

1 Rapid evolution of the human mutation spectrum.

Kelley Harris¹ and Jonathan K. Pritchard^{1,2,3}

¹Department of Genetics, Stanford University

²Department of Biology, Stanford University

³Howard Hughes Medical Institute, Stanford University

2 Correspondence: kellelyh@stanford.edu, pritch@stanford.edu.

3 March 24, 2017

4

Abstract

5 DNA is a remarkably precise medium for copying and storing biological information.
6 This high fidelity results from the action of hundreds of genes involved in replication,
7 proofreading, and damage repair. Evolutionary theory suggests that in such a system,
8 selection has limited ability to remove genetic variants that change mutation rates
9 by small amounts or in specific sequence contexts. Consistent with this, using SNV
10 variation as a proxy for mutational input, we report here that mutational spectra differ
11 substantially among species, human continental groups and even some closely-related
12 populations. Close examination of one signal, an increased TCC→TTC mutation rate
13 in Europeans, indicates a burst of mutations from about 15,000 to 2,000 years ago,
14 perhaps due to the appearance, drift, and ultimate elimination of a genetic modifier of
15 mutation rate. Our results suggest that mutation rates can evolve markedly over short
16 evolutionary timescales and suggest the possibility of mapping mutational modifiers.

17 Main Text

18 Germline mutations provide the raw material for evolution, but also generate genetic load
19 and inherited disease. Indeed, the vast majority of mutations that affect fitness are deleterious,
20 and hence biological systems have evolved elaborate mechanisms for accurate DNA
21 replication and repair of diverse types of spontaneous damage. Due to the combined action
22 of hundreds of genes, mutation rates are extremely low—in humans, about 1 point mutation
23 per 100MB or about 60 genome-wide per generation [1, 2].

24 While the precise roles of most of the relevant genes have not been fully elucidated,
25 research on somatic mutations in cancer has shown that defects in particular genes can
26 lead to increased mutation rates within very specific sequence contexts [3, 4]. For example,
27 mutations in the proofreading exonuclease domain of DNA polymerase ϵ cause TCT→TAT
28 and TCG→TTG mutations on the leading DNA strand [5]. Mutational shifts of this kind
29 have been referred to as “mutational signatures”. Specific signatures may also be caused by
30 nongenetic factors such as chemical mutagens, UV damage, or guanine oxidation [6].

31 Together, these observations imply a high degree of specialization of individual genes
32 involved in DNA proofreading and repair. While the repair system has evolved to be ex-
33 tremely accurate overall, theory suggests that in such a system, natural selection may have
34 limited ability to fine-tune the efficacy of individual genes [7, 8]. If a variant in a repair gene
35 increases or decreases the overall mutation rate by a small amount—for example, only in a
36 very specific sequence context—then the net effect on fitness may fall below the threshold at
37 which natural selection is effective. (Drift tends to dominate selection when the change in
38 fitness is less than the inverse of effective population size). The limits of selection on muta-
39 tion rate modifiers are especially acute in recombining organisms such as humans because a
40 variant that increases the mutation rate can recombine away from deleterious mutations it
41 generates elsewhere in the genome.

42 Given these theoretical predictions, we hypothesized that there may be substantial scope
43 for modifiers of mutation rates to segregate within human populations, or between closely
44 related species. Most triplet sequence contexts have mutation rates that vary across the
45 evolutionary tree of mammals [9], but evolution of the mutation spectrum over short time
46 scales has been less well described. Weak natural mutators have recently been observed
47 in yeast [10] and inferred from human haplotype data [11]; if such mutators affect specific
48 pathways of proofreading or repair, then we may expect shifts in the abundance of mutations

49 within particular sequence contexts. Indeed, one of us has recently identified a candidate
50 signal of this type, namely an increase in TCC→TTC transitions in Europeans, relative
51 to other populations [12]; this was recently replicated [13]. Here we show that mutation
52 spectrum change is much more widespread than these initial studies suggested: although the
53 TCC→TTC rate increase in Europeans was unusually dramatic, smaller-scale changes are
54 so commonplace that almost every great ape species and human continental group has its
55 own distinctive mutational spectrum.

56 Results

57 To investigate the mutational processes in different human populations, we classified each
58 single nucleotide variants (SNV) in the 1000 Genomes Phase 3 data [14] in terms of its
59 ancestral allele, derived allele, and 5' and 3' flanking nucleotides. We collapsed strand com-
60 plements together to obtain 96 SNV categories. Since the detection of singletons may vary
61 across samples, and because some singletons may result from cell-line or somatic mutations,
62 we only considered variants seen in more than one copy. We further excluded variants in
63 annotated repeats (since read mapping error rates may be higher in such regions) and in
64 PhyloP conserved regions (to avoid selectively constrained regions) [15]. From the remaining
65 sites, we calculated the distribution of derived SNVs carried by each Phase 3 individual. We
66 used this as a proxy for the mutational input spectrum in the ancestors of each individual.

67 To explore global patterns of the mutation spectrum, we performed principal component
68 analysis (PCA) in which each individual was characterized simply by the fraction of their
69 derived alleles in each of the 96 SNV categories (Fig. 1A). PCA is commonly applied to
70 individual-level genotypes, in which case the PCs are usually highly correlated with geog-
71 raphy [16]. Although the triplet mutation spectrum is an extremely compressed summary
72 statistic compared to typical genotype arrays, we found that it contains sufficient information
73 to reliably classify individuals by continent of origin. The first principal component sepa-
74 rated Africans from non-Africans, and the second separated Europeans from East Asians,
75 with South Asians and admixed native Americans (Figure 1–Figure Supplement 2) appearing
76 intermediate between the two.

77 Remarkably, we found that the mutation spectrum differences among continental groups
78 are composed of small shifts in the abundance of many different mutation types (Fig. 1B).
79 For example, comparing Africans and Europeans, 43 of the 96 mutation types are signifi-

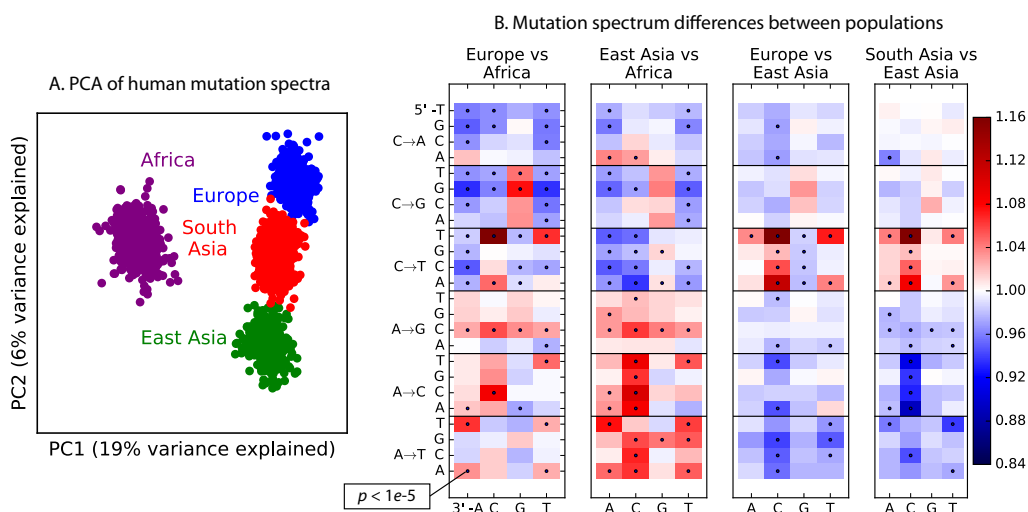


Figure 1: Global patterns of variation in SNV spectra. **A.** *Principal Component Analysis of individuals according to the fraction of derived alleles that each individual carries in each of 96 mutational types.* **B.** *Heatmaps showing, for pairs of continental groups, the ratio of the proportions of SNVs in each of the 96 mutational types. Each block corresponds to one mutation type; within blocks, rows indicate the 5' nucleotide, and columns indicate the 3' nucleotide. Red colors indicate a greater fraction of a given mutation type in the first-listed group relative to the second. Points indicate significant contrasts at $p < 10^{-5}$. See Figure Supplements 1, 2, and 3 for heatmap comparisons between additional population pairs as well as a description of PCA loadings and the p -values of all mutation class enrichments. Figure Supplement 4 demonstrates that these patterns are unlikely to be driven by biased gene conversion. In Figure Supplement 5, we see that this mutation spectrum structure replicates on both strands of the transcribed genome as well as the non-transcribed portion of the genome. Figure Supplements 6, 7, and 8 show that most of this structure replicates across multiple chromatin states and varies little with replication timing.*

80 cant at a $p < 10^{-5}$ threshold using a forward variable selection procedure. The previously
81 described TCC→TTC signature partially drives the difference between Europeans and the
82 other groups, but most other shifts are smaller in magnitude and appear to be spread over
83 more diffuse sets of related mutation types. East Asians have excess A→T transversions in
84 most sequence contexts, as well as about 10% more *AC→*CC mutations than any other
85 group. Compared to Africans, all Eurasians have proportionally fewer C→* mutations rela-
86 tive to A→* mutations.

87 **Replication of mutation spectrum shifts.** One possible concern is that batch effects or
88 other sequencing artifacts might contribute to differences in mutational spectra. Therefore
89 we replicated our analysis using 201 genomes from the Simons Genome Diversity Project [17].
90 The SGDP genomes were sequenced at high coverage, independently from 1000 Genomes,
91 using an almost non-overlapping panel of samples. We found extremely strong agreement
92 between the mutational shifts in the two data sets (Fig. 2). For example, all of the 43
93 mutation types with a significant difference between Africa and Europe (at $p < 10^{-5}$) in
94 1000 Genomes also show a frequency difference in the same direction in SGDP (comparing
95 Africa and West Eurasia). In both 1000 Genomes and SGDP, the enrichment of *AC→*CC
96 mutations in East Asia is larger in magnitude than any other signal aside from the previously
97 described TCC→TTC imbalance.

98 The greatest discrepancies between 1000 Genomes and SGDP involve transversions at
99 CpG sites, which are among the rarest mutational classes. These discrepancies might re-
100 sult from data processing differences or random sampling variation, but might also reflect
101 differences in the fine-scale ethnic composition of the two panels.

102 **Evidence for a pulse of TCC→TTC mutations in Europe and South Asia.** To
103 investigate the timescale over which the mutation spectrum change occurred, we analyzed
104 the allele frequency distribution of TCC→TTC mutations, which are highly enriched in
105 Europeans (Fig. 3A; $p < 1 \times 10^{-300}$ for Europe vs. Africa) and to a lesser extent in South
106 Asians. We calculated allele frequencies both in 1000 Genomes and in the larger UK10K
107 genome panel [18]. As expected for a signal that is primarily European, we found particular
108 enrichment of these mutations at low frequencies. But surprisingly, the enrichment peaks
109 around 0.6% frequency in UK10K, and there is practically no enrichment among the very
110 lowest frequency variants (Figure 3B and Figure 3–Figure Supplement 1). C→T mutations

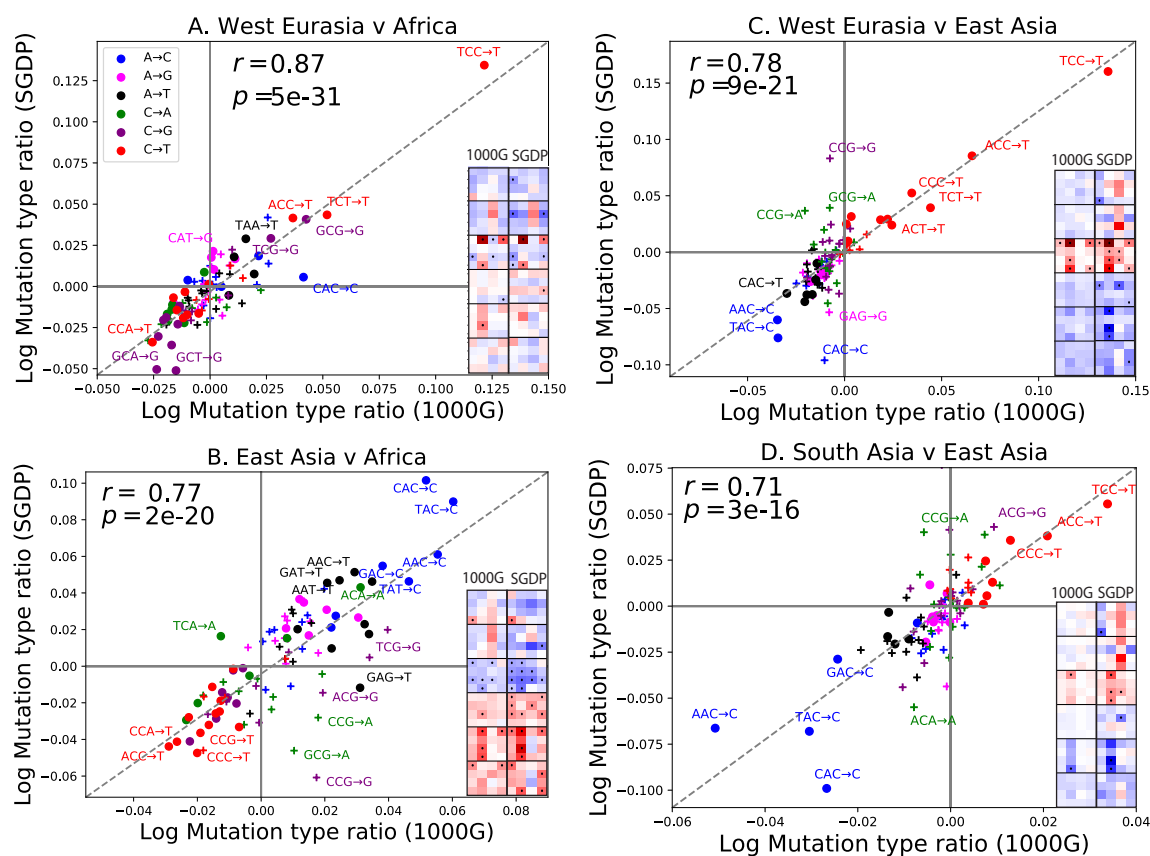


Figure 2: **Concordance of mutational shifts in 1000 Genomes vs. SGDP.** Each panel shows natural-log mutation spectrum ratios between a pair of continental groups, based on 1000 Genomes (x-axis) and SGDP (y-axis) data. Data points encoded by (+) symbols denote mutation types that are not significantly enriched in either population in the Figure 1 1000 Genomes analysis ($p < 10^{-5}$). These heatmaps use the same labeling and color scale as in Figure 1. All 1000 Genomes ratios in this figure were estimated after projecting the 1000 Genomes site frequency spectrum down to the sample size of SGDP. See Figure Supplements 1 and 2 for a complete set of SGDP heatmaps and regressions versus 1000 Genomes.

111 on other backgrounds, namely within TCT, CCC and ACC contexts, are also enriched in
112 Europe and South Asia and show a similar enrichment around 0.6% frequency that declines
113 among rarer variants (Fig. 3C). This suggests that these four mutation types comprise the
114 signature of a single mutational pulse that is no longer active. No other mutation types show
115 such a pulse-like distribution in UK10K, though several types show evidence of monotonic
116 rate change over time (Figure 3–Figure Supplements 3,4 and 5).

117 We used the enrichment of TCC→TTC mutations as a function of allele frequency to
118 estimate when this mutation pulse was active. Assuming a simple piecewise-constant model,
119 we infer that the rate of TCC→TTC mutations increased dramatically ~15,000 years ago
120 and decreased again ~2,000 years ago. This time-range is consistent with results showing this
121 signal in a pair of prehistoric European samples from 7,000 and 8,000 years ago, respectively
122 [13]. We hypothesize that this mutation pulse may have been caused by a mutator allele
123 that drifted up in frequency starting 15,000 years ago, but that is now rare or absent from
124 present day populations.

125 Although low frequency allele calls often contain a higher proportion of base calling er-
126 rors than higher frequency allele calls do, it is not plausible that base-calling errors could
127 be responsible for the pulse we have described. In the UK10K data, a minor allele present
128 at 0.6% frequency corresponds to a derived allele that is present in 23 out of 3854 sampled
129 haplotypes and supported by 80 short reads on average (assuming 7x coverage per individ-
130 ual). When independently generated datasets of different sizes are projected down to the
131 same sample size, the TCC→TTC pulse spans the same range of allele frequencies in both
132 datasets (Figure 3–Figure Supplements 1 and 2), which would not be the case if the shape
133 of the curve were a function of low frequency errors.

134 **Fine-scale mutation spectrum variation in other populations.** Encouraged by these
135 results, we sought to find other signatures of recent mutation pulses. We generated heatmaps
136 and PCA plots of mutation spectrum variation within each continental group, looking for
137 fine-scale differences between closely related populations (Figure 4 and Figure 4–Figure Sup-
138 plements 1 through 6). In some cases mutational spectra differ even between very closely
139 related populations. For example, the *AC→*CC mutations with elevated rates in East
140 Asia appear to be distributed heterogeneously within that group, with most of the load
141 carried by a subset of the Japanese individuals. These individuals also have elevated rates
142 of ACA→AAA and TAT→TTT mutations (Figure 4A and Figure Supplement 4). This

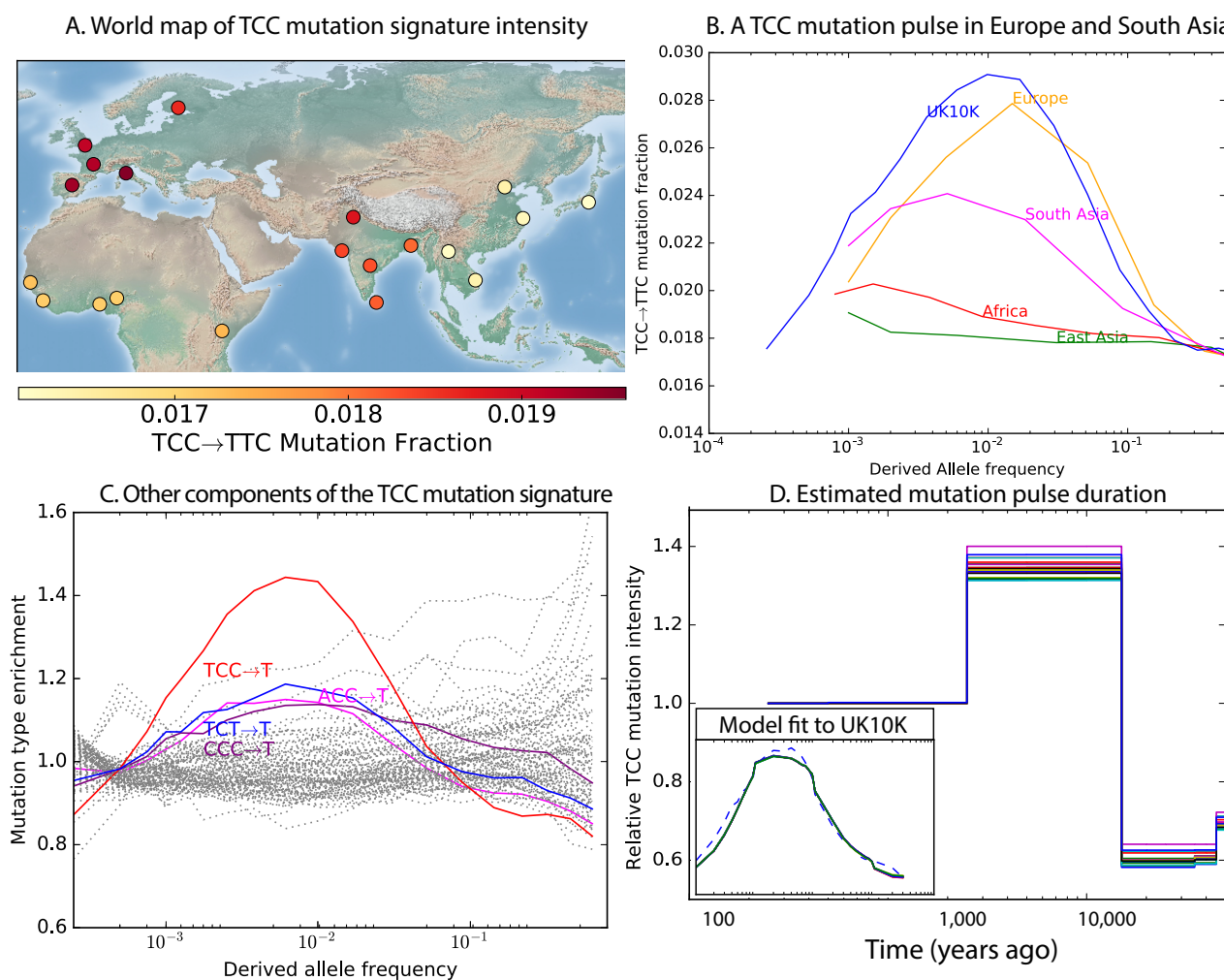


Figure 3: Geographic distribution and age of the TCC mutation pulse. (A) Observed frequencies of TCC→TTC variants in 1000 Genomes populations. (B) Fraction of TCC→TTC variants as a function of allele frequency in different samples indicates that these peak around 1%. See Figure Supplement 1 for distributions of TCC→TTC allele frequency within all 1000 Genomes populations, and see Figure Supplement 2 for the replication of this result in the Exome Aggregation Consortium Data. In the UK10K data, which has the largest sample size, the peak occurs at 0.6% allele frequency. (C) Other enriched C→T mutations with similar context also peak at 0.6% frequency in UK10K. See Figure Supplements 3, 4 and 5 for labeled allele frequency distributions of all 96 mutation types (most represented here as unlabeled grey lines). See Figure Supplement 6 for heatmap comparisons of the 1000 Genomes populations partitioned by allele frequency, which provide a different view of these patterns. (D) A population genetic model supports a pulse of TCC→TTC mutations from 15,000–2,000 years ago. Inset shows the observed and predicted frequency distributions of this mutation under the inferred model.

143 signature appears to be present in only a handful of Chinese individuals, and no Kinh or
144 Dai individuals. As seen for the European TCC mutation, the enrichment of these muta-
145 tion types peaks at low frequencies, i.e., $\sim 1\%$. Given the availability of only 200 Japanese
146 individuals in 1000 Genomes, it is hard to say whether the true peak is at a frequency much
147 lower than 1%.

148 PCA reveals relatively little fine-scale structure within the mutational spectra of Euro-
149 peans or South Asians (Figure 4–Figure Supplements 5 and 6). However, Africans exhibit
150 some substructure (Figure 4–Figure Supplement 3), with the Luhya exhibiting the most
151 distinctive mutational spectrum. Unexpectedly, a closer examination of PC loadings reveals
152 that the Luhya outliers are enriched for the same mutational signature identified in the
153 Japanese. Even in Europeans and South Asians, the first PC is heavily weighted toward
154 $*AC \rightarrow *CC$, $ACA \rightarrow AAA$, and $TAT \rightarrow TTT$, although this signature explains less of the mu-
155 tation spectrum variance within these more homogeneous populations. The sharing of this
156 signature may suggest either parallel increases of a shared mutation modifier, or a shared
157 aspect of environment or life history that affects the mutation spectrum.

158 **Mutation spectrum variation among the great apes.** Finally, given our finding of
159 extensive fine-scale variation in mutational spectra between human populations, we hypoth-
160 esized that mutational variation between species is likely to be even greater. To compare
161 the mutation spectra of the great apes in more detail, we obtained SNV data from the Great
162 Ape Diversity Panel, which includes 78 whole genome sequences from six great ape species
163 including human [19]. Overall, we find dramatic variation in mutational spectra among the
164 great ape species (Figure 5 and Figure 5–Figure Supplement 1).

165 As noted previously [20], one major trend is a higher proportion of CpG mutations among
166 the species closest to human, possibly reflecting lengthening generation time along the human
167 lineage, consistent with previous indications that species closely related to humans have lower
168 mutation rates than more distant species [21, 22, 23]. However, most other differences are
169 not obviously related to known processes such as biased gene conversion and generation
170 time change. The $A \rightarrow T$ mutation rate appears to have sped up in the common ancestor
171 of humans, chimpanzees, and bonobos, a change that appears consistent with a mutator
172 variant that was fixed instead of lost. It is unclear whether this ancient $A \rightarrow T$ speedup is
173 related to the $A \rightarrow T$ speedup in East Asians. Other mutational signatures appear on only a
174 single branch of the great ape tree, such as a slowdown of $A \rightarrow C$ mutations in gorillas.

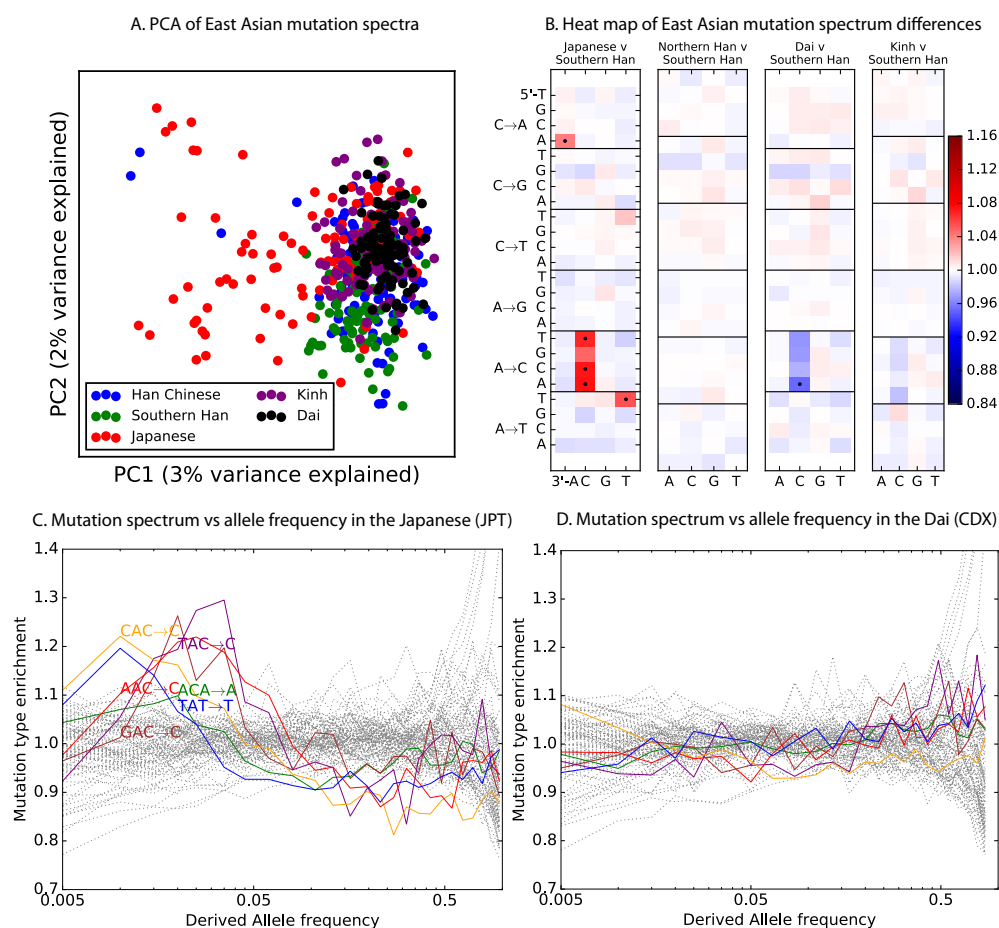


Figure 4: Mutational variation among east Asian populations. (A) *PCA of east Asian samples from 1000 Genomes, based on the relative proportions of each of the 96 mutational types. See Figure Supplements 2 through 6 for other finescale population PCAs.* (B) *Heatmaps showing, for pairs of east Asian samples, the ratio of the proportions of SNVs in each of the 96 mutational types. Points indicate significant contrasts at $p < 10^{-5}$. See Figure Supplement 1 for additional finescale heatmaps.* (C) and (D) *Relative enrichment of each mutational type in Japanese and Dai, respectively as a function of allele frequency. Six mutation types that are enriched in JPT are indicated. Populations: CDX=Dai, CHB=Han (Beijing); CHS=Han (south China); KHV=Kinh; JPT=Japanese.*

175 Discussion

176 The widespread differences captured in Figures 1 and 2 may be footprints of allele frequency
177 shifts affecting different mutator alleles. But in principle, other genetic and non-genetic
178 processes may also impact the observed mutational spectrum. First, biased gene conver-
179 sion (BGC) tends to favor C/G alleles over A/T, and BGC is potentially more efficient in
180 populations of large effective size compared to populations of smaller effective size [24]. How-
181 ever, despite the bottlenecks that are known to have affected Eurasian diversity, there is no
182 clear trend of an increased fraction of C/G→A/T relative to A/T→C/G in non-Africans vs.
183 Africans, or with distance from Africa (Figure 1–Figure Supplement 7), and previous studies
184 have also found little evidence for a strong genome-wide effect of BGC on the mutational
185 spectrum in humans and great apes [26, 20]. For these reasons, we think that evolution of
186 the mutational process is a better explanation than BGC or selection for differences that
187 have been observed between the spectra of ultra-rare singleton variants and older human
188 genetic variation [25];

189 It is also known that shifts in generation time or other life-history traits may affect
190 mutational spectra, particularly for CpG transitions [27, 28]. Most CpG transitions result
191 from spontaneous methyl-cytosine deamination as opposed to errors in DNA replication.
192 Hence the rate of CpG transitions is less affected by generation time than other mutations
193 [9, 29, 30]. We observe that Europeans have a lower fraction of CpG variants compared to
194 Africans, East Asians and South Asians (Fig. 1B), consistent with a recent report of a lower
195 rate of *de novo* CCG→CTG mutations in European individuals compared to Pakistanis [31].
196 Such a pattern may be consistent with a shorter average generation time in Europeans [29],
197 though it is unclear that a plausible shift in generation-time could produce such a large
198 effect. Apart from this, the other patterns evident in Figure 1 do not seem explainable by
199 known processes.

200 In summary, we report here that, mutational spectra differ significantly among closely
201 related human populations, and that they differ greatly among the great ape species. Our
202 work shows that subtle, concerted shifts in the frequencies of many different mutation types
203 are more widespread than dramatic jumps in the rate of single mutation types, although the
204 existence of the European TCC→TTC pulse shows that both modes of evolution do occur
205 [12, 29, 13].

206 At this time, we cannot exclude a role for nongenetic factors such as changes in life history

207 or mutagen exposure in driving these signals. However given the sheer diversity of the effects
208 reported here, it seems parsimonious to us to propose that most of this variation is driven
209 by the appearance and drift of genetic modifiers of mutation rate. This situation is perhaps
210 reminiscent of the earlier observation that genome-wide recombination patterns are variable
211 among individuals [32], and ultimate discovery of PRDM9 [33]; although in this case it is
212 unlikely that a single gene is responsible for all signals seen here. As large datasets of *de*
213 *novo* mutations become available, it should be possible to map mutator loci genome-wide.
214 In summary, our results suggest the likelihood that mutational modifiers are an important
215 part of the landscape of human genetic variation.

216 **Acknowledgements**

217 This work was funded by NIH grants GM116381 and HG008140, and by the Howard Hughes
218 Medical Institute. We thank Jeffrey Spence and Yun S. Song for technical assistance. We also
219 thank Ziyue Gao, Arbel Harpak, Molly Przeworski, Joshua Schraiber, and Aylwyn Scally
220 for comments and discussion, as well as two anonymous reviewers.

221 **Methods**

222 **Data Availability.**

223 All datasets analyzed here are publicly available at the following websites:

224	1000 Genomes Phase 3 UK10K Simons Genome Diversity Panel	http://www.1000genomes.org/category/phase-3/ http://www.uk10k.org/data-access.html https://www.simonsfoundation.org/life-sciences/simons-genome-diversity-project-dataset/
-----	--	---

225 **Human Mutation Spectrum Processing.**

226 Mutation spectra were computed using 1000 Genomes Phase 3 SNPs [14] that are biallelic,
227 pass all 1000 Genomes quality filters, and are not adjacent to any N's in the hg19 reference se-
228 quence. Ancestral states were assigned using the UCSC Genome Browser alignment of hg19
229 to the PanTro2 chimpanzee reference genome; SNPs were discarded if neither the reference
230 nor alternate allele matched the chimpanzee reference. To minimize the potential impact of
231 ancestral misidentification errors, SNPs with derived allele frequency higher than 0.98 were
232 discarded. We also filtered out regions annotated as “conserved” based on the 100-way Phy-
233 loP conservation score [15], download from <http://hgdownload.cse.ucsc.edu/goldenPath/hg19/phastCons100way/>, as well as regions annotated as repeats by RepeatMasker [34],
234 downloaded from
235 <http://hgdownload.cse.ucsc.edu/goldenpath/hg19/database/nestedRepeats.txt.gz> . To be
236 counted as part of the mutation spectrum of population P (which can be either a continen-
237 tal group or a finer-scale population from one city), a SNP should not be a singleton within
238 population P —at least two copies of the ancestral and derived alleles must be present within
239 that population.
240

241 An identical approach was used to extract the mutation spectrum of the UK10K ALSPAC
242 panel [18], which is not subdivided into smaller populations. The data were filtered as
243 described in [35]. The filtering procedure performed by Field, et al. reduces the ALSPAC
244 sample size to 1927 individuals.

245 We also computed mutation spectra of the Simons Genome Diversity Panel (SGDP)
246 populations [17]. Four of the SGDP populations, West Eurasia, East Asia, South Asia,
247 and Africa, were compared to their direct counterparts in the 1000 Genomes data. Three
248 additional SGDP populations, Central Asia and Siberia, Oceania, and America, had no close

249 1000 Genomes counterparts and were not analyzed here (although each project contained a
250 panel of people from the Americans, the composition of the American panels was extremely
251 different, with the 1000 Genomes populations being much more admixed with Europeans
252 and Africans). SGDP sites with more than 20% missing data were not utilized. All other
253 data processing was done the same way described for the 1000 Genomes data.

254 The following table gives the same size of each population panel, as well as the total
255 number of SNPs segregating in the panel that are used to compute mutation type ratios:

Dataset	Population	Number of individuals	Number of SNPs
1kg	Africa	504	16,870,400
1kg	Europe	503	8,508,040
1kg	East Asia	504	7,895,925
1kg	South Asia	489	9,552,781
SGDP	Africa	45	6,569,658
SGDP	West Eurasia	69	4,201,571
SGDP	East Asia	49	3,312,645
SGDP	South Asia	38	3,449,624

257 **Great Ape Diversity Panel Data Processing.**

258 Biallelic great ape SNPs were extracted from the Great Ape Diversity Panel VCF [19],
259 which is aligned to the hg18 human reference sequence. Ancestral states were assigned using
260 the Great Ape Genetic Diversity project annotation, which used the Felsenstein pruning
261 algorithm to assign allelic states to internal nodes in the great ape tree. In the Great Ape
262 Diversity Panel, the most recent common ancestor (MRCA) of the human species is labeled
263 as node 18; the MRCAs of chimpanzees, bonobos, gorillas, and orangutans, respectively, are
264 labeled as node 16, node 17, node 19, and node 15. We extracted the state of each MRCA
265 at each SNP in the alignment and used it to polarize the ancestral and derived allele at
266 that site; a SNP was discarded whenever the ancestral node was assigned an uncertain or
267 polymorphic ancestral state. As with the human data, SNPs with derived allele frequency
268 higher than 0.98 were not used, and both repeats and PhyloP-annotated conserved regions
269 were filtered away.

270 **Visual representation of mutation spectra.**

271 The mutation type of a SNP is defined in terms of its ancestral allele, its derived allele, and
272 its two immediate 5' and 3' neighbors. Two mutation types are considered equivalent if they

273 are strand-complementary to each other (e.g. ACG→ATG is equivalent to CGT→CAT).
274 This scheme classifies SNPs into 96 different mutation types, each that can be represented
275 with an A or C ancestral allele.

276 To compute the frequency $f_P(m)$ of SNP m in population P , we count up all SNPs of type
277 m where the derived allele is present in at least one representative of population P (which can
278 be either a specific population such as YRI or a broader continental group such as AFR).
279 After obtaining this count $C_P(m)$, we define $f_P(m)$ to be the ratio $C_P(m)/\sum_{m'} C_P(m')$,
280 where the sum in the denominator ranges over all 96 mutation types m' . The enrichment of
281 mutation type m in population P_1 relative to population P_2 is defined to be $f_{P_1}(m)/f_{P_2}(m)$;
282 these enrichments are visualized as heat maps in Figures 1B, 3B, and 4A.

283 To track changes in the mutational spectrum over time, we compute $f_P(m)$ in bins of
284 restricted allele frequency. This involves counting the number of SNPs of type m that are
285 present at frequency ϕ in population P to obtain counts $C_P(m, \phi)$ and frequencies $f_P(m, \phi) =$
286 $C_P(m, \phi)/\sum_{m'} C_P(m', \phi)$. Deviation of the ratio $f_P(m, \phi)/f_P(m)$ from 1 indicates that the
287 rate of m has fluctuated recently in the history of population P . To make the sampling
288 noise approximately uniform across alleles of different frequencies, alleles of derived count
289 greater than 5 were grouped into approximately log-spaced bins that each contained similar
290 numbers of UK10K SNPs. More precisely, we defined a set of bin endpoints b_1, b_2, \dots such
291 that the total number of SNPs ranging in derived allele count between b_i and $b_{i+1} - 1$ is
292 greater than or equal to the number of 5-ton SNPs, while the total number of SNPs ranging
293 in derived allele count from b_i to $b_{i+1} - 2$ is less than the number of 5-ton SNPs.

294 In some cases, e.g. Figures 2, Figure 2–Figure Supplement 1B, and Figure 3–Figure
295 Supplement 1, site frequency spectra were projected down to a smaller sample size before
296 counting SNPs in order to more accurately compare datasets of different sample sizes. A
297 binomial sampling approach was used to project a sample of N haplotypes does to a smaller
298 sample size n . Letting $C_P^{(N)}(m, \phi)$ denote the SNP counts in the large sample of N haplotypes,
299 effective SNP counts $C_P^{(n)}(m, \phi)$ in a sample of n haplotypes are computed as follows:

$$C_P^{(n)}(m, k/n) = \binom{n}{k} \sum_{\ell=1}^{N-1} (\ell/N)^k (1 - \ell/N)^{n-k} C_P^{(N)}(m, \ell/N)$$

300 Significance testing.

301 One central goal of this paper is to test whether many mutation types differ in rate between
 302 human populations or whether mutation spectrum shifts have been rare events affecting only
 303 a small proportion of mutation types. A simple statistical method for answering this question
 304 would be to perform 96 separate chi-square tests, one for each triplet-context-dependent
 305 mutation type, as follows:

306 Let S_i denote the total number of SNPs segregating in population P_i , and let $S_i^{(m)}$
 307 denote the number of SNPs of mutation type m . If mutation type m is more prevalent in
 308 population P_1 than in population P_2 , a chi-square test provides a natural way of assessing
 309 the significance of this difference. As described in [12], this test is performed on the following
 310 2-by-2 contingency table:

311	$S_1^{(m)}$	$P_1 - S_1^{(m)}$
	$S_2^{(m)}$	$P_2 - S_2^{(m)}$

312 It would be appealing to conclude that every mutation type “passing” this chi-square test
 313 is a mutation type that has changed in rate during recent human history. However, if we
 314 were to perform the full set of 96 tests, they would not be independent. A sufficiently large
 315 increase in the rate of one mutation type m_1 in population P_1 after divergence from P_2 could
 316 cause another mutation type m_2 , whose rate has remained constant, to comprise significantly
 317 different fractions of the SNPs from P_1 and P_2 . To minimize this effect, we formulate the
 318 following iterative procedure of conditionally independent tests: first, compute a chi-square
 319 significance value $p_{\text{unordered}}(m)$ for each mutation type m using the 2-by-2 chi-square table
 320 above. We then use these values to order the SNPs from lowest p value to highest and
 321 compute a set of ordered p values $p_{\text{ordered}}(m)$. For the mutation type m_0 with the lowest
 322 unordered p value, $p_{\text{unordered}}(m_0) = p_{\text{ordered}}(m_0)$. For mutation type m_i , which has the i th
 323 lowest unordered p value and $i < 96$, $p_{\text{ordered}}(m_i)$ is computed from the following contingency
 324 table:

325	$S_1^{(m_i)}$	$\sum_{j=i+1}^{96} S_1^{(m_j)}$
	$S_2^{(m_i)}$	$\sum_{j=i+1}^{96} S_2^{(m_j)}$

326 For mutation type m_{96} , which has the highest unordered p value, the ordered p value is
 327 computed from the contingency table

328	$S_1^{(m_{96})}$	$S_1^{(m_{95})}$
	$S_2^{(m_{96})}$	$S_2^{(m_{95})}$

329 This procedure is guaranteed to find fewer mutation types to differ significantly in rate
330 between populations compared to separate chi-square tests.

331 **Principal component analysis (PCA).**

332 The python package `matplotlib.mlab.PCA` was used to perform PCA on the complete set of
333 1000 Genomes haplotypes, each haplotype h represented by a 96-element vector encoding the
334 mutation frequencies $(f_h(m))_m$ of the non-singleton derived alleles present on that haplotype.
335 In the same way, a separate PCA was performed on each of the 5 continental groups to reveal
336 finescale components of mutation spectrum variation.

337 **Dating of the TCC→T mutation pulse.**

338 We estimated the duration and intensity of TCC→T rate acceleration in Europe by fitting a
339 simple piecewise-constant rate model to the UK10K frequency data. To specify the param-
340 eters of the model, we divide time into discrete log-spaced intervals bounded by time points
341 t_1, \dots, t_d , assigning each interval a TCC→T mutation rate r_0, \dots, r_d . In units of generations
342 before the present, the time discretization points were chosen to be: 20, 40, 200, 400, 800,
343 1200, 1600, 2000, 2400, 2800, 3200, 3600, 4000, 8000, 12000, 16000, 20000, 24000, 28000,
344 32000, 36000, 40000. We assume that the total rate r of mutations other than TCC→T
345 stays constant over time (a first-order approximation).

346 In terms of these rate variables, we can calculate the expected shape of the TCC→T
347 pulse shown in Figure 2B of the main text. The shape of this curve depends on both the
348 mutation rate parameters r_i and the demographic history of the European population, which
349 determines the joint distribution of allele frequency and allele age. To account for the effects
350 of demography, we use Hudson's `ms` program to simulate 10,000 random coalescent trees
351 under a realistic European demographic history inferred from allele frequency data [36] and
352 condition our inference upon this collection of trees as follows:

Let $A(m, t)$ be the function for which $\int_{t_i}^{t_{i+1}} A(m, t) dt$ equals the coalescent tree branch length, averaged over the sample of simulated trees, that is ancestral to exactly m lineages and falls between time t_i and t_{i+1} . Given this function, which can be empirically estimated from a sample of simulated trees, the expected frequency spectrum entry k/n is

$$E(k/n) = \frac{\sum_{i=1}^d \int_{t_{i-1}}^{t_i} A(k, t) dt}{\sum_{j=1}^n \sum_{i=1}^d \int_{t_{i-1}}^{t_i} A(j, t) dt}$$

and the expected fraction of TCC→T mutations in allele frequency bin k/n is

$$E(f_{\text{TCC} \rightarrow \text{T}}(k/n)) = \frac{\sum_{i=1}^d r_i \int_{t_{i-1}}^{t_i} A(k, t) dt}{r \sum_{i=1}^d \int_{t_{i-1}}^{t_i} A(k, t) dt}.$$

The expected value of the TCC→T enrichment ratio being plotted in Figure 2B is

$$E(r_{\text{TCC} \rightarrow \text{T}}(k/n)) = \frac{\sum_{i=1}^d r_i \int_{t_{i-1}}^{t_i} A(k, t) dt \cdot \sum_{j=1}^n \sum_{i=1}^d \int_{t_{i-1}}^{t_i} A(j, t) dt}{\sum_{i=1}^d \int_{t_{i-1}}^{t_i} A(k, t) dt \cdot \sum_{j=1}^n \sum_{i=1}^d r_i \int_{t_{i-1}}^{t_i} A(j, t) dt}$$

In Figure 2B, enrichment ratios are not computed for every allele frequency in isolation, but for allele frequency bins that each contain similar numbers of SNPs. Given integers $1 \leq k_m < k_{m+1} \leq n$, the expected TCC→T enrichment ratio averaged over all SNPs with allele frequency between k_m/n and k_{m+1}/n is:

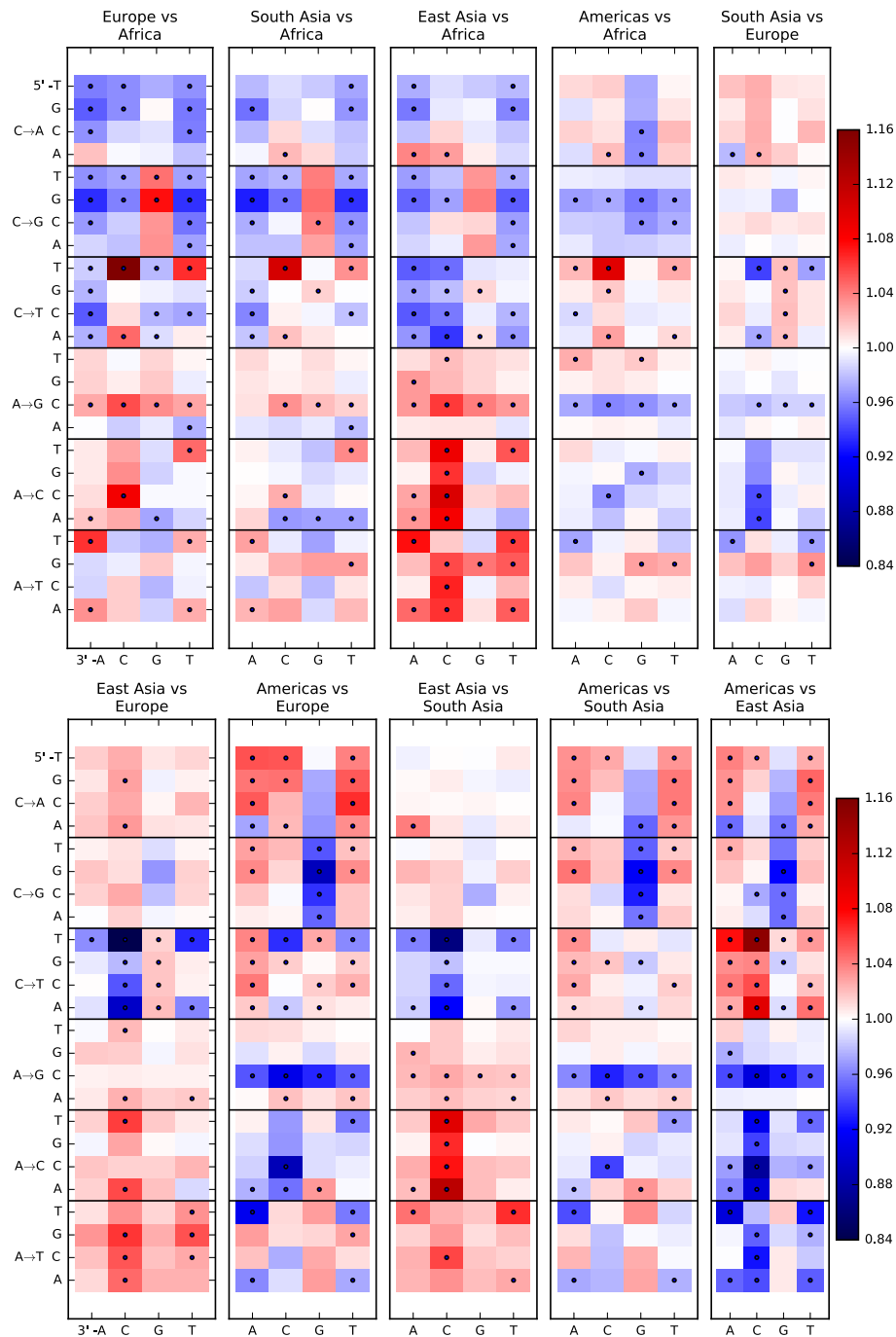
$$E(r_{\text{TCC} \rightarrow \text{T}}(k_m/n)) = \frac{\sum_{i=1}^d r_i \int_{t_{i-1}}^{t_i} \sum_{k=k_m}^{k_{m+1}} A(k, t) dt \cdot \sum_{j=1}^n \sum_{i=1}^d \int_{t_{i-1}}^{t_i} A(j, t) dt}{\sum_{i=1}^d \int_{t_{i-1}}^{t_i} \sum_{k=k_m}^{k_{m+1}} A(k, t) dt \cdot \sum_{j=1}^n \sum_{i=1}^d r_i \int_{t_{i-1}}^{t_i} A(j, t) dt}$$

We optimize the mutation rates r_1, \dots, r_d using a log-spaced quantization of allele frequencies $k_1/n, \dots, k_m/n$ defined such that all bins contain similar numbers of SNPs. The chosen allele count endpoints k_1, \dots, k_m are: 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 20, 30, 40, 50, 60, 70, 80, 90, 100, 200, 300, 400, 500, 600, 700, 800, 900, 1000, 2000, 3000, 4000. Given this quantization of allele frequencies, we optimize r_1, \dots, r_d by using the BFGS algorithm to minimize the least squares distance $D(r_0, \dots, r_d)$ between $E(r_{\text{TCC} \rightarrow \text{T}}(k_m/n))$ and the empirical ratio $r_{\text{TCC} \rightarrow \text{T}}(k_m/n)$ computed from the UK10K data. This optimization is subject to a regularization penalty that minimizes the jumps between adjacent mutation rates r_i and r_{i+1} :

$$D(r_0, \dots, r_d) = \sum_{m=1}^d (E(r_{\text{TCC} \rightarrow \text{T}}(k_m/n)) - r_{\text{TCC} \rightarrow \text{T}}(k_m/n))^2 + 0.25 \sqrt{\sum_{i=1}^d (r_{i-1} - r_i)^2}$$

353 Although the underlying model of mutation rate change assumed here is very simple, it
 354 still represents an advance over the method used in [12] to estimate of the timing of the
 355 TCC→TTC mutation rate increase. That method relied upon explicit estimates of allele
 356 age from a dataset of less than 100 individuals, which are much noisier than integration of
 357 a joint distribution of allele age and frequency across a sample of thousands of haplotypes.

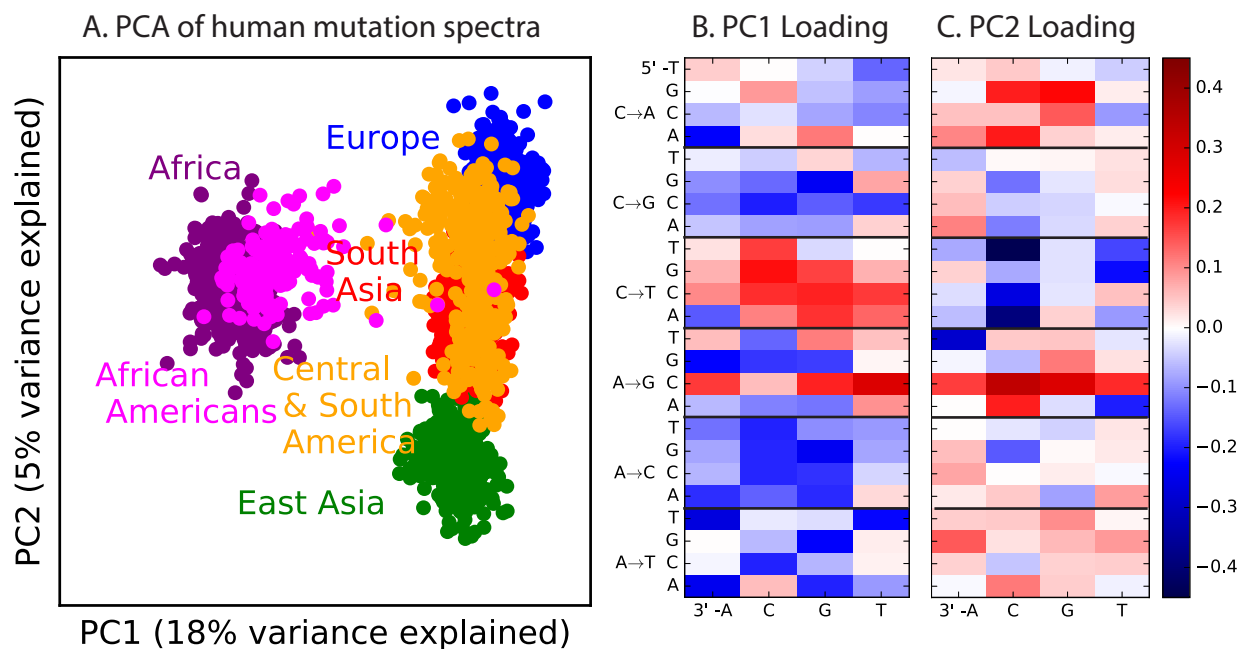
358 **Supplementary Figures**



359

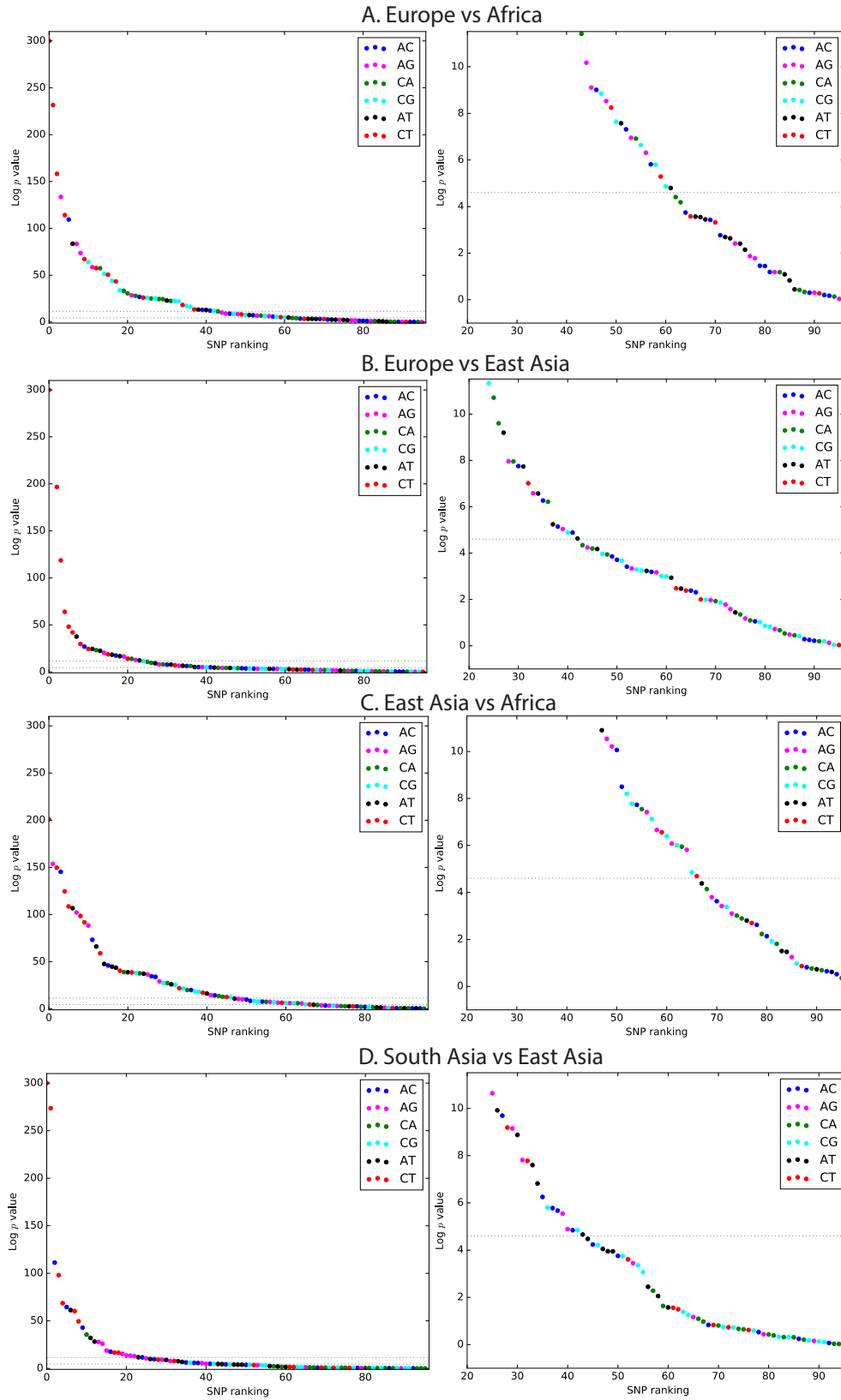
360 **Figure 1– Figure Supplement 1: Pairwise mutation spectrum comparisons**

361 **among continental groups.** Each of these plots compares the mutation spectra of two
 362 populations P_1 and P_2 . Letting f_i denote the fraction of SNVs in population P_i that have a
 363 given triplet context, ancestral allele, and derived allele, the corresponding heat map
 364 square visualizes the enrichment ratio f_1/f_2 . Black dots mark mutation types for which the
 365 difference between populations has a χ^2 p -value less than 10^{-5} .

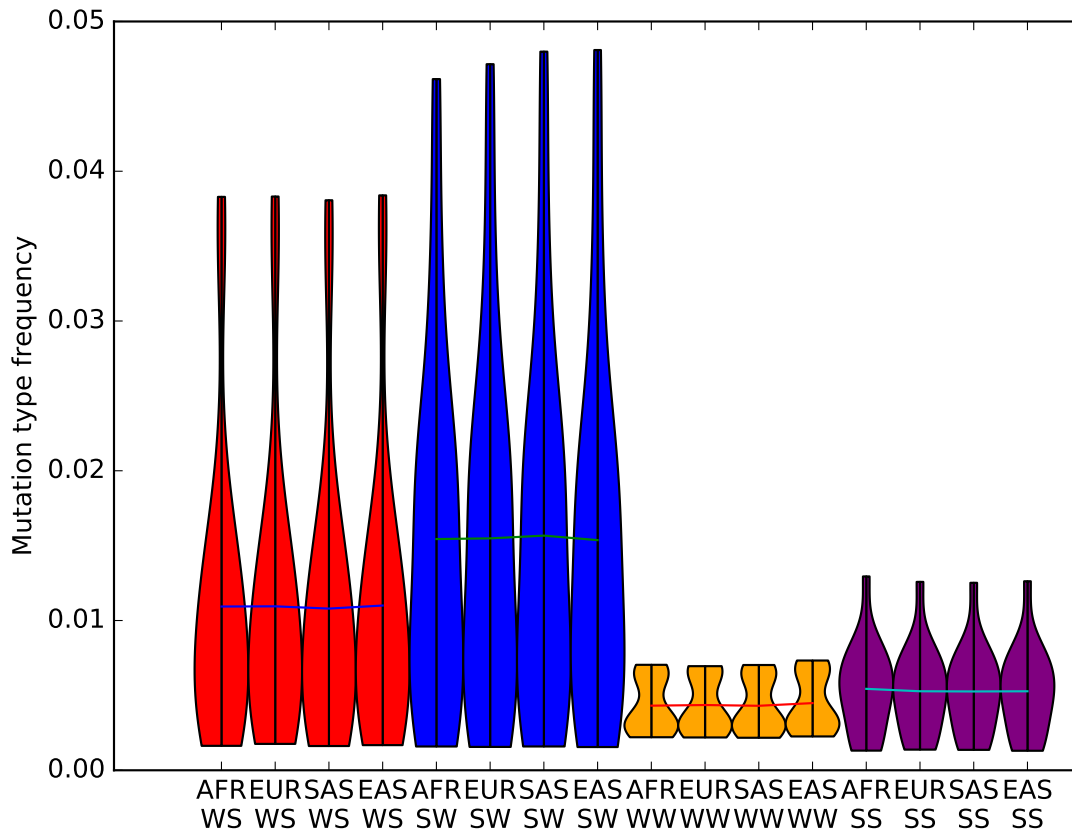


366

367 **Figure 1–Figure Supplement 2: PCA of all 1000 Genomes continental groups.**
 368 All admixed North and South American individuals were omitted from Figure 1 in the
 369 main text to clarify the separation of other populations along an African vs non-African
 370 axis and an East vs West Eurasian axis. Here, admixed Americans are added in black. As
 371 expected, some African-Americans group with the Africans, while other admixed
 372 Americans fall within the variation of other East and West Eurasians. The accompanying
 373 heat maps show the mutation type loadings of the first two principal components, the
 374 second of which is heavily weighted toward the European TCC→TTC signature.



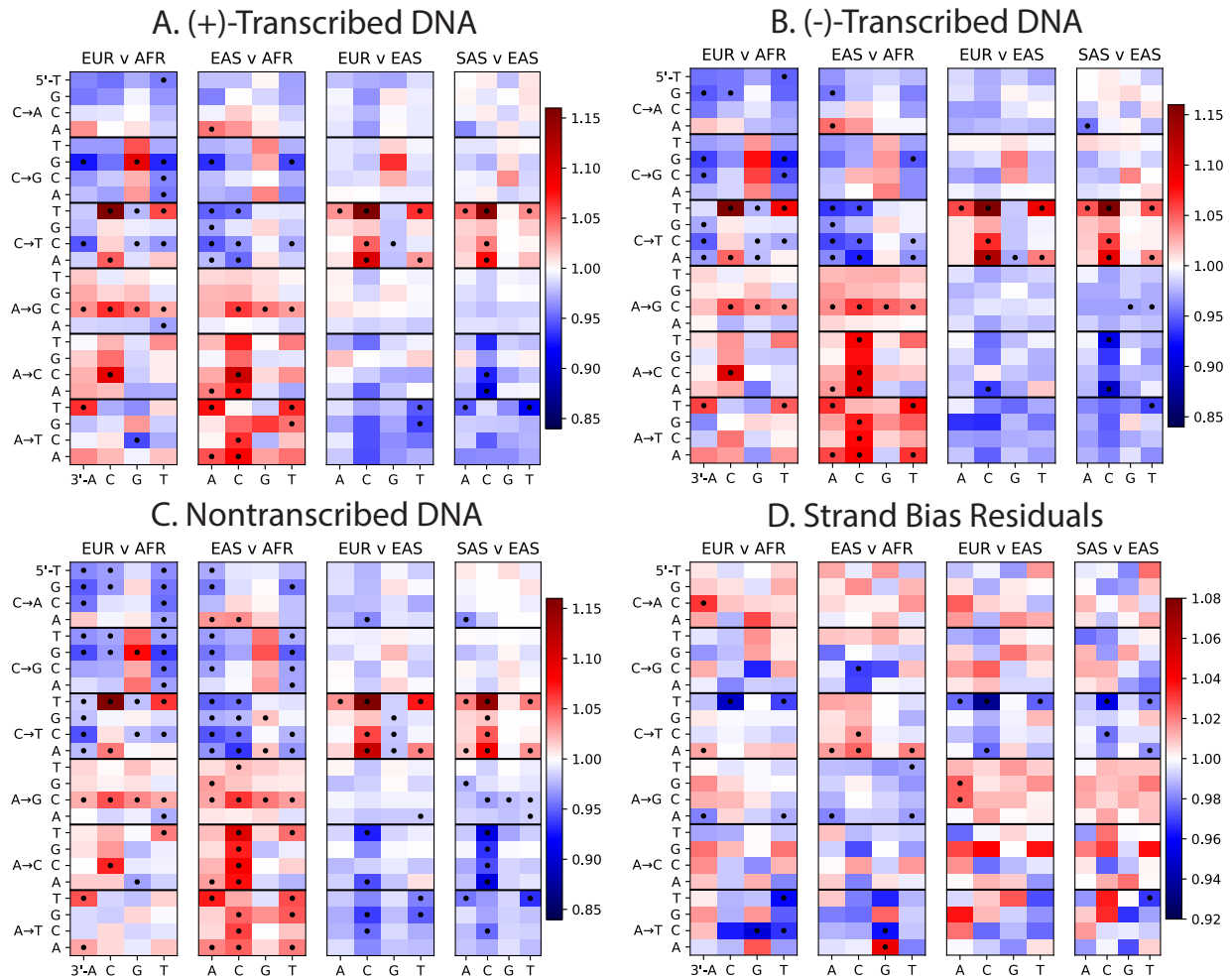
376 **Figure 1–Figure Supplement 3: Mutation spectrum comparison p -values.** Each
377 left-hand plot shows all chi-squared p -values corresponding to the ratios from Figure 1A. In
378 the absence of recent mutation spectrum evolution, only one out of 96 SNP categories is
379 expected to have a p -value below 0.01 (lower dotted line). In contrast, the majority of p
380 values meet the more stringent threshold $p < 1e - 5$. The corresponding right hand panel
381 shows a closeup of the distribution of p -values greater than $1e - 5$.



382

383 **Figure 1–Figure Supplement 4: The effects of biased gene conversion on**
384 **mutation spectra.** When using segregating variation to study the mutation spectrum,
385 one potential source of bias is that strong-to-weak mutations, where the ancestral allele is
386 G or C and the derived allele is A or T, have a lower fixation probability than
387 weak-to-strong mutations due to biased gene conversion (BGC). If this effect were
388 sufficiently strong, it would inflate the apparent mutation fractions of weak-to-strong
389 mutations, especially in populations with large effective sizes where natural selection is
390 particularly efficient. Within humans, Africans have the largest long-term effective
391 population size, while East Asians and Native Americans have the lowest. Therefore, if
392 BGC has created differences in mutation spectra between populations, the fraction of
393 weak-to-strong SNVs should be highest in Africans, intermediate in Europeans and South
394 Asians, and lowest in East Asians and Native Americans. This violin plot reveals no such
395 pattern, suggesting that BGC is not a strong driver of mutation spectrum differences

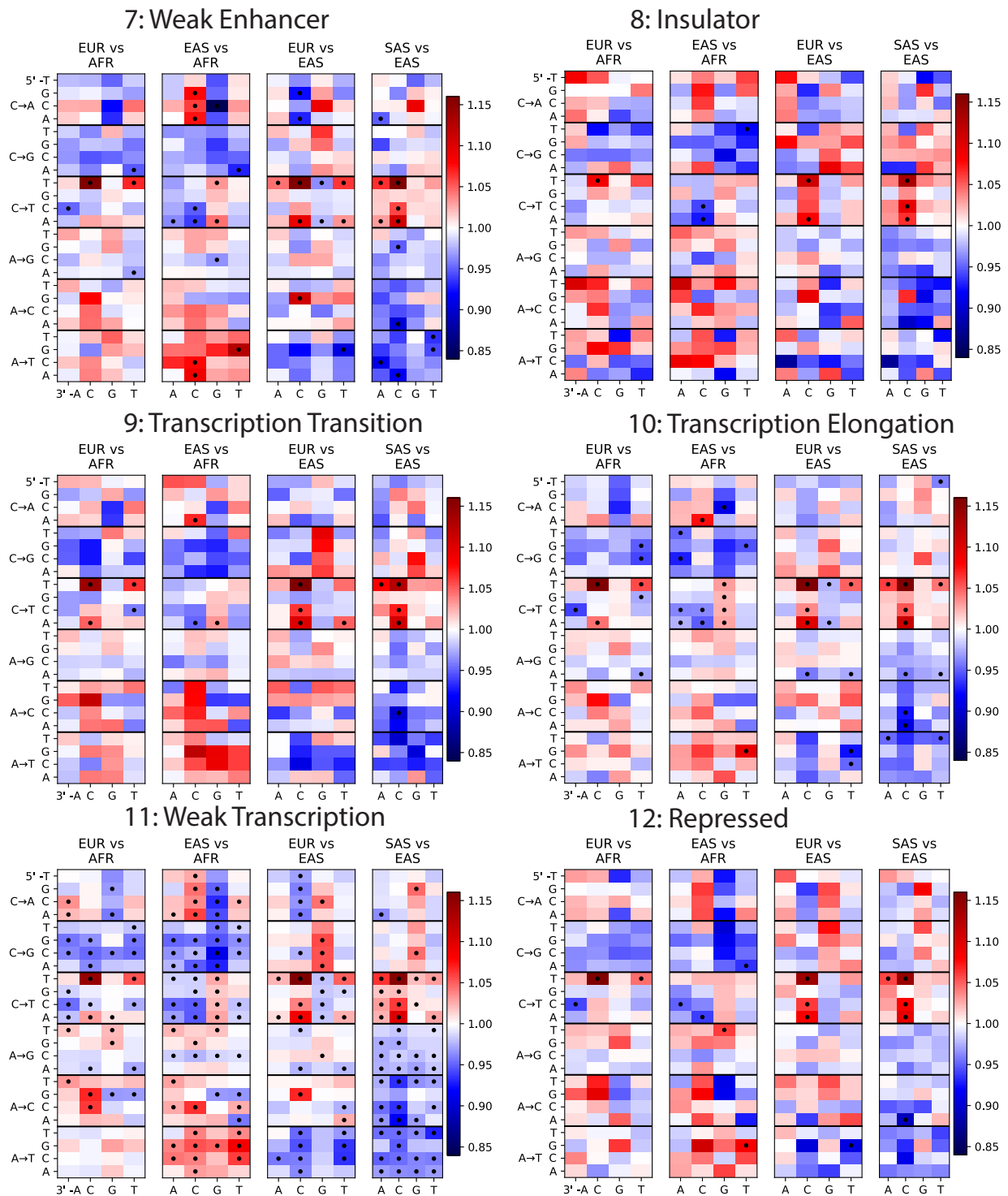
396 between human populations. We do not observe either a direct correlation between in
 397 strong-to-weak mutation fraction and distance from Africa or an inverse correlation
 398 between weak-to-strong mutation fraction and distance from Africa.



399

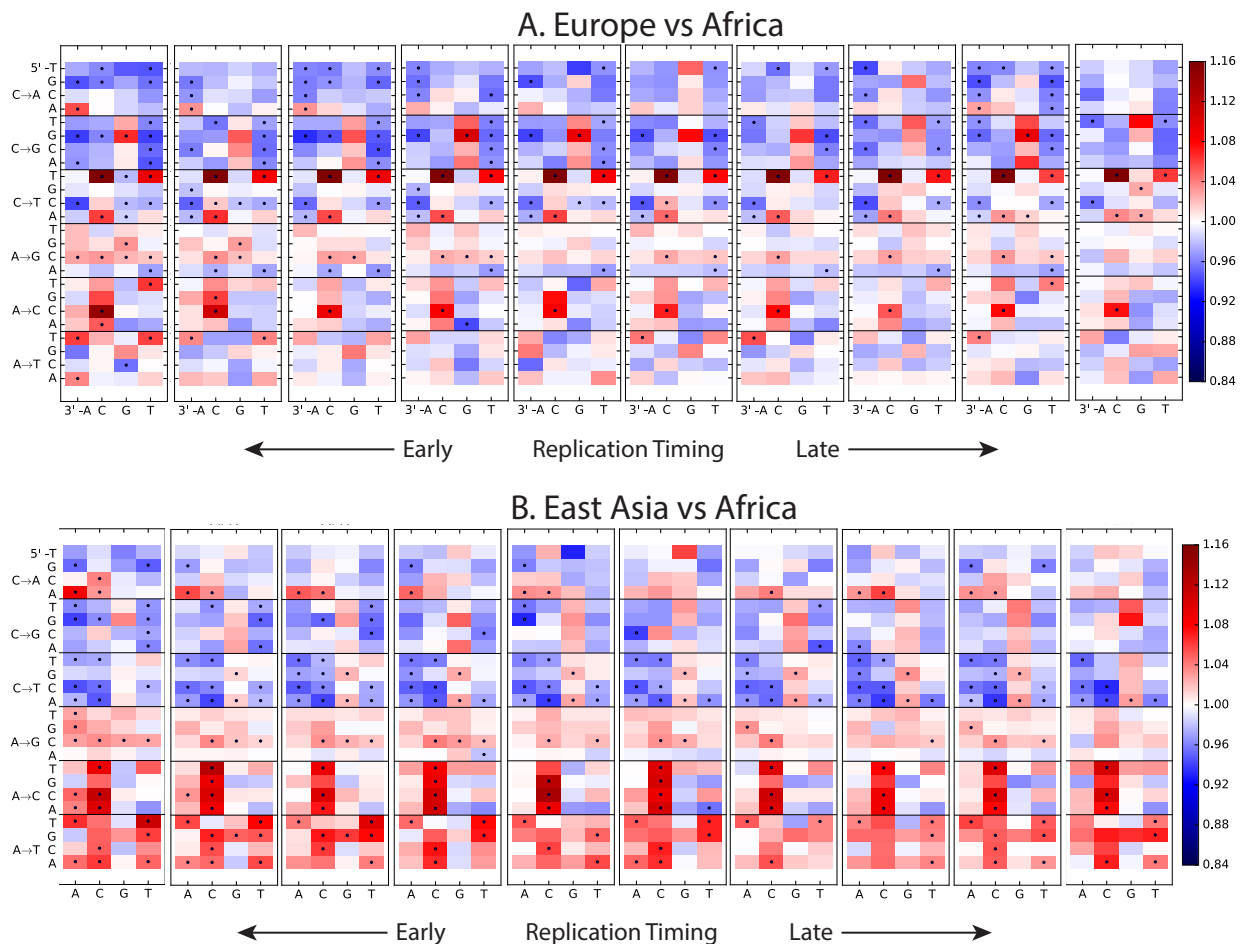
400 **Figure 1–Figure Supplement 5: Mutation Spectra of Transcribed vs**
 401 **Non-Transcribed DNA.** Using the UCSC Genome Browser annotations of the human
 402 reference hg19, we determined whether each SNP occurs in a transcribed or non
 403 transcribed region. We further divided SNPs occurring in transcribed regions according to
 404 whether the ancestral A or C allele occurs on the (+)-strand or the (-)-strand. Panels A,
 405 B, and C all show the same population-specific mutation type enrichments that are
 406 observed in Figure 1B. Panel D plots the residuals between panels A and B, highlighting
 407 mutation types that show a modest difference in strand bias between populations.

414 dots mark mutation types that show a significant enrichment in one population at the level
 415 $p < 0.01$. Every chromatin state shows enrichment of the TCC→TTC signature in Europe
 416 and South Asia. Some heat maps are noisy due to the small sample size of SNPs contained
 417 within these regions, but all showcase the same general patterns as Figure 1B.



418

419 **Figure 1–Figure Supplement 7: Mutation Spectra of ChromHMM chromatin**
 420 **states (Part II of II).**



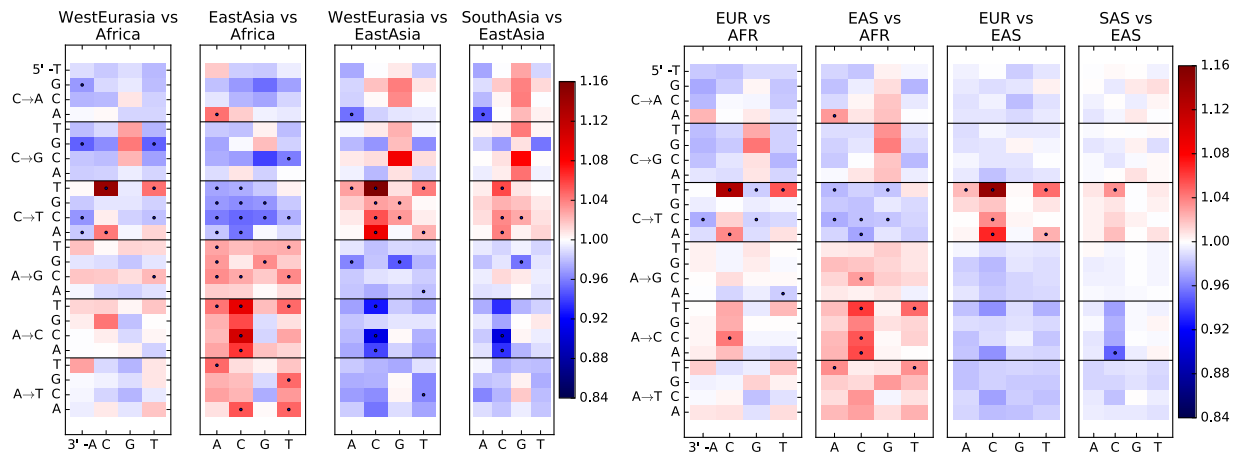
421

422 **Figure 1–Figure Supplement 8: Variation of the mutation spectrum with DNA**
 423 **replication timing.** We partitioned the genome into 10 equal replication timing quantiles
 424 using data obtained from [38], then computed mutation spectrum differences within each
 425 quantile. Although most patterns from Figure 1B replicate within each replication timing
 426 bin, there are a few exceptions. CpG transitions, which occur most often in
 427 early-replicating regions, vary in population bias depending on replication timing. In
 428 addition, the deficit of ACA→AAA and AAA→ATA mutations in Africa compared to
 429 Europe and Asia is observed mainly in early-replicating regions.

430 **Figure 1–Source Data 1.** This text file shows the number of SNPs in each of the 96
 431 mutational categories that passed all filters in each 1000 Genomes continental group.

A. Simons Genomic Diversity Panel (SGDP)

B. 1000 Genomes Phase 3 Panel (1000G)

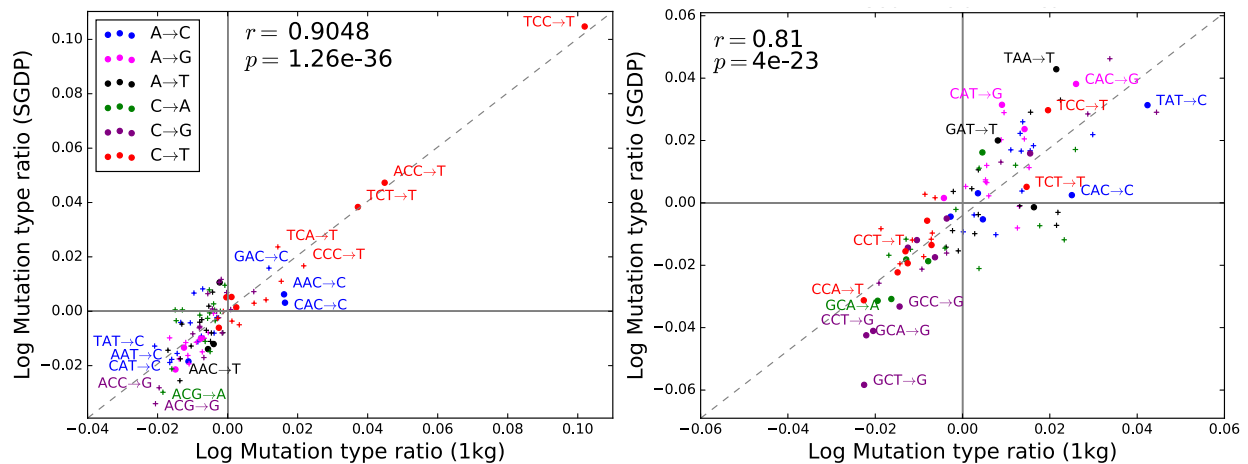


432

433 **Figure 2–Figure Supplement 1: Heatmap comparisons between continental**
 434 **groups in 1000 Genomes and the SGDP.** Here, each 1000 Genomes population is
 435 projected down to the sample size of the corresponding SGDP population in order to
 436 sample alleles with a similar distribution of ages and frequencies.

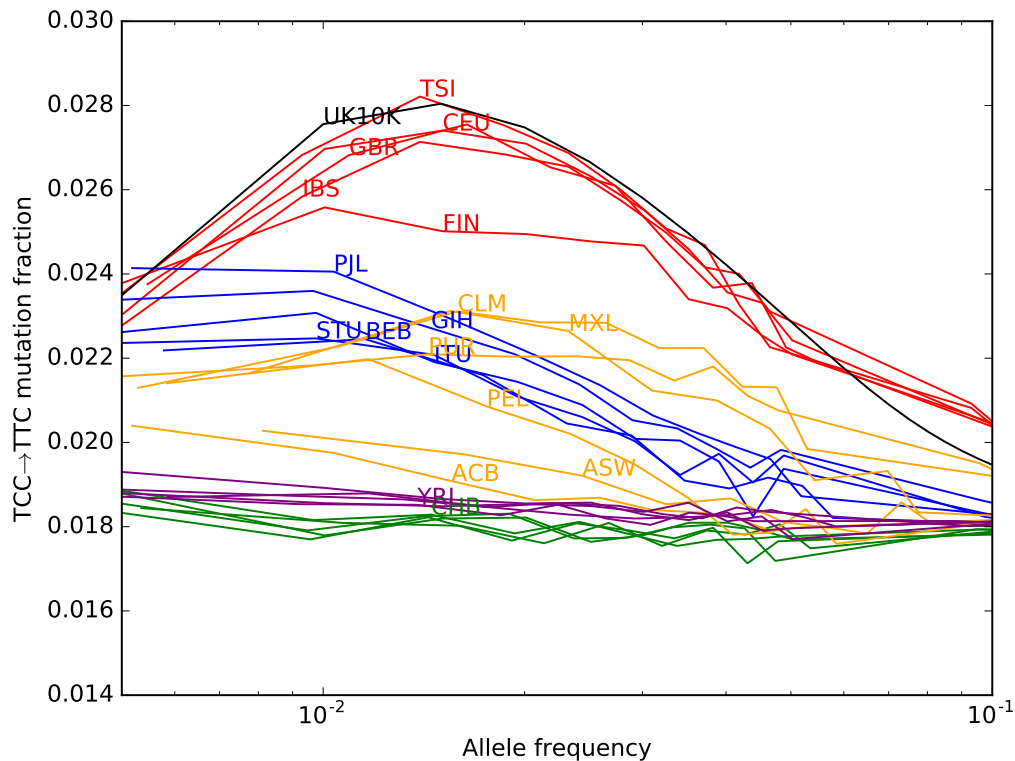
A. West Eurasia vs South Asia

B. South Asia vs Africa



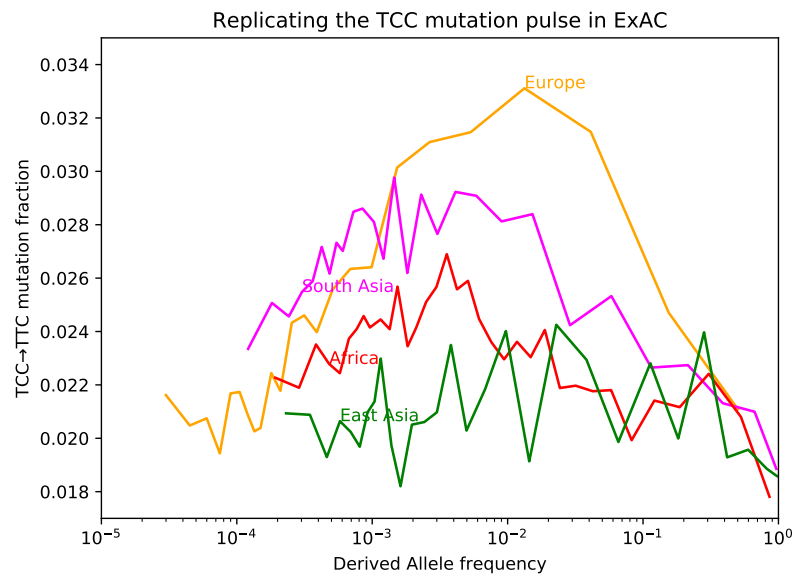
437

438 **Figure 2–Figure Supplement 2: Regression of the SGDP heatmap coefficients**
 439 **versus the corresponding 1000 Genomes heatmap coefficients.**



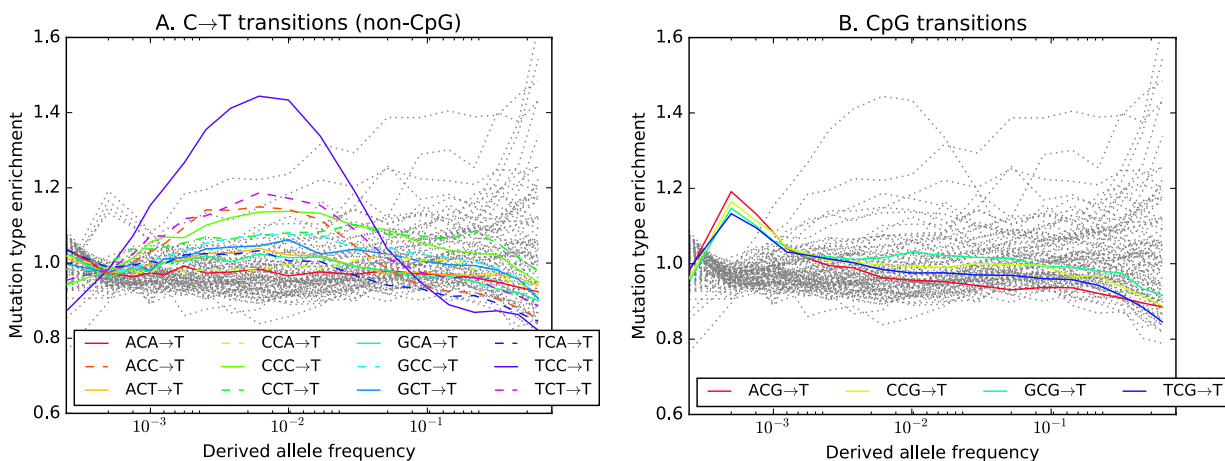
440

441 **Figure 3—Figure Supplement 1: TCC→TTC mutation fraction as a function of**
442 **allele frequency in all 1000 Genomes populations.** To enable better comparison with
443 the 1000 Genomes data, the UK10K SNPs have been downsampled to 200 individuals. The
444 age distribution of alleles of a given frequency varies as a function of the number of lineages
445 being sampled—this is why the UK10K pulse peaks around 0.6% frequency when measured
446 in a dataset of thousands of lineages, but peaks around 2% in a subsample of only 400
447 lineages. Some African and East Asian population names have been omitted for clarity
448 since the TCC→TTC mutation fraction is so uniform within these continental groups. Red
449 = European populations; Blue = South Asian; Orange = Americas; Purple = Africa;
450 Green = East Asia.



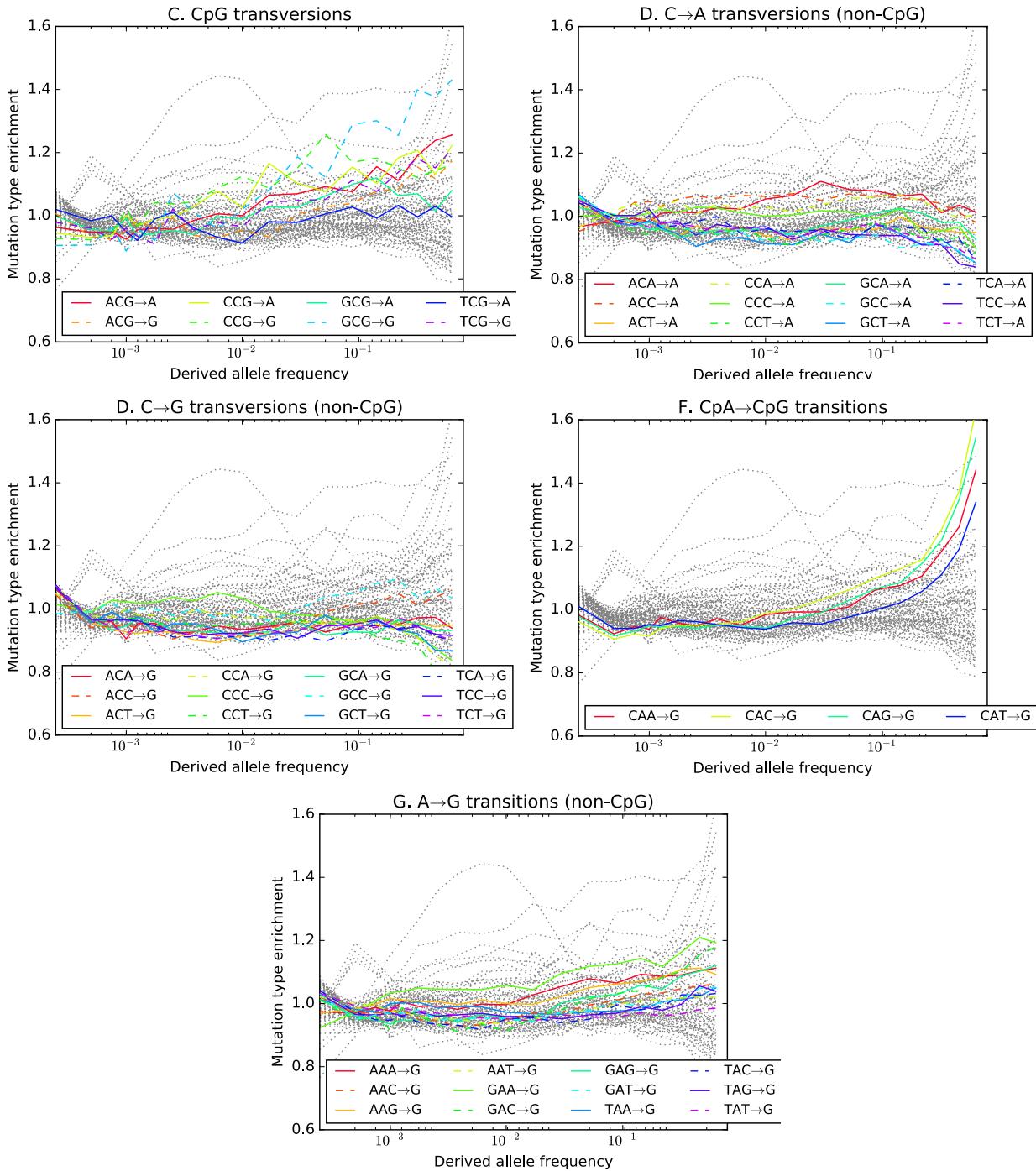
451

452 **Figure 3–Figure Supplement 2: Fraction of TCC→TTC mutations as a function**
453 **of allele frequency in ExAC.** Lek, et al. compiled data from 60,706 exomes to create
454 the Exome Aggregation Consortium dataset, which enables the analysis of ultra-rare
455 human variation [39]. The overall fraction of TCC→TTC mutations is slightly higher in
456 exome data than in whole genome data because exons contain a skewed distribution of
457 triplet contexts, but the pulse pattern from Figure 3B reproduces unmistakably.



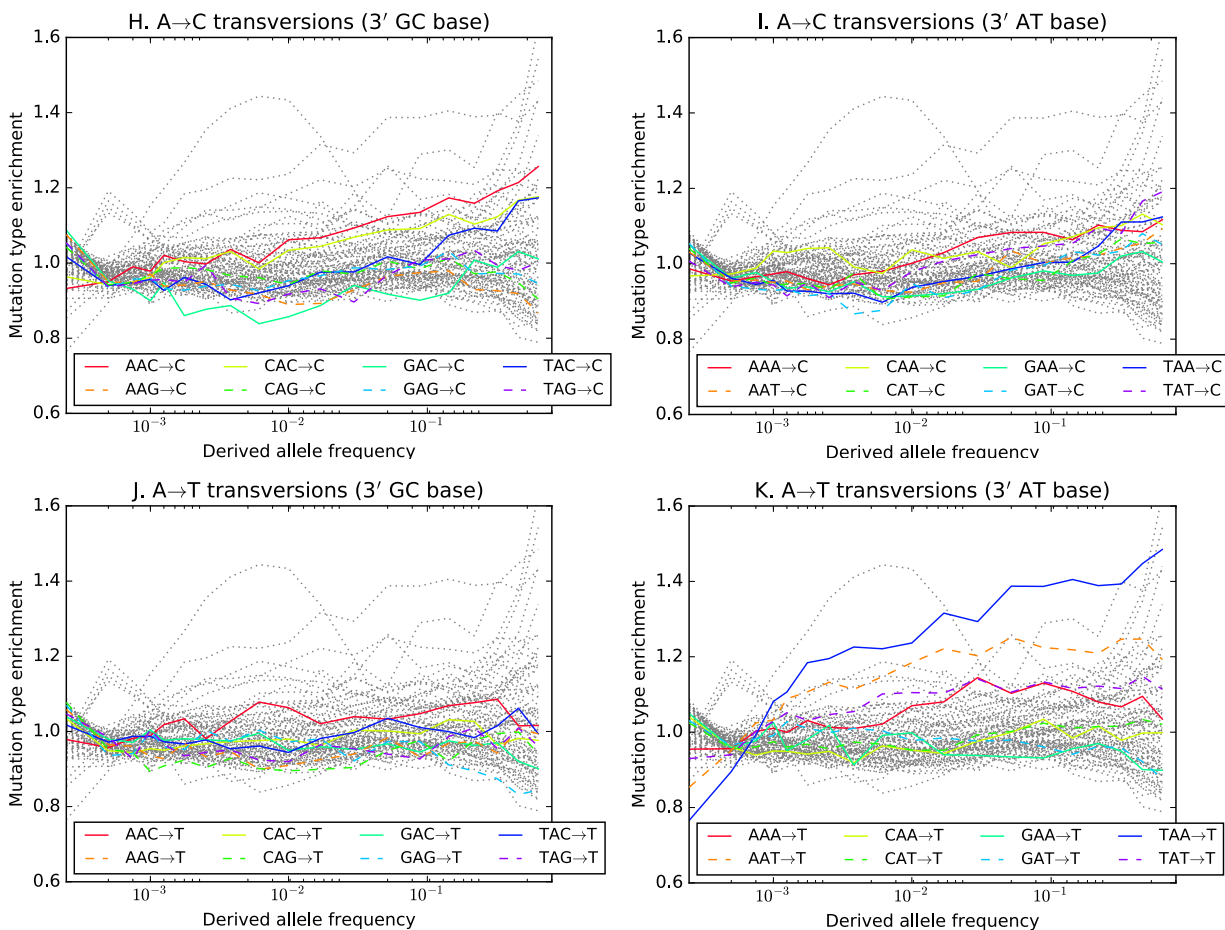
458

459 **Figure 3—Figure Supplement 3: Mutation type enrichment as a function of**
 460 **allele frequency in UK10K (Part I of III).** The eleven panels in Figure Supplements
 461 2, 3, and 4 show the full dependence of mutation spectrum on allele frequency in the
 462 UK10K data. If we let $F(f, m)$ denote the fraction of SNVs of frequency f that are of type
 463 m and let $F(m)$ denote the fraction of all mutations that are of type m , the enrichment of
 464 mutation type m as a function of frequency is $F(f, m)/F(m)$. This function is expected to
 465 fluctuate around $y = 1$ unless the rate of m has recently increased or decreased. All 96
 466 mutation types are visualized in every panel, but most corresponding lines are greyed out
 467 to enhance readability. Some lines deviate from $y = 1$ due to the effects of biased gene
 468 conversion (BGC)—this occurs when one of the ancestral or derived alleles is a weak base
 469 (A or T, abbreviated W) and the other allele is a strong base (G or C, abbreviated S).
 470 W→S mutations are more abundant at high allele frequencies, while S→W mutations are
 471 more abundant at low frequencies. These effects are visible but modest in panels D, G, H,
 472 and I, but much more pronounced in panels B, C, and F, which focus on mutations in the
 473 CpG context. Transitions of the type CpA→CpG, which create CpG motifs, are extremely
 474 enriched at high frequencies, and this pattern may be an artifact of ancestral
 475 misidentification [40]. CpG motifs have such high mutation rates that CpG→CpT
 476 transitions often happen at the same site in humans and chimps, and these low-frequency
 477 double mutations are misclassified as high-frequency CpT→CpG mutations. Although it is
 478 not surprising to see a peak of CpT→CpG transitions at high frequencies in panel F, it is
 479 somewhat surprising to see CpG→GpG transversions peak in abundance at high
 480 frequencies in panel C. This might be a signature of recent declines in the rates of these
 481 mutations, since neither ancestral misidentification nor biased gene conversion is thought
 482 to produce such a pattern. In addition, neither of these processes can explain the strong
 483 enrichment of certain A→T mutations at high frequencies that is observed in panel K.

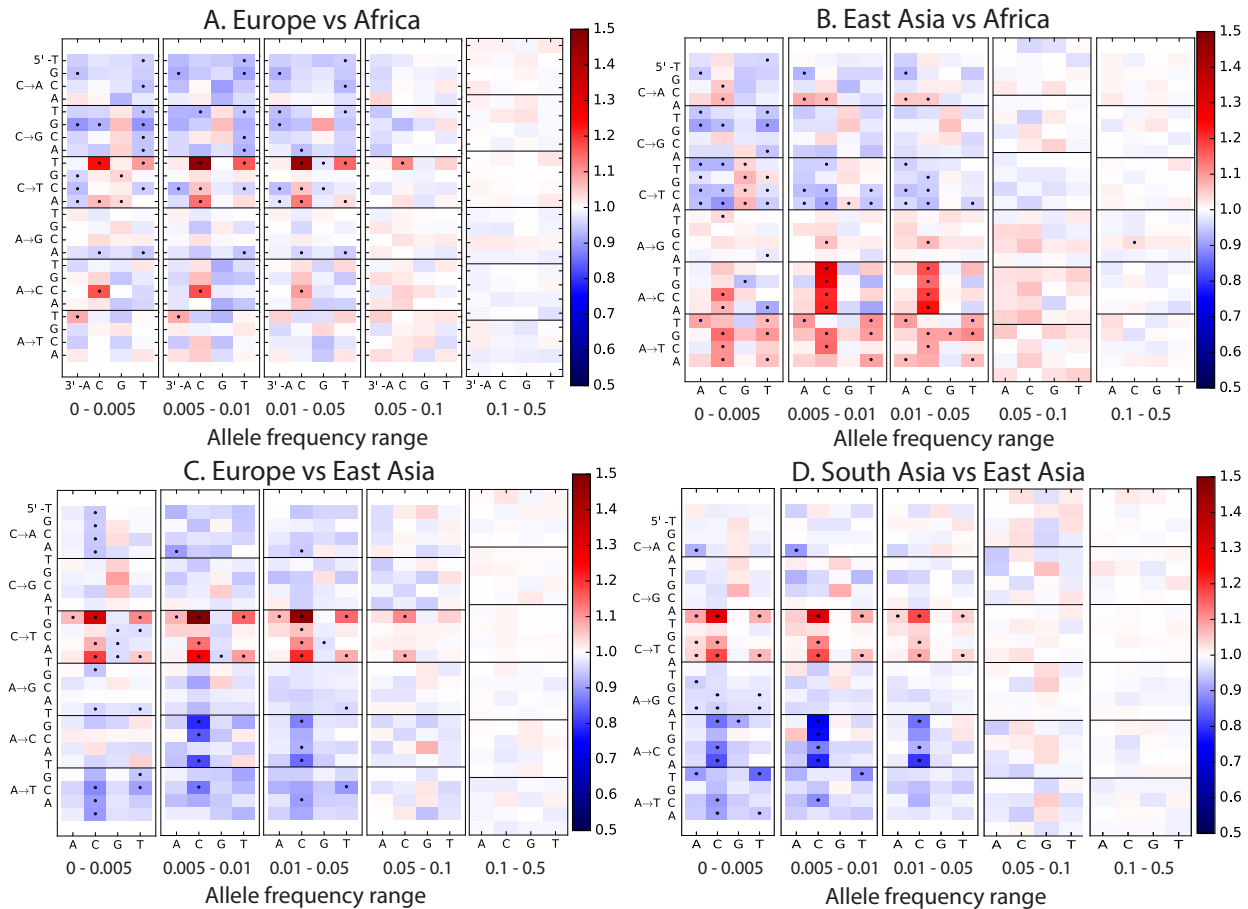


484

485 **Figure 3–Figure Supplement 4: Mutation type enrichment as a function of**
486 **allele frequency in UK10K (Part II of III).** The eleven panels in this 3-part figure
487 show the full dependence of mutation spectrum on allele frequency in the UK10K data.

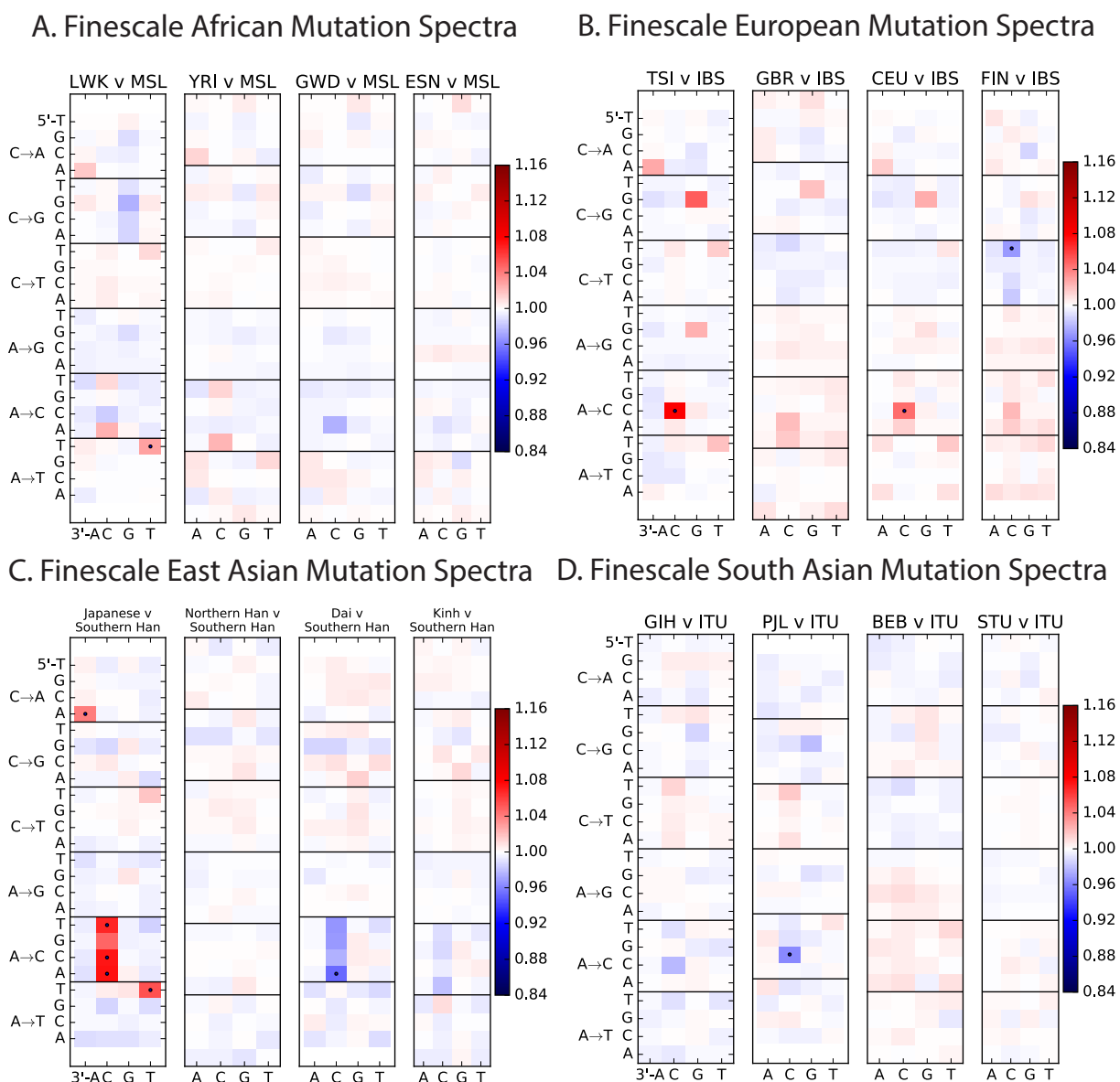


489 **Figure 3–Figure Supplement 5: Mutation type enrichment as a function of**
490 **allele frequency in UK10K (Part III of III).** The eleven panels in this 3-part figure
491 show the full dependence of mutation spectrum on allele frequency in the UK10K data.



492

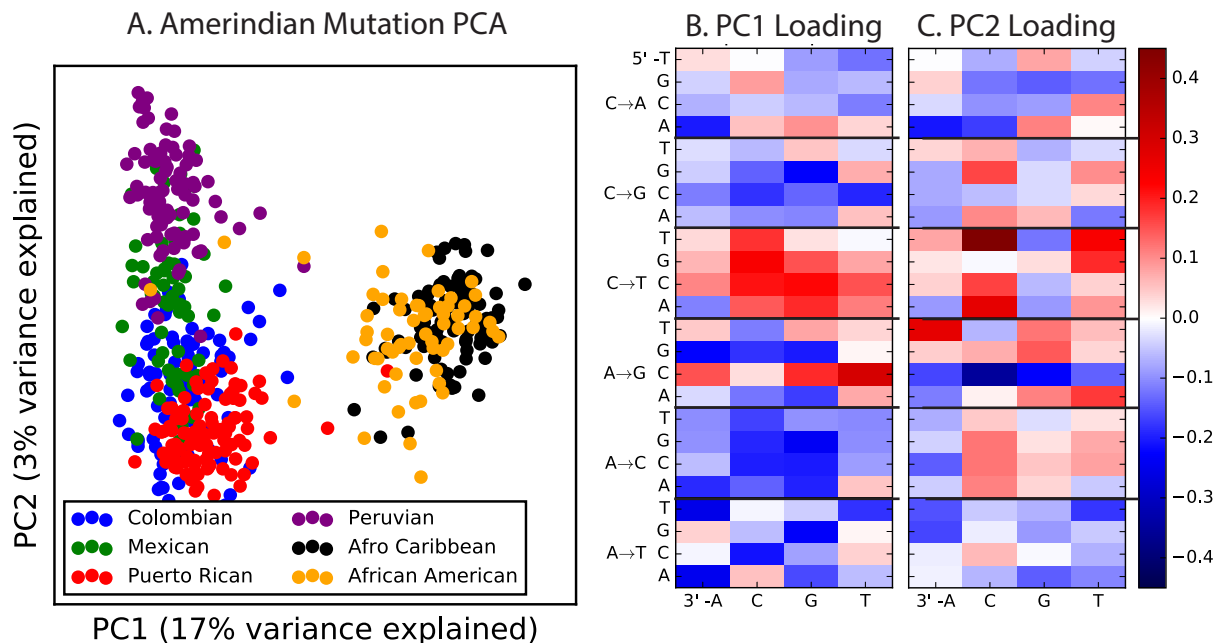
493 **Figure 3–Figure Supplement 6: Mutation spectrum comparisons partitioned by**
494 **allele frequency.** Each of these heatmaps shows a subset of the data used to construct
495 Figure 1B, partitioned by allele frequency to show how rare variants are the most highly
496 differentiated between populations. Black dots highlight mutation types that are
497 significantly different in abundance between two populations in a particular frequency class
498 at the $p < 10^{-5}$ level according to a chi-square test.



499

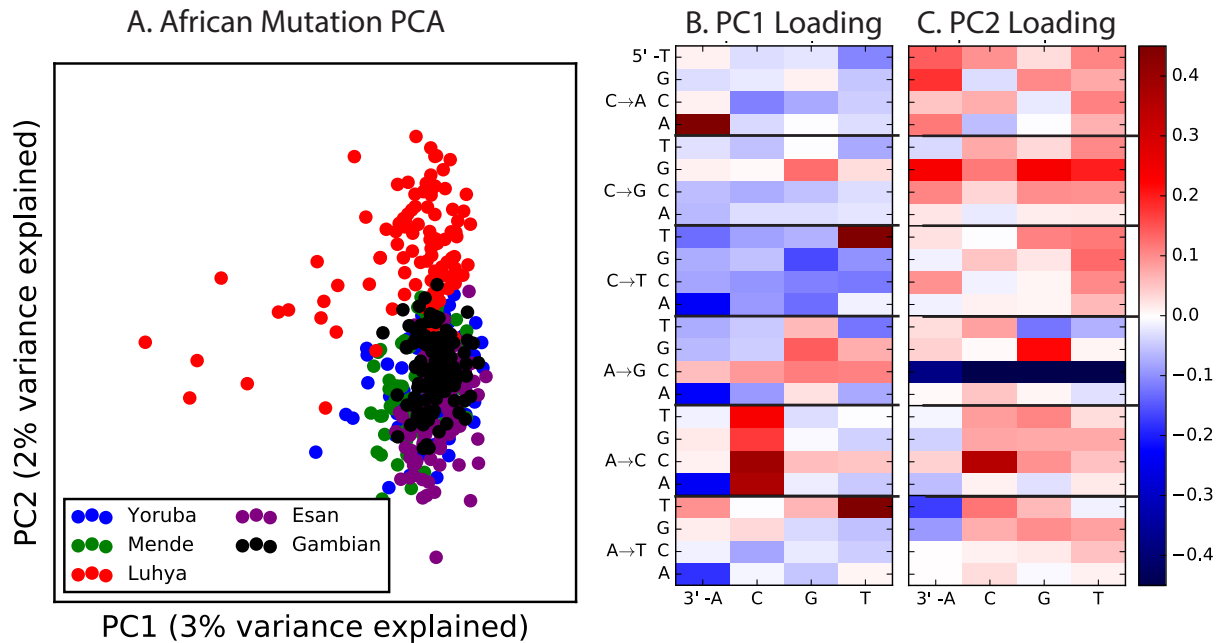
500 **Figure 4–Figure Supplement 1: Mutation spectrum differences within Africa,**
 501 **Europe, East Asia, and South Asia.** Figure 4B of the main text shows heat map
 502 comparisons between East Asian populations, which display fine-scale differences that are
 503 exceptionally well defined. For completeness, this figure shows finescale heatmap
 504 comparisons within all 1kG continental groups. We can see that CAC→CCC and
 505 TAT→TTT are heterogeneously distributed within multiple continents, but to the greatest
 506 extent in East Asia. In addition the TCC→TTC signature is somewhat heterogeneously
 507 distributed within Europe and South Asia, being depleted in Finns and enriched in the
 508 Punjabi and Gujarati. Each continental group in the 1000 Genomes data is divided into 5
 509 sub-populations. These heat maps compare the mutation spectra of these fine-scale
 510 populations to each other. African populations are: MSL = Mende in Sierra Leone; LWK =
 511 Luhya in Webuye, Kenya; YRI = Yoruba in Ibadan, Nigeria; GWD = Gambian in Western
 512 Divisions; ESN = Esan in Nigeria. European populations are: IBS = Iberian Population in

513 Spain; TSI = Toscani in Italia; GBR = British in England and Scotland; CEU = Utah
 514 Residents (CEPH) with Northern and Western Ancestry; FIN = Finnish in Finland. East
 515 Asian populations are: CDX = Chinese Dai in Xishuangbanna, China; JPT = Japanese in
 516 Tokyo, Japan; CHB = Han Chinese in Beijing, China; CHS = Southern Han Chinese; KHV
 517 = Kinh in Ho Chi Minh City, Vietnam. South Asian populations are: ITU = Indian Telugu
 518 from the UK; GIH = Gujarati Indian from Houston, Texas; PJJ = Punjabi from Lahore,
 519 Pakistan; BEB = Bengali from Bangladesh; STU = Sri Lankan Tamil from the UK.



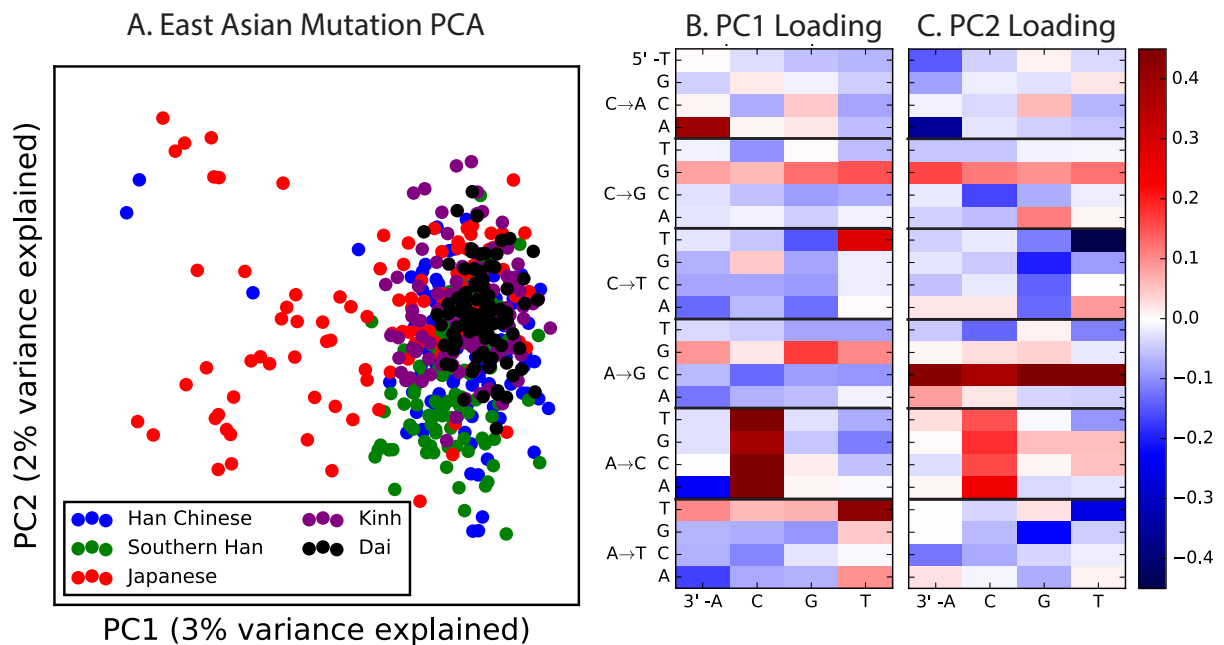
520

521 **Figure 4-Figure Supplement 2: PCA of American populations.** Population
 522 abbreviations are: CLM = Colombians from Medellin, Colombia; MXL = Mexican
 523 Ancestry from Los Angeles, USA; PUR = Puerto Ricans from Puerto Rico; PEL =
 524 Peruvians from Lima, Peru; ACB = African Caribbeans in Barbados; ASW = Americans
 525 of African Ancestry in SW USA. Admixed populations from the Americans show structure
 526 that mirrors the continental groups, with PC1 essentially measuring the ratio between
 527 African and non-African ancestry and PC2 measuring the ratio between European and
 528 Native American ancestry. The accompanying heat maps show the loadings of the first two
 529 principal components.



530

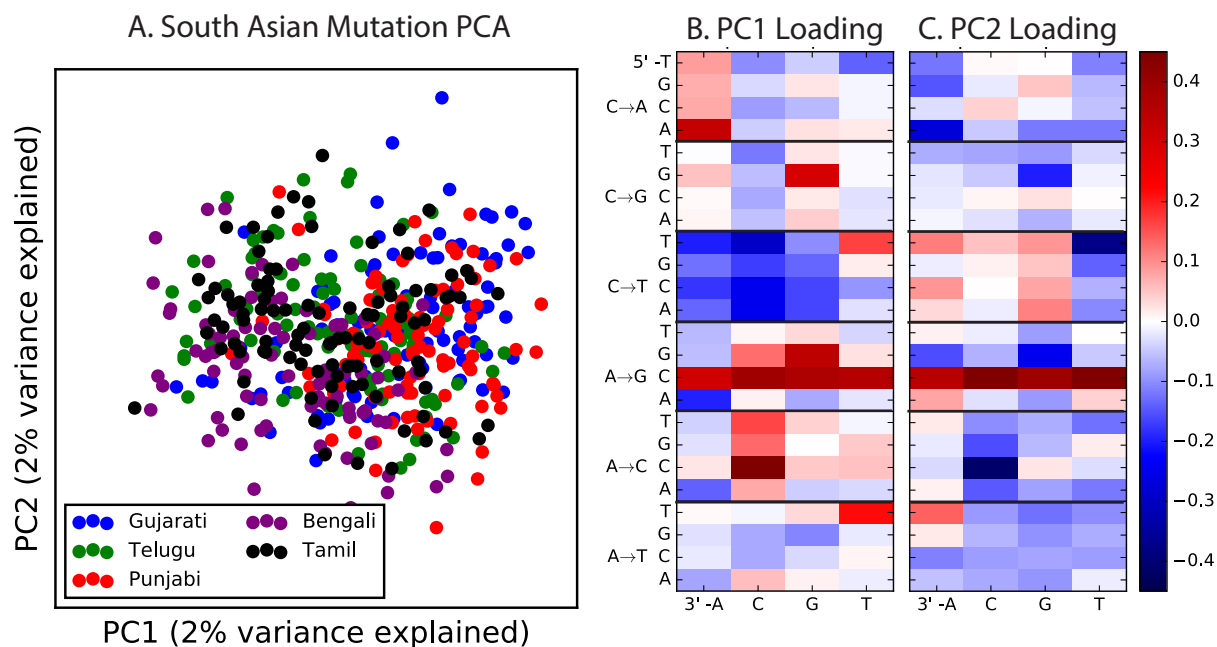
531 **Figure 4–Figure Supplement 3: PCA of African populations.** Population
 532 abbreviations are: MSL = Mende in Sierra Leone; LWK = Luhya in Webuye, Kenya; YRI
 533 = Yoruba in Ibadan, Nigeria; GWD = Gambian in Western Divisions; ESN = Esan in
 534 Nigeria.



535

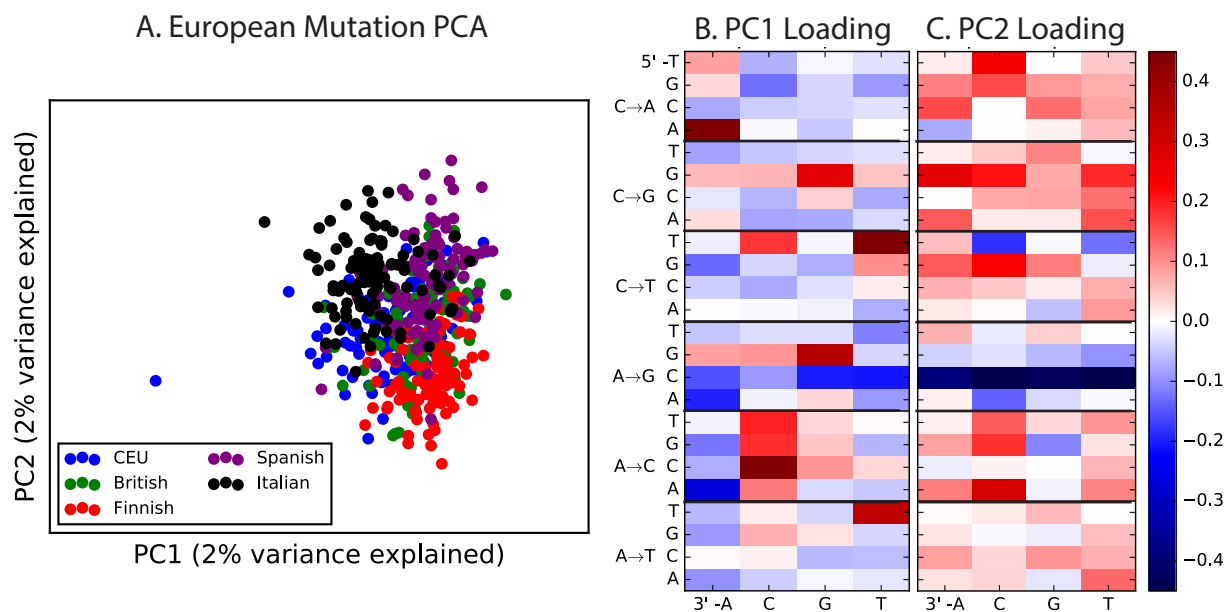
536 **Figure 4–Figure Supplement 4: PCA of East Asian populations.** Population
 537 abbreviations are: CDX = Chinese Dai in Xishuangbanna, China; JPT = Japanese in

538 Tokyo, Japan; CHB = Han Chinese in Beijing, China; CHS = Southern Han Chinese; KHV
 539 = Kinh in Ho Chi Minh City, Vietnam.



540

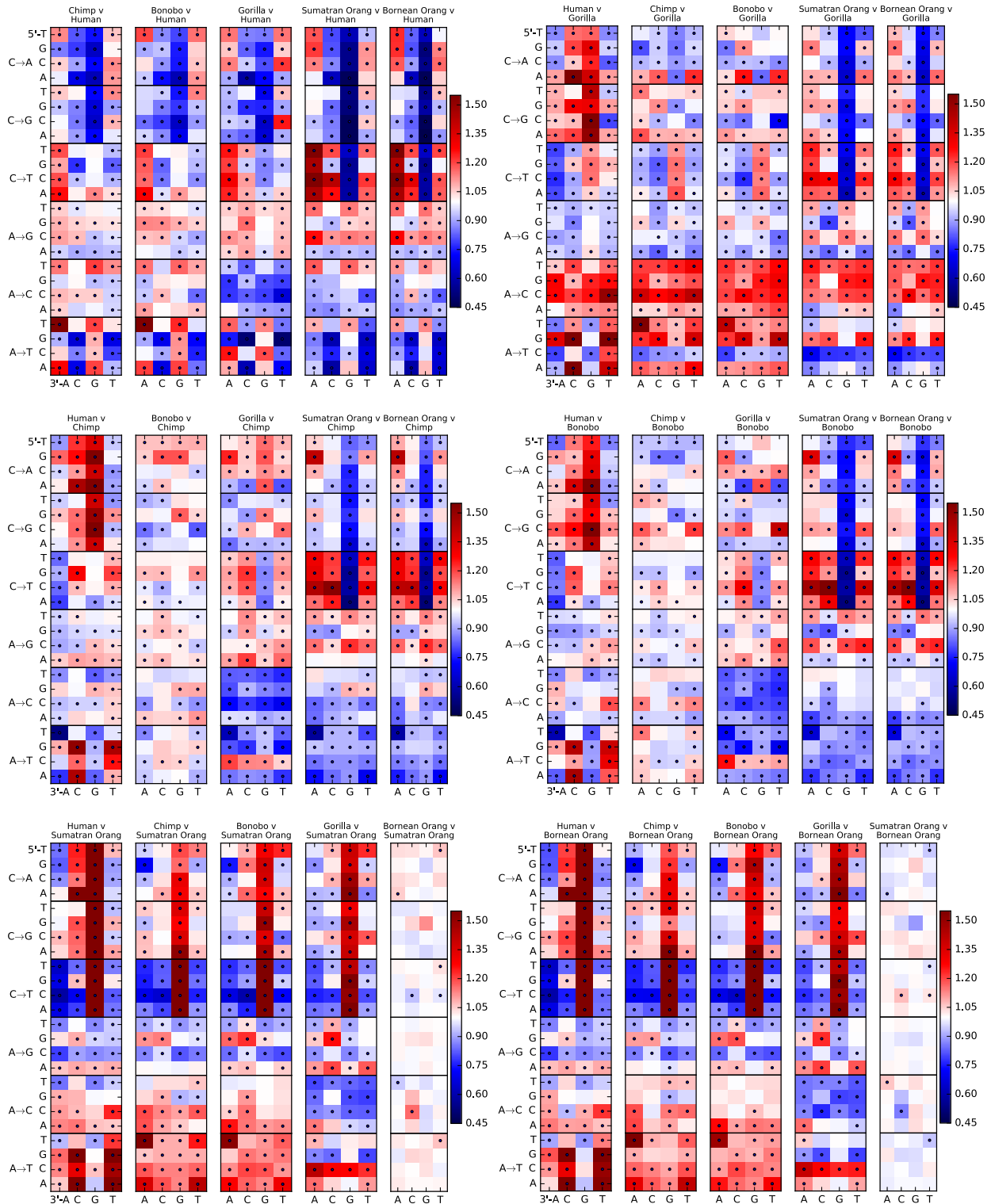
541 **Figure 4–Figure Supplement 5: PCA of South Asian populations.** Population
 542 abbreviations are: ITU = Indian Telugu from the UK; GIH = Gujarati Indian from
 543 Houston, Texas; PJJ = Punjabi from Lahore, Pakistan; BEB = Bengali from Bangladesh;
 544 STU = Sri Lankan Tamil from the UK.



545

546 **Figure 4–Figure Supplement 6: PCA of European populations.** Population
547 abbreviations are: IBS = Iberian Population in Spain; TSI = Toscani in Italia; GBR =
548 British in England and Scotland; CEU = Utah Residents (CEPH) with Northern and
549 Western Ancestry; FIN = Finnish in Finland.

550 **Figure 4–Source Data 1.** This text file shows the number of SNPs in each of the 96
551 mutational categories that passed all filters in each finescale 1000 Genomes population.



552

553 **Figure 5—Figure Supplement 1: Mutation spectra of great apes.** These heatmap
 554 comparisons demonstrate that closely related great apes such as Chimpanzees and Bonobos
 555 have more similar mutation spectra than more distantly related apes do.

References

- 556
- 557 [1] Kong, A. *et al.* Rate of *de novo* mutations and the importance of father's age to
558 disease risk. *Nature* **488**, 471–475 (2012).
- 559 [2] Ségurel, L., Wyman, M. & Przeworski, M. Determinants of mutation rate variation in
560 the human germline. *Annu Rev Genomics Hum Genet* **15**, 19.1–19.24 (2014).
- 561 [3] Alexandrov, L. *et al.* Signatures of mutational processes in human cancer. *Nature*
562 **500**, 415–421 (2013).
- 563 [4] Helleday, T., Eshtad, S. & Nik-Zainal, S. Mechanisms underlying mutational
564 signatures in human cancers. *Nature Reviews Genetics* **15**, 585–598 (2014).
- 565 [5] Shinbrot, E. *et al.* Exonuclease mutations in DNA polymerase epsilon reveal
566 replication strand specific mutation patterns and human origins of replication.
567 *Genome research* **24**, 1740–1750 (2014).
- 568 [6] Ohno, M. *et al.* 8-oxoguanine causes spontaneous *de novo* germline mutations in mice.
569 *Scientific Reports* **4**, 10.1038/srep04689 (2014).
- 570 [7] Lynch, M. The lower bound to the evolution of mutation rates. *Genome biology and
571 evolution* **3**, 1107–1118 (2011).
- 572 [8] Sung, W., Ackerman, M., Miller, S., Doak, T. & Lynch, M. Drift-barrier hypothesis
573 and mutation-rate evolution. *Proc Natl Acad Sci USA* **109**, 18488–18492 (2012).
- 574 [9] Hwang, D. & Green, P. Bayesian Markov chain Monte Carlo sequence analysis reveals
575 varying neutral substitution patterns in mammalian evolution. *Proc Natl Acad Sci
576 USA* **101**, 13994–14001 (2004).
- 577 [10] But, D. *et al.* Mismatch repair incompatibilities in diverse yeast populations. *Genetics*
578 **205**, 10.1534/genetics.116.199513 (2017).
- 579 [11] Seoighe, C. & Scally, A. Inference of candidate germline mutator loci in humans from
580 genome-wide haplotype data. *PLoS Genetics* **13**, e1006549 (2017).
- 581 [12] Harris, K. Evidence for recent, population-specific evolution of the human mutation
582 rate. *Proc Natl Acad Sci USA*, **112**, 3439–3444 (2015).

- 583 [13] Mathieson, I. & Reich, D. E. Variation in mutation rates among human populations.
584 *bioRxiv* 063578 (2016).
- 585 [14] 1000 Genomes Project Consortium *et al.* A global reference for human genetic
586 variation. *Nature* **526**, 68–74 (2015).
- 587 [15] Pollard, K., Hubisz, M., Rosenbloom, K. & Siepel, A. Detection of nonneutral
588 substitution rates on mammalian phylogenies. *Genome Research* **20**, 110–121 (2010).
- 589 [16] Novembre, J. *et al.* Genes mirror geography within Europe. *Nature* **456**, 98–101
590 (2008).
- 591 [17] Mallick, S. *et al.* The Simons Genome Diversity Project: 300 genomes from 142
592 diverse populations. *Nature* **538**, 201–206 (2016).
- 593 [18] UK10K Consortium. The UK10K project identifies rare variants in health and disease.
594 *Nature* **526**, 82–90 (2015).
- 595 [19] Prado-Martinez, J. *et al.* Great ape genetic diversity and population history. *Nature*
596 **499**, 471–475 (2013).
- 597 [20] Moorjani, P., Amorim, C., Arndt, P. & Przeworski, M. Variation in the molecular
598 clock of primates. *Proc Natl Acad Sci USA* **113**, 10607–10612 (2016).
- 599 [21] Goodman, M. The role of immunochemical differences in the phyletic development of
600 human behavior. *Human Biol* **33**, 131–162 (1961).
- 601 [22] Li, W. & Tanimura, M. The molecular clock runs more slowly in man than in apes
602 and monkeys. *Nature* **326**, 93–96 (1987).
- 603 [23] Scally, A. & Durbin, R. Revising the human mutation rate: implications for
604 understanding human evolution. *Nature Rev Genetics* **13**, 745–753 (2012).
- 605 [24] Galtier, N., Piganeau, G., Mouchiroud, D. & Duret, L. GC-content evolution in
606 mammalian genomes: the biased gene conversion hypothesis. *Genetics* **159**, 907–911
607 (2001).
- 608 [25] Carlson, J. *et al.* Extremely rare variants reveal patterns of germline mutation rate
609 heterogeneity in humans. *bioRxiv preprint* <https://doi.org/10.1101/108290> (2017).

- 610 [26] Do, R. *et al.* No evidence that selection has been less effective at removing deleterious
611 mutations in Europeans than in Africans. *Nature Genetics* **47**, 126–131 (2015).
- 612 [27] Martin, A. P. & Palumbi, S. R. Body size, metabolic rate, generation time, and the
613 molecular clock. *Proceedings of the National Academy of Sciences* **90**, 4087–4091
614 (1993).
- 615 [28] Amster, G. & Sella, G. Life history effects on the molecular clock of autosomes and
616 sex chromosomes. *Proceedings of the National Academy of Sciences* **113**, 1588–1593
617 (2016).
- 618 [29] Moorjani, P., Gao, Z. & Przeworski, M. Human germline mutation and the erratic
619 molecular clock. *bioRxiv* 058024 (2016).
- 620 [30] Gao, Z., Wyman, M., Sella, G. & Przeworski, M. Interpreting the dependence of
621 mutation rates on age and time. *PLoS Biology* **14**, e1002355 (2016).
- 622 [31] Narasimhan, V. *et al.* A direct multi-generational estimate of the human mutation rate
623 from autozygous segments seen in thousands of parentally related individuals. *bioRxiv*
624 *preprint* <http://dx.doi.org/10.1101/059436> (2016).
- 625 [32] Coop, G., Wen, X., Ober, C., Pritchard, J. K. & Przeworski, M. High-resolution
626 mapping of crossovers reveals extensive variation in fine-scale recombination patterns
627 among humans. *Science* **319**, 1395–1398 (2008).
- 628 [33] Baudat, F. *et al.* PRDM9 is a major determinant of meiotic recombination hotspots in
629 humans and mice. *Science* **327**, 836–840 (2010).
- 630 [34] Smit, A., Hubley, R. & Green, P. RepeatMasker Open-4.0.
631 <http://www.repeatmasker.org> (2013-15).
- 632 [35] Field, Y. *et al.* Detection of human adaptation during the past 2,000 years. *Science*
633 **353**, 10.1126/science.aag0776 (2016).
- 634 [36] Tennessen, J. *et al.* Evolution and functional impact of rare coding variation from
635 deep sequencing of human exomes. *Science* **338**, 64–69 (2012).
- 636 [37] Hoffman, M. *et al.* Integrative annotation of regulatory elements from ENCODE data.
637 *Nucleic Acids Research* **41**, 827–841 (2013).

- 638 [38] Woodfine, K. *et al.* Replication timing of the human genome. *Hum Mol Genet* **13**,
639 191–202 (2004).
- 640 [39] Lek, M. *et al.* Analysis of protein-coding genetic variation in 60,706 humans. *Nature*
641 **536**, 285–291 (2016).
- 642 [40] Hernandez, R., Williamson, S. & Bustamante, C. Context dependence, ancestral
643 misidentification, and spurious signatures of natural selection. *Mol Biol Evol* **24**,
644 1792–1800 (2007).