

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

2 **Somatic Mutations: A Comparative Analysis**

3

4 Zhan Zhou<sup>1,4,†</sup>, Yangyun Zou<sup>2,†</sup>, Gangbiao Liu<sup>1</sup>, Jingqi Zhou<sup>1</sup>, Jingcheng Wu<sup>4</sup>,

5 Shimin Zhao<sup>1</sup>, Zhixi Su<sup>2,\*</sup>, Xun Gu<sup>3,2,\*</sup>

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

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

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

12 <sup>4</sup>College of Pharmaceutical Sciences, Zhejiang University, Hangzhou, Zhejiang, China

13

14 †These authors contributed equally to this work

15 \*Corresponding authors

16 Zhixi Su, **e-mail:** zxsu@fudan.edu.cn, Tel: +86-21-51630616, Fax: +86-21-51630616

17 Xun Gu, **e-mail:** xgu@iastate.edu, Tel: +1-515-294- 8075, Fax: +1-515-294-8457

18

19 Note: This manuscript is currently under revision for *Oncotarget*.

20

21 **ABSTRACT**

22 Human genes exhibit different effects on fitness in cancer and normal cells. Here,  
23 we present an evolutionary approach to measure the selection pressure on human genes,  
24 using the well-known ratio of the nonsynonymous to synonymous substitution rate in  
25 both cancer genomes ( $C_N/C_S$ ) and normal populations ( $p_N/p_S$ ). A new mutation-profile-  
26 based method that adopts sample-specific mutation rate profiles instead of conventional  
27 substitution models was developed. We found that cancer-specific selection pressure is  
28 quite different from the selection pressure at the species and population levels. Both the  
29 relaxation of purifying selection on passenger mutations and the positive selection of  
30 driver mutations may contribute to the increased  $C_N/C_S$  values of human genes in cancer  
31 genomes compared with the  $p_N/p_S$  values in human populations. The  $C_N/C_S$  values also  
32 contribute to the improved classification of cancer genes and a better understanding of  
33 the onco-functionalization of cancer genes during oncogenesis. The use of our  
34 computational pipeline to identify cancer-specific positively and negatively selected  
35 genes may provide useful information for understanding the evolution of cancers and  
36 identifying possible targets for therapeutic intervention.

37

38 **Keywords:** Cancer somatic mutations, Mutation profile, Natural selection, Cancer-  
39 associated genes, Evolution

40

41

## 42 INTRODUCTION

43 Since the pioneering work of Cairns and Nowell [1, 2], the evolutionary concept  
44 of cancer progression has been widely accepted [3-7]. In this model, cancer cells evolve  
45 through random somatic mutations and epigenetic changes that may alter several crucial  
46 pathways, a process that is followed by clonal selection of the resulting cells.  
47 Consequently, cancer cells can survive and proliferate under deleterious circumstances  
48 [8, 9]. Therefore, knowledge of evolutionary dynamics will benefit our understanding  
49 of cancer initiation and progression. For example, there are two types of somatic  
50 mutations in cancer genomes: driver mutations and passenger mutations [10, 11].  
51 Driver mutations are those that confer a selective advantage on cancer cells, as indicated  
52 by statistical evidence of positive selection. Passenger mutations do not confer a clonal  
53 growth advantage and are usually considered neutral in cancer. However, some  
54 passenger mutations in protein-coding regions that would have potentially deleterious  
55 effects on cancer cells may be under negative selection in cancer [12, 13].

56 Cancer somatic mutations, especially driver mutations, promote the cancer specific  
57 functionalization of cancer-associated genes, i.e., onco-functionalization. Onco-  
58 functionalization of cancer-associated genes would promote cancer initiation and  
59 progression. For example, oncogenes may gain new functions during carcinogenesis,  
60 which could be considered cancer-specific neo-functionalization [14]. By contrast, the  
61 mutation of tumor suppressor genes to cause a loss or reduction of their function could  
62 be considered cancer-specific non-functionalization [15].

63 Analyses of large-scale cancer somatic mutation data have revealed that the effects

64 of positive selection are much stronger on cancer cells than on germline cells [16, 17].  
65 Given that many of the positively selected genes in tumor development act as the  
66 driving force behind tumor initiation and development and are thus considered “driver  
67 genes”, it is understandable that almost all previous studies have focused on positively  
68 selected genes in cancer genomes [3, 18-21]. Nevertheless, we have realized that an  
69 alternative approach, i.e., identifying cancer-constrained genes that are highly  
70 conserved in tumor cell populations (under purifying selection), is also valuable. For  
71 example, TP73, a homolog of TP53, is rarely mutated but frequently overexpressed in  
72 tumor cells. TP73 has been reported to activate the expression of glucose-6-phosphate  
73 dehydrogenase and support the proliferation of human cancer cells [22]. As essential  
74 genes are crucial for carcinogenesis, progression and metastasis, this idea may be  
75 advantageous in addressing issues related to drug resistance in cancer therapies,  
76 especially in cancers with high intratumor heterogeneity.

77 Many previous studies have used the ratio of nonsynonymous to synonymous  
78 substitution rates to identify genes that might be under strong positive selection both in  
79 organismal evolution and carcinogenesis [11, 16, 17, 23-26]. However, most of these  
80 studies applied conventional methods, which are usually based on simple nucleotide  
81 mutation/substitution models, e.g., the simplest equal-rate model assuming that every  
82 mutation or substitution pattern has the same probability [27]. Unfortunately, this may  
83 not be a realistic biological model because many recent cancer genomics studies have  
84 shown that mutation profiles vary greatly between different cancer samples [17, 28]. In  
85 addition, context-dependent mutation bias (i.e., base-substitution profiles that are

86 influenced by the flanking 5' and 3' bases of each mutated base) should also be  
87 considered [28, 29].

88 In this study, we describe a mutation-profile-based method to estimate the selective  
89 constraint for each gene in pan-cancer samples and human populations. In brief, the  
90 new method discards an unrealistic assumption inherent in the equal-rate model that  
91 every mutation or substitution pattern has the same probability [27]. This assumption  
92 can lead to nontrivial biased estimations when it is significantly violated. By contrast,  
93 our method implements an empirical nucleotide mutation model that simultaneously  
94 considers account several factors, including single-base mutation patterns, local-  
95 specific effects of surrounding DNA regions, and tissue/cancer types. Using simple  
96 somatic mutations from 9,155 tumor-normal paired whole-exome/genome sequences  
97 (ICGC Release 20), as well as rare germline substitutions from 6,500 exome sequences  
98 from the National Heart, Lung, and Blood Institute (NHLBI) Grant Opportunity (GO)  
99 Exome Sequencing Project (ESP), as references, we used this mutation-profile-based  
100 method to identify selective constraints on human genes, especially cancer-associated  
101 genes, in cancer cells. Our results may provide useful information for the precise  
102 classification of known cancer-associated genes and for an improved understanding of  
103 the evolution of cancers.

104

## 105 **RESULTS**

### 106 **The mutation rate profiles in cancer genomes and human populations differ**

107 Estimating evolutionary selective pressure on human genes is a practical method

108 for inferring the functional importance of genes in a specific population. By comparing  
109 selective pressures, we may be able to identify different functional and fitness effects  
110 of human genes in cancer and normal cells. The conventional method for measuring  
111 selective pressure is to calculate the ratio of nonsynonymous to synonymous  
112 substitution rates using the equal-rate method [27], which assumes equal substitution  
113 rates among different nucleotides. In this study, we used the cancer somatic mutations  
114 from 9,155 tumor-normal pairs from ICGC (Release 20) as well as rare variants (minor  
115 allele frequency <0.01%) from 6,500 exome sequences from ESP as a reference. We  
116 used these data to compare the empirical mutation rate profiles of cancer somatic  
117 mutations and germline substitutions using 96 substitution classifications [28, 29]. The  
118 empirical mutation rate profiles reveal the prevalence of each substitution pattern for  
119 point mutations and present not only the substitution types but also the sequence context  
120 (see Methods). The exonic mutation profiles of cancer somatic mutations and germline  
121 substitutions are both enriched in C-to-T transitions (Figure 1). The mutation rates for  
122 each trinucleotide context differ from each other, and the ratio of transition to  
123 transversion for each trinucleotide context is much greater than 1:2 for both cancer  
124 somatic mutations (ratio=2.70±0.47) and germline rare variants (ratio=3.28±0.53)  
125 (Supplementary Figure S1). These different mutation profiles may lead to different  
126 biological progressions in carcinogenesis, as depicted in several publications [19, 28].  
127 For example, the mutation profiles of melanoma are highly enriched in C-to-T  
128 transitions, indicating a direct mutagenic role of ultraviolet (UV) light in melanoma  
129 pathogenesis [30]. Thus, it is inappropriate to use conventional methods such as the

130 equal-rate model to measure selective pressure because this approach ignores the  
131 mutation bias of different nucleotide substitution types.

132

133 **Measuring selective pressure on human genes in cancer and germline cells using**  
134 **the mutation-profile-based method**

135 We therefore formulated an evolutionary approach that was designed specifically  
136 to estimate the selective pressure imposed on human genes in cancer cells and then  
137 identify genes that had undergone positive and purifying selection in cancer cells  
138 compared with in normal cells (see Figure 2 for an illustration). In cancer genomics,  
139 distinguishing synonymous from nonsynonymous somatic mutations is straightforward.  
140 We developed the mutation-profile-based method to estimate the  $C_N/C_S$  ratio of each  
141 human gene based on the mutation profiles of cancer somatic mutations and the  $p_N/p_S$   
142 ratio for germline substitutions. In contrast to the equal-rate method [27], our method  
143 considers differences in substitution rates and uses the overall mutation rate profile as  
144 the weight matrix (Figure 1).

145 We calculated the expected number of nonsynonymous and synonymous sites  
146 based on the exonic mutation rate profiles. We then counted the number of  
147 nonsynonymous and synonymous substitutions in the protein-coding region of each  
148 human gene for all cancer somatic mutations or germline substitutions. A  $\chi^2$  test was  
149 performed to identify the genes whose  $C_N/C_S$  values were either significantly greater  
150 than one or less than one, which indicates positive or negative (purifying) selection,  
151 respectively. Of the 16,953 genes with at least one germline substitution and cancer

152 somatic mutation, the overall  $C_N/C_S$  value for cancer somatic mutations  
153 (mean±s.e.=1.199±0.008) was much greater than the overall  $p_N/p_S$  of germline  
154 substitutions (mean±s.e.=0.738±0.005) (Wilcoxon test,  $p<2.2\times 10^{-16}$ ) (Table 1A,  
155 Supplementary Table S1). In the cancer genomes, 365 genes had  $C_N/C_S$  values  
156 significantly greater than one, and 923 genes had  $C_N/C_S$  values significantly less than  
157 one ( $\chi^2$  test,  $p<0.01$ , FDR<0.1). By contrast, germline substitutions included only 24  
158 genes with  $p_N/p_S$  values significantly greater than one, whereas 4,897 genes had  $p_N/p_S$   
159 values significantly less than one ( $\chi^2$  test,  $p<0.01$ , FDR<0.1). Of these 365 cancer  
160 positively selected genes, only one gene (*RSRC1*) also exhibited positive selection  
161 whereas 117 genes exhibited negative selection in germline substitutions. Additionally,  
162 500 cancer negatively selected genes did not exhibit significant negative selection in  
163 germline substitutions. These genes may therefore be under different selective pressure  
164 in cancer and germline genomes.

165 Previous studies have attributed elevated  $C_N/C_S$  values to the relaxation of  
166 purifying selection [16] or increased positive selection of globally expressed genes [17].  
167 Our results show that the number of genes under positive selection increased, whereas  
168 the number of genes under negative selection decreased, in cancer genomes compared  
169 with germline genomes. This result indicates that both the relaxation of purifying  
170 selection on passenger mutations and the positive selection of driver mutations may  
171 contribute to the increased  $C_N/C_S$  values of human genes in cancer genomes.

172

173 **Selection pressures on cancer-associated genes**

174 The Cancer Gene Census (CGC) [31, 32] contains more than 500 cancer-associated  
175 genes that have been reported in the literature to exhibit mutations and that are causally  
176 implicated in cancer development. Of those genes, 553 were included in the 16,953  
177 genes that we tested. These known cancer genes have significantly greater  $C_N/C_S$  values  
178 (Wilcoxon test,  $p=2.9\times 10^{-10}$ ) for cancer somatic mutations but significantly lower  $p_N/p_S$   
179 values for germline substitutions (Wilcoxon test,  $p<2.2\times 10^{-16}$ ) than other genes (Table  
180 1A). For selection over longer evolutionary time scales, we extracted the  $d_N/d_S$  values  
181 between human-mouse orthologs from the Ensembl database (Release 75) [33]. The  
182 known cancer genes have significantly lower human-mouse  $d_N/d_S$  values than other  
183 human genes (Wilcoxon test,  $p<2.2\times 10^{-16}$ ). These results support the work of Thomas  
184 *et al.* [34], who showed that known cancer genes may be more constrained and more  
185 important than other genes at the species and population levels, especially for  
186 oncogenes. By contrast, known cancer genes are more likely to gain onco-functional  
187 somatic mutations in cancer than other genes.

188 Among the 365 cancer positively selected genes, 45 (12.3%) genes are known  
189 cancer genes, indicating that cancer genes are significantly enriched in cancer positively  
190 selected genes (Fisher's Exact Test,  $p=6.7\times 10^{-15}$ ). When we choose a more stringent  
191 cut-off of  $p<10^{-5}$ , 17 of the 29 (58.6%) positively selected genes are known cancer genes,  
192 according to the CGC, and the work of Lawrence *et al.* [20] and Kandath *et al.* [35],  
193 such as the well-known cancer drivers *TP53*, *KRAS*, *PIK3CA*, and *BRAF*.  
194 (Supplementary Table S2). In addition, the 29 strong positively selected genes are  
195 significantly enriched in biological processes related to cancer, according to the

196 functional analysis using DAVID v6.7 [36] (Supplementary Table S3). Some cancer  
197 genes also show negative selection in cancer genomes, such as the oncogene *MLLT3*  
198 ( $C_N/C_S=0.11$ ,  $p=3.14\times 10^{-44}$ ,  $FDR=5.52\times 10^{-41}$ ). The *MLL-MLLT3* gene fusion is the  
199 main mutation type of *MLLT3* that drives tumorigenesis in acute leukemia [37].  
200 Interestingly, *MLLT3* has recurrent synonymous mutations at amino acid positions 166  
201 to 168 (S166S, 8/9155; S167S, 33/9155; S168S, 23/9155).

202 Using the  $C_N/C_S$  values, we classified known cancer genes according to the  
203 selection pressure on these genes in cancer cells, as well as their onco-functionalization  
204 in oncogenesis (Table 2). The most important two classes are oncogenes and tumor  
205 suppressor genes that are under strong positive selection, such as *TP53*, the most  
206 famous tumor suppressor gene [38], which shows strong positive selection pressure  
207 ( $C_N/C_S=32.57$ ,  $p=1.06\times 10^{-159}$ ,  $FDR=6.55\times 10^{-156}$ ). The non-synonymous mutations of  
208 *TP53* with onco-nonfunctionalization are distributed in a wide range of cancers. The  
209 oncogene *KRAS* [39] also showed a strong positive selection pressure ( $C_N/C_S=45.88$ ,  
210  $p=4.25\times 10^{-87}$ ,  $FDR=1.74\times 10^{-83}$ ). Recurrent non-synonymous mutation with onco-  
211 neofunctionalization of *KRAS* are highly enriched in codons 12 and 13; mutations in  
212 these codons represent 79.4% and 8.0% of all non-synonymous mutations of *KRAS*.

213 We also observed 12 cancer positively selected genes ( $p<10^{-5}$ ) that have not been  
214 reported as cancer-associated genes. These genes are recurrently mutated in several  
215 tumor types and are potential cancer driver genes. According to the mouse insertional  
216 mutagenesis experiments [40], three of these genes (*DMD*, *MYO9A*, and *COL5A2*) have  
217 been identified as cancer-causing genes [41-44].

218 When we chose a more stringent cut-off of  $p < 10^{-5}$  for cancer negatively selected  
219 genes, we found 112 genes that showed an enrichment in the Notch signaling pathway  
220 (Supplementary Table S3). Forty-seven of the 112 negatively selected genes showed  
221 more stringent selective constraint in cancer cells than in normal cells ( $p_N/p_S > C_N/C_S$ ,  
222  $p > 0.05$  for  $p_N/p_S$ ). It would be quite valuable to uncover the roles of these evolutionarily  
223 conserved genes in cancer cells. Out of the 47 genes, 14 genes showed a significantly  
224 increased expression level in cancers than in normal tissues (fold change  $> 2$ ,  $p < 10^{-4}$ )  
225 (Supplementary Table S4). For example, *SPRR3*, a member of the small proline-rich  
226 protein family, is under purifying selection in cancer cells ( $C_N/C_S = 0.27$ ,  $p = 5.73 \times 10^{-11}$ ,  
227  $FDR = 1.91 \times 10^{-8}$ ) and neutral selection in germline cells ( $p_N/p_S = 0.88$ ,  $p = 0.75$ ,  
228  $FDR = 0.37$ ). It has been reported that *SPRR3* is overexpressed in several tumor types,  
229 and is associated with tumor cell proliferation and invasion. Therefore, *SPRR3* could  
230 be a potential biomarker and novel therapeutic target [45-47].

231 We also examined essential genes during human development and cancer  
232 development. We extracted 2,452 human orthologs of mouse essential genes from  
233 DEG10 (the Database of Essential Genes) [48]. These genes, which are human  
234 orthologs of known essential genes in mice [49], are critical for cell survival and are  
235 therefore more conserved than other genes at the species and population levels. Here,  
236 we found that human orthologs of mouse essential genes have significantly lower  $d_N/d_S$   
237 values (measured between human-mouse orthologs) and lower  $p_N/p_S$  values for  
238 germline substitutions but similar  $C_N/C_S$  values for cancer somatic mutations compared  
239 with the values for non-essential genes (Table 1A). Human orthologs of mouse essential

240 genes are also enriched among cancer positively selected genes. Eighteen of the twenty-  
241 nine (62.1%) positively selected genes ( $p < 10^{-5}$ ) are human orthologs of mouse essential  
242 genes (Supplementary Table S2). We also used the human orthologs of mouse essential  
243 genes from OGEE (the database of Online GENE Essentiality) [50] to confirm these  
244 results (Supplementary Table S2).

245 Cancer essential genes were identified by performing genome-scale pooled RNAi  
246 screens. RNAi screens with the 45k shRNA pool in 12 cancer cell lines, including  
247 small-cell lung cancer, non-small-cell lung cancer, glioblastoma, chronic myelogenous  
248 leukemia, and lymphocytic leukemia, revealed 268 common essential genes [51].  
249 Compared to other human genes, these cancer essential genes have significantly lower  
250  $d_N/d_S$  values and lower  $p_N/p_S$  values for germline substitutions and greater  $C_N/C_S$  values  
251 for cancer somatic mutations (Table 1A), suggesting a functional shift of these genes in  
252 human populations and cancer cells.

253 We further tested the correlations of the  $d_N/d_S$ ,  $p_N/p_S$  and  $C_N/C_S$  values of human  
254 genes for human-mouse orthologs, germline substitutions and cancer somatic mutations  
255 to compare selective pressures among species, populations and cancer cells (Table 1B).  
256 For different gene sets, the  $d_N/d_S$  values show a weak positive correlation with the  $p_N/p_S$   
257 values, but no significant correlation with  $C_N/C_S$  values. The  $p_N/p_S$  values and  $C_N/C_S$   
258 values also do not have significant correlation for different gene sets. These results  
259 indicate that the cancer-specific selection pressure is quite different from the selection  
260 pressure at the species and population levels.

261

## 262 **Selection pressure among different cancer types**

263 As cancer is highly heterogeneous, we further analyzed the selection pressure of  
264 human genes in different cancer types. The 9,155 tumor samples from the ICGC  
265 database could be classified as 20 cancer types according to the primary site. The  
266 overall  $C_N/C_S$  values for the cancer somatic mutations in the different cancer types  
267 ranged from  $1.078 \pm 0.022$  to  $1.827 \pm 0.013$  (mean  $\pm$  s.e., Table 3). The detected positively  
268 and negatively selected genes ( $\chi^2$  test,  $p < 0.01$ ) varied in the different cancer types  
269 (Supplementary Table S5). Due to the limited number of tumor samples and somatic  
270 mutations for each cancer type, particularly in the cancer types with low mutation rates,  
271 our method might not be sensitive enough to detect the selection pressure for each gene.  
272 For example, only one positively selected gene was detected in bone cancer (IDH1) and  
273 nervous system cancer (ALK), respectively. There were also three genes (TP53,  
274 PIK3CA and KRAS) that showed positive selection in more than five cancer types. In  
275 particular, TP53 showed positive selection in 15 cancer types. On the other hand, more  
276 genes (164/188, 87.2%) were under positive selection in only one cancer type. We also  
277 found that six genes (TBP, EP400, DSPP, MUC21, MLLT3, and MUC2) were under  
278 negative selection in more than five cancer types. These genes also showed negative  
279 selection at the species and population levels. Furthermore, 85.8% (2,417/2,817) of  
280 genes showed negative selection in only one cancer type. These results indicate the  
281 divergence of selection pressure in different cancer types.

282

## 283 **Comparison of the equal-rate model and empirical mutation profile model**

284        Considering that different nucleotide substitution models might provide varying  
285 estimates, we used the equal-rate method [27] as the simplest model to calculate the  
286 expected numbers of nonsynonymous and synonymous sites. The overall  $C_N/C_S$  value  
287 for cancer somatic mutations (mean $\pm$ s.e.=0.892 $\pm$ 0.006) is greater than the  $p_N/p_S$  value  
288 for germline substitutions (mean $\pm$ s.e.=0.633 $\pm$ 0.004) for the 16,953 genes  
289 (Supplementary Table S1) but lower than that calculated using the mutation-profile-  
290 based method (Wilcoxon test,  $p < 2.2 \times 10^{-16}$ ) (Figure 3A). Consequently, the number of  
291 genes with  $C_N/C_S$  values  $> 1$  ( $\chi^2$  test,  $p < 0.01$ , FDR $< 0.1$ ) is much lower than those  
292 calculated using the exonic mutation profiles (37 versus 365), whereas the number of  
293 genes with  $C_N/C_S$  values  $< 1$  ( $\chi^2$  test,  $p < 0.01$ , FDR $< 0.1$ ) is much greater (2851 versus  
294 923) (Figure 3B and 3C).

295        We also used the intergenic mutation rate profile from 2,900 tumor-normal whole  
296 genome sequences, which are included in the 9,155 cancer samples of ICGC database,  
297 to calculate the  $C_N/C_S$  value for cancer somatic mutations. The overall  $C_N/C_S$  value  
298 (mean $\pm$ s.e.=1.503 $\pm$ 0.010) is greater than that calculated from the exonic mutation rate  
299 profile (mean $\pm$ s.e.=1.199 $\pm$ 0.008) (Wilcoxon test,  $p < 2.2 \times 10^{-16}$ ), resulting in more  
300 positively selected genes (1526 versus 365) and fewer negatively selected genes (298  
301 versus 923) (Figure 3B and 3C).

302        The equal-rate method ignores the mutation rate bias between different substitution  
303 types, especially the ratio of transition to transversion, leading to underestimation of  
304 the  $C_N/C_S$  ratio. Therefore, the equal-rate method is strict for positive selection detection  
305 but relaxed for the detection of negative selection [52]. In contrast, the mutation-profile-

306 based method considers the mutation bias, which can be depicted as the internal  
307 variance between mutation rates of different substitution types. Thus, the mutation-  
308 profile-based method can correct the underestimation of the  $C_N/C_S$  ratio estimated by  
309 the equal-rate method. Furthermore, the mutation-profile-based method would also  
310 increase the false-positive results for detecting positively selected genes but be more  
311 conservative in detecting negatively selected genes. The mutation bias may simulate  
312 the detection of genes under strong selection pressure but may suppress the detection  
313 of genes under weak selection pressure.

314

## 315 **DISCUSSION**

316 A key goal of cancer research is to identify cancer-associated genes, such as  
317 oncogenes and tumor suppressor genes, that might promote tumor occurrence and  
318 progression when mutated [28]. Instead of searching for cancer-causing genes with  
319 multiple driver mutations, an alternative approach is to identify cancer essential genes  
320 in tumor cell populations because they are crucial for carcinogenesis, progression and  
321 metastasis. Cancer essential genes are important for the growth and survival of cancer  
322 cells [51] and are expected to be highly conserved in cancer cells. In this study, we  
323 aimed to detect both cancer-specific positively and negatively selected genes using a  
324 molecular evolution approach.

325 Based on analyses of large-scale cancer somatic mutation data derived from The  
326 Cancer Genome Atlas (TCGA) or International Cancer Genome Consortium (ICGC),  
327 previous studies identified important differences between the evolutionary dynamics of

328 cancer somatic cells and whole organisms [6, 16, 18]. However, these studies applied  
329 canonical nucleotide substitution models to identify the molecular signatures of natural  
330 selection in cancer cells or human populations and neglected the apparently different  
331 mutation profiles of these cell types. Here, we developed a new mutation-profile-based  
332 method to calculate the  $C_N/C_S$  values of human genes for cancer somatic mutations. In  
333 our results, a large number of known cancer genes did not show significant positive  
334 selection according to our analysis. One possible reason for this finding suggests that  
335 positive selection for driver mutations is obscured by the relaxed purifying selection of  
336 passenger mutations. Additionally, among the strong positively selected genes, more  
337 than half are known cancer genes. Another possible reason might be that the main  
338 mutation type of more than 300 cancer-associated genes is translocation or copy number  
339 variation, rather than point mutation. Furthermore, some of the positively selected  
340 genes might also be related to cancer, such as *DMD*, *MYO9A*, and *COL5A2*, which have  
341 been identified as cancer-causing genes based on mouse insertional mutagenesis  
342 experiments [40].

343 Two prerequisites are crucial to properly apply the mutation-profile-based method.  
344 First, a large number of samples with similar mutation profiles are necessary to increase  
345 the power of selection pressure detection. Second, a subset of nucleotide substitutions  
346 should be chosen to represent the background neutral mutation profiles of the samples.  
347 In this study, because of the limited number of cancer samples, especially the number  
348 of whole-genome sequenced tumor-normal tissue pairs, we pooled all samples to  
349 analyze pan-cancer-level selection pressures. However, cancer somatic mutation

350 profiles are well known to be heterogeneous among different cancer types, even for  
351 samples with the same tissue origin [19, 20, 28, 35]. As the number of sequenced cancer  
352 genomes increases, we will be able to classify cancer samples by their specific mutation  
353 profiles and infer evolutionarily selective pressures more precisely using the mutation-  
354 profile-based method.

355 Background neutral mutation profiles can be calculated based on intergenic regions  
356 from the corresponding samples. In this study, we assumed that most of the exonic  
357 somatic mutations in the cancer samples do not have significant effects on the fitness  
358 of cancer cells. Under this assumption, we can apply the mutation profiles of coding  
359 regions to approximate the background. The exonic mutation profiles used in our  
360 mutation-profile-based method considered the weight of the 96 substitution  
361 classifications within the cancer exomes, which may reflect the mutation bias of  
362 different substitution types within the protein-coding regions. This method would  
363 correct the underestimation of the  $C_N/C_S$  value that occurs with the equal-rate method  
364 [52]. The mutation-profile-based method is more sensitive for the detection of positive  
365 selection but more conservative for the detection of negative selection compared with  
366 the equal-rate method. As more tumor-normal whole genome sequence data become  
367 available, it would be better to choose suitable mutation profiles for the mutation-  
368 profile-based method. With the expansion of these data in the future, we may apply  
369 more precise methods to identify neutral background mutation properties.

370

## 371 MATERIALS AND METHODS

## 372 **Datasets**

373 Cancer somatic mutation data from 9,155 cancer samples corresponding to 20  
374 primary sites were extracted from the ICGC Data Portal (<http://dcc.icgc.org>, Release  
375 20), which includes 36,985,985 somatic mutations and small insertions/deletions. Data  
376 on rare human protein-coding variants (minor allele frequency <0.01%) from 6,500  
377 human exomes (ESP6500) were extracted from the NHLBI GO Exome Sequencing  
378 Project (<http://evs.gs.washington.edu/EVS>, Exome Variant Server NGESPE, Seattle,  
379 WA). A total of 572 known cancer genes were extracted from the Cancer Gene Census  
380 (<http://cancer.sanger.ac.uk/cancergenome/projects/census/>, COSMIC v72) [31, 32].

381 Human gene sequences and annotations were extracted from the Ensembl database  
382 (GRCh37, Release 75) [33]. For each gene, only the longest transcript was selected for  
383 the subsequent analyses. The  $d_N/d_S$  values between human-mouse orthologs were  
384 extracted from the Ensembl database. The HGNC (HUGO Gene Nomenclature  
385 Committee) database [53] (<http://www.genenames.org/>) and the Genecards database  
386 [54] (<http://www.genecards.org>) were used to map the gene IDs from different datasets.  
387 DAVID v6.7 was utilized for the functional annotation analysis [36]. ANNOVAR was  
388 utilized to perform biological and functional annotations of the cancer somatic  
389 mutations and germline substitutions [55]. The OncoPrint database [56]  
390 (<https://www.oncoPrint.org>) was used to compare the gene expression level of  
391 negatively selected genes between cancer and normal tissues. The human orthologs of  
392 mouse essential genes were extracted from the DEG 10 [48] and the OGEE v2 [50]  
393 databases.

394

395 **Statistical measure for gene-specific selection pressure in cancer evolution ( $C_N/C_S$ )**

396 In cancer genomics, distinguishing synonymous from nonsynonymous somatic  
397 mutations is straightforward. Thus, given a set of independent cancer samples, the ratio  
398 of nonsynonymous counts (N) to synonymous counts (S) of a gene, denoted by  $N/S$ , is  
399 simply given by the sum over all samples, under the assumption of no double mutations  
400 at the same nucleotide site (e.g., for the observed mutation A>C, the mutation path  
401 A>G>C is almost impossible in cancers). To further explore the underlying mechanisms,  
402 the  $N/S$  ratio must be normalized by  $L_N/L_S$ , that is,

403 
$$C_N/C_S=(N/L_N)/(S/L_S)=q_N/q_S \quad (1)$$

404 where  $L_N$  is the number of expected nonsynonymous sites and  $L_S$  is the number of  
405 expected synonymous sites. Note that  $C_N/C_S$  is specific for cancer somatic mutations,  
406 to avoid notation confusions with  $d_N/d_S$  in molecular evolution and  $p_N/p_S$  in population  
407 genetics. To avoid a calculation error for the small sample size, 0.5 was added to each  
408 parameter for the calculation of  $C_N/C_S$  if N or S was equal to zero.

409 The calculation of  $L_N$  and  $L_S$  from the nucleotide sequence is not a trivial task. For  
410 instance, in the codon TTT (coding for amino acid Phe), the first two positions are  
411 counted as nonsynonymous sites because no synonymous changes can occur at these  
412 positions. At the third position, the transition change (T>C) is synonymous, whereas  
413 the remaining two transversion changes (T>A and T>G) are nonsynonymous.  
414 Apparently, the weight of the third position of codon TTT as synonymous ( $w_S$ ) or  
415 nonsynonymous ( $w_N$ ) depends on the pattern of somatic mutations. At one extreme, if

416 the transition mutation is dominant, this position should nearly be counted as a  
417 synonymous site ( $w_S=1$ ); at the other extreme (transversion dominant), this position  
418 would be counted as a nonsynonymous site ( $w_S=0$ ).

419

#### 420 **Equal-rate model**

421 The weight of a nucleotide as synonymous ( $w_S$ ) is simple when the rate of base  
422 change is the same. Let  $I_S$  be the number of possible synonymous changes at a site. This  
423 is counted as  $w_S=I_S/3$  synonymous and  $(1- I_S/3)$  nonsynonymous. For instance, in the  
424 codon TTT (Phe), the first two positions are counted as nonsynonymous sites because  
425 no synonymous changes can occur at these positions ( $w_S=0$ ). The third position of  
426 codon TTT is then counted as one third of a synonymous site ( $w_S=1/3$ ) and two-thirds  
427 of a nonsynonymous site ( $w_N=2/3$ ) because only one of the three possible changes is  
428 synonymous. It is then straightforward to calculate the numbers of synonymous and  
429 nonsynonymous sites.

430

#### 431 **Empirical mutation profile model**

432 Substantial evidence has demonstrated that the rate of somatic mutations in cancer  
433 depends on not only the nucleotide site (e.g., synonymous or nonsynonymous sites) and  
434 the mutation type (e.g., transition or transversion) but also on the sequence context of  
435 each mutated site, i.e., the effects of near-by nucleotides on somatic mutations are  
436 nontrivial. Recent studies [28, 29, 57] proposed an empirical mutation profile of any  
437 position with base P, considering two immediate neighbor nucleotides (x, y) of a

438 trinucleotide string denoted by  $xPy$ . Since base P has six base-change patterns (under  
439 Watson-Crick pairing) and both x and y have four types of bases, there are a total of  
440  $4 \times 6 \times 4 = 96$  substitution classifications, with the empirical profile denoted by  
441  $M(xPy \rightarrow xP_i^*y)$ , where  $P_i^*$  ( $i=1,2,3$ ) for the other three bases instead of P. To determine  
442 the probability of the mutation type ( $xPy \rightarrow xP_i^*y$ ), we divided the number of mutations  
443 in that trinucleotide context ( $xPy \rightarrow xP_i^*y$ ) by the number of occurrences of the  
444 trinucleotide ( $xPy$ ). Our computational pipeline is illustrated by the following example.

445 In the encoding sequence with two codons ... TTT-ATG..., we consider the third  
446 position of codon TTT (Phe). Under the trinucleotide TTA for the mutation profile (not  
447 the codon), the corresponding three substitution configurations are given by  
448  $M(TTA \rightarrow TCA)$ ,  $M(TTA \rightarrow TAA)$  and  $M(TTA \rightarrow TGA)$ , respectively, and the number of  
449 occurrences of TTA is  $M(TTA)$ . Next, we consider codon TTT. Because TTT and TTC  
450 are synonymous codons but TTA and TTG are not, the probabilities that this site will  
451 be synonymous and nonsynonymous are simply given by the following:

$$452 \quad w_S = M(TTA \rightarrow TCA) / M(TTA)$$

$$453 \quad w_N = (M(TTA \rightarrow TGA) + M(TTA \rightarrow TAA)) / M(TTA) \quad (2)$$

454 We counted all somatic mutations in the protein-coding regions of the 9,155 tumor-  
455 normal paired cancer samples, as well as all the rare protein-coding variants of the  
456 ESP6500 dataset. The mutation profiles were depicted as the mutation rate of each  
457 mutation type according to the 96 substitution classifications.

458 The ratio of transition to transversion for each trinucleotide context was calculated  
459 based on the mutation rate of transitions and transversions. For example, the ratio of

460 transition to transversion for  $ACA = M(ACA > ATA) / (M(ACA > AAA) + M(ACA > AGA))$ .

461

## 462 **Detection of positive and negative selections**

463 The  $\chi^2$  test was used to compare the number of nonsynonymous and synonymous  
464 substitutions to the number of nonsynonymous and synonymous sites for each gene to  
465 test the statistical significance of the difference between the  $C_N/C_S$  values and one.  
466 Genes with  $C_N/C_S$  values significantly greater than one were classified as under positive  
467 selection in tumors, whereas genes with  $C_N/C_S$  values significantly less than one were  
468 classified as under negative, or purifying, selection. The false-discovery rate was  
469 estimated using the qvalue package from Bioconductor [58]. The software tool R was  
470 used for statistical analysis (<http://www.r-project.org/>).

471

## 472 **ACKNOWLEDGEMENTS**

473 We are grateful to Xiaopu Wang for his help with the manuscript preparation. We  
474 would like to thank the NHLBI GO Exome Sequencing Project and its ongoing studies  
475 which produced and provided the exome variant calls for comparison: the Lung GO  
476 Sequencing Project (HL-102923), the WHI Sequencing Project (HL-102924), the  
477 Broad GO Sequencing Project (HL-102925), the Seattle GO Sequencing Project (HL-  
478 102926) and the Heart GO Sequencing Project (HL-103010). We also gratefully  
479 acknowledge the clinical contributors and data producers from the International Cancer  
480 Genome Consortium (ICGC) for referencing the ICGC datasets.

481

482 **CONFLICTS OF INTEREST**

483 The authors declare that they have no conflicts of interest.

484

485 **GRANT SUPPORT**

486 This work was supported by grants from the Ministry of Science and Technology  
487 China (2012CB910101), the National Natural Science Foundation of China (31272299,  
488 31301034, 31501021), the Zhejiang Provincial Natural Sciences Foundation of China  
489 (LY15C060001), the Shanghai Pujiang Program (13PJD005), the China Postdoctoral  
490 Science Foundation (2013M531117), the Fundamental Research Funds for the Central  
491 Universities, and the Open Research Funds of the State Key Laboratory of Genetic  
492 Engineering, Fudan University.

493

494 **REFERENCES**

- 495 1. Cairns J. Mutation selection and the natural history of cancer. *Nature*. 1975; 255(5505):197-200.
- 496 2. Nowell PC. The clonal evolution of tumor cell populations. *Science*. 1976; 194(4260):23-28.
- 497 3. Crespi BJ and Summers K. Positive selection in the evolution of cancer. *Biol Rev Camb Philos Soc*.  
498 2006; 81(3):407-424.
- 499 4. Merlo LM, Pepper JW, Reid BJ and Maley CC. Cancer as an evolutionary and ecological process.  
500 *Nat Rev Cancer*. 2006; 6(12):924-935.
- 501 5. Podlaha O, Riester M, De S and Michor F. Evolution of the cancer genome. *Trends Genet*. 2012;  
502 28(4):155-163.
- 503 6. Yates LR and Campbell PJ. Evolution of the cancer genome. *Nat Rev Genet*. 2012; 13(11):795-806.
- 504 7. Greaves M and Maley CC. Clonal evolution in cancer. *Nature*. 2012; 481(7381):306-313.
- 505 8. Luo J, Solimini NL and Elledge SJ. Principles of cancer therapy: oncogene and non-oncogene  
506 addiction. *Cell*. 2009; 136(5):823-837.
- 507 9. Hanahan D and Weinberg RA. Hallmarks of Cancer: The Next Generation. *Cell*. 2011; 144(5):646-  
508 674.
- 509 10. Stratton MR, Campbell PJ and Futreal PA. The cancer genome. *Nature*. 2009; 458(7239):719-724.
- 510 11. Greenman C, Stephens P, Smith R, Dalgliesh GL, Hunter C, Bignell G, Davies H, Teague J, Butler  
511 A, Stevens C, Edkins S, O'Meara S, Vastrik I, Schmidt EE, Avis T, Barthorpe S, et al. Patterns of  
512 somatic mutation in human cancer genomes. *Nature*. 2007; 446(7132):153-158.
- 513 12. Beckman RA and Loeb LA. Negative clonal selection in tumor evolution. *Genetics*. 2005;  
514 171(4):2123-2131.

- 515 13. McFarland CD, Korolev KS, Kryukov GV, Sunyaev SR and Mirny LA. Impact of deleterious  
516 passenger mutations on cancer progression. *Proc Natl Acad Sci USA*. 2013; 110(8):2910-2915.
- 517 14. Croce CM. Molecular origins of cancer: Oncogenes and cancer. *New Engl J Med*. 2008; 358(5):502-  
518 511.
- 519 15. Sherr CJ. Principles of tumor suppression. *Cell*. 2004; 116(2):235-246.
- 520 16. Woo YH and Li WH. DNA replication timing and selection shape the landscape of nucleotide  
521 variation in cancer genomes. *Nat Commun*. 2012; 3:1004.
- 522 17. Ostrow SL, Barshir R, DeGregori J, Yeger-Lotem E and Hershberg R. Cancer Evolution Is  
523 Associated with Pervasive Positive Selection on Globally Expressed Genes. *PLoS Genet*. 2014;  
524 10(3):e1004239.
- 525 18. Davoli T, Xu AW, Mengwasser KE, Sack LM, Yoon JC, Park PJ and Elledge SJ. Cumulative  
526 haploinsufficiency and triplosensitivity drive aneuploidy patterns and shape the cancer genome. *Cell*.  
527 2013; 155(4):948-962.
- 528 19. Lawrence MS, Stojanov P, Polak P, Kryukov GV, Cibulskis K, Sivachenko A, Carter SL, Stewart C,  
529 Mermel CH, Roberts SA, Kiezun A, Hammerman PS, McKenna A, Drier Y, Zou L, Ramos AH, et  
530 al. Mutational heterogeneity in cancer and the search for new cancer-associated genes. *Nature*. 2013;  
531 499(7457):214-218.
- 532 20. Lawrence MS, Stojanov P, Mermel CH, Robinson JT, Garraway LA, Golub TR, Meyerson M,  
533 Gabriel SB, Lander ES and Getz G. Discovery and saturation analysis of cancer genes across 21  
534 tumour types. *Nature*. 2014; 505(7484):495-501.
- 535 21. Chen H, Xing K and He X. The dJ/dS Ratio Test Reveals Hundreds of Novel Putative Cancer Drivers.  
536 *Mol Biol Evol*. 2015; 32(8):2181-2185.
- 537 22. Du W, Jiang P, Mancuso A, Stonestrom A, Brewer MD, Minn AJ, Mak TW, Wu M and Yang X.  
538 TAp73 enhances the pentose phosphate pathway and supports cell proliferation. *Nat Cell Biol*. 2013;  
539 15(8):991-1000.
- 540 23. Endo T, Ikeo K and Gojobori T. Large-scale search for genes on which positive selection may  
541 operate. *Mol Biol Evol*. 1996; 13(5):685-690.
- 542 24. Arbiza L, Dopazo J and Dopazo H. Positive selection, relaxation, and acceleration in the evolution  
543 of the human and chimp genome. *PLoS Comput Biol*. 2006; 2(4):e38.
- 544 25. Ovens K and Naugler C. Preliminary evidence of different selection pressures on cancer cells as  
545 compared to normal tissues. *Theor Biol Med Model*. 2012; 9:44.
- 546 26. Martincorena I, Roshan A, Gerstung M, Ellis P, Van Loo P, McLaren S, Wedge DC, Fullam A,  
547 Alexandrov LB, Tubio JM, Stebbings L, Menzies A, Widaa S, Stratton MR, Jones PH and Campbell  
548 PJ. Tumor evolution. High burden and pervasive positive selection of somatic mutations in normal  
549 human skin. *Science*. 2015; 348(6237):880-886.
- 550 27. Nei M and Gojobori T. Simple methods for estimating the numbers of synonymous and  
551 nonsynonymous nucleotide substitutions. *Mol Biol Evol*. 1986; 3(5):418-426.
- 552 28. Alexandrov LB, Nik-Zainal S, Wedge DC, Aparicio SAJR, Behjati S, Biankin AV, Bignell GR, Bolli  
553 N, Borg A, Børresen-Dale A-L, Boyault S, Burkhardt B, Butler AP, Caldas C, Davies HR, Desmedt  
554 C, et al. Signatures of mutational processes in human cancer. *Nature*. 2013; 500(7463):415-421.
- 555 29. Samocha KE, Robinson EB, Sanders SJ, Stevens C, Sabo A, McGrath LM, Kosmicki JA, Rehnstrom  
556 K, Mallick S, Kirby A, Wall DP, MacArthur DG, Gabriel SB, DePristo M, Purcell SM, Palotie A, et  
557 al. A framework for the interpretation of de novo mutation in human disease. *Nat Genet*. 2014;  
558 46(9):944-950.

- 559 30. Hodis E, Watson IR, Kryukov GV, Arold ST, Imielinski M, Theurillat JP, Nickerson E, Auclair D,  
560 Li L, Place C, Dicara D, Ramos AH, Lawrence MS, Cibulskis K, Sivachenko A, Voet D, et al. A  
561 landscape of driver mutations in melanoma. *Cell*. 2012; 150(2):251-263.
- 562 31. Futreal PA, Coin L, Marshall M, Down T, Hubbard T, Wooster R, Rahman N and Stratton MR. A  
563 census of human cancer genes. *Nat Rev Cancer*. 2004; 4(3):177-183.
- 564 32. Forbes SA, Bindal N, Bamford S, Cole C, Kok CY, Beare D, Jia M, Shepherd R, Leung K, Menzies  
565 A, Teague JW, Campbell PJ, Stratton MR and Futreal PA. COSMIC: mining complete cancer  
566 genomes in the Catalogue of Somatic Mutations in Cancer. *Nucleic Acids Res*. 2011; 39(Database  
567 issue):D945-950.
- 568 33. Flicek P, Amode MR, Barrell D, Beal K, Billis K, Brent S, Carvalho-Silva D, Clapham P, Coates G,  
569 Fitzgerald S, Gil L, Giron CG, Gordon L, Hourlier T, Hunt S, Johnson N, et al. Ensembl 2014.  
570 *Nucleic Acids Res*. 2014; 42(Database issue):D749-755.
- 571 34. Thomas MA, Weston B, Joseph M, Wu W, Nekrutenko A and Tonellato PJ. Evolutionary dynamics  
572 of oncogenes and tumor suppressor genes: higher intensities of purifying selection than other genes.  
573 *Mol Biol Evol*. 2003; 20(6):964-968.
- 574 35. Kandath C, McLellan MD, Vandin F, Ye K, Niu B, Lu C, Xie M, Zhang Q, McMichael JF,  
575 Wyczalkowski MA, Leiserson MDM, Miller CA, Welch JS, Walter MJ, Wendl MC, Ley TJ, et al.  
576 Mutational landscape and significance across 12 major cancer types. *Nature*. 2013; 502(7471):333-  
577 339.
- 578 36. Huang da W, Sherman BT and Lempicki RA. Systematic and integrative analysis of large gene lists  
579 using DAVID bioinformatics resources. *Nat Protoc*. 2009; 4(1):44-57.
- 580 37. Meyer C, Hofmann J, Burmeister T, Groger D, Park TS, Emerenciano M, Pombo de Oliveira M,  
581 Renneville A, Villarese P, Macintyre E, Cave H, Clappier E, Mass-Malo K, Zuna J, Trka J, De  
582 Braekeleer E, et al. The MLL recombinome of acute leukemias in 2013. *Leukemia*. 2013;  
583 27(11):2165-2176.
- 584 38. Aubrey BJ, Strasser A and Kelly GL. Tumor-Suppressor Functions of the TP53 Pathway. *Cold  
585 Spring Harb Perspect Med*. 2016; 6(5).
- 586 39. Morris JPt, Wang SC and Hebrok M. KRAS, Hedgehog, Wnt and the twisted developmental biology  
587 of pancreatic ductal adenocarcinoma. *Nat Rev Cancer*. 2010; 10(10):683-695.
- 588 40. Ranzani M, Annunziato S, Adams DJ and Montini E. Cancer gene discovery: exploiting insertional  
589 mutagenesis. *Mol Cancer Res*. 2013; 11(10):1141-1158.
- 590 41. Abbott KL, Nyre ET, Abrahante J, Ho YY, Isaksson Vogel R and Starr TK. The Candidate Cancer  
591 Gene Database: a database of cancer driver genes from forward genetic screens in mice. *Nucleic  
592 Acids Res*. 2015; 43(Database issue):D844-848.
- 593 42. Takeda H, Wei Z, Koso H, Rust AG, Yew CC and Mann MB. Transposon mutagenesis identifies  
594 genes and evolutionary forces driving gastrointestinal tract tumor progression. *Nat Genet*. 2015;  
595 47(2):142-150.
- 596 43. Rahrmann EP, Watson AL, Keng VW, Choi K, Moriarity BS, Beckmann DA, Wolf NK, Sarver A,  
597 Collins MH, Moertel CL, Wallace MR, Gel B, Serra E, Ratner N and Largaespada DA. Forward  
598 genetic screen for malignant peripheral nerve sheath tumor formation identifies new genes and  
599 pathways driving tumorigenesis. *Nat Genet*. 2013; 45(7):756-766.
- 600 44. Perez-Mancera PA, Rust AG, van der Weyden L, Kristiansen G, Li A, Sarver AL, Silverstein KA,  
601 Grutzmann R, Aust D, Rummele P, Knosel T, Herd C, Stemple DL, Kettleborough R, Brosnan JA,  
602 Li A, et al. The deubiquitinase USP9X suppresses pancreatic ductal adenocarcinoma. *Nature*. 2012;

- 603 486(7402):266-270.
- 604 45. Liu Q, Zhang C, Ma G and Zhang Q. Expression of SPRR3 is associated with tumor cell  
605 proliferation and invasion in glioblastoma multiforme. *Oncol Lett.* 2014; 7(2):427-432.
- 606 46. Kim JC, Yu JH, Cho YK, Jung CS, Ahn SH, Gong G, Kim YS and Cho DH. Expression of SPRR3  
607 is associated with tumor cell proliferation in less advanced stages of breast cancer. *Breast Cancer*  
608 *Res Treat.* 2012; 133(3):909-916.
- 609 47. de AST, Souza-Santos PT, de Oliveira DS, Bernardo V, Lima SC, Rapozo DC, Krueel CD, Faria PA,  
610 Ribeiro Pinto LF and Albano RM. Quantitative evaluation of SPRR3 expression in esophageal  
611 squamous cell carcinoma by qPCR and its potential use as a biomarker. *Exp Mol Pathol.* 2011;  
612 91(2):584-589.
- 613 48. Luo H, Lin Y, Gao F, Zhang CT and Zhang R. DEG 10, an update of the database of essential genes  
614 that includes both protein-coding genes and noncoding genomic elements. *Nucleic Acids Res.* 2014;  
615 42(Database issue):D574-580.
- 616 49. Georgi B, Voight BF and Bucan M. From mouse to human: evolutionary genomics analysis of  
617 human orthologs of essential genes. *PLoS Genet.* 2013; 9(5):e1003484.
- 618 50. Chen WH, Lu G, Chen X, Zhao XM and Bork P. OGEE v2: an update of the online gene essentiality  
619 database with special focus on differentially essential genes in human cancer cell lines. *Nucleic*  
620 *Acids Res.* 2017; 45(D1):D940-D944.
- 621 51. Luo B, Cheung HW, Subramanian A, Sharifnia T, Okamoto M, Yang X, Hinkle G, Boehm JS,  
622 Beroukhim R, Weir BA, Mermel C, Barbie DA, Awad T, Zhou X, Nguyen T, Piqani B, et al. Highly  
623 parallel identification of essential genes in cancer cells. *Proc Natl Acad Sci USA.* 2008;  
624 105(51):20380-20385.
- 625 52. Yang Z and Bielawski JP. Statistical methods for detecting molecular adaptation. *Trends Ecol Evol.*  
626 2000; 15(12):496-503.
- 627 53. Gray KA, Daugherty LC, Gordon SM, Seal RL, Wright MW and Bruford EA. Genenames.org: the  
628 HGNC resources in 2013. *Nucleic Acids Res.* 2013; 41(Database issue):D545-552.
- 629 54. Safran M, Dalah I, Alexander J, Rosen N, Iny Stein T, Shmoish M, Nativ N, Bahir I, Doniger T,  
630 Krug H, Sirota-Madi A, Olender T, Golan Y, Stelzer G, Harel A and Lancet D. GeneCards Version  
631 3: the human gene integrator. *Database (Oxford).* 2010; 2010:baq020.
- 632 55. Wang K, Li M and Hakonarson H. ANNOVAR: functional annotation of genetic variants from high-  
633 throughput sequencing data. *Nucleic Acids Res.* 2010; 38(16):e164.
- 634 56. Rhodes DR, Kalyana-Sundaram S, Mahavisno V, Varambally R, Yu J, Briggs BB, Barrette TR,  
635 Anstet MJ, Kincaid-Beal C, Kulkarni P, Varambally S, Ghosh D and Chinnaiyan AM. OncoPrint  
636 3.0: genes, pathways, and networks in a collection of 18,000 cancer gene expression profiles.  
637 *Neoplasia.* 2007; 9(2):166-180.
- 638 57. Krawczak M, Ball EV and Cooper DN. Neighboring-nucleotide effects on the rates of germ-line  
639 single-base-pair substitution in human genes. *Am J Hum Genet.* 1998; 63(2):474-488.
- 640 58. Storey JD and Tibshirani R. Statistical significance for genomewide studies. *Proc Natl Acad Sci*  
641 *USA.* 2003; 100(16):9440-9445.

642

643

644 **FIGURE LEGENDS**

645 **Figure 1.** Mutation profiles of cancer somatic mutations and germline substitutions,  
646 including the exonic mutation profile of 9,155 cancer samples and the exonic mutation  
647 profile of ESP6500.

648 **Figure 2.** The pipeline used to identify positively and negatively selected cancer genes  
649 with the mutation-profile-based method.

650 **Figure 3.** The overall omega ratio (A) and overlap of cancer positively selected (B) and  
651 negatively selected (C) genes based on different models.

652

653 **TABLES**

654 **Table 1.** The  $\omega$  ratios ( $d_N/d_S$ ,  $p_N/p_S$ ,  $C_N/C_S$  values) (A) and the correlations of the  $\omega$   
 655 ratios (B) for the different gene sets for the human-mouse orthologs and for germline  
 656 and cancer somatic mutations. The positively and negatively selected genes indicates  
 657 the genes that are under positive and negative selection in cancer cells, respectively ( $\chi^2$   
 658 test,  $p < 0.01$ ,  $FDR < 0.1$ ).

659 (A)

|                             | $d_N/d_S$     | $p_N/p_S$     | $C_N/C_S$     |
|-----------------------------|---------------|---------------|---------------|
| All genes                   | 0.154 ± 0.006 | 0.738 ± 0.005 | 1.199 ± 0.008 |
| Cancer genes                | 0.106 ± 0.005 | 0.537 ± 0.014 | 1.550 ± 0.116 |
| Oncogenes                   | 0.088 ± 0.009 | 0.473 ± 0.029 | 1.940 ± 0.508 |
| Tumor suppressor genes      | 0.121 ± 0.017 | 0.545 ± 0.037 | 1.994 ± 0.497 |
| Human essential genes(DEGs) | 0.092 ± 0.002 | 0.559 ± 0.007 | 1.217 ± 0.032 |
| Cancer essential genes      | 0.090 ± 0.008 | 0.587 ± 0.030 | 1.465 ± 0.190 |
| Positively selected genes   | 0.137 ± 0.007 | 0.757 ± 0.029 | 3.264 ± 0.198 |
| Negatively selected genes   | 0.129 ± 0.004 | 0.600 ± 0.012 | 0.471 ± 0.005 |

660 (B)

|                            | $d_N/d_S$ vs $p_N/p_S$ |                         | $d_N/d_S$ vs $C_N/C_S$ |         | $p_N/p_S$ vs $C_N/C_S$ |                         |
|----------------------------|------------------------|-------------------------|------------------------|---------|------------------------|-------------------------|
|                            | r                      | p-value                 | r                      | p-value | r                      | p-value                 |
| All genes                  | 0.03                   | $6.6 \times 10^{-5}$    | 0.00                   | 0.80    | 0.09                   | $< 2.2 \times 10^{-16}$ |
| Known cancer genes         | 0.38                   | $< 2.2 \times 10^{-16}$ | -0.01                  | 0.87    | -0.04                  | 0.36                    |
| Oncogenes                  | 0.12                   | 0.23                    | -0.01                  | 0.94    | -0.05                  | 0.65                    |
| Tumor suppressor genes     | 0.34                   | $3.5 \times 10^{-3}$    | 0.02                   | 0.86    | -0.07                  | 0.56                    |
| Human essential genes(DEG) | 0.32                   | $< 2.2 \times 10^{-16}$ | 0.01                   | 0.58    | 0.04                   | 0.06                    |
| Cancer essential genes     | 0.20                   | $1.4 \times 10^{-3}$    | -0.02                  | 0.73    | -0.04                  | 0.50                    |
| Positively selected genes  | 0.34                   | $1.0 \times 10^{-10}$   | -0.02                  | 0.67    | 0.13                   | 0.01                    |
| Negatively selected genes  | 0.35                   | $< 2.2 \times 10^{-16}$ | -0.04                  | 0.22    | -0.07                  | 0.03                    |

661

662

663

664 **Table 2.** Classification of cancer genes according to cancer-specific selection pressures

|                        | #Positive selection | #Negative selection | #Neutral |
|------------------------|---------------------|---------------------|----------|
| Known cancer genes     | 45                  | 29                  | 479      |
| Oncogenes              | 11                  | 7                   | 79       |
| Tumor suppressor genes | 10                  | 6                   | 54       |

665

666

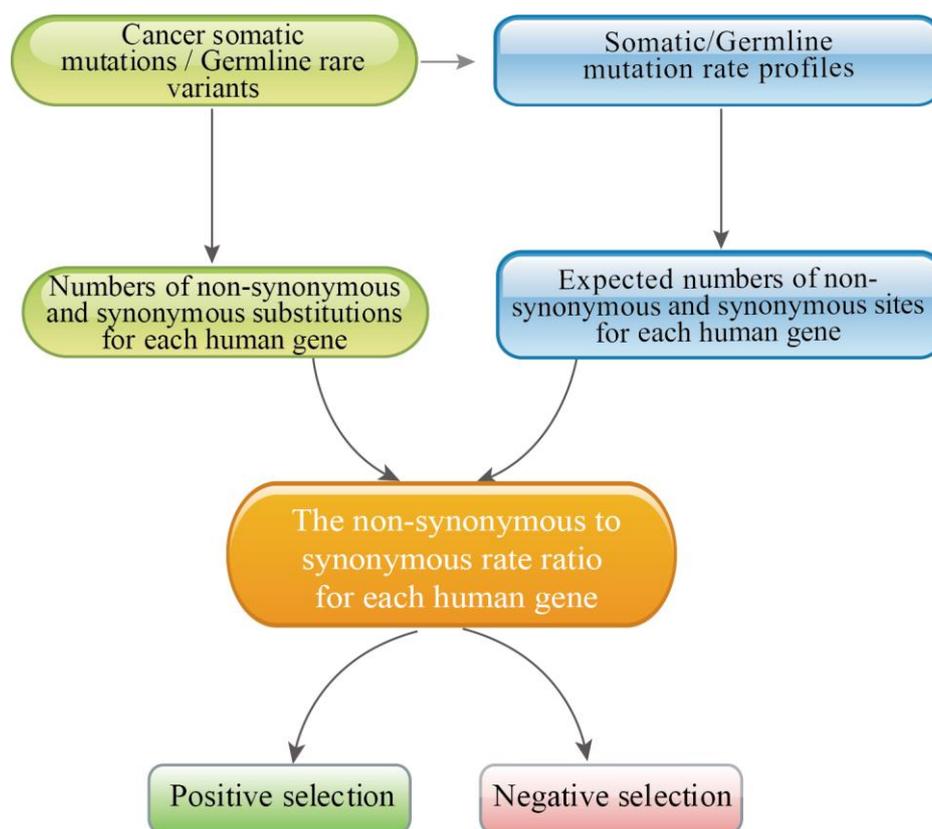
667 **Table 3.** The selection pressure in different cancer types

| Cancer type           | #Samples | $C_N/C_S$   | #Positive selection | #Negative selection |
|-----------------------|----------|-------------|---------------------|---------------------|
| Bladder cancer        | 233      | 1.389±0.010 | 5                   | 99                  |
| Blood cancer          | 686      | 1.145±0.014 | 4                   | 86                  |
| Bone cancer           | 164      | 1.078±0.022 | 1                   | 0                   |
| Brain cancer          | 797      | 1.392±0.021 | 10                  | 57                  |
| Breast cancer         | 1072     | 1.589±0.012 | 15                  | 105                 |
| Cervix cancer         | 194      | 1.402±0.011 | 3                   | 67                  |
| Colorectal cancer     | 443      | 1.563±0.014 | 41                  | 472                 |
| Esophagus cancer      | 347      | 1.350±0.012 | 4                   | 67                  |
| Gall bladder cancer   | 239      | 1.251±0.010 | 2                   | 57                  |
| Head and neck cancer  | 521      | 1.315±0.012 | 10                  | 256                 |
| Kidney cancer         | 668      | 1.381±0.010 | 2                   | 70                  |
| Liver cancer          | 966      | 1.551±0.011 | 10                  | 125                 |
| Lung cancer           | 224      | 1.410±0.011 | 11                  | 141                 |
| Nervous system cancer | 108      | 1.585±0.134 | 1                   | 0                   |
| Ovary cancer          | 181      | 1.244±0.010 | 1                   | 13                  |
| Pancreas cancer       | 685      | 1.333±0.014 | 5                   | 81                  |
| Prostate cancer       | 499      | 1.232±0.011 | 3                   | 41                  |
| Skin cancer           | 584      | 1.148±0.011 | 45                  | 1303                |
| Stomach cancer        | 298      | 1.560±0.013 | 20                  | 163                 |
| Uterus cancer         | 246      | 1.827±0.013 | 56                  | 135                 |

668

669



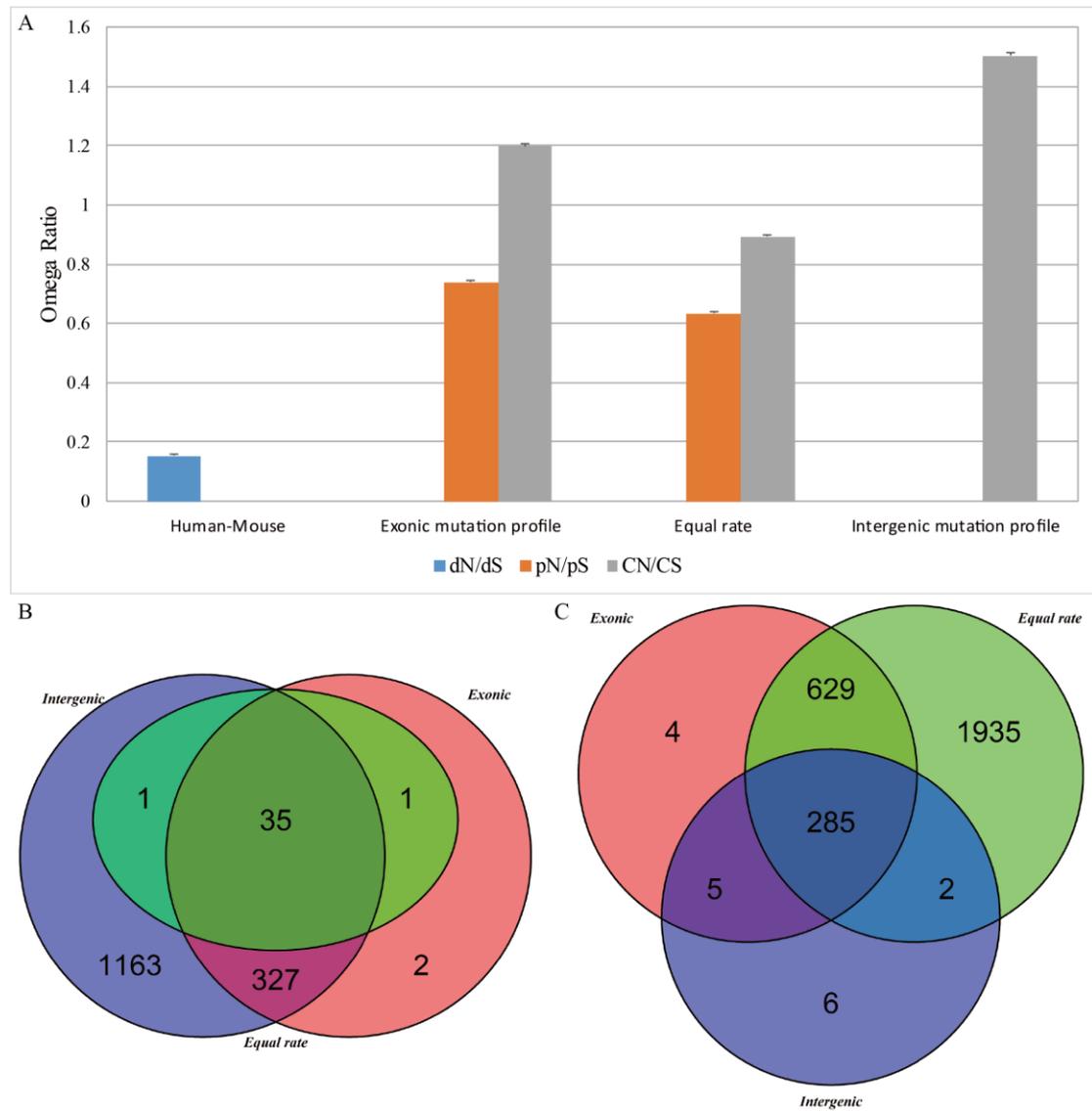


675

676

677

**Figure 2**



678

679

680

**Figure 3**