342. Available from:
754      https://dx.plos.org/10.1371/journal.pgen.1001342

755  Pfeifer SP. 2017. Direct estimate of the spontaneous germ line mutation rate in African green monkeys.
756      Evolution (N. Y). [Internet] 71:2858–2870. Available from:
757      http://doi.wiley.com/10.1111/evo.13383

758  Poplin R, Ruano-Rubio V, DePristo MA, Fennell TJ, Carneiro MO, Auwera GA Van der, Kling DE,
759      Gauthier LD, Levy-Moonshine A, Roazen D, et al. 2018. Scaling accurate genetic variant discovery
760      to tens of thousands of samples. bioRxiv [Internet]:201178. Available from:
761      https://www.biorxiv.org/content/10.1101/201178v2

762  Pozzi L, Hodgson JA, Burrell AS, Sterner KN, Raaum RL, Disotell TR. 2014. Primate phylogenetic
763      relationships and divergence dates inferred from complete mitochondrial genomes. Mol.
764      Phylogenet. Evol. 75:165–183.

765  Pulquério MJF, Nichols RA. 2007. Dates from the molecular clock: how wrong can we be? Trends Ecol.
766      Evol. 22:180–184.

767  Rahbari R, Wuster A, Lindsay SJ, Hardwick RJ, Alexandrov LB, Turki S Al, Dominiczak A, Morris A,
768      Porteous D, Smith B, et al. 2016. Timing, rates and spectra of human germline mutation. Nat.
769      Genet. [Internet] 48:126–133. Available from:
770      http://www.nature.com/authors/editorial_policies/license.html#terms

771  Rawlins RG, Kessler MJ. 1986. The Cayo Santiago Macaques: History, Behavior, and Biology. Available
772      from:
773      https://books.google.fr/books?hl=en&lr=&id=3fMCzTve890C&oi=fnd&pg=PR9&dq=The+Cayo+S
774      antiago+macaques:+History,+behavior,+and+biology&ots=wwfUkTHWoC&sig=SUQ3CN6nFJo3h
775      521NL6uKEzrRD4#v=onepage&q=The Cayo Santiago macaques%3A History%2C behavior%2C
776      and biology

777  Roach JC, Glusman G, Smit AFA, Huff CD, Hubley R, Shannon PT, Rowen L, Pant KP, Goodman N,
778      Bamshad M, et al. 2010. Analysis of Genetic Inheritance in a Family Quartet by Whole Genome
779      Sequencing. Science (80-. ). [Internet] 328:636–639. Available from:
780      https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3037280/pdf/nihms247436.pdf

781  Rosenberg N. A., & Nordborg M. 2002. Genealogical trees, coalescent theory and the analysis of genetic
782      polymorphisms. Nature Reviews Genetics, Vol. 3, pp. 380–390. https://doi.org/10.1038/nrg795

783  Scally A. 2016. Mutation rates and the evolution of germline structure. Philos. Trans. R. Soc. B Biol. Sci.
784      [Internet] 371:20150137. Available from: http://dx.doi.org/10.1098/rstb.2015.0137

785  Schrago CG. 2014. The effective population sizes of the anthropoid ancestors of the human-chimpanzee
786      lineage provide insights on the historical biogeography of the great apes. Mol. Biol. Evol. 31:37–47.

787  Ségurel L, Wyman MJ, Przeworski M. 2014. Determinants of Mutation Rate Variation in the Human
788      Germline. Annu. Rev. Genomics Hum. Genet. [Internet] 15:47–70. Available from:
789      http://www.annualreviews.org/doi/10.1146/annurev-genom-031714-125740

26

790   Song K, Li L, Zhang G. 2016. Coverage recommendation for genotyping analysis of highly heterologous
791       species using next-generation sequencing technology. Sci. Rep. 6:35736.

792   Steiper ME, Young NM. 2008. Timing primate evolution: Lessons from the discordance between
793       molecular and paleontological estimates. Evol. Anthropol. 17:179–188.

794   Stewart C-B, Disotell TR. 1998. Primate evolution – in and out of Africa. Curr. Biol. [Internet] 8:R582–
795       R588. Available from: https://linkinghub.elsevier.com/retrieve/pii/S0960982207003673

796   Tatsumoto S, Go Y, Fukuta K, Noguchi H, Hayakawa T, Tomonaga M, Hirai H, Matsuzawa T, Agata K,
797       Fujiyama A. 2017. Direct estimation of de novo mutation rates in a chimpanzee parent-offspring trio
798       by ultra-deep whole genome sequencing. Sci. Rep. 7.

799   Teeling EC, Springer MS, Madsen O, Bates P, O'Brien SJ, Murphy WJ. 2005. A molecular phylogeny
800       for bats illuminates biogeography and the fossil record. Science (80-. ). 307:580–584.

801   Thomas GWC, Wang RJ, Puri A, Rogers J, Radivojac P, Hahn MW, Thomas GWC, Wang RJ, Puri A,
802       Harris RA, et al. 2018. Reproductive Longevity Predicts Mutation Rates in Primates. Curr. Biol.
803       [Internet] 28:1–5. Available from: https://doi.org/10.1016/j.cub.2018.08.050

804   Venn O, Turner I, Mathieson I, De Groot N, Bontrop R, McVean G. 2014. Strong male bias drives
805       germline mutation in chimpanzees. Science (80-. ). 344:1272–1275.

806   Wang H, Zhu X. 2014. De novo mutations discovered in 8 Mexican American families through whole
807       genome sequencing. BMC Proc. [Internet] 8:S24. Available from:
808       https://www.ncbi.nlm.nih.gov/pmc/articles/PMC4143763/pdf/1753-6561-8-S1-S24.pdf

809   Wang RJ, Thomas GWC, Raveendran M, Harris RA, Doddapaneni H, Muzny DM, Capitanio JP,
810       Radivojac P, Rogers J, Hahn MW. 2020. Paternal age in rhesus macaques is positively associated
811       with germline mutation accumulation but not with measures of offspring sociability. Genome
812       Research, gr-255174.

813   Watterson GA. 1975. On the number of segregating sites in genetical models without recombination.
814       Theor. Popul. Biol. 7:256–276.

815   Wu C-I, Lit W-H. 1985. Evolution evidence for higher rates of nucleotide substitution in rodents than in
816       man. Proc. Nati. Acad. Sci. USA 82:1741–1745.

817   Wu FL, Strand A, Ober C, Wall JD, Moorjani P, Przeworski M. 2019. A comparison of humans and
818       baboons suggests germline mutation rates do not track cell divisions. bioRxiv:844910.

819   Xue C, Raveendran M, Harris RA, Fawcett GL, Liu X, White S, Dahdouli M, Rio Deiros D, Below JE,
820       Salerno W, et al. 2016. The population genomics of rhesus macaques (Macaca mulatta) based on
821       whole-genome sequences. Genome Res. [Internet] 26:1651–1662. Available from:
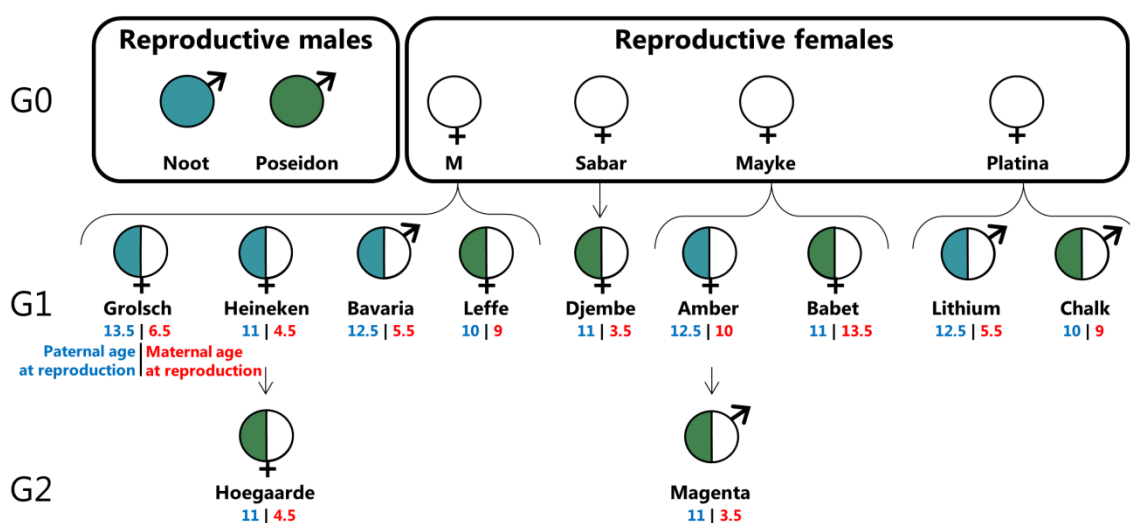822       http://www.ncbi.nlm.nih.gov/pubmed/27934697

823   Yuan Q, Zhou Z, Lindell SG, Higley JD, Ferguson B, Thompson RC, Lopez JF, Suomi SJ, Baghal B,
824       Baker M, et al. 2012. The rhesus macaque is three times as diverse but more closely equivalent in
825       damaging coding variation as compared to the human. BMC Genet. [Internet] 13:52. Available
826       from: http://www.ncbi.nlm.nih.gov/pubmed/22747632

27

827 Yuen RKC, Merico D, Cao H, Pellecchia G, Alipanahi B, Thiruvahindrapuram B, et al. 2016. Genome-
828   wide characteristics of de novo mutations in autism. Npj Genomic Medicine. 1:1–10.
829   https://doi.org/10.1038/npjgenmed.2016.27

830 Zeng K, Jackson BC, Barton HJ. 2018. Methods for estimating demography and detecting between-locus
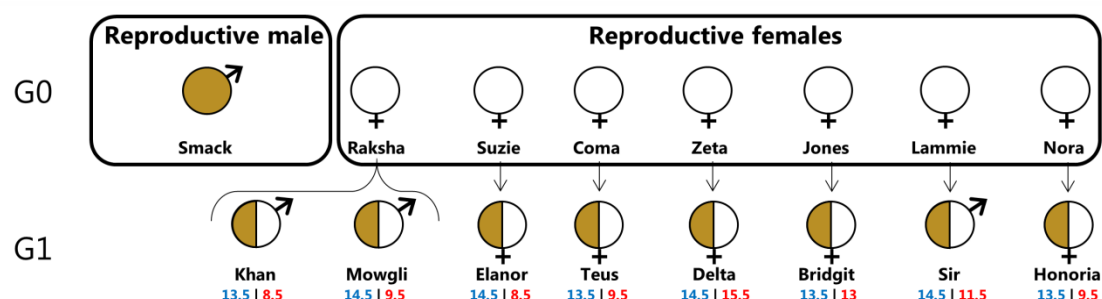831   differences in the effective population size and mutation rate. Mol. Biol. Evol. 36:423–433.
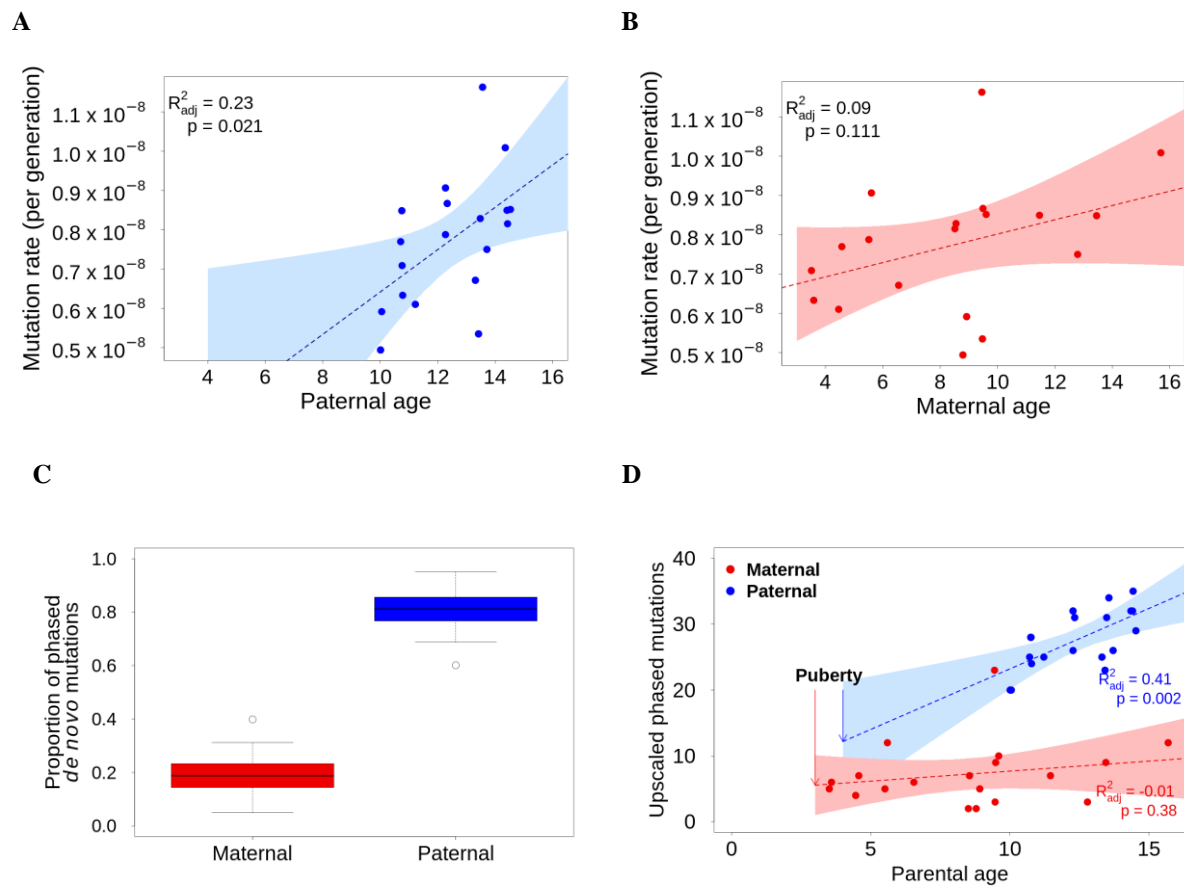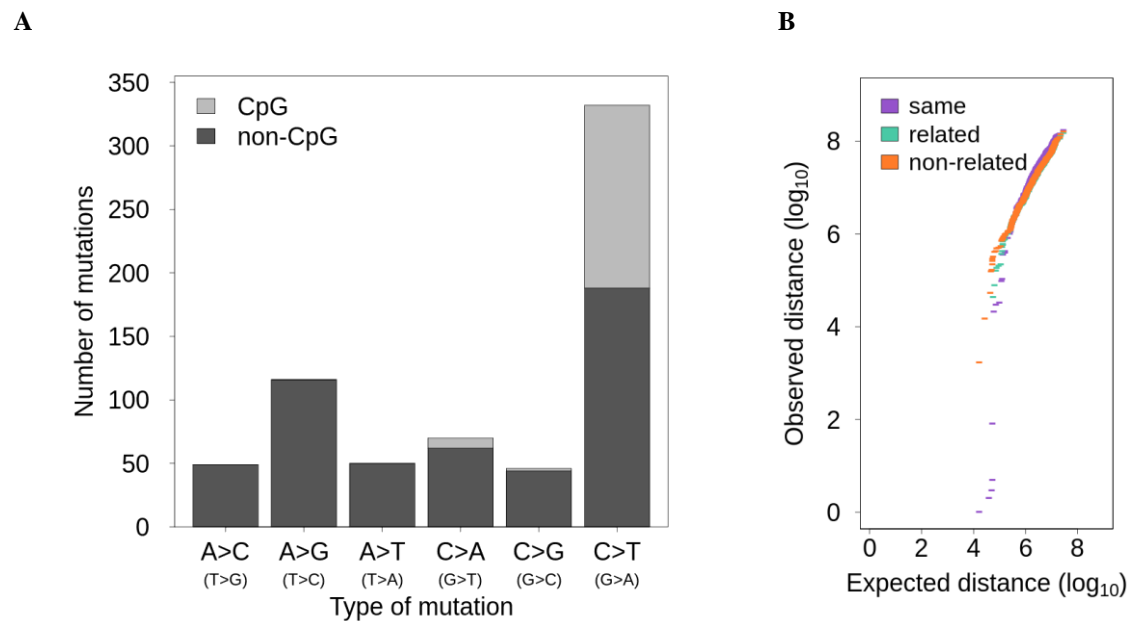
832

# Figures

**A**



**B**



**Fig 1. Pedigree of the 19 trios used for the direct estimation of mutation rate**. A: The first group is composed of two reproductive males and four reproductive females. B: The second group contained one reproductive male and seven reproductive females. In each offspring, the color on the left corresponds to the paternal lineage and under the name are the age of the father (in blue) and mother (in red) at the time of reproduction. The reproductive ranges are 4.5 years for males and 12.2 years for females.
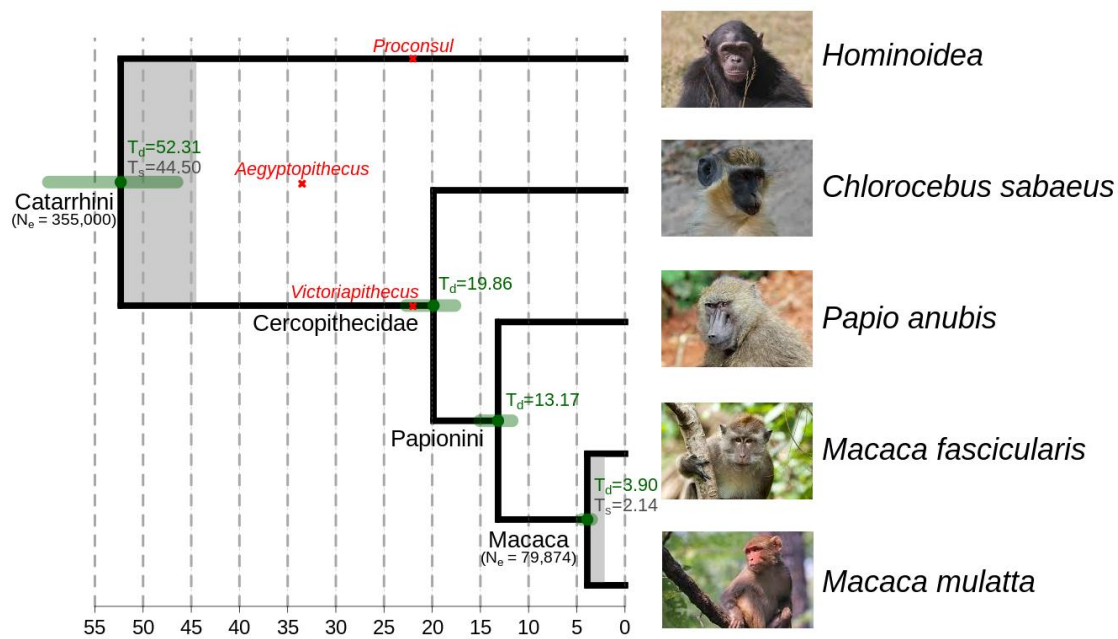
29

**Fig 2. Parental contribution and age effect to the *de novo* mutation rate.** A: There is a positive correlation between the mutation rate and the paternal age. B: The correlation between maternal age and mutation rate is not significant. C: Males contribute to 80.6 % of the *de novo* mutations while females contribute to 19.4 % of them. D: Upscaled number of *de novo* mutations given by each parent shows a similar contribution at the age of sexual maturation and a substantial increase with male age.
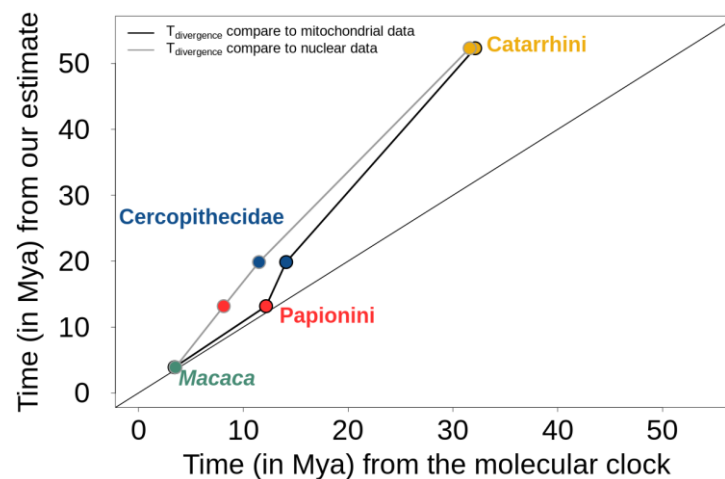
**Fig 3. Characterizations of the *de novo* mutations.** A: The type of *de novo* mutations in CpG and non-CpG sites. B: QQ-plot of the distance between *de novo* mutations compared to a uniform distribution within individuals (purple), between related individuals (green), and between non-related individuals (orange).

**Fig 4. Molecular dating with pedigree-based mutation rate.** A: Primates phylogeny based on the yearly mutation rate ($0.62 \times 10^{-9}$ per site per year). In green are the confidence interval of our divergence time estimates (Td) and grey shades represent the time of speciation (Ts). The effective population sizes are indicated under the nodes ($N_e$ Macaca ancestor is our estimate of $N_e$ *Macaca mulatta* and $N_e$ Catarrhini from the literature (Schrago 2014)). B: Comparison of our divergence time and speciation time with the previous estimation using the molecular clock from mitochondrial (Pozzi et al. 2014) and nuclear data (Perelman et al. 2011) calibrated with fossils records.

32