1 **TITLE:**

2 **Mutation distribution density in tumors reconstructs human's lost diversity**

3

4

5 **Authors:**

6 José María Heredia-Genestar[1], Tomàs Marquès-Bonet[1,2,3,4], David Juan*[1], Arcadi

7 Navarro*[1,2,3,5]

8

9 1- Institute of Evolutionary Biology (CSIC-UPF), Department of Experimental and

10 Health Sciences, Universitat Pompeu Fabra, Dr. Aiguader 88, 08003, Barcelona,

11 Spain.

12 2- CRG-CNAG, Centre for Genomic Regulation (CRG), Barcelona Institute of Science

13 and Technology (BIST), Baldiri i Reixac 4, 08028, Barcelona, Spain.

14 3- Institució Catalana de Recerca i Estudis Avançats (ICREA), Passeig de Lluís

15 Companys, 23, 08010, Barcelona, Spain.

16 4- Institut Català de Paleontologia Miquel Crusafont, Universitat Autònoma de

17 Barcelona, Edifici ICTA-ICP, c/ Columnes s/n, 08193, Cerdanyola del Vallès,

18 Barcelona, Spain

19 5- National Institute for Bioinformatics (INB), Barcelona Biomedical Research Park

20 (PRBB), Dr. Aiguader 88, 08003, Barcelona, Spain

21 * Corresponding authors

22

23

24

25    **Introductory Paragraph:**

26    Mutations do not accumulate uniformly across the genome. Human germline and tumor

27    mutation density correlate poorly, and each is associated with different genomic

28    features. Here, we analyze the genome-wide distribution of mutation densities in

29    human and non-human Great Ape (NHGA) germlines as well as human tumors.

30    Strikingly, non-human Great Ape germlines present higher correlation with tumors than

31    the human germline does. This situation is mediated by a different distribution in the

32    human germline of mutations at non-CpG sites, but not of CpG>T transitions. We

33    propose that the impact of ancestral and historical human demographic events on

34    human mutation density leads to this specific disruption in its expected genome-wide

35    distribution. Tumors partially recover this distribution by the accumulation of pre-

36    neoplastic-like somatic mutations. Our results highlight the potential utility of using

37    Great Ape population data, rather than human controls, to establish the expected

38    mutational background of healthy somatic cells.

39

2

## Introduction

Mutation density, at different scales, has been shown to correlate with different genomic features, such as regional GC-content or recombination rate[1–5]. In cancer, mutation density has been linked to chromatin states[6], with higher mutation accumulation in closed chromatin. It has been suggested that the tumor's higher mutation accumulation in closed chromatin is due to poorer accessibility or recruitment of the mismatch repair machinery to late-replicating, closed-chromatin regions[7,8]. Recent studies have shown that the correlation between tumor mutation density and chromatin state is highly tissue-dependent, allowing the identification of the tissue of origin of metastatic tumor samples[9,10].

At a smaller scale, sequence context is a good predictor of the mutation rate[11], beyond hypermutable CpG sites[12–15]. Sequence context has been widely used in cancer analyses to detect signatures of mutation associated with mutagens such as UV-light, tobacco smoke, or APOBEC activity[16,17]. These effects have also been detected in healthy somatic tissues[18,19].

*De novo* mutations are also affected by sequence context[20–22]. The rates of some particular mutation types have changed recently across ancestries[23–26]. Mutation rates seem to have been under selection in the human lineage. Sequence context studies have shown differences in the relative proportion of certain mutation types between Great Ape species[25]. Furthermore, studies of *de novo* mutations in Great Ape samples revealed a slowdown of the overall mutation rate in humans relative to chimpanzees and gorillas[27].

Here we study mutation rate evolution, through the differences in mutation distribution (at the 1Mbp scale) between human tumoral tissues and healthy populations in the Great Ape lineage.

68    **Results**

69

70    We compared the mutation density distribution in human (1kGP[28], sgdp_50[29]), non-

71    human Great Apes (NHGA: chimpanzee[30,31], gorilla[30,32]), and human cancer[33] datasets.

72    We focused on high-quality orthologous regions shared between human, chimpanzee

73    and gorilla genomes, measuring the number of variants per 1Mbp independently of the

74    frequency of each variant (see Methods).

75

76    In agreement with previous reports[1,3,4,6], we observe a variable distribution of the

77    mutation density across the genome in all datasets (**Figure 1a**). Mutation densities

78    correlate weakly between the human germline and tumors[1,6] (**Table 1**). Strikingly, the

79    NHGA-tumor correlations are much stronger than the human-tumor correlation and are

80    similar to the human-NHGA germline correlations (**Table 1** & **Supplementary Table 1)**.

81

82    We compared the distribution of mutation density between pairs of datasets

83    (**Supplementary Figure 1**). Interestingly, we observed that mutation density in tumors

84    is higher in windows where NHGAs have higher mutation density than humans

85    (**Figures 1b,c**). To control for differences in the shapes of distributions, we ranked

86    each set of windows according to their mutation density (**Figures 1d,e**). These ranked

87    distributions show a clear pattern: tumor mutation densities are higher in windows with

88    higher ranks in NHGAs than in human (two-sided Mann-Whitney U test p-value human-

89    chimpanzee= 3.7e-216; human-gorilla = 2.8e-161). This behavior is exclusive to

90    human-NHGA comparisons, as it cannot be observed when comparing chimpanzee to

91    gorilla (**Supplementary Figure 1**), and can be replicated under different conditions and

92    datasets (**Supplementary Notes, Supplementary Tables 2-6 & Supplementary**

93    **Figure 1**).

94

4

95    High-diversity NHGA subspecies have stronger correlations with both human and

96    tumor than the low-diversity subspecies (**Supplementary Table 3**). Furthermore, the

97    diagonal pattern is only characteristic of comparisons between the germlines of

98    humans and high-diversity NHGA subspecies. A comparison of high and low-diversity

99    chimpanzee and gorilla subspecies showed a clear horizontal split (**Supplementary**

100   **Figure 2**). Mutation density in tumors co-localizes with the most diverse NHGA

101   subspecies, regardless of the mutation density in the least diverse. In other words;

102   while a lack of diversity distorts the distribution of the genome-wide mutation densities,

103   the diagonal pattern is caused by effects intrinsic to the human lineage. We observed a

104   weak intermediate pattern when comparing NHGA to three archaic hominid genomes

105   (**Supplementary Figure 1**; **Supplementary Note**). This suggests that at least part of

106   the differentiation process in the distribution of mutation densities was already

107   established before the human-Neanderthal split.

108

109   Interestingly, correlations between a variety of genomic features and tumor mutation

110   density are consistently more similar to the correlations with NHGAs than with humans

111   (**Figure 2a**). Mutation densities in NHGAs have, like in humans, strong correlations

112   with sequence conservation and recombination rate (**Supplementary Figure 3**).

113   However, and strikingly, NHGAs show strong positive correlations with epigenomic

114   features associated with closed chromatin, just as tumors do (**Figure 2a,**

115   **Supplementary Table 7**). We also observe consistent associations with human

116   chromatin states[34] (**Figures 2b,c**). GC-content, H3K36me1, and CpG-content show a

117   clear positive correlation with human but negative with NHGAs and tumors, suggesting

118   that they might be contributing to the diagonal pattern (**Figure 2d,e and 3a,b**).

119   Interestingly, H3K36me1 has been shown to be specifically recruited in the gene

120   bodies of genes regulated by CpG islands although its role remains unclear[35].

121

122    Intrigued by the connection of several CpG-related features with the diagonal pattern

123    that implies stronger correlation between mutation densities in tumors and NHGA than

124    with the human germline (**Figure 3a,b**), we analyzed separately CpG>T transitions and

125    mutations at non-CpG sites. (**Figure 3c-f**). CpG>T transitions present very strong

126    correlations between all germline datasets and very poor correlations with tumor

127    (**Figure 3c,d**). The relationship between CpG-content and mutation density at non-

128    CpG sites is different in humans compared to NHGAs and tumors. Moreover, their

129    correlations are similar to those using all sites (**Figure 3e,f**). Correcting the mutation

130    density of CpG>T transitions by the regional CpG content homogenizes the directions

131    of the correlations with genomic features in all datasets (**Supplementary Notes,**

132    **Supplementary Figure 2**). Interestingly, this correlation is weaker in human than in

133    NHGA and in tumors (**Supplementary Notes, Supplementary Table 8,**

134    **Supplementary Figure 3**). This suggests that the differences in correlations with

135    genomic features are caused by differences in the relative contribution of non-

136    CpG/CpG>T mutation density in each dataset. The distribution of human *de novo*

137    mutations[36] at both non-CpG and CpG sites replicates the behavior of human germline

138    mutations showing very low correlations with tumor (**Supplementary Notes,**

139    **Supplementary Tables 3&9**). When comparing the distribution of non-CpG mutations,

140    we detect a horizontal pattern (**Supplementary Figure 3**) similar to those observed in

141    comparisons of high- and low-diversity subspecies. Therefore, the combination of the

142    behaviors of both non-CpG and CpG>T mutations causes the diagonal pattern

143    observed when comparing all SNVs.

144

145    To explore the contribution of different mechanisms to the observed mutation densities,

146    we analyzed their trinucleotide context. The triplet mutation spectra of human,

147    chimpanzee, and gorilla are very similar (**Supplementary Figure 4, Supplementary**

148    **Table 10**). It has been shown that the human mutation spectrum can be recapitulated

149    by a combination the cancer signatures SBS1 and SBS5[20,37]. We were able to replicate

150    this association in NHGA and another primate species (Vervet monkey)

151    (**Supplementary Notes, Supplementary Table 10**), suggesting its conservation in the

152    primate lineage.

153

154    A subset of trinucleotides is significantly enriched in one of the species (Chi-Squared

155    test p-value <10e-5). We detected no association between these trinucleotides and

156    known mutation mechanisms (**Figure 4a, see Methods, Supplementary Note,**

157    **Supplementary Figure 4, Supplementary Table 11**). However, linear regression

158    models show a positive and significant (p-value <10e-4) effect of the triplet's GC-

159    content and its fold-enrichment in the human-chimpanzee comparison

160    (**Supplementary Figure 4**). Only trinucleotides with similar enrichment between

161    species (non-CpG, mainly C>G and T>C) show differences in their distribution across

162    the genome between human, NHGA, and tumor (trinucleotide-difference test p-value

163    <10e-5, see Methods, **Supplementary Note, Figure 4a**).

164

165

166    We compared the association of the number of mutations caused by each cancer

167    signature[17,38] in each individual tumor type to the human-NHGA-tumor pattern

168    (**Supplementary Table 12**). Signatures SBS5 and SBS40 showed a significant

169    association (signature-difference test p-value <10e-4, see Methods) of the pattern with

170    the tumor's signature mutation load (**Figure 4b**). Both SBS5 and SBS40 are flat

171    signatures whose mutation load is associated with the age of the sample and with pre-

172    neoplastic mutations in tumors[17,38] This suggests that the strong correlation between

173    NHGA and tumor mutation densities is driven by conserved mechanisms in healthy

174    cells in the Great Ape lineage, while the genome-wide distribution of mutations has

175    been altered in the human germline.

176

177

178   **Discussion**

179   We analyzed the mutation density distribution at the 1Mbp scale in the human and

180   NHGA germlines, as well as in human tumors. We observed a moderate similitude

181   between human and NHGA germlines and, surprisingly, a higher resemblance

182   between human tumors with the germlines of NHGAs than with humans

183

184   These discrepancies in mutation density in the human and NHGA germlines are

185   differently associated with genomic and epigenomic features. Regions more densely

186   mutated in humans than in NHGAs tend to be GC-rich, exon-rich, promoter and

187   enhancer-rich, open chromatin and early replicating. Particularly, CpG-related features

188   show a positive correlation with human and a negative correlation with NHGA and

189   tumor mutation densities. The possible functional implications in human evolution

190   require further study.

191

192   These observations are driven by the different behavior of mutation density at CpG>T

193   transitions (very similar in all germlines and very different in tumors) and at non-CpG

194   sites (more similar in NHGAs and human tumors than in human germline). This is

195   exclusive of the human germline and, thus, must have been caused by human-specific

196   conditions.

197

198   We observed that human and other primates showed a very similar global triplet

199   mutation spectrum. We detected an enrichment of certain trinucleotide mutations in

200   humans and NHGAs consistent with previous results (non-CpG, GC-rich mutations are

201   enriched in humans)[25]. The enriched trinucleotides are not associated with mutation

202   signatures with known causes, nor do they contribute significantly to the higher

203   similitude of human tumors to NHGA germlines. This suggests the absence of strong

204   mechanistic changes biasing the accumulation of mutations in any of the studied

205   germlines.

8

206

207    As previously described for human[20,37], we observed that mutation rates of three non-

208    human Primates are explained by mutation signatures SBS1 (mostly CpG>T

209    mutations) and SBS5 (associated with "normal" accumulation of mutation in healthy

210    somatic and germline cells[16,39]). Moreover, the lower human-tumor than NHGA-tumor

211    correlation is driven by the accumulation of mutations associated with signatures SBS5

212    and SBS40 (similar to SBS5 and recently discovered[17]). These results suggest that the

213    poor human-tumor correlation is caused by the fact that human (but not NHGAs)

214    germline (and *de novo* mutations) do not currently reflect the expected mutation

215    densities of healthy (and pre-neoplastic-like) human somatic cells. One possible

216    explanation of this effect, would be if the recent slowdown in mutation rates in

217    humans[27] affected differently the different types of mutations.

218

219    We observed that the moderate human-NHGA and the low human-tumor correlations

220    of mutation densities at non-CpG sites could be caused by losses in population

221    diversity (as observed in low-diversity NHGA subspecies). We propose that successive

222    bottlenecks during human evolution removed a substantial part of nucleotide variation

223    that still remains to be recovered as a whole. In contrast, the hypermutability of CpG

224    sites and its concentration in specific regions caused CpG>T transitions to have

225    already recovered diversity levels similar to those of high-diversity NHGAs. Moreover,

226    the recent human-exclusive population expansions[30,40] are expected to cause an

227    increase of clock-like CpG>T mutations in the population[41,42], leaving signatures akin to

228    positive selection, as it has been described in Native Americans[24]. These effects

229    caused a decoupling of the CpG>T/non-CpG mutation rates within the same region,

230    stronger in humans than in NHGA and tumors. We cannot disregard an additional

231    contribution of human-specific shifts in CpG>T transitions mutation rates, although they

232    have been suggested to be similar across all Great Apes[42]. We propose that the

233    combination of population bottlenecks and expansions, together with the specific

234    nature of the different mutation types, drives the differences observed in the

235    distributions of human mutation densities.

236

237    Our results imply that accumulated mutations in human populations are a poor proxy of

238    the expected mutational background in healthy somatic cells. In fact, accumulated

239    mutations in NHGAs (at least at non-CpG sites) or even in tumors happen to be more

240    informative about the normal occurrence of mutations in healthy somatic cells.

241

**Methods:**

**Datasets used:**

For the human datasets we used the release variant calling of 2,504 humans from the 1000 Genomes Project[28] (1kGP), our own calling of 50 additional human samples from the Simons Genome Diversity Panel[29] (sgdp_50), and *de novo* mutations from 1,548 trios[36] that were mapped to the human reference hg19 using the liftOver tool[43]. We used our own mapping and calling of 69 chimpanzees and bonobos (59 chimpanzees and 10 bonobos, referred as chimpanzees in short)[30,31] and 43 gorillas[30,32]. We used the release variant calling of 3 archaic samples: Altai and Vindija 33.19 Neanderthals[44,45], and Denisova[46]. Finally, for the tumor dataset, we used the release variant calling of 2,583 human tumors from the Pan-Cancer Analysis of Whole Genomes Consortium[33].

**Definition of high-quality orthologous regions shared between human, chimpanzee and gorilla genomes**

We mapped and called chimpanzee and bonobo, gorilla, and human (sgdp_50) samples to the human reference hg19 using BWA MEM[47] and GATK[48] following the best practices protocols[49,50] and additional quality filters (**Supplementary Notes**).

To avoid missmappings to the human reference and erroneous estimates of mutation density in the NHGA samples (too low density caused by lack of mapping reads or deletions or too high density caused by collapsed duplications[51]) we filtered out any region of the human reference genome hg19 failing one of the following criteria: poor mappability of the human reference split into 35bp k-mers, poor callability in ≥25% of the chimpanzee or gorilla samples, or, matching a known Copy Number Variable region in NHGA samples[52] (**Supplementary Notes**). 2,052Mbp of autosomal sequence passed this filtering (76.54% of the non-N human reference autosomes). We divided

11

270    the autosomes into 1Mbp overlapping (500kb) windows and kept all windows where

271    ≥50% of its bases passed our filtering. This left 5,040 1Mbp windows to analyze

272    (**Supplementary Figure 1, Supplementary Table 2**).

273

274    These filters were applied to all datasets used, including both our callings and external

275    datasets used as released. All SNV counts, trinucleotide counts, and genomic features

276    measurements through this study used only regions passing this filtering. For the

277    analysis of archaic samples, we combined this filtering with the intersection of the

278    callability mask of all 3 archaic samples. This specific filtering was applied to all

279    datasets when compared with the archaic samples.

280

281

282    **Mutation density:**

283    We measured mutation density of each window in each dataset by counting either the

284    number of non-fixed segregating sites (in the human, chimpanzee and gorilla datasets)

285    or the number of somatic mutations (in the tumor and human *de novo* datasets,

286    accounting repeated mutations as independent mutational events). We divided this

287    count of Single Nucleotide Variants (SNV) by the fraction of the window passing our

288    filtering. This results in a measure of mutations per Megabasepair (Mbp) of sequence

289    for each window. We standardized the resulting distribution within each dataset

290    deeming it as the mutation density. We ranked all windows within a dataset by their

291    distribution of mutation density to control for the different shapes of the datasets

292    distributions.

293

294    **Correlations between distributions:**

295    All correlations used in this analysis are Pearson's correlation (using the R function

296    cor.test) between the standardized mutation densities (unranked) of the two datasets

12

297    unless otherwise specified. Partial correlations, when used, were calculated using the

298    pcor function from the ppcor R package.

299

300    **Significance of the diagonal split:**

301    To measure the significance of the diagonal split pattern observed when comparing the

302    human and NHGA datasets, we divided all windows into two groups depending on if

303    the ranked mutation density is higher in human than NHGAs or vice-versa. We

304    calculated the two-sided Mann-Whitney U test on the variable of interest (usually, the

305    tumor mutation density) on both groups using the R function wilcox.test.

306

307    **Genomic Features:**

308    The genomic features used were filtered using the same mappability, callability and

309    copy-number filters used for the mutation density data. The features used were either

310    the overlap of the feature's genomic coordinates with the fraction of the 1Mbp window

311    passing our filtering (e.g. GC-content, CpG-content), or the average value or intensity

312    of the feature in the passing fraction of the window (e.g. histone marks), depending on

313    the original data (**Supplementary Table 7**).

314

315    **Trinucleotides**:

316    We classified each SNV into the 96 possible combinations of trinucleotides (12

317    different mutation types, by 16 combinations of the adjacent nucleotides, divided by

318    two when folding them). We determined the adjacent reference sequence of each SNV

319    using the getfasta option of bedtools[53]. We filtered out any variant where the liftOver

320    tool[43] could not map them to the chimpanzee panTro5 or the gorilla gorGor5 reference

321    genomes, or the trinucleotide sequence differed in one of the three reference

322    genomes. This filter was applied to all windows and we used for our analysis only

323    windows where ≥50% of it passed both the original high-quality orthologous regions

324    and this 3-reference filter, leaving 4,920 windows to use. We applied additional filters

325     requiring the trinucleotide to be species-exclusive and to not overlap variants in other

326     species (**Supplementary Note**). This resulted in a high-confidence set of species-

327     exclusive trinucleotides where the ancestral and derived alleles could be reliably

328     inferred. This filtering affected more CpG>T than non-CpG sites, due to the recurrent

329     nature of CpG>T transitions (**Supplementary Table 10**).

330

331     **Mutation spectra:**

332     We calculated each species' mutation spectra as the fraction of all trinucleotides in a

333     dataset belonging to one of the 96 trinucleotides. We calculated correlations between

334     datasets using Pearson's correlation (cor.test function in R). We measured the

335     correlation of the mutation spectrum of each species and the combined effect of cancer

336     mutation signatures SBS1 and SBS5[17,38] by the formula: 0.1*SBS1+0.9*SBS5, as

337     CpG>T transitions are the main components of signature SBS1 and they represent

338     ~10% of the trinucleotides in both the human and NHGA datasets.

339

340     **Whole-genome enrichment of trinucleotides:**

341     We calculated the enrichment and its significance in each germline dataset pair

342     (human-chimpanzee, human-gorilla, chimpanzee-gorilla) using the method described

343     in Harris, 2017[25]. We calculated the enrichment of trinucleotide T between species A

344     and B by dividing *fraction of T in species A / fraction of T in species B*. We calculated a

345     chi-squared test using a contingency table with: the trinucleotide count in species A, in

346     species B, the count of the rest of trinucleotides in species A, and in species B. As the

347     counts of trinucleotides are not independent from each other, we sorted all

348     trinucleotides from most to least significant, and rerun the test by decreasing

349     significance order, while removing the previously used trinucleotides from the count of

350     total trinucleotides.

351

352    CpG>T transitions are highly affected by the sample size of the datasets. We ran all

353    the tests using both 1kGP and sgdp_50 as the human dataset. We detected

354    incoherences on the significance and direction of the results in two CpG>T

355    trinucleotides. We report the results using 1kGP where tests using both 1kGP and

356    sgdp_50 are coherent in both significance and direction of the enrichment.

357    The top 10% most enriched trinucleotides in each species pairwise comparison were

358    compared with cancer mutation signatures[38], and reported when the trinucleotide

359    represented ≥5% of the mutations within a signature.

360

361    **Trinucleotide distribution difference test (trinucleotide-difference test):**

362    We developed a method to determine which trinucleotides contribute significantly to the

363    difference between NHGAs-tumors and human-tumors mutation density correlations:

364    For each trinucleotide T and each pair of species (human-chimpanzee, human-gorilla,

365    and, chimpanzee-gorilla) we, subtract the ranked mutation density of T in species A

366    minus the ranking in tumor, and in species B minus tumor. We calculate the two-sided

367    Kolmogorov-Smirnov test (using the R function ks.test) of the two resulting

368    distributions. We use the p-value of the ks-test as the significance of the test and the

369    difference between the standard deviation of both distributions (as both have a mean of

370    0) as the test's effect size. The results when using 1kGP or sgdp_50 as the human

371    datasets are concordant in the direction of the association, but we discarded the

372    sgdp_50 results because the smaller number of SNV (and of each trinucleotide type)

373    results in lower power when using sgdp_50.

374

375    **Association of GC-content in the trinucleotide sequence:**

376    We counted the number of Cytosine and Guanine bases in each trinucleotide and built

377    a linear regression (using the R function glm). The GC-content of the triplet acted as a

378    predictor of the result of the test (the log10 fold-enrichment in the whole-genome

379    enrichment analysis or the difference between the standard deviation of both

380    distributions in the trinucleotide-difference test).

381

382    **Mutation load-difference test per mutation signature (signature-difference test):**

383    In order to determine the contribution of each mutation signature to the difference

384    between NHGAs-tumors and human-tumors mutation density correlations, we rerun

385    the trinucleotide-difference test using the 1kGP and chimpanzee datasets, while using

386    the different individual tumor types (**Supplementary Table 12**). For each trinucleotide,

387    tumor type and mutation signature, we built a linear regression (using R's glm function)

388    where the mutation load of that signature in that tumor type[17] predicted the effect size

389    in the trinucleotide-difference test for that tumor type (**Supplementary Note**). For each

390    signature, we built a contingency table where all 96 trinucleotides where classified by

391    whether being significant or not (p-value <0.05) in the trinucleotide-difference test, and

392    the significance of the mutation load in the linear regression model. We ran a chi-

393    squared test on that contingency table and obtained its significance.

16

**Author Contributions:**

J.M.H.G. performed all the analysis. J.M.H.G and D.J wrote the manuscript. T.M.B., D.J. and A.N. conceived and supervised this work. All the authors read and approved the final manuscript.

**Competing interests statement:**

All authors declare no competing interests

**Data availability statement:**

No new data was generated for this study. All the analyses were performed using publicly available data obtained from their original publications, as referenced.

# References

1. Hodgkinson, A., Chen, Y. & Eyre-Walker, A. The large-scale distribution of somatic mutations in cancer genomes. *Hum. Mutat.* **33**, 136–143 (2012).

2. Tyekucheva, S. *et al.* Human-macaque comparisons illuminate variation in neutral substitution rates. *Genome Biol.* **9**, R76 (2008).

3. Ananda, G., Chiaromonte, F. & Makova, K. D. A genome-wide view of mutation rate co-variation using multivariate analyses. *Genome Biol.* **12**, R27 (2011).

4. Hodgkinson, A. & Eyre-Walker, A. Variation in the mutation rate across mammalian genomes. *Nat. Rev. Genet.* **12**, 756–766 (2011).

5. Makova, K. D. & Hardison, R. C. The effects of chromatin organization on variation in mutation rates in the genome. *Nat. Rev. Genet.* **16**, 213–223 (2015).

6. Schuster-Böckler, B. & Lehner, B. Chromatin organization is a major influence on regional mutation rates in human cancer cells. *Nature* **488**, 504–507 (2012).

7. Supek, F. & Lehner, B. Differential DNA mismatch repair underlies mutation rate variation across the human genome. *Nature* **521**, 81–84 (2015).

8. Supek, F. & Lehner, B. Scales and mechanisms of somatic mutation rate variation across the human genome. *DNA Repair (Amst).* 102647 (2019). doi:10.1016/j.dnarep.2019.102647

9. Polak, P. *et al.* Cell-of-origin chromatin organization shapes the mutational landscape of cancer. *Nature* **518**, 360–364 (2015).

10. Kübler, K. *et al.* Tumor mutational landscape is a record of the pre-malignant state. *bioRxiv* 517565 (2019). doi:10.1101/517565

11. Aggarwala, V. & Voight, B. F. An expanded sequence context model broadly explains variability in polymorphism levels across the human genome. *Nat. Genet.* **48**, 349–55 (2016).

12. Hwang, D. G. & Green, P. Bayesian Markov chain Monte Carlo sequence

18

445        analysis reveals varying neutral substitution patterns in mammalian evolution.

446        *Proc. Natl. Acad. Sci. U. S. A.* **101**, 13994–4001 (2004).

447   13.   Keightley, P. D., Eöry, L., Halligan, D. L. & Kirkpatrick, M. Inference of Mutation

448        Parameters and Selective Constraint in Mammalian Coding Sequences by

449        Approximate Bayesian Computation. *Genetics* **187**, 1153–1161 (2011).

450   14.   Siepel, A. & Haussler, D. Phylogenetic Estimation of Context-Dependent

451        Substitution Rates by Maximum Likelihood. *Mol. Biol. Evol.* **21**, 468–488 (2003).

452   15.   Nachman, M. W. & Crowell, S. L. Estimate of the mutation rate per nucleotide in

453        humans. *Genetics* **156**, 297–304 (2000).

454   16.   Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer.

455        *Nature* **500**, 415–421 (2013).

456   17.   Alexandrov, L. *et al.* The Repertoire of Mutational Signatures in Human Cancer.

457        *bioRxiv* 322859 (2018). doi:10.1101/322859

458   18.   Martincorena, I. *et al.* Tumor evolution. High burden and pervasive positive

459        selection of somatic mutations in normal human skin. *Science* **348**, 880–6

460        (2015).

461   19.   Martincorena, I. *et al.* Somatic mutant clones colonize the human esophagus

462        with age. *Science* **362**, 911–917 (2018).

463   20.   Rahbari, R. *et al.* Timing, rates and spectra of human germline mutation. *Nat.*

464        *Genet.* **48**, 126–133 (2016).

465   21.   Smith, T. C. A., Arndt, P. F. & Eyre-Walker, A. Large scale variation in the rate of

466        germ-line de novo mutation, base composition, divergence and diversity in

467        humans. *PLOS Genet.* **14**, e1007254 (2018).

468   22.   Carlson, J. *et al.* Extremely rare variants reveal patterns of germline mutation

469        rate heterogeneity in humans. *Nat. Commun.* **9**, 3753 (2018).

470   23.   Harris, K. Evidence for recent, population-specific evolution of the human

471        mutation rate. *Proc. Natl. Acad. Sci. U. S. A.* **112**, 3439–44 (2015).

472   24.   Mathieson, I. & Reich, D. Differences in the rare variant spectrum among human

473        populations. *PLOS Genet.* **13**, e1006581 (2017).

474    25.   Harris, K. & Pritchard, J. K. Rapid evolution of the human mutation spectrum.

475        *Elife* **6**, (2017).

476    26.   Narasimhan, V. M. *et al.* Estimating the human mutation rate from autozygous

477        segments reveals population differences in human mutational processes. *Nat.*

478        *Commun.* **8**, 303 (2017).

479    27.   Besenbacher, S., Hvilsom, C., Marques-Bonet, T., Mailund, T. & Schierup, M. H.

480        Direct estimation of mutations in great apes reconciles phylogenetic dating. *Nat.*

481        *Ecol. Evol.* **3**, 286–292 (2019).

482    28.   The 1000 Genomes Project Consortium. A global reference for human genetic

483        variation. *Nature* **526**, 68–74 (2015).

484    29.   Mallick, S. *et al.* The Simons Genome Diversity Project: 300 genomes from 142

485        diverse populations. *Nature* **538**, (2016).

486    30.   Prado-Martinez, J. *et al.* Great ape genetic diversity and population history.

487        *Nature* **499**, 471–5 (2013).

488    31.   de Manuel, M. *et al.* Chimpanzee genomic diversity reveals ancient admixture

489        with bonobos. *Science* **354**, 477–481 (2016).

490    32.   Xue, Y. *et al.* Mountain gorilla genomes reveal the impact of long-term

491        population decline and inbreeding. *Science (80-. ).* **348**, 242–245 (2015).

492    33.   The ICGC/TCGA Pan-Cancer Analysis of Whole Genomes Network. Pan-cancer

493        analysis of whole genomes. *Prep.* (2019).

494    34.   Ernst, J. *et al.* Mapping and analysis of chromatin state dynamics in nine human

495        cell types. *Nature* **473**, 43–49 (2011).

496    35.   Vavouri, T. & Lehner, B. Human genes with CpG island promoters have a

497        distinct transcription-associated chromatin organization. *Genome Biol.* **13**, R110

498        (2012).

499    36.   Jónsson, H. *et al.* Parental influence on human germline de novo mutations in

500        1,548 trios from Iceland. *Nature* **549**, 519–522 (2017).

501    37.    Alexandrov, L. B. *et al.* Clock-like mutational processes in human somatic cells.

502           *Nat. Genet.* **47**, 1402–1407 (2015).

503    38.    Tate, J. G. *et al.* COSMIC: the Catalogue Of Somatic Mutations In Cancer.

504           *Nucleic Acids Res.* **47**, D941–D947 (2019).

505    39.    Blokzijl, F. *et al.* Tissue-specific mutation accumulation in human adult stem

506           cells during life. *Nature* **538**, 260–264 (2016).

507    40.    Li, H. & Durbin, R. Inference of human population history from individual whole-

508           genome sequences. *Nature* **475**, 493–496 (2011).

509    41.    Harpak, A., Bhaskar, A. & Pritchard, J. K. Mutation Rate Variation is a Primary

510           Determinant of the Distribution of Allele Frequencies in Humans. *PLoS Genet.*

511           **12**, e1006489 (2016).

512    42.    Moorjani, P., Amorim, C. E. G., Arndt, P. F. & Przeworski, M. Variation in the

513           molecular clock of primates. *Proc. Natl. Acad. Sci.* **113**, 10607–10612 (2016).

514    43.    Hinrichs, A. S. *et al.* The UCSC Genome Browser Database: update 2006.

515           *Nucleic Acids Res.* **34**, D590-8 (2006).

516    44.    Prüfer, K. *et al.* The complete genome sequence of a Neanderthal from the Altai

517           Mountains. *Nature* **505**, 43–9 (2014).

518    45.    Prüfer, K. *et al.* A high-coverage Neandertal genome from Vindija Cave in

519           Croatia. *Science* **358**, 655–658 (2017).

520    46.    Meyer, M. *et al.* A high-coverage genome sequence from an archaic Denisovan

521           individual. *Science* **338**, 222–6 (2012).

522    47.    Li, H. Aligning sequence reads, clone sequences and assembly contigs with

523           BWA-MEM. (2013).

524    48.    McKenna, A. *et al.* The Genome Analysis Toolkit: A MapReduce framework for

525           analyzing next-generation DNA sequencing data. *Genome Res.* **20**, 1297–1303

526           (2010).

527    49.    DePristo, M. A. *et al.* A framework for variation discovery and genotyping using

528           next-generation DNA sequencing data. *Nat. Genet.* **43**, 491–498 (2011).

529    50.    Van der Auwera, G. A. *et al.* From FastQ data to high confidence variant calls:

530          the Genome Analysis Toolkit best practices pipeline. *Curr. Protoc. Bioinforma.*

531          **43**, 11.10.1-33 (2013).

532    51.    Hartasánchez, D. A., Brasó-Vives, M., Heredia-Genestar, J. M., Pybus, M. &

533          Navarro, A. Effect of Collapsed Duplications on Diversity Estimates: What to

534          Expect. *Genome Biol. Evol.* **10**, 2899–2905 (2018).

535    52.    Sudmant, P. H. *et al.* Evolution and diversity of copy number variation in the

536          great ape lineage. *Genome Res.* **23**, 1373–82 (2013).

537    53.    Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing

538          genomic features. *Bioinformatics* **26**, 841–842 (2010).

539    54.    Barski, A. *et al.* High-Resolution Profiling of Histone Methylations in the Human

540          Genome. *Cell* **129**, 823–837 (2007).

541

542

543

544 **Figure legends:**

545 **Figure 1:** Distribution of mutation density across datasets. **a)** Distribution of the

546 standardized mutation density in 1Mbp windows in human, NHGA, and tumor datasets.

547 The numbers next to the legend represent the fold-enrichment between the 95th and 5th

548 quantiles. **b)** Distribution of the standardized mutation density in humans, chimpanzee

549 and tumor. Each point represents a 1Mbp window. The x-axis represents the human

550 mutation density, the y-axis the chimpanzee mutation density, and the point color, the

551 tumor mutation density. The black line represents the diagonal where the mutation

552 density is equal in human and chimpanzee. **c)** Same as b but comparing human and

553 gorilla. **d)** Distribution of the ranked mutation density in humans, chimpanzee and

554 tumor. Each point represents a 1Mbp window. The x and y axis represent the ranking

555 in mutation density in human and chimpanzee, respectively. Color of points represents

556 the ranked mutation density in the tumor dataset. The solid black line represents the

557 diagonal where the ranked mutation density is equal in human and chimpanzee. The

558 dashed lines represent 25% difference in ranking in both species. **e)** Same as d,

559 comparing human and gorilla.

560

23

Figure 1



Distribution of mutations

561     **Figure 2:** Genomic Features. **a)** Pearson's correlation R of different datasets with

562     human genomic features (**Supplementary Table 7**). **b)** Overlap of heterochromatin in

563     human lymphoblastoid cell lines (LCLs) measured by chromHMM states[34] compared

564     with the human and chimpanzee ranked mutation density distribution. **c)** Same as b but

565     using the aggregate chromHMM states associated with the presence of promoters. **d)**

566     same as b and c but color denotes the window's GC-content, **e)** density of H3K36me1

567     histone mark ChIP-seq reads[54].

568

**Figure 2**

a. Correlation with genomic features

Legend:
- tumor
- 1kGP
- sgdp_50
- chimpanzee
- gorilla

b. Heterochromatin
c. Promoters
d. GC-content
e. H3K36me1

569    **Figure 3:** CpG-content. **a)** Distribution of the CpG-content in the human reference

570    hg19 compared with the ranked mutation density in human and chimpanzee, **b)** loess

571    smoothers of mutation density rank and CpG-content for the different datasets. **c)**

572    CpG>T transitions corrected by the whole window size; loess smoothers same as in b;

573    **d)** correlation of the standardized mutation density of CpG>T transitions in different

574    species; **e)** same as in b,c, but using only mutations at non-CpG sites; **f)** correlation of

575    the standardized mutation density of mutations at non-CpG sites in different species.


576

25
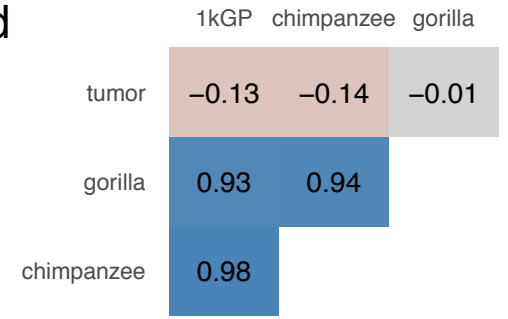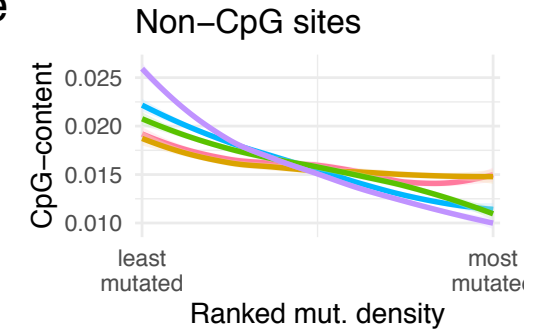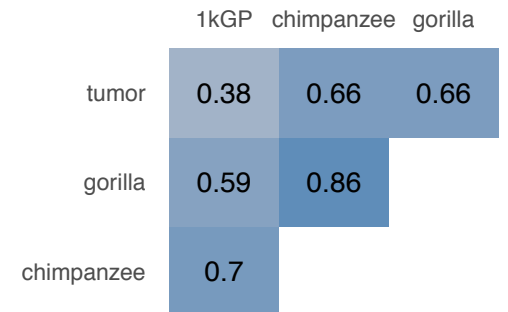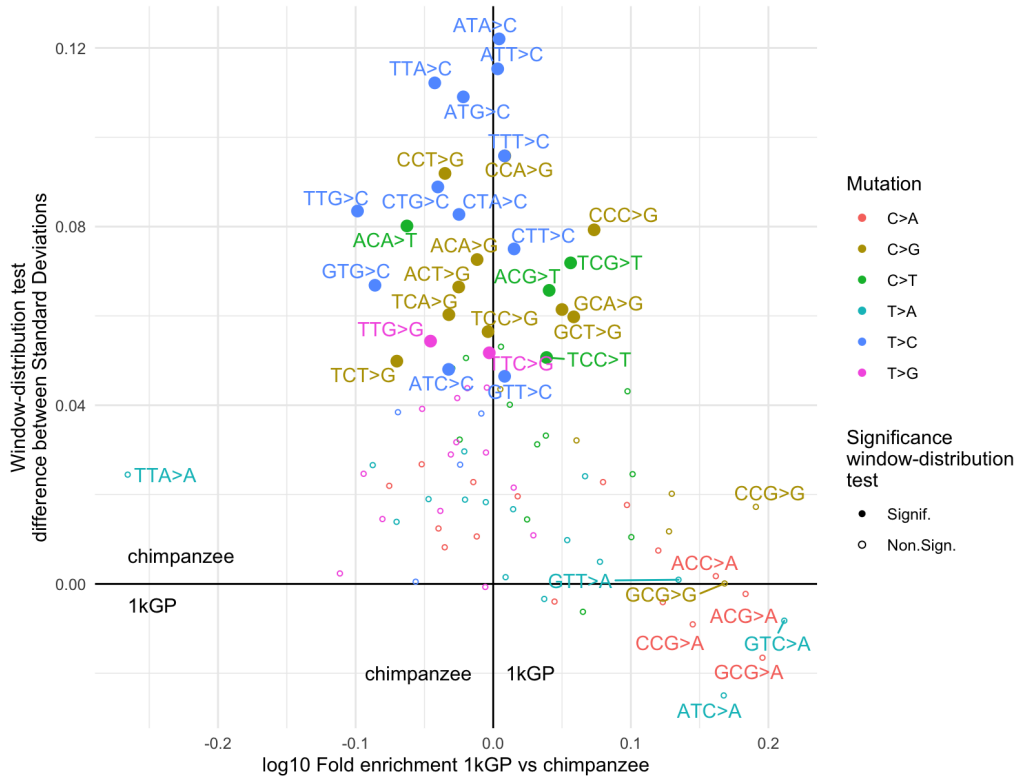
Figure 3

**Association with CpG sites**

577    **Figure 4:** Trinucleotide analysis. **a)** Contribution to the higher chimpanzee-tumors

578    mutation distribution similarity Vs. genome-wide enrichment in human compared to

579    chimpanzee. X-axis: log10 of the enrichment of trinucleotides comparing human and

580    chimpanzee. Left: enriched in chimpanzee; right: enriched in human. Y-axis: effect size

581    (difference between the standard deviations of human-tumor and chimpanzee-tumor)

582    of the trinucleotide-difference test (see Methods). Positive values: tumor distribution

583    more similar to chimpanzee; negative values: tumor distribution more similar to human.

584    Color represents the central nucleotide mutation type. Filled dots represent mutation

585    types significant (p-value <1e-5) in the trinucleotide-difference test. **b)** -Log10 p-values

586    of the association of each cancer signature mutation load to the trinucleotide-difference

587    test (signature-difference test; see Methods). Color represents the number of mutations

588    associated with each signature in the whole dataset. Dot size represents the number of

589    tumor types with two or more samples showing the signature. Only non-artifact

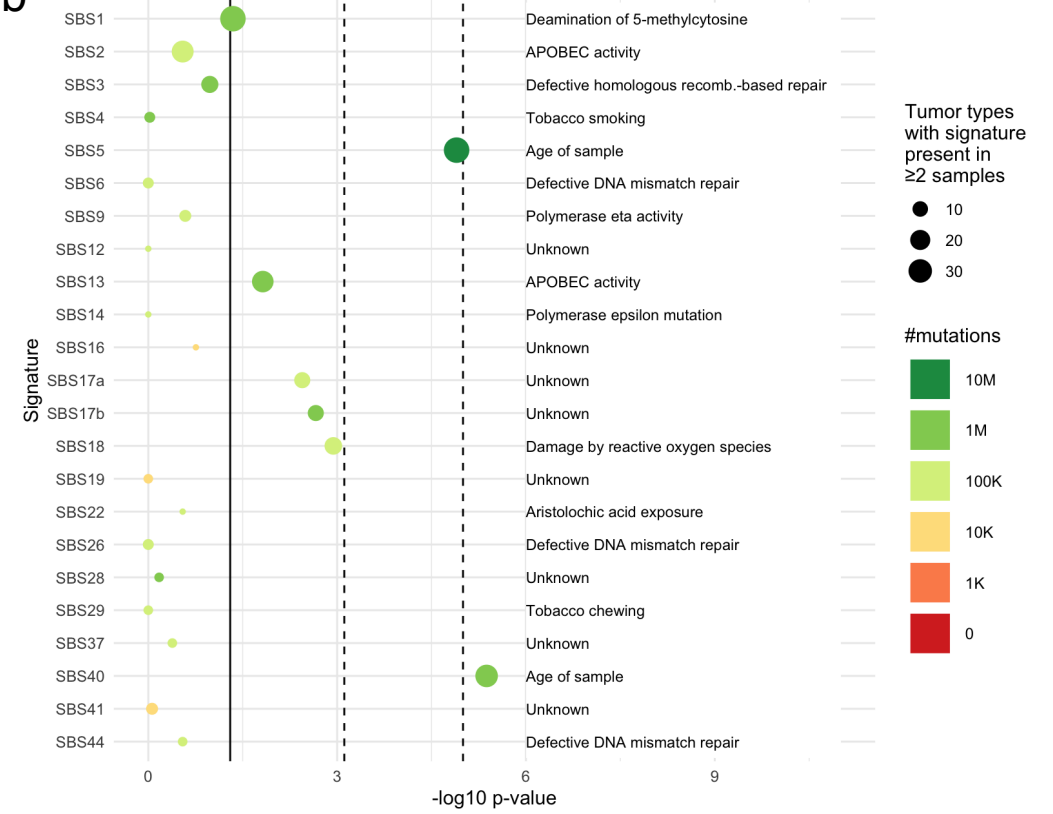590    signatures present in 2 or more tumor types are shown.

591

# Figure 4

## Analysis of trinucleotide mutations

592 **Table legends:**

593

594 **Table 1:** Correlations. Pairwise Pearson's correlation R of the standardized mutation

595 density of 5,040 1Mbp windows between datasets.

596

Correlation between distributions of mutation density

|  | 1kGP | Chimpanzee | Gorilla |
|---|---|---|---|
| **Chimpanzee** | 0.65 | - | - |
| **Gorilla** | 0.53 | 0.84 | - |
| **Tumor** | 0.16 | 0.55 | 0.58 |

597