

1 **Mutation-Profile-Based Methods for Understanding Selection Forces in Cancer**

2 **Somatic Mutations: A Comparative Analysis**

3

4 Zhan Zhou^{1,4,†}, Yangyun Zou^{2,†}, Gangbiao Liu¹, Jingqi Zhou¹, Jingcheng Wu⁴,

5 Shimin Zhao¹, Zhixi Su^{2,*}, Xun Gu^{3,2,*}

6 ¹State Key Laboratory of Genetic Engineering, School of Life Sciences, Fudan
7 University, Shanghai, China,

8 ²Ministry of Education Key Laboratory of Contemporary Anthropology, Center for
9 Evolutionary Biology, School of Life Sciences, Fudan University, Shanghai, China,

10 ³Department of Genetics, Development and Cell Biology, Program of Bioinformatics
11 and Computational Biology, Iowa State University, Ames, Iowa, USA,

12 ⁴College of Pharmaceutical Sciences, Zhejiang University, Hangzhou, Zhejiang,
13 China

14

15 [†]These authors contributed equally to this work

16 *Corresponding authors

17 Xun Gu: xgu@iastate.edu

18 Zhixi Su: zxsu@fudan.edu.cn

19

20

21 **Abstract**

22 Human genes perform different functions and exhibit different effects on fitness
23 in cancer and normal cell populations. Here, we present an evolutionary approach to
24 measuring the selective pressure on human genes in both cancer and normal cell
25 genomes using the well-known dN/dS (nonsynonymous to synonymous substitution
26 rate) ratio. We develop a new method called the mutation-profile-based Nei-Gojobori
27 (mpNG) method, which applies sample-specific nucleotide substitution profiles
28 instead of conventional substitution models to calculating dN/dS ratios in cancer and
29 normal populations. Compared with previous studies that focused on positively
30 selected genes in cancer genomes, which potentially represent the driving force
31 behind tumor initiation and development, we employed an alternative approach to
32 identifying cancer-constrained genes that strengthen negative selection pressure in
33 tumor cells. In cancer cells, we found a conservative estimate of 45 genes under
34 intensified positive selection and 16 genes under strengthened purifying selection
35 relative to germline cells. The cancer-specific positively selected genes were enriched
36 for cancer genes and human essential genes, while several cancer-specific negatively
37 selected genes were previously reported as prognostic biomarkers for cancers. Thus,
38 our computation pipeline used to identify positively and negatively selected genes in
39 cancer may provide useful information for understanding the evolution of cancer
40 somatic mutations.

41

42 **Keywords:** Cancer somatic mutations, Natural selection, Cancer essential gene,
43 Evolution

44 **Introduction**

45 Since the pioneering work of Cairns and Nowell^{1,2}, the evolutionary concept of
 46 cancer progression has been widely accepted³⁻⁷. In this model, cancer cells evolve
 47 through random somatic mutations and epigenetic changes that may alter several
 48 crucial pathways, a process that is followed by the clonal selection of the resulting
 49 cells. Consequently, cancer cells can survive and proliferate under deleterious
 50 circumstances^{8,9}. Therefore, knowledge of evolutionary dynamics will benefit our
 51 understanding of cancer initiation and progression. For example, there are two types
 52 of somatic mutations in cancer genomes: driver mutations and passenger
 53 mutations^{10,11}. Driver mutations are those that confer a selective advantage on cancer
 54 cells, as indicated by statistical evidence of positive selection. However, some
 55 passenger mutations undergo purifying selection because they would have potentially
 56 deleterious effects on cancer cells^{12,13}. Between these two extremes are passenger
 57 mutations that are usually considered to be neutral in cancer.

58 Analyses of large-scale cancer somatic mutation data have revealed that the effects
 59 of positive selection are much stronger on cancer cells than on germline cells^{14,15}.
 60 Given that many of the genes positively selected for in tumor development act as the
 61 driving force behind tumor initiation and progression, it is understandable that almost
 62 all previous studies focused on the positively selected genes in cancer genomes^{3,16-19}.
 63 We have realized that an alternative approach, i.e., identifying cancer-constrained
 64 genes that are highly conserved in tumor cell populations (under purifying selection),
 65 is also valuable. Because essential genes are more evolutionarily conserved²⁰, it

would be feasible to identify cancer essential genes from the genes that are evolutionarily conserved in cancer cells. Because cancer essential genes may not be the driver genes for carcinogenesis but are crucial for cancer cell proliferation and survival²¹, using evolutionary conservation to identify relevant genes may be advantageous in addressing therapeutic issues related to drug resistance, especially in cancers with high intratumor heterogeneity.

Many previous studies have used the ratio of nonsynonymous to synonymous substitution rates (dN/dS) to identify genes that might be under strong positive selection both in organismal evolution and tumorigenesis^{11,14,15,22-24}. However, most of these studies applied well-known methods that are usually based on simple nucleotide mutation/substitution models, where every mutation or substitution pattern has the same probability²⁵. Unfortunately, this may not be a realistic biological model because many recent cancer genomics studies have shown that mutation profiles are quite different between different cancer samples^{15,26}. In addition, context-dependent mutation bias (i.e., base-substitution profiles that consider the flanking 5' and 3' bases of each mutated base) should be taken into consideration^{26,27}.

In this study, we describe a new method, called the mutation-profile-based Nei-Gojobori (mpNG) method, to estimate the selective constraint in cancer somatic mutations. Simply stated, the mpNG method removes an unrealistic assumption inherent in the original NG method (named NG86), wherein each type of nucleotide change has the same mutation rate²⁵. This assumption can lead to nontrivial biased estimations when it is significantly violated. In contrast, mpNG implements an

empirical nucleotide mutation model that simultaneously takes into account several factors, including single-base mutation patterns, local-specific effects of surrounding DNA regions, and tissue/cancer types. Using 7,042 tumor-normal paired whole-exome sequences (WESs), as well as rare germline variations from 6,500 exome sequences (ESP6500) as references, we used the mpNG method to identify the selective constraint of human genes in cancer cells. The potential for our computational pipeline to identify cancer-constrained genes may provide useful information for identifying promising drug targets or prognostic biomarkers.

Results

The mutation profiles of cancer genomes and human populations are different

Estimating evolutionary selective pressure on human genes is a practical method for inferring the functional importance of genes to a specific population. By comparing selective pressures on genes in cancer cell populations with those in normal cell populations, we may be able to identify different functional and fitness effects of human genes in cancer and normal cells. The conventional method for measuring selective pressure is to calculate the dN/dS ratio using the NG86 method²⁵, which assumes equal substitution rates among different nucleotides. In our study, we used the cancer somatic mutations from 7,042 tumor-normal pairs, as well as rare variations from 6,500 exome sequences from the National Heart, Lung, and Blood Institute (NHLBI) Grant Opportunity (GO) Exome Sequencing Project (ESP6500), as a reference. We used these data to compare the relative mutation probabilities from

110 cancer somatic mutations and germline substitutions for all possible base substitutions,
 111 considering the identities of the bases immediately 5' and 3' of each mutated base. We
 112 then depicted the mutation profiles as 96 substitution classifications^{26,27}. The mutation
 113 profiles exhibit the prevalence of each substitution pattern for somatic point mutations,
 114 which present not only the substitution types but also the sequence context (see
 115 Methods). The exonic mutation profiles of cancer somatic substitutions and germline
 116 substitutions differed from one another, and the intronic and intergenic mutation
 117 profiles were quite different from the exonic mutation profile of cancer cells (Fig. 1).
 118 We also calculated the exonic mutation profiles of four different cancer types: colon
 119 adenocarcinoma (COAD), lung adenocarcinoma (LUAD), skin cutaneous melanoma
 120 (SKCM), and breast carcinoma (BRCA). These cancer types varied considerably not
 121 only in their mutation rates but also in their mutation patterns. Specifically, the
 122 mutation rate of SKCM was much higher than that of the other three types.
 123 Additionally, the mutation profiles of SKCM were highly enriched in the C-to-T
 124 substitution pattern (Fig. 1). These data indicated a direct mutagenic role for
 125 ultraviolet (UV) light in SKCM pathogenesis²⁸. The different mutation profiles may
 126 lead to different biological progressions in carcinogenesis, which have been depicted
 127 in several publications^{17,26}. Thus, it is inappropriate to use conventional methods such
 128 as NG86²⁵ to measure selective pressure by means of dN/dS calculation because this
 129 approach ignores the mutation bias of different nucleotide substitution types.

130

131 **Measuring selective pressure on human genes in cancer and germline cells using**

132 **the mpNG method**

133 We therefore formulated an evolutionary approach that was designed specifically
134 to estimate the selective pressure imposed on human genes in cancer cells and then to
135 identify genes that had undergone positive and purifying selection in cancer cells
136 rather than in normal cells (see Fig. 2 for illustration). We developed the mpNG
137 method to estimate the dN/dS ratio of each human gene based on the mutation profiles
138 of cancer somatic mutations and germline substitutions. In contrast to the NG86
139 method²⁵, our method considered the difference in substitution rate and took the
140 overall mutation profile as the weight matrix (Fig. 1).

141 We calculated the expected number of nonsynonymous and synonymous sites
142 based on the exonic mutation profiles. We then counted the number of
143 nonsynonymous and synonymous substitutions in the protein-coding region of each
144 human gene for all cancer somatic mutations or germline substitutions. A χ^2 test was
145 performed to identify the genes whose dN/dS values were either significantly greater
146 than one or less than one, which indicates positive or negative (purifying) selection,
147 respectively. Of the 18,602 genes with at least one germline substitution and cancer
148 somatic substitution, the overall dN/dS value for cancer somatic substitutions
149 (mean \pm s.e.=1.367 \pm 0.009) was much greater than that of germline substitutions
150 (mean \pm s.e.=0.903 \pm 0.006) (Wilcoxon test, $P < 10^{-16}$) (Table 1, Supplementary Table S1).
151 In the cancer genomes, 1,230 genes had dN/dS values significantly greater than one,
152 and 326 genes had dN/dS values significantly less than one (χ^2 test, $P < 0.05$). In
153 contrast, the germline substitutions included only 306 genes with dN/dS values

154 significantly greater than one, whereas 4,357 genes had dN/dS values significantly
155 less than one (χ^2 test, $P < 0.05$) (Table 1). Of these cancer positively selected genes,
156 1,191 genes exhibited positive selection in cancer genomes but non-positive selection
157 in germline genomes. Additionally, 275 genes exhibited negative selection in cancer
158 genomes but non-negative selection in germline genomes. These genes may therefore
159 be under different selective pressure in cancer and normal genomes.

160 Considering that different models might provide varying estimates, we used the
161 NG86 method²⁵ as the simplest model to calculate the numbers of nonsynonymous
162 and synonymous sites. The overall dN/dS value for cancer somatic substitutions
163 (mean \pm s.e.=0.990 \pm 0.006) was greater than that for germline substitutions
164 (mean \pm s.e.=0.624 \pm 0.004) for the 18,602 genes (Supplementary Table S1), whereas
165 the ratio was less than that calculated using mpNG method (Wilcoxon test, $P < 10^{-16}$)
166 (Table 1). Consequently, for both germline and cancer somatic substitutions, the
167 number of genes with dN/dS values > 1 (χ^2 test, $P < 0.05$) was much lower than those
168 calculated using the exonic mutation profiles, whereas the number of genes with
169 dN/dS values < 1 (χ^2 test, $P < 0.05$) was much greater (Table 1). We further used the
170 intergenic and intronic somatic mutation profiles of 507 cancer samples with
171 whole-genome sequences (WGSs) within the 7,042 tumor-normal pairs as a contrast.
172 The overall dN/dS values calculated using these mutation profiles were between the
173 values obtained using the NG86 method and the exonic mutation profiles, as was the
174 number of genes under positive and negative selection (Table 1, Supplementary Table
175 S1). Different models show different single-nucleotide substitution properties, which

176 resulted in a different list of candidate genes under positive and negative selection.
177 However, the genes under positive and negative selection calculated using different
178 models almost overlapped (Fig. 3A,B). The NG86 method ignores the mutation rate
179 bias between different substitution types, leading to underestimation of the dN/dS
180 ratio. Therefore, the NG86 method is strict with regard to detecting positive selection,
181 but it is relaxed about detection of negative selection²⁹. In contrast, the mpNG method
182 takes the mutation bias, which can be depicted as the internal variance between
183 mutation rates of different substitution types, into consideration. Thus, the mpNG
184 method could recover the underestimation of the true dN/dS ratio estimated by the
185 NG86 method, which would increase the sampling errors and false discovery rates
186 (FDRs). It would also increase the false positive results for detecting positively
187 selected genes, but be more conservative in detecting negatively selected genes. The
188 mutation bias does not affect the detection of genes under strong selection pressure,
189 while it may affect the detection of genes under weak selection pressure. The
190 mutation bias could be depicted by the internal variance of different substitution types.
191 The exonic mutation profile had greater internal variance ($\sigma=0.015$) than that of
192 intronic ($\sigma=0.008$) and intergenic ($\sigma=0.008$) mutation profiles, leading to the
193 maximum estimation of dN/dS ratios.

194 Regardless of the method used to calculate the dN/dS values for germline and
195 cancer somatic substitutions, we found that the dN/dS value for cancer somatic
196 substitutions was much greater than that for germline substitutions. Previous studies
197 have attributed the elevated dN/dS values to the relaxation of purifying selection¹⁴ or

the increased positive selection of globally expressed genes¹⁵. Our results show that the number of genes under positive selection increased, whereas the number of genes under negative selection decreased, in cancer genomes compared with germline genomes. This result indicates that both the relaxation of purifying selection on passenger mutations and the positive selection of driver mutations may contribute to the increased dN/dS values of human genes in cancer genomes.

Relaxation of purifying selection for human genes in cancer cells

In this study, we used the mpNG method with exonic mutation profiles to estimate the dN/dS values for germline substitutions and cancer somatic mutations. The Cancer Gene Census^{30,31} contains more than 500 cancer genes that have been reported in the literature to exhibit mutations and that are causally implicated in cancer development. Of those, 503 genes were included in the 18,602 genes we tested. These known cancer genes had significantly lower dN/dS values for germline substitutions (Wilcoxon test, $P < 10^{-16}$), but slightly greater dN/dS values (Wilcoxon test, $P = 0.01$) for cancer somatic mutations than those of other genes (Table 2A). For selection over longer time scales, we extracted the dN/dS values between human-mouse orthologs from the Ensembl database (Release 73)^{32,33}. The known cancer genes had significantly lower human-mouse dN/dS values than other human genes. Among the cancer genes, oncogenes (OGs) had significantly lower dN/dS values than non-cancer genes (Wilcoxon test, $P < 10^{-15}$), whereas the mean dN/dS values of tumor suppressor genes (TSGs) were not significantly different from those

220 of non-cancer genes (Wilcoxon test, $P=0.89$). These results support the work of
 221 Thomas *et al.*³⁴, who showed that known cancer genes may be more constrained and
 222 more important than other genes at the species and population levels, especially for
 223 oncogenes. In contrast, known cancer genes are more likely to gain functional somatic
 224 mutations in cancer relative to all other genes. However, within the known cancer
 225 genes, only 53 genes exhibited positive selection (χ^2 test, $P<0.05$) for cancer somatic
 226 substitutions, which suggests that positive selection for driver mutations is obscured
 227 by the relaxed purifying selection of passenger mutations.

228 We also examined human essential genes³⁵ and cancer common essential genes²¹.
 229 We extracted 2,452 human essential genes from DEG10 (the Database of Essential
 230 Genes)³⁵. These genes are critical for cell survival and are therefore more conserved
 231 than other genes at species and population levels. Here, we found that human essential
 232 genes had significantly lower dN/dS values of human-mouse orthologs and germline
 233 substitution, and similar dN/dS values for cancer somatic mutations, relative to the
 234 values for non-essential genes (Table 2A). Cancer essential genes were identified by
 235 performing genome-scale pooled RNAi screens. RNAi screens with the 45k shRNA
 236 pool in 12 cancer cell lines, including small-cell lung cancer, non-small-cell lung
 237 cancer, glioblastoma, chronic myelogenous leukemia, and lymphocytic leukemia,
 238 revealed 268 common essential genes²¹. Compared to other human genes, these
 239 cancer essential genes had significantly lower dN/dS values of human-mouse
 240 orthologs and germline substitutions, and similar dN/dS values for cancer somatic
 241 mutations (Table 2A).

242 The cancer positively selected genes displayed a pattern similar to that of the
243 cancer genes, cancer common essential genes, and human essential genes. These
244 genes had lower dN/dS values for human-mouse orthologs (Wilcoxon test, $P=4.5 \times 10^{-4}$)
245 and germline substitutions (Wilcoxon test, $P=0.01$), but significantly greater dN/dS
246 values for cancer somatic mutations (Wilcoxon test, $P < 10^{-16}$). However, the
247 negatively selected cancer genes displayed a different pattern, with greater dN/dS
248 values for human-mouse orthologs (Wilcoxon test, $P=7.3 \times 10^{-4}$) and germline
249 substitutions (Wilcoxon test, $P=2.3 \times 10^{-4}$), and significantly lower dN/dS values for
250 cancer somatic mutations (Wilcoxon test, $P < 10^{-16}$). These results indicate that the
251 positively selected genes may include the cancer associated genes or human essential
252 genes, while the negatively selected genes may include genes under greater selective
253 constraints in cancer cells than in normal cells.

254 We further tested the correlation of dN/dS values of human genes for
255 human-mouse orthologs, germline substitutions and cancer somatic mutations, in
256 order to compare selective pressures among species, populations and cancers (Table
257 2B). For different gene sets, the dN/dS values between human-mouse orthologs
258 showed a weak positive correlation with those of germline substitutions, but no
259 correlation with the values for cancer somatic substitutions. The dN/dS values for
260 human germline and cancer somatic substitutions displayed different correlation
261 patterns between different gene sets. The tumor suppressor genes and positively
262 selected cancer genes showed weak positive correlation, while other gene sets had no
263 correlation.

264

265 **Roles of cancer positively and negatively selected genes in cancer cells**

266 We next tested the genes under positive or purifying selection for their roles in
267 cancer. Functional annotation analysis based on the Database for Annotation,
268 Visualization and Integrated Discovery (DAVID) v6.7^{36,37} showed an enrichment of
269 genes involved in cell morphogenesis and pathways in cancer for positively selected
270 genes (Table 3A). Additionally, we found an enrichment of genes involved in sensory
271 perception for cancer negatively selected cancer genes (Table 3B). It is important to
272 note that we only used a relaxed filter ($P < 0.05$) for detecting cancer positively or
273 negatively selected genes, which would lead to high FDRs. We further calculated the
274 FDR for each P-value, using the qvalue (Supplementary Table S1)³⁸. We set the
275 strengthened filter for detecting positively and negatively selected genes at $P < 10^{-3}$ and
276 $FDR < 0.25$. Only 61 genes met this requirement, which included 45 cancer positively
277 selected genes and 16 cancer negatively selected genes (Supplementary Table S2).

278 Among the 45 cancer positively selected genes, there were three oncogenes
279 (GANP, NFE2L2, RHOA) and five tumor suppressor genes (TP53, CSMD1,
280 CDKN2A and SPOP), according to the Cancer Gene Census³⁰. Fourteen of the 61
281 genes are human essential genes, and seven are orthologs of mouse or yeast essential
282 genes, according to the DEG10³⁵. In addition, four positively selected genes (IKBIP,
283 TEX13A, FZD10 and PGAP2) also had dN/dS values significantly greater than one
284 ($P < 0.01$, $FDR < 0.05$) for germline substitutions. Six genes showed negative selection
285 ($P < 0.01$, $FDR < 0.5$) in germline substitutions. Among those six genes, CAMD2,

286 CSMD1 and CSMD3 have been reported as candidate tumor suppressor genes³⁹⁻⁴².
 287 Additionally, ACTG1 is associated with cancer cell migration⁴³. There were also 13
 288 cancer positively selected genes displaying neutral selection for germline substitutions.
 289 It would be interesting to investigate the roles in cancer of these cancer-specific
 290 positively selected genes and the four human essential genes that are not known to be
 291 cancer-related.

292 Among the 16 cancer negatively selected genes, there were two human essential
 293 genes: an oncogene (FUS) and a tumor suppressor gene (APC). Both were also under
 294 negative selection for germline substitutions ($P < 0.02$, $FDR < 0.06$). BRCA1 mutations,
 295 which would increase cancer risk for breast and ovarian cancer, can be germline
 296 mutations as well as somatic mutations⁴⁴. The other 13 genes showed greater selective
 297 constraint in cancer cells than in normal cells. It would be quite valuable to uncover
 298 the roles of these evolutionarily conserved genes in cancer cells. Several of these
 299 genes were reported to be required for the survival and proliferation of cancer cells
 300 and might therefore serve as potential drug targets or prognostic biomarkers. For
 301 example, BCL2L12 is a member of the BCL2 family and is an anti-apoptotic factor
 302 that can inhibit the p53 tumor suppressor and caspases 3 and 7^{45,46}. Overexpression of
 303 BCL2L12 has been detected in several cancer types, and BCL2L12 is therefore
 304 considered a molecular prognostic biomarker in these cancers⁴⁷⁻⁵⁰. MAP4 is a major
 305 non-neuronal microtubule-associated protein that promotes microtubule assembly. Ou
 306 *et al.* have reported that the protein level of MAP4 is positively correlated with
 307 bladder cancer grade. Additionally, silencing MAP4 can efficiently disrupt the

308 microtubule cytoskeleton, inhibiting the invasion and migration of bladder cancer
309 cells⁵¹. EPPK1 is a member of the plakin family, which plays a role in the
310 organization of cytoskeletal architecture. Guo *et al.* used proteomics to identify
311 EPPK1 as a predictive plasma biomarker for cervical cancer⁵². These negatively
312 selected, cancer-specific genes are more conserved in cancer cells than in normal cells,
313 indicating they may be crucial for the basic cellular processes of cancer cells.

314

315 **Discussion**

316 A key goal of cancer research is to identify cancer-related genes, such as OGs
317 and TSGs, the mutation of which might promote the occurrence and progression of
318 tumors²⁶. There are also cancer essential genes that are important for the growth and
319 survival of cancer cells²¹. Different methods are needed to identify different types of
320 cancer-related genes. In contrast to recent studies focused on the detection of driver
321 mutations^{16-18,53}, we aimed to detect cancer essential genes using a molecular
322 evolution approach. Advances in the understanding of positively selected cancer
323 drivers, as well as the severe side effects of classical chemotherapy and radiation
324 therapies that target DNA integrity and cell division, have fueled efforts to develop
325 anticancer drugs with more precise molecular targeting and fewer side effects.
326 Although personalized therapeutic approaches that target genetically activated drivers
327 have greatly improved patient outcomes in a number of common and rare cancers, the
328 rapid acquisition of drug resistance due to high intra-tumor heterogeneity is becoming
329 a challenging problem⁵⁴. In other words, driver mutations may differ considerably

330 between tumor sub-clones. Instead of looking for cancer-causing genes with multiple
331 driver mutations, an alternative approach is to identify cancer essential genes that are
332 highly conserved in tumor cell populations because they are crucial for carcinogenesis,
333 progression and metastasis. To some extent, this idea may overcome drug resistance
334 in targeted cancer therapies, as mutations in cancer essential genes are deleterious in
335 tumor populations.

336 Several approaches can be utilized to identify cancer essential genes suitable for
337 targeting with drugs, including siRNA-mediated knockdown of specific components
338 and genetic tumor models. The genome-wide pooled shRNA screens promoted by the
339 RNAi Consortium⁵⁵, however, can only be performed in cell lines *in vitro* and are
340 limited to the analysis of genes important for proliferation and survival^{21,56}. Thus,
341 these screens will miss certain classes of genes that may function only in the proper *in*
342 *vivo* tumor environment. Furthermore, siRNA screens may not be sensitive to target
343 genes whose products are components of the cellular machinery. These types of
344 targets may be frequently stabilized by their participation in complexes with a long
345 biological half-life. Indeed, this longevity may be the reason why not all such targets
346 seem to be essential for cancer cells in standard short-term siRNA screens⁸. Genetic
347 tumor models can also enable screening strategies within an entire organism to
348 identify cancer essential genes. However, this method is not suitable for large-scale
349 screening. With the explosive increase in cancer somatic mutation data from cancer
350 genome sequencing, it is now possible to investigate the natural selection of each
351 human gene in cancer genomes using evolutionary genomics methods⁸. One major

352 aim is to identify genes that are under significantly increased purifying selective
353 constraint in cancer cells relative to normal cells, which would suggest that they are
354 cancer-specific essential genes.

355 Through analyses of large-scale cancer somatic mutation data derived from The
356 Cancer Genome Atlas (TCGA) or International Cancer Genome Consortium (ICGC),
357 previous studies found important differences between the evolutionary dynamics of
358 cancer somatic cells and whole organisms^{6,14,16}. However, these studies applied
359 canonical nucleotide substitution models to identify the molecular signatures of
360 natural selection in cancer cells or human populations, which neglected the apparently
361 different mutation profiles between these cell types. Here, we developed a new
362 mutation-profile-based Nei-Gojobori method (mpNG) to calculate the dN/dS values
363 of 18,602 human genes for both cancer somatic and normal human germline
364 substitutions.

365 Two prerequisites are crucial to properly apply the mpNG method. First, a large
366 number of samples with similar mutation profiles is necessary to increase the power
367 of selection pressure detection. Second, a subset of nucleotide substitutions should be
368 chosen to represent the background neutral mutation profiles of the samples. In this
369 study, because of the limited number of cancer samples, especially the number of
370 whole-genome sequenced tumor-normal tissue pairs, we pooled all of the samples to
371 analyze pan-cancer-level selection pressures. Mutation profiles are well known to be
372 heterogeneous, even for samples with the same tissue origin^{17,26}. As an increasing
373 number of cancer genomes are sequenced in the near future, we will be able to

374 classify cancer samples by their specific mutation profiles and infer evolutionarily
 375 selective pressures using the mpNG method. With respect to background neutral
 376 mutation profiles, it will be appropriate to calculate them based on intergenic regions
 377 from the corresponding samples. However, only a small number of tumor-normal
 378 paired WGSs are currently available. Therefore, in this study, we assume that most
 379 exonic somatic mutations in the cancer samples do not have significant effects on the
 380 fitness of cancer cells. Under this assumption, we can apply the mutation profiles of
 381 WESs to approximate the background. The exonic mutation profiles used in our
 382 mpNG method consider the weight of the 96 substitution classifications within the
 383 cancer exomes, which may reflect the mutation bias of different substitution types
 384 within the protein-coding regions. This method would recover the underestimation of
 385 the dN/dS value that occurs with the NG86 method²⁹. Using the mpNG method, the
 386 detection of positive selection would be relaxed, whereas the detection of negative
 387 selection would be conservative relative to the NG86 method. Were more
 388 tumor-normal WGSs available, it would be better to choose suitable mutation profiles
 389 for the mpNG method. With the expansion of these data in the future, we may apply
 390 more precise methods to identify neutral background mutation properties.

391 As a conservative estimate of positively and negatively selected genes in cancer,
 392 we found 45 genes under intensified positive selection and 16 genes under
 393 strengthened purifying selection in cancer cells compared with germline cells. The set
 394 of cancer-specific positively selected genes was enriched for known cancer genes
 395 and/or human essential genes, while several of the cancer-specific negatively selected

genes have previously been reported as prognostic biomarkers for cancers. Because cancer-specific negatively selected genes are more evolutionarily constrained in cancer cells than in normal cells, identification of cancer-specific negatively selected genes would inform the potential options for cancer therapeutic targets or diagnostic biomarkers. However, cancer somatic mutations vary greatly among different cancer types and even among individual cancer genomes^{17,18,26,57}, therefore, further studies will be needed to better understand the evolution of human cancer.

Methods

Datasets

Cancer somatic mutation data from 7,042 primary cancers corresponding to 30 different classes were extracted from the work of Alexandrov *et al.*²⁶, which includes 4,938,362 somatic substitutions and small insertions/deletions from 507 WGSs and 6,535 WESs. Data on rare human protein-coding variants (minor allele frequency <0.01%) from 6,500 human WESs (ESP6500) were extracted from the ANNOVAR database⁵⁸ based on the NHLBI GO Exome Sequencing Project. A total of 522 known cancer genes were extracted from the Cancer Gene Census (<http://cancer.sanger.ac.uk/cancergenome/projects/census/>, COSMIC v68)^{30,31}.

Human gene sequences and annotations were extracted from the Ensembl database (Release 73)^{32,33}. For each gene, we only chose the longest sequence to avoid duplicate records of each single substitution. The HGNC (HUGO Gene Nomenclature Committee) database⁵⁹ (<http://www.genenames.org/>) and the Genecards database⁶⁰

418 (<http://www.genecards.org>) were also used to map the gene IDs from different
419 datasets. DAVID v6.7 was utilized for the functional annotation analysis^{36,37}.

420

421 **Calculating mutation rate profiles**

422 We calculated the mutation rate profiles using the 96 substitution
423 classifications^{26,27}, which not only show the base substitution but also include
424 information on the sequence context of each mutated base. We counted all somatic
425 substitutions in the protein-coding regions of the 7,042 tumor-normal paired WESs, as
426 well as all the protein-coding variants of the ESP6500 dataset. We also counted the
427 total number of each trinucleotide type for the exonic, intronic, and intergenic regions
428 in the human genome. We calculated the mutation rate of each substitution type as the
429 number of substitutions per trinucleotide type per patient. The mutation profiles were
430 depicted as the mutation rate of each mutation type according to the 96 substitution
431 classifications.

432

433 **Detection of positive and negative selections**

434 ANNOVAR was utilized to perform biological and functional annotations of the
435 cancer somatic mutations and germline substitutions⁵⁸. Substitutions within
436 protein-coding genes were classified as either nonsynonymous or synonymous. We
437 counted the number of nonsynonymous (n) and synonymous (s) substitutions for each
438 gene across all somatic mutations in the 7,042 tumor-normal pairs. Somatic mutations
439 at the same site and with the same mutation type that occurred in different patients

were counted as different substitutions because they, unlike germline evolution, occurred independently.

We further calculated the number of nonsynonymous (N) and synonymous (S) sites in each human protein-coding gene utilizing different models. The simple method of Nei and Gojobori was used²⁵. We also considered cancer somatic mutation profiles, which were depicted as the percentage of each mutation type according to the 96 substitution classifications. For each gene, we calculated the proportion of substitutions that would be nonsynonymous or synonymous for each protein-coding site, as the probability of mutation types for each site was determined according to the mutation profiles. Then, we added up the proportions to calculate the total number of nonsynonymous (N) and synonymous (S) sites for each gene.

After counting the number of nonsynonymous (n) and synonymous (s) substitutions, as well as the number of nonsynonymous (N) and synonymous (S) sites for each gene, we calculated the ratio of the rates of nonsynonymous and synonymous substitutions (dN/dS) for each human gene as follows: $\frac{dN}{dS} = \frac{n / N}{(s + 0.5) / (S + 0.5)}$.

The dN/dS for germline substitutions was calculated using the same approach.

A χ^2 test was used to compare the number of nonsynonymous and synonymous substitutions to the number of nonsynonymous and synonymous sites for each gene in order to test the statistical significance of the difference between the dN/dS values and one. The genes with dN/dS values significantly greater than one were classified as being under positive selection in tumors, whereas the genes with dN/dS values significantly less than one were classified as being under negative, or purifying,

462 selection. The false discovery rate was estimated using the qvalue package from
463 Bioconductor³⁸. A Wilcoxon test was performed to compare dN/dS values between
464 cancer somatic substitutions and germline substitutions, as well as between known
465 cancer genes and all other genes. The software tool R was used for statistical analysis
466 (<http://www.r-project.org/>).

467

468 References

- 469 1 Cairns, J. Mutation selection and the natural history of cancer. *Nature* **255**, 197-200 (1975).
- 470 2 Nowell, P. C. The clonal evolution of tumor cell populations. *Science* **194**, 23-28 (1976).
- 471 3 Crespi, B. J. & Summers, K. Positive selection in the evolution of cancer. *Biol. Rev. Camb. Philos.*
472 *Soc.* **81**, 407-424, doi:10.1017/S1464793106007056 (2006).
- 473 4 Merlo, L. M., Pepper, J. W., Reid, B. J. & Maley, C. C. Cancer as an evolutionary and ecological
474 process. *Nat. Rev. Cancer* **6**, 924-935, doi:10.1038/nrc2013 (2006).
- 475 5 Podlaha, O., Riester, M., De, S. & Michor, F. Evolution of the cancer genome. *Trends Genet.* **28**,
476 155-163, doi:10.1016/j.tig.2012.01.003 (2012).
- 477 6 Yates, L. R. & Campbell, P. J. Evolution of the cancer genome. *Nat. Rev. Genet.* **13**, 795-806,
478 doi:10.1038/nrg3317 (2012).
- 479 7 Greaves, M. & Maley, C. C. Clonal evolution in cancer. *Nature* **481**, 306-313,
480 doi:10.1038/nature10762 (2012).
- 481 8 Luo, J., Solimini, N. L. & Elledge, S. J. Principles of cancer therapy: oncogene and non-oncogene
482 addiction. *Cell* **136**, 823-837, doi:10.1016/j.cell.2009.02.024 (2009).
- 483 9 Hanahan, D. & Weinberg, R. A. Hallmarks of Cancer: The Next Generation. *Cell* **144**, 646-674
484 (2011).
- 485 10 Stratton, M. R., Campbell, P. J. & Futreal, P. A. The cancer genome. *Nature* **458**, 719-724,
486 doi:10.1038/nature07943 (2009).
- 487 11 Greenman, C. *et al.* Patterns of somatic mutation in human cancer genomes. *Nature* **446**, 153-158,
488 doi:10.1038/nature05610 (2007).
- 489 12 Beckman, R. A. & Loeb, L. A. Negative clonal selection in tumor evolution. *Genetics* **171**,
490 2123-2131, doi:10.1534/genetics.105.040840 (2005).
- 491 13 McFarland, C. D., Korolev, K. S., Kryukov, G. V., Sunyaev, S. R. & Mirny, L. A. Impact of
492 deleterious passenger mutations on cancer progression. *Proc. Natl. Acad. Sci. USA* **110**, 2910-2915,
493 doi:10.1073/pnas.1213968110 (2013).
- 494 14 Woo, Y. H. & Li, W. H. DNA replication timing and selection shape the landscape of nucleotide
495 variation in cancer genomes. *Nat. Commun.* **3**, 1004, doi:10.1038/ncomms1982 (2012).
- 496 15 Ostrow, S. L., Barshir, R., DeGregori, J., Yeger-Lotem, E. & Hershberg, R. Cancer Evolution Is
497 Associated with Pervasive Positive Selection on Globally Expressed Genes. *PLoS Genet.* **10**,
498 e1004239, doi:10.1371/journal.pgen.1004239 (2014).
- 499 16 Davoli, T. *et al.* Cumulative haploinsufficiency and triplosensitivity drive aneuploidy patterns and

- 500 shape the cancer genome. *Cell* **155**, 948-962, doi:10.1016/j.cell.2013.10.011 (2013).
- 501 17 Lawrence, M. S. *et al.* Mutational heterogeneity in cancer and the search for new
502 cancer-associated genes. *Nature* **499**, 214-218, doi:10.1038/nature12213 (2013).
- 503 18 Lawrence, M. S. *et al.* Discovery and saturation analysis of cancer genes across 21 tumour types.
504 *Nature* **505**, 495-501, doi:10.1038/nature12912 (2014).
- 505 19 Chen, H., Xing, K. & He, X. The dJ/dS Ratio Test Reveals Hundreds of Novel Putative Cancer
506 Drivers. *Mol. Biol. Evol.*, doi:10.1093/molbev/msv083 (2015).
- 507 20 Jordan, I. K., Rogozin, I. B., Wolf, Y. I. & Koonin, E. V. Essential genes are more evolutionarily
508 conserved than are nonessential genes in bacteria. *Genome Res.* **12**, 962-968, doi:10.1101/gr.87702
509 (2002).
- 510 21 Luo, B. *et al.* Highly parallel identification of essential genes in cancer cells. *Proc. Natl. Acad. Sci.*
511 *USA* **105**, 20380-20385, doi:10.1073/pnas.0810485105 (2008).
- 512 22 Endo, T., Ikeo, K. & Gojobori, T. Large-scale search for genes on which positive selection may
513 operate. *Mol. Biol. Evol.* **13**, 685-690 (1996).
- 514 23 Arbiza, L., Dopazo, J. & Dopazo, H. Positive selection, relaxation, and acceleration in the
515 evolution of the human and chimp genome. *PLoS Comput. Biol.* **2**, e38,
516 doi:10.1371/journal.pcbi.0020038 (2006).
- 517 24 Ovens, K. & Naugler, C. Preliminary evidence of different selection pressures on cancer cells as
518 compared to normal tissues. *Theor. Biol. Med. Model* **9**, 44, doi:10.1186/1742-4682-9-44 (2012).
- 519 25 Nei, M. & Gojobori, T. Simple methods for estimating the numbers of synonymous and
520 nonsynonymous nucleotide substitutions. *Mol. Biol. Evol.* **3**, 418-426 (1986).
- 521 26 Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature* **500**, 415-421,
522 doi:10.1038/nature12477 (2013).
- 523 27 Samocha, K. E. *et al.* A framework for the interpretation of de novo mutation in human disease.
524 *Nat. Genet.* **46**, 944-950, doi:10.1038/ng.3050 (2014).
- 525 28 Hodis, E. *et al.* A landscape of driver mutations in melanoma. *Cell* **150**, 251-263,
526 doi:10.1016/j.cell.2012.06.024 (2012).
- 527 29 Yang, Z. & Bielawski, J. P. Statistical methods for detecting molecular adaptation. *Trends Ecol.*
528 *Evol.* **15**, 496-503, doi:10.1016/S0169-5347(00)01994-7 (2000).
- 529 30 Futreal, P. A. *et al.* A census of human cancer genes. *Nat. Rev. Cancer* **4**, 177-183,
530 doi:10.1038/nrc1299 (2004).
- 531 31 Forbes, S. A. *et al.* COSMIC: mining complete cancer genomes in the Catalogue of Somatic
532 Mutations in Cancer. *Nucleic Acids Res.* **39**, D945-950, doi:10.1093/nar/gkq929 (2011).
- 533 32 Kinsella, R. J. *et al.* Ensembl BioMarts: a hub for data retrieval across taxonomic space. *Database*
534 *(Oxford)* **2011**, bar030, doi:10.1093/database/bar030 (2011).
- 535 33 Flicek, P. *et al.* Ensembl 2014. *Nucleic Acids Res.* **42**, D749-755, doi:10.1093/nar/gkt1196 (2014).
- 536 34 Thomas, M. A. *et al.* Evolutionary dynamics of oncogenes and tumor suppressor genes: higher
537 intensities of purifying selection than other genes. *Mol. Biol. Evol.* **20**, 964-968,
538 doi:10.1093/molbev/msg110 (2003).
- 539 35 Luo, H., Lin, Y., Gao, F., Zhang, C. T. & Zhang, R. DEG 10, an update of the database of essential
540 genes that includes both protein-coding genes and noncoding genomic elements. *Nucleic Acids*
541 *Res.* **42**, D574-580, doi:10.1093/nar/gkt1131 (2014).
- 542 36 Huang da, W., Sherman, B. T. & Lempicki, R. A. Systematic and integrative analysis of large gene
543 lists using DAVID bioinformatics resources. *Nat. Protoc.* **4**, 44-57, doi:10.1038/nprot.2008.211

(2009).

37 Huang da, W., Sherman, B. T. & Lempicki, R. A. Bioinformatics enrichment tools: paths toward the comprehensive functional analysis of large gene lists. *Nucleic Acids Res.* **37**, 1-13, doi:10.1093/nar/gkn923 (2009).

38 Storey, J. D. & Tibshirani, R. Statistical significance for genomewide studies. *Proc. Natl. Acad. Sci. USA* **100**, 9440-9445, doi:10.1073/pnas.1530509100 (2003).

39 Kamal, M. *et al.* Loss of CSMD1 expression is associated with high tumour grade and poor survival in invasive ductal breast carcinoma. *Breast Cancer Res. Treat.* **121**, 555-563, doi:10.1007/s10549-009-0500-4 (2010).

40 Zhang, R. & Song, C. Loss of CSMD1 or 2 may contribute to the poor prognosis of colorectal cancer patients. *Tumour Biol.* **35**, 4419-4423, doi:10.1007/s13277-013-1581-6 (2014).

41 Chang, G. *et al.* Hypoexpression and epigenetic regulation of candidate tumor suppressor gene CADM-2 in human prostate cancer. *Clin. Cancer Res.* **16**, 5390-5401, doi:10.1158/1078-0432.CCR-10-1461 (2010).

42 He, W. *et al.* Aberrant methylation and loss of CADM2 tumor suppressor expression is associated with human renal cell carcinoma tumor progression. *Biochem. Biophys. Res. Commun.* **435**, 526-532, doi:10.1016/j.bbrc.2013.04.074 (2013).

43 Luo, Y. *et al.* Loss of ASAP3 destabilizes cytoskeletal protein ACTG1 to suppress cancer cell migration. *Mol. Med. Rep.* **9**, 387-394, doi:10.3892/mmr.2013.1831 (2014).

44 Roy, R., Chun, J. & Powell, S. N. BRCA1 and BRCA2: different roles in a common pathway of genome protection. *Nat. Rev. Cancer* **12**, 68-78, doi:10.1038/nrc3181 (2012).

45 Stegh, A. H. *et al.* Bcl2L12-mediated inhibition of effector caspase-3 and caspase-7 via distinct mechanisms in glioblastoma. *Proc. Natl. Acad. Sci. USA* **105**, 10703-10708, doi:10.1073/pnas.0712034105 (2008).

46 Stegh, A. H. *et al.* Glioma oncoprotein Bcl2L12 inhibits the p53 tumor suppressor. *Genes Dev.* **24**, 2194-2204, doi:10.1101/gad.1924710 (2010).

47 Thomadaki, H. *et al.* Overexpression of the novel member of the BCL2 gene family, BCL2L12, is associated with the disease outcome in patients with acute myeloid leukemia. *Clin. Biochem.* **45**, 1362-1367, doi:10.1016/j.clinbiochem.2012.06.012 (2012).

48 Karan-Djurasevic, T. *et al.* Expression of Bcl2L12 in chronic lymphocytic leukemia patients: association with clinical and molecular prognostic markers. *Med. Oncol.* **30**, 405, doi:10.1007/s12032-012-0405-7 (2013).

49 Foutadakis, S., Avgeris, M., Tokas, T., Stravodimos, K. & Scorilas, A. Increased BCL2L12 expression predicts the short-term relapse of patients with TaT1 bladder cancer following transurethral resection of bladder tumors. *Urol. Oncol.* **32**, 39 e29-36, doi:10.1016/j.urolonc.2013.04.005 (2014).

50 Tzovaras, A. *et al.* BCL2L12: a promising molecular prognostic biomarker in breast cancer. *Clin. Biochem.* **47**, 257-262, doi:10.1016/j.clinbiochem.2014.09.008 (2014).

51 Ou, Y. *et al.* Activation of cyclic AMP/PKA pathway inhibits bladder cancer cell invasion by targeting MAP4-dependent microtubule dynamics. *Urol. Oncol.* **32**, 47.e21-e28, doi:10.1016/j.urolonc.2013.06.017 (2014).

52 Guo, X. *et al.* Potential predictive plasma biomarkers for cervical cancer by 2D-DIGE proteomics and Ingenuity Pathway Analysis. *Tumour Biol.* **36**, 1711-1720, doi:10.1007/s13277-014-2772-5 (2015).

588 53 Supek, F., Miñana, B., Valcárcel, J., Gabaldón, T. & Lehner, B. Synonymous Mutations Frequently
589 Act as Driver Mutations in Human Cancers. *Cell* **156**, 1324-1335, doi:10.1016/j.cell.2014.01.051
590 (2014).

591 54 Dobbstein, M. & Moll, U. Targeting tumour-supportive cellular machineries in anticancer drug
592 development. *Nat. Rev. Drug Discov.* **13**, 179-196, doi:10.1038/nrd4201 (2014).

593 55 Root, D. E., Hacohen, N., Hahn, W. C., Lander, E. S. & Sabatini, D. M. Genome-scale
594 loss-of-function screening with a lentiviral RNAi library. *Nat. Methods* **3**, 715-719,
595 doi:10.1038/nmeth924 (2006).

596 56 Marcotte, R. *et al.* Essential gene profiles in breast, pancreatic, and ovarian cancer cells. *Cancer*
597 *Discov.* **2**, 172-189, doi:10.1158/2159-8290.CD-11-0224 (2012).

598 57 Kandoth, C. *et al.* Mutational landscape and significance across 12 major cancer types. *Nature* **502**,
599 333-339, doi:10.1038/nature12634 (2013).

600 58 Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic variants from
601 high-throughput sequencing data. *Nucleic Acids Res.* **38**, e164, doi:10.1093/nar/gkq603 (2010).

602 59 Gray, K. A. *et al.* Genenames.org: the HGNC resources in 2013. *Nucleic Acids Res.* **41**, D545-552,
603 doi:10.1093/nar/gks1066 (2013).

604 60 Safran, M. *et al.* GeneCards Version 3: the human gene integrator. *Database (Oxford)* **2010**,
605 baq020, doi:10.1093/database/baq020 (2010).

606

607

608 **Acknowledgements**

609 We are grateful to Xiaopu Wang for his help with the manuscript preparation. We
610 would like to thank the NHLBI GO Exome Sequencing Project and its ongoing
611 studies which produced and provided exome variant calls for comparison: the Lung
612 GO Sequencing Project (HL-102923), the WHI Sequencing Project (HL-102924), the
613 Broad GO Sequencing Project (HL-102925), the Seattle GO Sequencing Project
614 (HL-102926) and the Heart GO Sequencing Project (HL-103010). We also gratefully
615 acknowledge the TCGA Research Network and the Broad Institute TCGA GDAC
616 Firehose for referencing the TCGA datasets.

617 This work was supported by a grant from the Ministry of Science and Technology
618 China (2012CB910101), grants from the National Natural Science Foundation of
619 China (31272299, 31301034), a grant from the Zhejiang Provincial Natural Sciences
620 Foundation of China (LY15C060001), grants from Fudan University and Iowa State
621 University to XG, the Shanghai Pujiang Program (13PJD005) to ZS, and a General
622 Financial Grant from China Postdoctoral Science Foundation (2013M531117) to ZZ.

623

624 **Authors Contributions**

625 ZZ, ZS and XG conceived of the study. ZZ, YZ, GL, JZ and ZS performed the data
626 analyses. SZ contributed ideas and tools for the analysis and edited the manuscript.
627 ZZ, YZ, ZS and XG wrote the manuscript. All authors read and approved the final
628 manuscript.

629

630 **Competing financial interests**

631 The authors declare that they have no competing financial interests.

632

633

634 **Figure Legends**

635 **Figure 1.** Mutation profiles of cancer somatic substitutions and germline substitutions,
636 including the exonic mutation profile of 7,042 cancer samples, the exonic mutation
637 profile of ESP6500, the intronic mutation profile of 507 whole cancer genomes, the
638 intergenic mutation profile of 507 whole cancer genomes, and the exonic mutation
639 profiles of breast carcinoma (BRCA), lung adenocarcinoma (LUAD), colon
640 adenocarcinoma (COAD), and skin cutaneous melanoma (SKCM).

641 **Figure 2.** The pipeline used to identify positively and negatively selected cancer
642 genes with the mpNG method.

643 **Figure 3.** The overlap of positively selected (A) and negatively selected (B) genes
644 based on different models.

645

646 Tables

647 **Table 1.** The dN/dS values and number of human genes under positive or negative
648 selection in germline and cancer based on the NG86 and mpNG methods with
649 different mutation profiles. P-values are according to a χ^2 test.

	dN/dS	# Positive selection*	# Negative selection*
Germline (NG86)	0.624 ± 0.004	42	9093
Cancer (NG86)	0.990 ± 0.006	306	2330
Cancer (intergenic)	1.240 ± 0.008	697	722
Cancer (intronic)	1.281 ± 0.008	822	624
Germline (exonic)	0.903 ± 0.006	264	4357
Cancer (exonic)	1.367 ± 0.009	1230	326

650 Note: *P<0.05

651

Table 2. The dN/dS values (A) and correlation of dN/dS values (B) of different gene sets for human-mouse orthologs, and for germline and cancer somatic substitutions.

(A)

	Human-Mouse	Germline	Cancer
All genes	0.155 ± 0.006	0.903 ± 0.006	1.367 ± 0.009
Known cancer genes	0.111 ± 0.005	0.675 ± 0.017	1.350 ± 0.033
Oncogenes	0.101 ± 0.006	0.665 ± 0.020	1.336 ± 0.038
Tumor suppressor genes	0.151 ± 0.014	0.732 ± 0.039	1.350 ± 0.066
Human essential genes	0.093 ± 0.002	0.704 ± 0.013	1.288 ± 0.015
Cancer essential genes	0.089 ± 0.007	0.698 ± 0.032	1.413 ± 0.067
Positively selected genes	0.136 ± 0.004	0.918 ± 0.029	3.216 ± 0.091
Negatively selected genes	0.172 ± 0.008	0.915 ± 0.023	0.479 ± 0.009

(B)

	Human-Mouse vs Germline		Human-Mouse vs Cancer		Germline vs Cancer	
	r	P-Value	r	P-Value	r	P-Value
All genes	0.04	3.3×10^{-7}	-0.01	0.47	0.10	$<10^{-16}$
Known cancer genes	0.45	$<10^{-16}$	-0.02	0.72	0.11	0.02
Oncogenes	0.43	$<10^{-16}$	-0.01	0.85	0.04	0.43
Tumor suppressor genes	0.52	1.0×10^{-8}	0.04	0.66	0.36	1.6×10^{-4}
Human essential genes	0.19	$<10^{-16}$	-0.05	0.01	0.06	1.4×10^{-3}
Cancer essential genes	0.30	5.7×10^{-7}	-0.07	0.29	0.03	0.65
Positively selected genes	0.17	2.4×10^{-9}	-0.02	0.57	0.23	$<10^{-16}$
Negatively selected genes	0.22	6.3×10^{-5}	0.10	0.07	0.04	0.60

Table 3. Functional enrichment of positively and negatively selected genes in cancer

genomes ($P < 0.01$, $FDR < 10\%$).

(A)

Category	Term	P-Value	FDR (%)
GOTERM_BP_FAT	GO:0032989~cellular component morphogenesis	7.42×10^{-4}	1.34
GOTERM_BP_FAT	GO:0043009~chordate embryonic development	2.40×10^{-3}	4.28
GOTERM_BP_FAT	GO:0009792~embryonic development ending in birth or egg hatching	2.89×10^{-3}	5.13
GOTERM_BP_FAT	GO:0000902~cell morphogenesis	3.28×10^{-3}	5.80
GOTERM_BP_FAT	GO:0030098~lymphocyte differentiation	4.90×10^{-3}	8.55
GOTERM_BP_FAT	GO:0051276~chromosome organization	5.19×10^{-3}	9.02
KEGG_PATHWAY	hsa05200:Pathways in cancer	4.23×10^{-3}	0.52
KEGG_PATHWAY	hsa05215:Prostate cancer	5.88×10^{-4}	0.72
KEGG_PATHWAY	hsa05213:Endometrial cancer	1.46×10^{-3}	1.78
KEGG_PATHWAY	hsa05210:Colorectal cancer	2.27×10^{-3}	2.75
KEGG_PATHWAY	hsa05216:Thyroid cancer	2.74×10^{-3}	3.32

(B)

Category	Term	P-Value	FDR (%)
GOTERM_BP_FAT	GO:0007600~sensory perception	1.35×10^{-3}	2.20
GOTERM_BP_FAT	GO:0050890~cognition	3.11×10^{-3}	5.00
GOTERM_BP_FAT	GO:0007608~sensory perception of smell	4.27×10^{-3}	6.80
GOTERM_BP_FAT	GO:0007606~sensory perception of chemical stimulus	4.67×10^{-3}	7.41
KEGG_PATHWAY	hsa04740:Olfactory transduction	1.03×10^{-3}	1.11