1 **Contrasting determinants of mutation rates in germline and soma**

2

3 Chen Chen[1,2,*], Hongjian Qi[3,4], Yufeng Shen[3,5], Joseph Pickrell[1,2,+], Molly

4 Przeworski[1,3,+]

5

6 [1]Department of Biological Sciences, Columbia University, New York, NY

7 [2]New York Genome Center, New York, NY

8 [3]Departments of Systems Biology, Columbia University Medical Center, New

9 York, NY

10 [4]Department of Applied Physics and Applied Mathematics, Columbia University,

11 New York, NY

12 [5]Department of Biomedical Informatics, Columbia University, New York, NY

13 [+]Co-supervised this project

14 [*]To whom correspondence should be addressed: cc3499@columbia.edu

15

16 Keywords:

17 human; mutation rate; germline mutations; somatic mutations; strand asymmetry

1

18      **Abstract**

19

20      Recent studies of somatic and germline mutations have led to the identification of

21      a number of factors that influence point mutation rates, including CpG

22      methylation, expression levels, replication timing and GC content. Intriguingly,

23      some of the effects appear to differ between soma and germline: in particular,

24      whereas mutation rates have been reported to decrease with expression levels in

25      tumors, no clear effect has been detected in the germline.  Distinct approaches

26      were taken to analyze the data, however, so it is hard to know whether these

27      apparent differences are real. To enable a cleaner comparison, we considered a

28      statistical model in which the mutation rate of a coding region is predicted by GC

29      content, expression levels, replication timing, and two histone repressive marks.

30      We applied this model to both a set of germline mutations identified in exomes

31      and to exonic somatic mutations in four types of tumors. Germline and soma

32      share most determinants of mutations; notably, we detected an effect of

33      expression levels on germline mutations as well as on somatic ones. However,

34      whereas in somatic tissues, increased expression levels are associated with

35      greater strand asymmetry and *decreased* mutation rates, in ovaries and testes,

36      increased expression leads to greater strand asymmetry but *increased* mutation

37      rates. This contrast points to differences in damage or repair rates during

38      transcription in soma and germline.

39

40 **Introduction**

41

42 Germline mutations are the source of all heritable variation, including in disease

43 susceptibility, and it is increasingly clear that somatic mutations also play

44 important roles in human diseases, notably cancers (Muller 1927; Stratton,

45 Campbell, and Futreal 2009). Understanding the rate and mechanisms by which

46 mutations occur is therefore of interest to both evolutionary biologists and to

47 human geneticists aiming to identify the underlying causes of genetic diseases

48 (Shendure and Akey 2015; Gao et al. 2016). In particular, an accurate estimate

49 of the local mutation rate is key to testing for an excess of disease mutations in

50 specific genes among cases (Lawrence et al. 2013; Samocha et al. 2014).

51 Characterization of the variation in mutation rate along the genome can also yield

52 important insights into DNA damage and repair mechanisms (Stratton 2011;

53 Ségurel, Wyman, and Przeworski 2014).

54

55 Until recently, our understanding of germline point mutations came mainly from

56 analysis of diversity along the genome or divergence among species (Green et al.

57 2003; Webster et al. 2004; Polak and Arndt 2008; Hodgkinson and Eyre-Walker

58 2011; Park, Qian, and Zhang 2012). In the past several years, analyses have

59 also been based on resequencing exomes or whole genomes from blood

60 samples of human pedigrees and identifying variants present in the offspring but

61 absent in the child (reviewed in Campbell and Eichler 2013 and Ségurel, Wyman,

62 and Przeworski 2014; Shendure and Akey 2015; Francioli et al. 2015; Rahbari et

63    al. 2016; Goldmann et al. 2016; Besenbacher et al. 2016). This approach is more

64    direct than analyzing divergence data and presents the advantage of being

65    almost unaffected by selection, but the analysis is technically challenging and,

66    with current study designs, some mutations may be missed, notably those that

67    occur in the early post-zygotic divisions (Rahbari et al. 2016; Moorjani, Gao, and

68    Przeworski 2016; Harland et al. 2016).

69

70    Our knowledge of somatic point mutations, in turn, relies primarily on

71    resequencing tumors. In these analyses, mutation calls are made by sequencing

72    tumor and non-cancerous tissue pairs from the same individual and then

73    excluding the variants shared between the two tissues (as the shared mutations

74    are likely to be germline). Because, in this approach, a large population of cells is

75    sequenced, the mutations identified tend to predate the tumorigenesis and thus

76    are mostly somatic mutations that occurred in normal tissues (see, e.g.,

77    Martincorena et al. 2015; Alexandrov et al. 2015).

78

79    Studies of both germline and soma reveal that the point mutation rate varies

80    across the genome, from the scale of a single base pair to much larger scales

81    (Hodgkinson and Eyre-Walker 2011; Hodgkinson, Chen, and Eyre-Walker 2012;

82    Ségurel, Wyman, and Przeworski 2014). At the single base pair level, the largest

83    source of variation in germline mutation rate is the identity of the adjacent base

84    pairs (Hwang and Green 2004; Hodgkinson and Eyre-Walker 2011). Notably, the

85    mutation rate of CpG transitions (henceforth CpG Ti) is an order of magnitude

4

86  higher than other mutation types (e.g., Kong et al. 2012). Most CpG dinucleotides

87  are methylated in the human genome; when the methylated cytosine undergoes

88  spontaneous deamination to generate thymine and is not corrected by the time of

89  replication, the damage leads to a mutation. Among other types of sites, rates of

90  mutation vary by 2 to 3 fold (Kong et al. 2012). In the soma, the mutation rate at

91  CpG sites is also elevated, although the extent of the increase differs across

92  tumor types (Pleasance, Stephens, et al. 2010; Pleasance, Cheetham, et al.

93  2010; Lee et al. 2010). More generally, tumors vary in their mutation spectrum:

94  analyses of mutations and their two neighboring base pairs (i.e., considering 96

95  mutation types) point to enrichment of distinct mutational signatures for different

96  types of cancers, a subset of which have been shown to reflect particular

97  mutagens or differences in the efficiency of repair (Alexandrov et al. 2013).

98

99  Over a larger scale of megabases, germline mutation rates have been

100  associated with a number of additional factors, including transcription level (in

101  testis), replication timing (in lymphoblastoid cell lines), chromatin state (both in

102  lymphoblastoid cells and in ovary), meiotic crossover rates and GC content

103  (Hodgkinson and Eyre-Walker 2011; Michaelson et al. 2012; Park, Qian, and

104  Zhang 2012; Francioli et al. 2015; Goldmann et al. 2016; Besenbacher et al.

105  2016). Somatic mutation rates have also been associated with replication timing

106  (in Hela cell lines) and with average transcription levels across 91 cell lines in

107  Cancer Cell Line Encyclopedia (Lawrence et al. 2013).

108

109   In many cases, little is known about the mechanistic basis for the association of a

110   given factor with mutation rates. However, the association of somatic mutation

111   rates with transcription levels appears to be a byproduct of transcription-coupled

112   repair (TCR), a sub-pathway of nucleotide excision repair (NER) (Hanawalt and

113   Spivak 2008; Nouspikel 2009). NER is a versatile repair pathway that senses

114   lesion-causing distortions to DNA structure and excises the lesion for repair.

115   Another subpathway of NER, global genome repair (GGR), can repair lesions on

116   both transcribed strand (henceforth TS) and non-transcribed strand (henceforth

117   NTS), including transcribed regions as well as transcriptionally-silent ones. In

118   contrast, TCR operates only within transcribed regions, triggered by lesions on

119   the TS, which it repairs off the NTS. This mechanism gives rise to a mutational

120   strand asymmetry as well as a compositional asymmetry between strands. For

121   example, TCR leads to more A to G mutations (A>G henceforth) on the NTS than

122   TS; acting over long periods of time, this phenomenon generates an excess of G

123   over A (and T over C) on the NTS (Green et al. 2003; McVicker and Green 2010).

124   Such mutational strand asymmetry has been found in both germline and soma

125   (Green et al. 2003; Polak and Arndt 2008; Rubin and Green 2009; Lawrence et al.

126   2013; Martincorena et al. 2015; Francioli et al. 2015).

127

128   While many of the same determinants appear to play important roles in both

129   germline and soma, there are hints of differences as well. For instance, studies of

130   pre-neoplastic somatic mutations indicate that, over a 100 kb scale, the mutation

131   rates in somatic tissues decrease with expression levels and increase with

132    replication timing (Lawrence et al. 2013). Similarly, two studies that focused on

133    somatic mutations in non-cancerous somatic tissues, normal eyelid tissue and

134    neurons, found mutations to be enriched in regions of low expression and

135    repressed chromatin (Martincorena et al. 2015; Lodato et al. 2015). A similar

136    effect of replication timing was identified in studies of germline mutation

137    (Stamatoyannopoulos et al. 2009; Francioli et al. 2015; Besenbacher et al. 2016;

138    Carlson et al. 2017). However, the effect of expression levels on germline

139    mutation rates remains unclear: one study reported increased divergence

140    between humans and macaques with greater germline expression (Park, Qian,

141    and Zhang 2012), but others found no discernable effect of expression levels on

142    mutation rates (Green et al. 2003; Webster et al. 2004; Polak and Arndt 2008;

143    Hodgkinson and Eyre-Walker 2011; Francioli et al. 2015). This difference

144    between germline and soma is particularly puzzling in light of the observation that

145    the strand asymmetry of mutation rates between TS and NTS is seen in the

146    germline as well as the soma (Pleasance, Cheetham, et al. 2010; Pleasance,

147    Stephens, et al. 2010; McVicker and Green 2010; Lawrence et al. 2013).

148    Together, these observations suggest that the determinants of mutation rates

149    may differ between germline and soma, raising the more general possibility that

150    the damage rate or the repair efficacy differs among cell types (Lynch 2010).

151

152    A limitation, however, is that studies have used different statistical approaches,

153    rendering the comparison hard to interpret. As an illustration, whereas some

154    studies binned the genome into windows of 100 kb (e.g., Lawrence et al. 2013)

155    or 1Mb regions (e.g., Polak et al. 2015), other studies have compared the mean

156    mutation rate in transcribed regions and non-transcribed regions or in genes

157    grouped by expression levels (Hodgkinson and Eyre-Walker 2011; Francioli et al.

158    2015; Lodato et al. 2015). Studies of somatic mutation also vary in whether they

159    group different tissues or distinguish among them (e.g., Pleasance, Stephens, et

160    al. 2010; Lawrence et al. 2013). An additional limitation of earlier studies of

161    germline mutation is that, by necessity, they relied on human-chimpanzee

162    divergence as a proxy for de novo mutation rates (Green et al. 2003; Webster et

163    al. 2004; Hodgkinson and Eyre-Walker 2011), even though divergence reflects

164    not only the mutation process but also effects of natural selection in the human-

165    chimpanzee ancestor and biased gene conversion (McVicker et al. 2009; Duret

166    and Galtier 2009).

167

168    To our knowledge, only one study has used a uniform approach to study

169    germline and soma. Their findings point to possible differences in their

170    determinants: for instance, the histone mark H3K9me3 accounts for more than

171    40% of mutation rate variation at 100 kb in tumors, when a much weaker

172    association is seen in the germline (Schuster-Böckler and Lehner 2012;

173    Goldmann et al. 2016). This analysis relied on pairwise correlations, however,

174    and therefore the results may be confounded by other factors that are correlated

175    to the histone marks and differ between tissues. Moreover, to our knowledge,

176    there has been no parallel treatment of strand asymmetry in germline and soma.

177

178    To overcome these limitations, we built a multivariable regression model, in

179    which the mutation rates of CpG Ti and other types of mutations in a coding

180    region are predicted by GC content, expression levels, replication timing and two

181    histone repressive marks. To this end, we used the expression levels, replication

182    timing and histone marker levels of matched cell types. We applied the model to

183    a large set of germline point mutations identified in exomes from recently

184    published studies on developmental disorders and to somatic point mutations in

185    exomes found in four types of tumors and reported by the Cancer Genome Atlas

186    (see Materials and Methods). In addition, we considered the mutational strand

187    asymmetry in the two sets of data.

188

189 **Materials and Methods**

190

191 **Datasets**. To study germline mutations, we relied on de novo mutation calls

192 made from 8681 trios surveyed by exome sequencing. We combined results from

193 two main sources: studies of neurodevelopmental disorders (NDD), which

194 considered 5542 cases and 1911 controls (unaffecteds), and studies of

195 congenital heart defect (CHD), conducted by the Pediatric Cardiac Genomics

196 Consortium, which included 1228 trios. The NDD cases include 3953 cases of

197 Autism Spectrum Disorder (ASD), 1133 cases of deciphering developmental

198 disorders (DDD), 264 cases of epileptic encephalopathies (EE), and 192 cases of

199 intellectual disability (ID). All these studies applied similar capture and

200 sequencing methods, and most samples were at >20X coverage (see Table 1).

201 We tested for an effect of the study, which could potentially arise from differences

202 in design or analysis pipeline, by adding a categorical variable (by an analogous

203 approach to the one described below to test for differences among tissues). We

204 found a marginally significant interaction between the study and the expression

205 level in testis (our proxy for expression levels in the germline), driven by one

206 study (CHD cases; Homsy et al. 2015), as well as for interactions between the

207 studies and the effects of H3K9me3 and GC content, driven by two small studies

208 (EE and ID) (see Figure S1). Given these very minor differences and in order to

209 increase our power, we combined all the germline mutation datasets in what

210 follows (see Supplementary Materials Table S1 for list of mutations).

211

| Datasets | Trios | References | Capture | Sequencing |
|---|---|---|---|---|
| Autism Spectrum Disorder (ASD) | 3953 | De Rubeis et al. 2014; Iossifov et al. 2014 | Exome | Illumina and SOLiD |
| Simons Simplex Collection, unaffected | 1911 | Iossifov et al. 2014 | Exome | Illumina |
| Congenital heart disease (CHD) | 1213 | Homsy et al. 2015 | Exome | Illumina |
| Deciphering Developmental Disorders Study (DDD) | 1133 | The Deciphering Developmental Disorders Study 2015 | Exome | Illumina |
| Epileptic Encephalopathies (EE) | 264 | Epi4K Consortium and Epilepsy Phenome/Genome Project 2013 | Exome | Illumina |
| Intellectual Disability (ID) | 192 | de Ligt et al. 2012; Rauch et al. 2012; Hamdan et al. 2014 | Exome | Illumina |

212 **Table 1**. Summary of germline datasets

213

11

214   To examine determinants of mutation rates in somatic tissues, we downloaded

215   somatic mutation calls identified in four types of cancer from the Cancer Genome

216   Atlas (TCGA) portal (in July 2015): breast invasive carcinoma (BRCA), cervical

217   squamous cell carcinoma and endocervical adenocarcinoma (CESC), brain lower

218   grade glioma (LGG), and liver hepatocellular carcinoma (LIHC). The numbers of

219   samples are listed below (Table 2). In all cases, both non-cancerous and tumor

220   tissues of patients were sampled and the exomes were sequenced using an

221   Illumina platform. In the studies, mutation calls shared by the normal and tumor

222   samples were removed (on the presumption that they are germline). What

223   remains are somatic mutations found at high enough frequency to be seen in a

224   large population of cells, which are therefore likely to predate the tumorigenesis,

225   i.e., mutations that occurred in the pre-neoplastic tissues (Martincorena et al.

226   2015).

227

228   For each type of cancer with more than one mutation annotation file available in

229   the TCGA data portal, we selected the file that included the largest number of

230   patient samples. We removed the ~7.6% of samples that had an unusually large

231   number of mutations per sample ($p < 0.05$ by Tukey's test), because they are

232   likely to reflect loss of some aspect of the DNA mismatch repair and hence their

233   mutational mechanisms likely differ (Supek and Lehner 2015).

234

| Datasets | Sample sizes |
|---|---|
| Breast Invasive carcinoma (BRCA) | 904 |

12

| Cervical squamous cell carcinoma and endocervical adenocarcinoma (CESC) | 181 |
|---|---|
| Low grade glioma (LGG) | 502 |
| Liver hepatocellular carcinoma (LIHC) | 171 |

235    **Table 2.** Sizes of TCGA datasets

236

237    **Possible determinants of mutation rates.** We considered the main factors

238    previously reported to be significantly correlated with mutation rates, namely

239    expression levels, replication timing, GC content and histone modification levels.

240    To quantify expression levels, we relied on gene expression data (measured as

241    RPKM) from the Genotype-Tissue Expression (GTEx) for breast, uterus, brain

242    cortex and liver tissues. We used gene expression levels of testis and ovary as

243    our proxy for germline expression.

244

245    The effect of the replication timing on somatic mutation rates was argued to be

246    cell-type specific (Supek and Lehner 2015). We therefore relied on Repli-Seq

247    measurements (provided per base pair) in ENCODE cell lines that match the four

248    types of cancer, namely MCF-7 (breast cancer), Hela-S3 (cervical cancer), SK-N-

249    SH (neuroblastoma), and HepG2 (liver hepatocellular carcinoma) cell lines.

250    These measurements were obtained from the UCSC Genome Browser.  In all

251    cases, the replication timing reported is a smooth measure of the relative

252    enrichment of early vs. late S-phase nascent strands, with high values indicating

253    early replication. For each gene, we computed the average replication timing by

254    taking the mean value of the data points that overlap with gene start-to-end

255    coordinates in UCSC Refseq gene database. For genes with multiple transcripts,

256    we took the union of all exons in all transcripts.  For germline mutations, there

257    are no data for the appropriate cell types, so we used replicating timing estimates

258    for lymphoblastoid cell lines (LCL) (provided in 10 kb windows) (Koren et al.

259    2012). We also tried using replication timing data from three somatic tissues

260    instead; the replication timing data are highly correlated among the tissues and

261    therefore the effects of mutation were estimated to be very similar (see Figure

262    S2).

263

264    We also considered the effects of chromatin marks that had been shown to

265    correlate individually with somatic and germline mutation rates (Schuster-Böckler

266    and Lehner 2012; Carlson et al. 2017): specifically, histone modification

267    H3K9me3 and H3K27me3, two repressive marks associated with constitutively

268    and facultatively repressed genes, respectively. Levels of these marks were

269    downloaded from roadmap epigenomics data browser (Dec 2015, hg19) and

270    converted to gene-based histone modification levels by averaging across the

271    gene. We used the histone modification levels of adult ovary, breast

272    myoepithelial cells, brain hippocampus and adult liver as proxies for germline,

273    breast, brain and liver, respectively. In the following regression analysis, we

274    considered only three of four somatic tissues, as we could not obtain histone

275    modification data for CESC. Finally, we computed exonic GC content as the

276    fraction of G or C residues in the union of exons in all isoforms of a given gene.

277

278    Germline mutation studies relied on the UCSC Refseq gene annotation, whereas

279    TCGA uses GENECODE annotation, which contains more transcripts (Larsson et

280    al. 2005; Zhao and Zhang 2015). To make the comparison cleaner, we focused

281    on exonic regions considered in both types of studies by using gene and exon

282    coordinates of Refseq database in build hg19 from UCSC genome browser.

283

284    **Statistical model.** Our main goal was to investigate possible relationships

285    between mutation rates and gene expression levels, while controlling for

286    replication timing, GC content and some histone modification levels. Because our

287    mutation counts are over-dispersed, with greater variance than mean, we used a

288    negative binomial regression model (instead of, e.g., a Poisson regression

289    model). Specifically, for every protein-coding gene, we counted the number of

290    CpG Ti or other types of mutations in the coding exons of a gene and treated it

291    as an outcome of a sequence of independent Bernoulli trials with probability $\lambda_i$,

292    where $\lambda_i$ is the probability of a mutation occurring in gene i.

293

294    Transitions at CpG sites are thought to primarily occur due to spontaneous

295    deamination at methylated cytosines, a distinct mutational source, and thus their

296    determinants may be distinct from other mutation types (reviewed in Ségurel,

297    Wyman, and Przeworski 2014). However, within CpG islands, most CpGs are

298    hypomethylated (Takai and Jones 2002). To focus on a more homogeneous set

299    of methylated CpGs, we therefore excluded CpG islands from the analyses of

15

300    CpG Ti. CpG island annotations were downloaded from UCSC browser (track:

301    CpG Islands).

302

303    We considered gene expression levels measured in RPKM ($X_1$), replication

304    timing ($X_2$), mean histone modification levels (H3K9me3 as $X_3$, H3K27me3 as $X_4$)

305    and GC content ($X_5$) as predictors. We also included L, the total number of CpG

306    sites (when considering CpG Ti) or all nucleotides (when considering all other

307    types of mutations) in the exons of the given gene, as an exposure variable, to

308    account for the variation in gene length. The logarithm of $\lambda_i$ is then modeled as a

309    linear combination of these features scores:

310    $$\log(\lambda_i) = \beta_0 + \sum_{j=1}^{5} \beta_j X_{ij} + \log(L) + \varepsilon$$

311    We used R function glm.nb to estimate the coefficients, where $\beta_0$ is an intercept

312    term, $\beta_j$ is the effect size of feature j, and $X_{ij}$ is the score for feature j in gene i. In

313    order to make the effect sizes of different features comparable within a model,

314    we normalized all the predictor variables to have a mean of 0 and a standard

315    deviation of 1. The gene expression levels measured in RPKM originally range

316    from 0 to a few hundred thousand. As is standard (e.g., Green et al. 2003;

317    Francioli et al. 2015), we added half of the smallest non-zero value in the

318    corresponding expression data sets and then log-transformed the expression

319    level before normalization.

320

321    We note that in this model, we are considering possible effects one at a time.

322    Including interaction terms affects the estimates and significance levels but

16

323    changes none of the qualitative results, with the exception of results for

324    H3K27me3, which become less significant (see Figure S3).

325

326    To examine whether the predictors have significantly different effects across

327    tissues, we combined the models into one by including a categorical variable C

328    for the tissue type (see Figure 2). In this approach:

329              C = 1 for somatic tissues, C = 0 for germline;

$$\log(\lambda_{ij}) = \beta_0 + \sum_{j=1}^{5} \beta_j X_{ij} + C\left( \beta_6 + \sum_{j=7}^{11} \beta_j X_{ij} \right) + \log(L) + \varepsilon$$

330    $X_1$, $X_2$, $X_3$, $X_4$ and $X_5$ are the same genomic or epigenomic features as in the

331    separate model, $\beta_1$, $\beta_2$, $\beta_3$, $\beta_4$, $\beta_5$ are the effect sizes of features $X_1$ to $X_5$ for testis,

332    and $\beta_7$, $\beta_8$, $\beta_9$, $\beta_{10}$, $\beta_{11}$ are the differences of effect size in the somatic tissue of

333    features $X_1$ to $X_5$ compared to those in testis. We used the R function glm.nb to

334    estimate the coefficients.

335

336    Similarly, in order to ask whether effects differ between CpG Ti and other type of

337    mutations in the same tissue, we included a binary variable C for the two

338    mutation types (see Figure S4).

339              C = 1 for CpG Ti, C = 0 for all other mutations;

$$\log(\lambda_{ij}) = \beta_0 + \sum_{j=1}^{5} \beta_j X_{ij} + C\left( \beta_6 + \sum_{j=7}^{11} \beta_j X_{ij} \right) + \log(L) + \varepsilon$$

340    All variables are set up the same way as in the combined model described

341    previously, except for that $\beta_7$, $\beta_8$, $\beta_9$, $\beta_{10}$, $\beta_{11}$ are now the differences of the effect

342    sizes for CpG Ti compared to those for all other mutation types.

343

344    **Mutation spectrum and strand asymmetry analysis**. We annotated the

345    direction of transcription using the UCSC CCDS track and filtered out genes that

346    are transcribed off both strands (1.7% of genes in Refseq), which left around

347    19,000 genes to consider. This annotation allowed us to classify mutations into

348    six types of mutation (A>C, A>G, A>T, G>A, G>C, G>T) on either TS or NTS.

349    There are thus 12 possible changes (each of the six on both strands). We then

350    calculated the mutation rate of any given type on NTS and TS separately, by

351    considering the number of corresponding mutations in the combined data sets,

352    divided by the total number of nucleotides that could give rise to such a mutation

353    in the exons. To obtain the confidence intervals on the mutation rates (reported in

354    Figure 3, 4 and Figure S5) as well as for the mutation asymmetry ratio (Figure 4

355    and Figure S5), we used bootstrap. Specifically, we created 100 samples, of the

356    same size as the original sample, by drawing randomly from the original sample

357    with replacement, and estimated the 95% CI from those 100 samples.

358    We tested for strand asymmetry by a Chi-squared test. Because A>G strand

359    asymmetry shows the greatest asymmetry (Green et al. 2003) and is the only

360    mutation type that we found in all tissues (Figure 3), we focused primarily on this

361    type, though we also considered A>T mutational patterns (see Figure S5). To

362    test if the extent of strand asymmetry changes with transcription levels, we

363    grouped genes into expression level quantiles and calculated A>G strand

364    asymmetry. Our measure of strand asymmetry is the ratio of the mutation rate on

365    NTS to that on TS.

366

367    **Data availability.** Germline mutations are provided in Table S1. TCGA somatic

368    mutations can be downloaded from GDC data portal (https://gdc-

369    portal.nci.nih.gov/search/s?facetTab=cases). The gene RPKM data are available

370    at GTEx website (http://www.gtexportal.org/home/datasets). The replication

371    timing data of LCL and other tissues are available from (Koren et al. 2012) and

372    ENCODE website

373    (https://www.encodeproject.org/search/?type=Experiment&assay_title=Repli-seq)

374    respectively. The histone modification data can be freely accessed at epigenome

375    roadmap website (http://www.roadmapepigenomics.org/data/tables/all).

376

**Results**

377

378

379     We began by applying our multivariable regression model (see Materials and

380     Methods) to compare the determinants of mutation rates per gene between the

381     two germline tissues and among the three somatic tissues (Figure 1). Results for

382     germline mutations are very similar using testis or ovary expression profiles. In

383     both, there is no discernable effect of replication timing, other than a marginally

384     significant negative effect for mutations other than CpG Ti. However, in contrast

385     to a previous study using de novo mutations (Francioli et al. 2013) and most

386     previous studies of divergence, we found a significant increase of germline

387     mutation rates with expression levels for both CpG Ti and other mutation types

388     (Figure 1; see also Figure S2 for similar results with replication timing for different

389     tissues). The difference with a previous analysis of de novo mutations may be

390     due to the scale of a gene considered here (rather than 100 kb windows).

391

392     The effect of expression levels is most clearly seen using testis expression (P =

393     0.03 for CpG Ti; P = $1.4 \times 10^{-5}$ for other mutation types) than using ovary

394     expression, possibly due to the fact that over three quarters of germline

395     mutations are of male origin (Kong et al. 2012; Rahbari et al. 2016; Goldmann et

396     al. 2016). Alternatively, the ovary expression profile may be a poorer proxy for

397     female germ cells than the testis expression profile is for male germ cells. In any

398     case, henceforth, we use testis expression profile for analysis of the germline

399     mutation rates.

20

400

401      We note that our analysis of germline mutation relies on calls made in exome

402      studies of blood samples from six sets, including five cases and unaffected

403      controls (see Table 1). A previous study reported that in one set of cases,

404      individuals with congenital heart disease (CHD), there is an increased number of

405      putatively damaging mutations in the genes most highly expressed in the

406      developing heart and brain (Homsy et al. 2015). Since the mutations are thought

407      to be germline mutations (rather than somatic mutations), this association cannot

408      be causal, instead reflecting an enrichment of damaging mutations in important

409      heart developmental genes in CHD patients. To evaluate whether our findings of

410      increased mutation rates with germline expression levels could be driven by a

411      similar ascertainment bias, we excluded the CHD set and obtained the same

412      results (see Figure S6). We also reran the analysis, comparing the effects in the

413      five cases compared to the controls; none of the qualitative results differed,

414      though as expected from the smaller size of the control sets, the estimated effect

415      sizes were more uncertain (see Figure S7). Thus, our results suggest that the

416      increase in mutation rates with expression levels in testes is not a result of

417      focusing primarily on cases.

418

419      Germline mutation rates are also associated with H3K27me3 levels. We also

420      found that, other than for CpG Ti, mutation rates in a gene increase with its GC

421      content. This observation is consistent with previous findings of a high rate of GC

422      to AT mutations relative to other types (e.g., Kong et al. 2012). Moreover, it is

21

423    thought that mis-incorporated bases during DNA replication in an AT rich regions

424    are more easily accessible and thus more easily repaired than GC rich regions

425    (Petruska and Goodman 1985; Bloom et al. 1994). In contrast, we found a

426    marginally negative effect of GC content on germline rates of CpG Ti. A possible

427    explanation for this observation is that spontaneous deamination, the likely

428    source of most CpG Ti, occurs more readily when DNA is single stranded, which

429    is more likely in AT-rich than GC-rich regions (Fryxell and Moon 2005; Elango et

430    al. 2008).

431

432    The effects of determinants on mutation rates are also concordant across

433    somatic tissues. Notably, mutation rates decrease with expression levels in all

434    three tissues, though the magnitudes of the effects differ. This finding is

435    consistent with previous studies and thought to be a result of TCR (Lawrence et

436    al. 2013). Intriguingly, in a model comparing the effects on CpG Ti and other

437    mutation types directly, in all three somatic tissues, the effect of expression levels

438    on mutation rates is most pronounced for CpG Ti (see Figure S4). This finding

439    suggests that damage or repair of CpG Ti is tightly coupled to transcription.

440

441    In all three somatic tissues, there is also a decrease in mutation rate with

442    replication timing and H3K27me3 levels, as well as an increase with H3K9me3

443    levels (Schuster-Böckler and Lehner 2012; Behjati et al. 2014; Blokzijl et al.

444    2016). The effect of replicating timing on mutation rate has been attributed to the

445    depletion of free nucleotides within later replicating regions, leading to the

22

446     accumulation of single-stranded DNA and thus rendering the DNA more

447     susceptible to endogenous DNA damage (Stamatoyannopoulos et al. 2009). An

448     alternative hypothesis is that DNA mismatch repair (MMR), which is coupled with

449     replication, is more effective in the early replicating regions of the genome; this

450     possibility is supported by the finding that this association is not detected in the

451     tissue of MMR-deficient patients (Supek and Lehner 2015). While on face value,

452     it may seem surprising that replication timing is a significant determinant for the

453     LGG samples, given that neurons are post-mitotic, glial cells still retain their

454     ability to divide and a substantial fraction of mutations detected in neuronal

455     samples may have occurred at earlier stages in development.

456

457     The only difference in the determinants of mutation rates across somatic tissues

458     appears to be the effect of GC content on CpG Ti rates: mutation rates decrease

459     with GC content in brain tissues and increase with GC content in liver and breast

460     tissues. This finding raises the possibility that damage or repair rates of CpG

461     sites differ in brain tissues (Lodato et al. 2015).

462

463     Figure 1 also hints at a difference between testes (also ovaries) and somatic

464     tissues in the directional effect of expression levels on mutation rates, with a

465     marginally significant positive effect for germline mutations (P = 0.03 for CpG Ti,

466     P = $1.4 \times 10^{-5}$ for other mutation types) and a significantly negative effect for

467     somatic tissues (e.g., BRCA: P = $8 \times 10^{-16}$ for CpG Ti; P<$2 \times 10^{-16}$ for other mutation

468     types). When we tested for this difference explicitly, by adding a binary variable

23

469    for soma and germline (see Materials and Methods), we found that expression

470    levels and replication timing differ in their effects, for both CpG Ti and other

471    mutation types (Figure 2).

472

473    Specifically, replication timing has a positive effect on both tissue types but its

474    effect is stronger in soma (Figure 2). The simplest explanation is that a larger

475    fraction of mutations in the soma are introduced by errors related to replication,

476    as opposed to other non-replicative sources. Another (not mutually-exclusive)

477    possibility is that the effect of early replication versus late replication differs to a

478    greater extent in the soma than in the germline. For example, if MMR is much

479    more efficient in early replicating regions (Supek and Lehner 2015) and more

480    efficient in soma than germline.

481

482    To examine this possibility further, we considered a signature of TCR—strand

483    asymmetry—in the different tissues, finding it among germline mutations as well

484    as in all four somatic tissues (Figure 3). Consistent with previous studies (Green

485    et al. 2003; Francioli et al. 2015), one type in particular, A > G, stands out. While

486    the asymmetry is significant in all five data sets, with more mutation on the NTS

487    than the TS, the degree of asymmetry is significantly different among the five

488    data sets ($\chi^2$ test, P = 3x10$^{-8}$), with the strongest seen in germline. Intriguingly,

489    other mutation types, notably G>C mutations, show even more pronounced

490    differences among tissues, with significant excess on the transcribed strand in

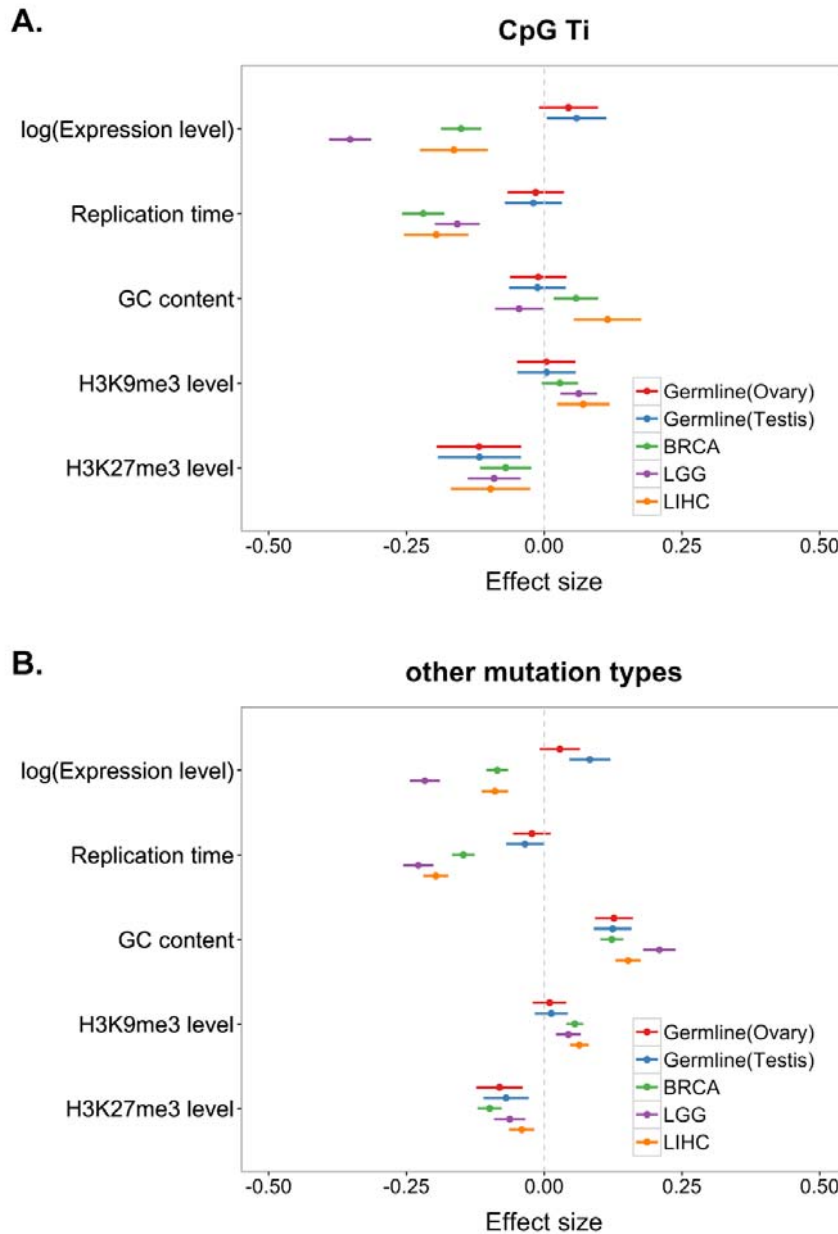491    the germline and LGG samples but a significant paucity on the NTS in BRCA and

24

492    CESC. These findings indicate a potential difference in either strand-biased

493    damage or in TCR (or both) among somatic tissues. In summary, the total

494    mutation rate appears to behave quite differently as a function of expression

495    levels in the germline and the soma (Figure 1 and 2), despite the fact that we

496    observed clear evidence for TCR in both types of tissues (Figure 3).

497

498    To examine this difference in more detail, we focused on A>G mutations and

499    considered how the mutation rate and degree of asymmetry covary with

500    expression (Figure 4). A striking contrast emerges: in the germline, as expression

501    levels increase, mutation rates and asymmetry increase, whereas in the soma,

502    asymmetry increases while mutation rates decrease. The same pattern is seen

503    when A>T mutation rate and asymmetry are considered (see Figure S5). This

504    difference in behavior with expression levels strongly suggests that the balance

505    between damage and repair rates during transcription differs between germline
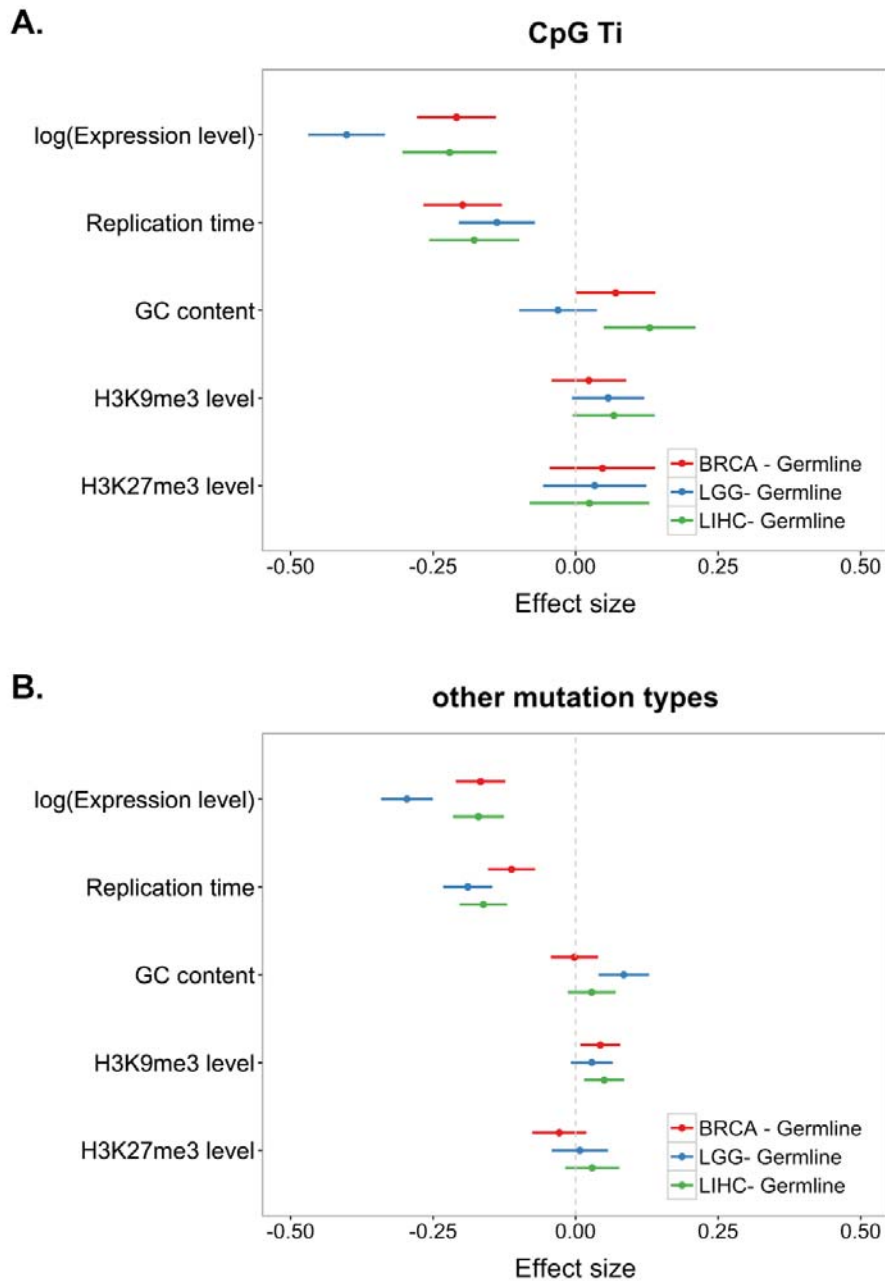
506    and soma.

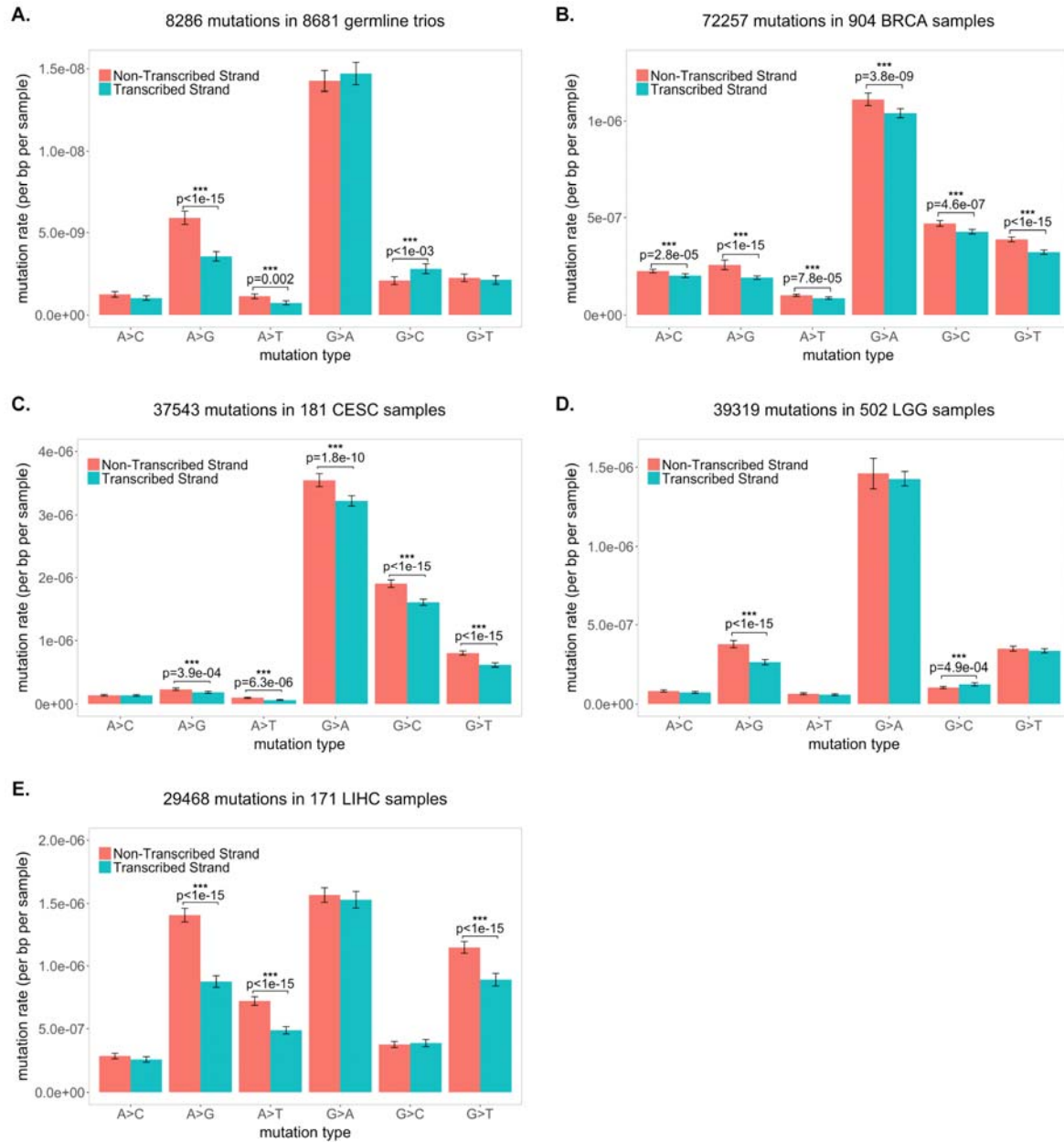507                                          **Figures**



508

509    **Figure 1.** Coefficients of multivariable binomial regression model fit to germline and somatic

510    mutation data. In panel A, are results for CpG Ti and in panel B, for other mutation types. Red,

511    blue and green, purple and orange bars represent the 95% CI for the estimate of the regression

512    coefficient in germline data set using ovary expression and testis expression, BRCA (breast

513    invasive carcinoma), LGG (brain lower grade glioma) and LIHC (liver hepatocellular carcinoma)

514    data sets respectively. For all replication timing data, high value means early.

515

**Figure 2.** Coefficients of combined model comparing each somatic data set to germline data set

using testis expression. In panel A, results for CpG Ti and in panel B, for other mutation types.

Red, blue and green bars represent the 95% CI of the deviation of the estimated coefficient from

the germline estimate; they are shown for BRCA (breast invasive carcinoma), LGG (brain lower

grade glioma) and LIHC (liver hepatocellular carcinoma) data sets respectively. For all replication

timing data, high value means early.

522

**Figure 3.** Strand asymmetry for six mutation types. In panel A are results for the germline; in panel B, for BRCA (breast invasive carcinoma); in panel C, for CESC (cervical squamous cell carcinoma and endocervic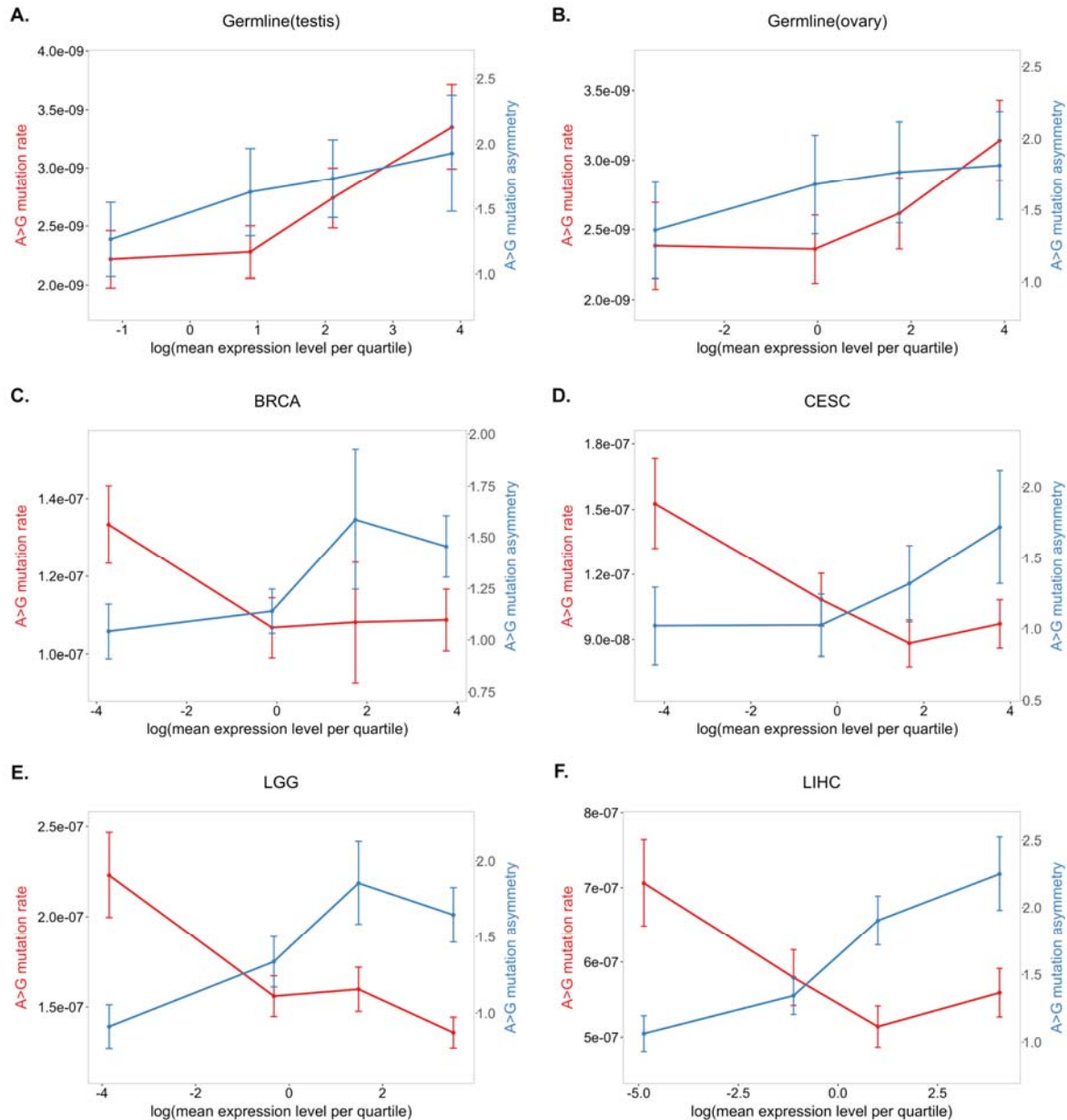al adenocarcinoma); in panel D, for LGG (brain lower grade glioma); and in panel E, for LIHC (liver hepatocellular carcinoma). The error bars of the mutation rate denote 95% confidence intervals estimated by bootstrapping (see Materials and Methods).

**Figure 4.** The degree of A>G strand asymmetry and the A>G mutation rate as a function of gene expression level quartiles. Shown are in panels A and B are results for the germline using testis expression levels and ovary expression levels, respectively; in panel C, for BRCA (breast invasive carcinoma); in panel D, for CESC (cervical squamous cell carcinoma and endocervical adenocarcinoma); in panel E, for LGG (brain lower grade glioma); and in panel F, for LIHC (liver hepatocellular carcinoma). The error bars for both the strand asymmetry and the mutation rate per quartile were estimated by bootstrapping (see Materials and Methods).

29

**Discussion**

537

538

539    We compared the determinants of mutation in the soma and the germline, using

540    the same unit of analysis (a coding region) and the same statistical model, and

541    applied it to similar exome data for germline de novo mutations and four types of

542    tumors, in which mutations largely predate tumorigenesis. We recapitulated

543    previous findings of the effects of GC content and of a histone mark indicative of

544    repression on germline and somatic mutations, as well as those of expression

545    levels and replicating time on somatic mutations (Schuster-Böckler and Lehner

546    2012; Lawrence et al. 2013). Strikingly, we also found clear differences in the

547    determinants of mutation rates between germline and soma, consistent with

548    earlier hints based on divergence data (Hodgkinson and Eyre-Walker 2011).

549    Notably, our results confirmed that somatic mutation rates decrease with

550    expression levels and reveal that, in sharp contrast, de novo germline mutation

551    rates increase with expression. This contrast suggests that transcription is

552    mutagenic in germline but not in soma, and that the DNA damage or repair

553    processes differ between them.

554

555    One limitation of our comparison—and of previous studies of germline and

556    somatic mutation—is the need to rely on proxies for determinants of interest,

557    such as replication timing data from cancer cell lines instead of normal cells. A

558    second limitation is that we considered only two types of mutations (CpG Ti and

559    other). Other work indicates that while these two types capture most of the

560    variation in mutation rates, the larger context (adjacent base pairs, but also

561    7mers) also impacts mutation rates (Hwang and Green 2004; Hodgkinson and

562    Eyre-Walker 2011; Aggarwala and Voight 2016). These different mutation

563    subtypes are likely affected somewhat differently by the determinants considered

564    here (Carlson et al. 2017). Despite these limitations, our work provides a

565    framework to contrast possible determinants of mutation rates in soma and

566    germline while controlling for some confounding effects, and results will only

567    improve as data sets increase and the measurements of salient genomic and

568    cellular features become more accurate. What is already clear is that the

569    divergent effect of expression on mutation rates in germline and soma is not

570    attributable to well-known covariates (included in our model). Moreover, the

571    differences that cannot readily be explained by the noise introduced by imperfect

572    proxies or limited data.

573

574    Notably, our results indicate that the tradeoff between damage and repair

575    associated with transcription must differ between germline and soma.

576    Transcription plausibly increases the rate of damage by opening up the DNA

577    helix, rendering the single strands more susceptible to mutagens (Polak and

578    Arndt 2008; Jinks-Robertson and Bhagwat 2014). One possibility is that, in the

579    germline, the rate of transcription-associated mutagenesis (TAM) swamps TCR,

580    leading to higher mutation rates with increased transcription, whereas in the

581    soma, TCR is relatively more efficient and the balance of TAM and TCR leads to

582    decreased mutagenesis with increased expression. Another possibility, which is

31

583    not mutually exclusive, is the presence of additional repair mechanisms in

584    somatic tissues. In support of this possibility, global genome repair (GGR) is

585    attenuated in differentiated cells, yet mutations on the NTS appear to

586    nonetheless be repaired efficiently (Nouspikel and Hanawalt 2000; Marteijn et al.

587    2014). This evidence led to the hypothesis of transcription-domain-associated

588    repair (DAR), which might repair damage on both strands in addition to TCR

589    (reviewed in Nouspikel 2007). From an evolutionary standpoint, the increased

590    efficiency of TCR relative to TAM in soma versus germline may be explained by

591    selection pressure on the repair of somatic tissues to prevent aging and cancer

592    (Lynch 2010).

593

594    Mounting evidence suggests that per cell division mutation rates differ across

595    tissues (Greenman et al. 2007; Lynch 2010; Alexandrov et al. 2013) and in

596    particular that they may be higher in early embryonic development than at other

597    stages of development (Ségurel, Wyman, and Przeworski 2014; Rahbari et al.

598    2016; Harland et al. 2016; Lindsay et al. 2016). This study further suggests that

599    at least part of the explanation may lie in the balance between damage and

600    repair, with TCR operating at different efficiencies relative to TAM or jointly with

601    other repair pathways, thereby maintaining low mutation rates in soma. As

602    mutation data from more tissues become available, it will be both feasible and

603    enlightening to examine tissue-specific differences in repair.

604

605                                    **Acknowledgments**

32

606    We thank Ziyue Gao and Priya Moorjani for comments on the manuscripts and

607    helpful discussions.

**References**

608

609     Aggarwala, Varun, and Benjamin F. Voight. 2016. "An Expanded Sequence

610          Context Model Broadly Explains Variability in Polymorphism Levels across

611          the Human Genome." *Nature Genetics* 48 (4): 349–55.

612          doi:10.1038/ng.3511.

613     Alexandrov, Ludmil B., Philip H. Jones, David C. Wedge, Julian E. Sale, Peter J.

614          Campbell, Serena Nik-Zainal, and Michael R. Stratton. 2015. "Clock-like

615          Mutational Processes in Human Somatic Cells." *Nature Genetics* 47 (12):

616          1402–7. doi:10.1038/ng.3441.

617     Alexandrov, Ludmil, Nik-Zainal Serena, David Wedge, Samuel Aparicio, Sam

618          Behjati, Andrew Biankin, Graham Bignell, et al. 2013. "Signatures of

619          Mutational Processes in Human Cancer." *Nature* 500 (7463): 415–21.

620          doi:10.1038/nature12477.

621     Behjati, Sam, Meritxell Huch, Ruben van Boxtel, Wouter Karthaus, David C.

622          Wedge, Asif U. Tamuri, Iñigo Martincorena, et al. 2014. "Genome

623          Sequencing of Normal Cells Reveals Developmental Lineages and

624          Mutational Processes." *Nature* 513 (7518): 422–25.

625          doi:10.1038/nature13448.

626     Besenbacher, Søren, Patrick Sulem, Agnar Helgason, Hannes Helgason, Helgi

627          Kristjansson, Aslaug Jonasdottir, Adalbjorg Jonasdottir, et al. 2016. "Multi-

628          Nucleotide de Novo Mutations in Humans." *PLOS Genetics* 12 (11):

629          e1006315. doi:10.1371/journal.pgen.1006315.

630   Blokzijl, Francis, Joep de Ligt, Myrthe Jager, Valentina Sasselli, Sophie Roerink,

631       Nobuo Sasaki, Meritxell Huch, et al. 2016. "Tissue-Specific Mutation

632       Accumulation in Human Adult Stem Cells during Life." *Nature* 538 (7624):

633       260–64. doi:10.1038/nature19768.

634   Bloom, L. B., M. R. Otto, R. Eritja, L. J. Reha-Krantz, M. F. Goodman, and J. M.

635       Beechem. 1994. "Pre-Steady-State Kinetic Analysis of Sequence-

636       Dependent Nucleotide Excision by the 3'-exonuclease Activity of

637       Bacteriophage T4 DNA Polymerase." *Biochemistry* 33 (24): 7576–86.

638   Brennan, Cameron W, Roel Verhaak, McKenna Aaron, Benito Campos, Houtan

639       Noushmehr, Sofie R Salama, Siyuan Zheng, et al. 2013. "The Somatic

640       Genomic Landscape of Glioblastoma." 155 (2): 462–77.

641       doi:10.1016/j.cell.2013.09.034.

642   Campbell, Catarina D, and Evan E Eichler. 2013. "Properties and Rates of

643       Germline Mutations in Humans" 29 (10): 575–84.

644       doi:10.1016/j.tig.2013.04.005.

645   Carlson, Jedidiah, Laura J. Scott, Adam E. Locke, Matthew Flickinger, Shawn

646       Levy, The BRIDGES Consortium, Richard M. Myers, et al. 2017.

647       "Extremely Rare Variants Reveal Patterns of Germline Mutation Rate

648       Heterogeneity in Humans." *bioRxiv*, February, 108290.

649       doi:10.1101/108290.

650   De Rubeis, Silvia, Xin He, Arthur P. Goldberg, Christopher S. Poultney, Kaitlin

651       Samocha, A. Ercument Cicek, Yan Kou, et al. 2014. "Synaptic,

652    Transcriptional and Chromatin Genes Disrupted in Autism." *Nature* 515

653    (7526): 209–15. doi:10.1038/nature13772.

654 Duret, Laurent, and Nicolas Galtier. 2009. "Biased Gene Conversion and the

655    Evolution of Mammalian Genomic Landscapes." *Annual Review of*

656    *Genomics and Human Genetics* 10 (1): 285–311. doi:10.1146/annurev-

657    genom-082908-150001.

658 Elango, Navin, Seong-Ho Kim, NISC Comparative Sequencing Program, Eric

659    Vigoda, and Soojin V. Yi. 2008. "Mutations of Different Molecular Origins

660    Exhibit Contrasting Patterns of Regional Substitution Rate Variation."

661    *PLOS Computational Biology* 4 (2): e1000015.

662    doi:10.1371/journal.pcbi.1000015.

663 Epi4K Consortium, and Epilepsy Phenome/Genome Project. 2013. "De Novo

664    Mutations in Epileptic Encephalopathies." *Nature* 501 (7466): 217–21.

665    doi:10.1038/nature12439.

666 Francioli, Laurent C., Paz P. Polak, Amnon Koren, Androniki Menelaou, Sung

667    Chun, Ivo Renkens, Genome of the Netherlands Consortium, et al. 2015.

668    "Genome-Wide Patterns and Properties of de Novo Mutations in Humans."

669    *Nature Genetics* 47 (7): 822–26. doi:10.1038/ng.3292.

670 Fryxell, Karl J., and Won-Jong Moon. 2005. "CpG Mutation Rates in the Human

671    Genome Are Highly Dependent on Local GC Content." *Molecular Biology*

672    *and Evolution* 22 (3): 650–58. doi:10.1093/molbev/msi043.

673    Gao, Ziyue, Minyoung J. Wyman, Guy Sella, and Molly Przeworski. 2016.

674         "Interpreting the Dependence of Mutation Rates on Age and Time." *PLOS*

675         *Biology* 14 (1): e1002355. doi:10.1371/journal.pbio.1002355.

676    Goldmann, Jakob M, Wendy SW Wong, Michele Pinelli, Terry Farrah, Dale

677         Bodian, Anna B Stittrich, Gustavo Glusman, et al. 2016. "Parent-of-Origin-

678         Specific Signatures of de Novo Mutations" 48 (8): 935–39.

679         doi:10.1038/ng.3597.

680    Green, Phil, Brent Ewing, Webb Miller, Pamela J. Thomas, NISC Comparative

681         Sequencing Program, and Eric D. Green. 2003. "Transcription-Associated

682         Mutational Asymmetry in Mammalian Evolution." *Nature Genetics* 33 (4):

683         514–17. doi:10.1038/ng1103.

684    Greenman, Christopher, Philip Stephens, Raffaella Smith, Gillian L Dalgliesh,

685         Christopher Hunter, Graham Bignell, Helen Davies, et al. 2007. "Patterns

686         of Somatic Mutation in Human Cancer Genomes." *Cah Rev The* 446

687         (7132): 153–58. doi:10.1038/nature05610.

688    Hamdan, Fadi F., Myriam Srour, Jose-Mario Capo-Chichi, Hussein Daoud,

689         Christina Nassif, Lysanne Patry, Christine Massicotte, et al. 2014. "De

690         Novo Mutations in Moderate or Severe Intellectual Disability." *PLOS*

691         *Genetics* 10 (10): e1004772. doi:10.1371/journal.pgen.1004772.

692    Hanawalt, Philip C, and Graciela Spivak. 2008. "Transcription-Coupled DNA

693         Repair: Two Decades of Progress and Surprises" 9 (12): 958–70.

694         doi:10.1038/nrm2549.

695    Harland, Chad, Carole Charlier, Latifa Karim, Nadine Cambisano, Manon

696        Deckers, Erik Mullaart, Wouter Coppieters, and Michel Georges. 2016.

697        "Frequency of Mosaicism Points towards Mutation-Prone Early Cleavage

698        Cell Divisions." *bioRxiv*, October, 79863. doi:10.1101/079863.

699    Hodgkinson, Alan, Ying Chen, and Adam Eyre-Walker. 2012. "The Large□scale

700        Distribution of Somatic Mutations in Cancer Genomes" 33 (1): 136–43.

701        doi:10.1002/humu.21616.

702    Hodgkinson, Alan, and Adam Eyre-Walker. 2011. "Variation in the Mutation Rate

703        across Mammalian Genomes." *Nat Rev Genetics* 12 (11): 756–66.

704        doi:10.1038/nrg3098.

705    Homsy, Jason, Samir Zaidi, Yufeng Shen, James S. Ware, Kaitlin E. Samocha,

706        Konrad J. Karczewski, Steven R. DePalma, et al. 2015. "De Novo

707        Mutations in Congenital Heart Disease with Neurodevelopmental and

708        Other Congenital Anomalies." *Science* 350 (6265): 1262–66.

709        doi:10.1126/science.aac9396.

710    Hwang, Dick G., and Phil Green. 2004. "Bayesian Markov Chain Monte Carlo

711        Sequence Analysis Reveals Varying Neutral Substitution Patterns in

712        Mammalian Evolution." *Proceedings of the National Academy of Sciences*

713        *of the United States of America* 101 (39): 13994–1.

714        doi:10.1073/pnas.0404142101.

715    Iossifov, Ivan, Brian J. O'Roak, Stephan J. Sanders, Michael Ronemus, Niklas

716        Krumm, Dan Levy, Holly A. Stessman, et al. 2014. "The Contribution of de

717    Novo Coding Mutations to Autism Spectrum Disorder." *Nature* 515 (7526):

718    216–21. doi:10.1038/nature13908.

719 Jinks-Robertson, Sue, and Ashok S. Bhagwat. 2014. "Transcription-Associated

720    Mutagenesis." *Annual Review of Genetics* 48 (1): 341–59.

721    doi:10.1146/annurev-genet-120213-092015.

722 Kong, Augustine, Michael L. Frigge, Gisli Masson, Soren Besenbacher, Patrick

723    Sulem, Gisli Magnusson, Sigurjon A. Gudjonsson, et al. 2012. "Rate of de

724    Novo Mutations and the Importance of Father/'s Age to Disease Risk."

725    *Nature* 488 (7412): 471–75. doi:10.1038/nature11396.

726 Koren, Amnon, Paz Polak, James Nemesh, Jacob J. Michaelson, Jonathan

727    Sebat, Shamil R. Sunyaev, and Steven A. McCarroll. 2012. "Differential

728    Relationship of DNA Replication Timing to Different Forms of Human

729    Mutation and Variation." *The American Journal of Human Genetics* 91 (6):

730    1033–40. doi:10.1016/j.ajhg.2012.10.018.

731 Larsson, Thomas P., Christian G. Murray, Tobias Hill, Robert Fredriksson, and

732    Helgi B. Schiöth. 2005. "Comparison of the Current RefSeq, Ensembl and

733    EST Databases for Counting Genes and Gene Discovery." *FEBS Letters*

734    579 (3): 690–98. doi:10.1016/j.febslet.2004.12.046.

735 Lawrence, Michael S, Petar Stojanov, Paz Polak, Gregory V Kryukov, Kristian

736    Cibulskis, Andrey Sivachenko, Scott L Carter, et al. 2013. "Mutational

737    Heterogeneity in Cancer and the Search for New Cancer-Associated

738    Genes." *Nature* 499 (7457): 214–18. doi:10.1038/nature12213.

739     Lee, William, Zhaoshi Jiang, Jinfeng Liu, Peter M. Haverty, Yinghui Guan,

740          Jeremy Stinson, Peng Yue, et al. 2010. "The Mutation Spectrum Revealed

741          by Paired Genome Sequences from a Lung Cancer Patient." *Nature* 465

742          (7297): 473–77. doi:10.1038/nature09004.

743     Ligt, Joep de, Marjolein H. Willemsen, Bregje W.M. van Bon, Tjitske Kleefstra,

744          Helger G. Yntema, Thessa Kroes, Anneke T. Vulto-van Silfhout, et al.

745          2012. "Diagnostic Exome Sequencing in Persons with Severe Intellectual

746          Disability." *New England Journal of Medicine* 367 (20): 1921–29.

747          doi:10.1056/NEJMoa1206524.

748     Lindsay, Sarah J., Raheleh Rahbari, Joanna Kaplanis, Thomas Keane, and

749          Matthew Hurles. 2016. "Striking Differences in Patterns of Germline

750          Mutation between Mice and Humans." *bioRxiv*, October, 82297.

751          doi:10.1101/082297.

752     Lodato, Michael A, Mollie B Woodworth, Semin Lee, Gilad D Evrony, Bhaven K

753          Mehta, Amir Karger, Soohyun Lee, et al. 2015. "Somatic Mutation in

754          Single Human Neurons Tracks Developmental and Transcriptional

755          History" 350 (6256): 94–98. doi:10.1126/science.aab1785.

756     Lynch, Michael. 2010. "Rate, Molecular Spectrum, and Consequences of Human

757          Mutation." *Proc National Acad Sci* 107 (3): 961–68.

758          doi:10.1073/pnas.0912629107.

759     Marteijn, Jurgen A., Hannes Lans, Wim Vermeulen, and Jan H. J. Hoeijmakers.

760          2014. "Understanding Nucleotide Excision Repair and Its Roles in Cancer

761 and Ageing." *Nature Reviews Molecular Cell Biology* 15 (7): 465–81.

762 doi:10.1038/nrm3822.

763 Martincorena, Iñigo, Amit Roshan, Moritz Gerstung, Peter Ellis, Peter Loo,

764 McLaren Stuart, David C Wedge, et al. 2015. "High Burden and Pervasive

765 Positive Selection of Somatic Mutations in Normal Human Skin" 348

766 (6237): 880–86. doi:10.1126/science.aaa6806.

767 McVicker, Graham, David Gordon, Colleen Davis, and Phil Green. 2009.

768 "Widespread Genomic Signatures of Natural Selection in Hominid

769 Evolution." *PLOS Genetics* 5 (5): e1000471.

770 doi:10.1371/journal.pgen.1000471.

771 McVicker, Graham, and Phil Green. 2010. "Genomic Signatures of Germline

772 Gene Expression." *Genome Research* 20 (11): 1503–11.

773 doi:10.1101/gr.106666.110.

774 Michaelson, Jacob J, Yujian Shi, Madhusudan Gujral, Hancheng Zheng, Dheeraj

775 Malhotra, Xin Jin, Minghan Jian, et al. 2012. "Whole-Genome Sequencing

776 in Autism Identifies Hot Spots for De Novo Germline Mutation" 151 (7):

777 1431–42. doi:10.1016/j.cell.2012.11.019.

778 Moorjani, Priya, Ziyue Gao, and Molly Przeworski. 2016. "Human Germline

779 Mutation and the Erratic Evolutionary Clock." *PLOS Biology* 14 (10):

780 e2000744. doi:10.1371/journal.pbio.2000744.

781 Muller, H. J. 1927. "ARTIFICIAL TRANSMUTATION OF THE GENE." *Science* 66

782 (1699): 84–87. doi:10.1126/science.66.1699.84.

783    Nouspikel, Thierry. 2007. "DNA Repair in Differentiated Cells: Some New

784          Answers to Old Questions." *Neuroscience*, Genome Dynamics and DNA

785          Repair in the CNS, 145 (4): 1213–21.

786          doi:10.1016/j.neuroscience.2006.07.006.

787    Nouspikel, Thierry. 2009. "DNA Repair in Mammalian Cells." *Cellular and*

788          *Molecular Life Sciences* 66 (6): 994–1009. doi:10.1007/s00018-009-8737-

789          y.

790    Nouspikel, Thierry, and Philip C. Hanawalt. 2000. "Terminally Differentiated

791          Human Neurons Repair Transcribed Genes but Display Attenuated Global

792          DNA Repair and Modulation of Repair Gene Expression." *Molecular and*

793          *Cellular Biology* 20 (5): 1562–70.

794    Park, Chungoo, Wenfeng Qian, and Jianzhi Zhang. 2012. "Genomic Evidence for

795          Elevated Mutation Rates in Highly Expressed Genes." *EMBO Reports* 13

796          (12): 1123–29. doi:10.1038/embor.2012.165.

797    Petruska, J., and M. F. Goodman. 1985. "Influence of Neighboring Bases on

798          DNA Polymerase Insertion and Proofreading Fidelity." *Journal of*

799          *Biological Chemistry* 260 (12): 7533–39.

800    Pleasance, Erin D., R. Keira Cheetham, Philip J. Stephens, David J. McBride,

801          Sean J. Humphray, Chris D. Greenman, Ignacio Varela, et al. 2010. "A

802          Comprehensive Catalogue of Somatic Mutations from a Human Cancer

803          Genome." *Nature* 463 (7278): 191–96. doi:10.1038/nature08658.

804    Pleasance, Erin D., Philip J. Stephens, Sarah O'Meara, David J. McBride, Alison

805          Meynert, David Jones, Meng-Lay Lin, et al. 2010. "A Small-Cell Lung

806    Cancer Genome with Complex Signatures of Tobacco Exposure." *Nature*

807    463 (7278): 184–90. doi:10.1038/nature08629.

808 Polak, Paz, and Peter F. Arndt. 2008. "Transcription Induces Strand-Specific

809    Mutations at the 5′ End of Human Genes." *Genome Research* 18 (8):

810    1216–23. doi:10.1101/gr.076570.108.

811 Polak, Paz, Rosa Karlić, Amnon Koren, Robert Thurman, Richard Sandstrom,

812    Michael Lawrence, Alex Reynolds, et al. 2015. "Cell-of-Origin Chromatin

813    Organization Shapes the Mutational Landscape of Cancer." *Nature* 518

814    (7539): 360–64. doi:10.1038/nature14221.

815 Rahbari, Raheleh, Arthur Wuster, Sarah Lindsay, Robert Hardwick, Ludmil

816    Alexandrov, Saeed Turki, Anna Dominiczak, et al. 2016. "Timing, Rates

817    and Spectra of Human Germline Mutation." *Nat Genet* 48 (2): 126–33.

818    doi:10.1038/ng.3469.

819 Rauch, Anita, Dagmar Wieczorek, Elisabeth Graf, Thomas Wieland, Sabine

820    Endele, Thomas Schwarzmayr, Beate Albrecht, et al. 2012. "Range of

821    Genetic Mutations Associated with Severe Non-Syndromic Sporadic

822    Intellectual Disability: An Exome Sequencing Study." *The Lancet* 380

823    (9854): 1674–82. doi:10.1016/S0140-6736(12)61480-9.

824 Rubin, Alan F., and Phil Green. 2009. "Mutation Patterns in Cancer Genomes."

825    *Proceedings of the National Academy of Sciences* 106 (51): 21766–70.

826    doi:10.1073/pnas.0912499106.

827 Samocha, Kaitlin E, Elise B Robinson, Stephan J Sanders, Christine Stevens,

828    Aniko Sabo, Lauren M McGrath, Jack A Kosmicki, et al. 2014. "A

829    Framework for the Interpretation of de Novo Mutation in Human Disease."

830    *Nature Genetics* 46 (9): 944–50. doi:10.1038/ng.3050.

831 Schuster-Böckler, Benjamin, and Ben Lehner. 2012. "Chromatin Organization Is

832    a Major Influence on Regional Mutation Rates in Human Cancer Cells"

833    488 (7412): 504–7. doi:10.1038/nature11273.

834 Ségurel, Laure, Minyoung J. Wyman, and Molly Przeworski. 2014. "Determinants

835    of Mutation Rate Variation in the Human Germline." *Annual Review of*

836    *Genomics and Human Genetics* 15: 47–70. doi:10.1146/annurev-genom-

837    031714-125740.

838 Shendure, Jay, and Joshua M. Akey. 2015. "The Origins, Determinants, and

839    Consequences of Human Mutations." *Science* 349 (6255): 1478–83.

840    doi:10.1126/science.aaa9119.

841 Stamatoyannopoulos, John A., Ivan Adzhubei, Robert E. Thurman, Gregory V.

842    Kryukov, Sergei M. Mirkin, and Shamil R. Sunyaev. 2009. "Human

843    Mutation Rate Associated with DNA Replication Timing." *Nature Genetics*

844    41 (4): 393–95. doi:10.1038/ng.363.

845 Stratton, Michael R. 2011. "Exploring the Genomes of Cancer Cells: Progress

846    and Promise." *Science* 331 (6024): 1553–58.

847    doi:10.1126/science.1204040.

848 Stratton, Michael R., Peter J. Campbell, and P. Andrew Futreal. 2009. "The

849    Cancer Genome." *Nature* 458 (7239): 719–24. doi:10.1038/nature07943.

850   Supek, Fran, and Ben Lehner. 2015. "Differential DNA Mismatch Repair

851        Underlies Mutation Rate Variation across the Human Genome" 521

852        (7550): 81–84. doi:10.1038/nature14173.

853   Takai, Daiya, and Peter A. Jones. 2002. "Comprehensive Analysis of CpG

854        Islands in Human Chromosomes 21 and 22." *Proceedings of the National*

855        *Academy of Sciences* 99 (6): 3740–45. doi:10.1073/pnas.052410099.

856   The Deciphering Developmental Disorders Study. 2015. "Large-Scale Discovery

857        of Novel Genetic Causes of Developmental Disorders." *Nature* 519 (7542):

858        223–28. doi:10.1038/nature14135.

859   Webster, Matthew T., Nick G. C. Smith, Martin J. Lercher, and Hans Ellegren.

860        2004. "Gene Expression, Synteny, and Local Similarity in Human

861        Noncoding Mutation Rates." *Molecular Biology and Evolution* 21 (10):

862        1820–30. doi:10.1093/molbev/msh181.

863   Zhao, Shanrong, and Baohong Zhang. 2015. "A Comprehensive Evaluation of

864        Ensembl, RefSeq, and UCSC Annotations in the Context of RNA-Seq

865        Read Mapping and Gene Quantification." *BMC Genomics* 16 (1): 97.

866        doi:10.1186/s12864-015-1308-8.