

1

2 **Quantification of neoantigen-mediated**

3 **immunoediting in cancer evolution**

4

5 Tao Wu¹⁻³, Guangshuai Wang¹, Xuan Wang¹, Shixiang Wang¹, Xiangyu Zhao¹,

6 Chenxu Wu¹, Wei Ning¹, Ziyu Tao¹, Fuxiang Chen⁴ and Xue-Song Liu¹

7

8 Affiliations of authors:

9 1 School of Life Science and Technology, ShanghaiTech University, Shanghai
10 201203, China;

11 2 Shanghai Institute of Biochemistry and Cell Biology, Chinese Academy of
12 Sciences, Shanghai, China;

13 3 University of Chinese Academy of Sciences, Beijing, China;

14 4 Department of Clinical Immunology, Ninth People's Hospital, Shanghai Jiao
15 Tong University School of Medicine, Shanghai, 200011, People's Republic of
16 China.

17

18 **Correspondence:** Xue-Song Liu, School of Life Science and Technology,
19 ShanghaiTech University, 230 Haik Road, Shanghai 201210, China. E-mail:
20 liuxs@shanghaitech.edu.cn.

21

22 **Key words:** immunoediting, immunogenicity, cancer biomarker, Neoantigen;
23 immunoediting-elimination; immunoediting-escape; Negative selection; Tumor
24 evolution;

25

26 **Disclosure of Potential Conflict of interests:** The authors have declared that
27 no competing interests exist.

28

29

30 **Abstract**

31 Immunoediting, which includes three temporally distinct stages, termed
32 elimination, equilibrium, and escape, has been proposed to explain the
33 interactions between cancer cells and the immune system during the evolution
34 of cancer. However the status of immunoediting in cancer remain unclear, and
35 the existence of neoantigen depletion signal in untreated cancer has been
36 debated. Here we developed a distribution pattern based method for quantifying
37 neoantigen mediated negative selection in cancer evolution. Our method
38 provides a robust and reliable quantification for immunoediting signal in
39 individual cancer patient. The prevalence of immunoediting signal in
40 immunotherapy untreated cancer genome has been demonstrated with this
41 method. Importantly, the elimination and escape stages of immunoediting can
42 be quantified separately, tumor types with strong immunoediting-elimination
43 tend to have weak immunoediting-escape signal, and vice versa. Quantified
44 immunoediting-elimination signal predicts cancer immunotherapy clinical
45 response. Immunoediting quantification provides an evolutional perspective for
46 evaluating the immunogenicity of neoantigen, and reveals potential biomarker
47 for cancer precision immunotherapy.

48

49

50 **Significance**

51 Neoantigen mediated negative selection has been demonstrated here, and the
52 quantified signal reveals distinct features of cancer immunoediting and also
53 potential biomarker for cancer immunotherapy response prediction.

54

55

56

57

58

59

60 **Introduction**

61 During cancer evolution, some genome DNA alterations can be positively
62 selected, such as driver mutations. Some genome alterations could posit a
63 deleterious effect, and consequently are negatively selected or depleted during
64 cancer evolution. Some (or the majority of) genome DNA alterations do not
65 have driving or deleterious effects on cancer, and follow neutral evolution
66 pattern. The interactions between immune cells and tumor cells are reflected
67 as immunoediting, which could mediate the negative selection of DNA
68 alterations encoding high immunogenicity (also known as neoantigenic
69 mutations) (1,2). Positively selected genome alterations can be readily detected
70 in the final mutation reservoir, however, the quantification of negative selection
71 in cancer evolution is still a significant challenge (3). The status of neoantigen
72 mediated negative selection in cancer evolution has been evaluated in several
73 studies, and controversial results have been reported (4-7).

74

75 Convincing cases of adaptive molecular evolution have been identified through
76 comparison of synonymous (silent; dS) and nonsynonymous (amino acid-
77 changing; dN) substitution rates in protein-coding DNA sequences. dN/dS is the
78 ratio between the rate of non-synonymous substitutions per non-synonymous
79 site and the rate of synonymous substitutions per synonymous site. dN/dS
80 method was originally developed to quantified the molecular evolution from
81 sequencing data (8,9). Recently, dN/dS method has been applied in cancer
82 evolution study (3). An important consideration in dN/dS analysis is the
83 selection of negative control regions. For example, a recent study reported that
84 neoantigen depletion signal is undetectable in pan-cancer dataset (5), however
85 the selection of negative control region is questionable (10). In addition, the
86 percentage of depleted neoantigen could be tiny, and this prohibit the accurate
87 detection of negative selection signal through dN/dS, especially in individual
88 cancer patient. Population genetics based method has been used to identify
89 the neutral pattern of cancer evolution, this method is based on the assumption

90 that the variant allele frequency (VAF) within a tumor follow a characteristic
91 power-law distribution in case of neutral evolution (11,12). The detection of
92 neutral evolution with this method has been questioned (13,14), and its
93 application in immune mediated negative selection has not been fully
94 established. In together, till now, the existence and the degree of neoantigen
95 mediated negative selection in human cancer remain unclear.

96

97 To address these challenges, we build a new distribution pattern based method
98 for quantifying neoantigen mediated negative selection in individual cancer
99 patient. With this new analysis framework, we demonstrate the pan-cancer
100 existence of neoantigen mediated negative selection signal. Importantly the
101 elimination and escape phases of immunoediting can be quantified separately.
102 In total, this study not only provides a novel method for quantifying negative
103 selection in cancer evolution, but also reveals potential biomarker for cancer
104 immunotherapy clinical response prediction.

105

106

107

108 **Materials and Methods**

109 **Pan-cancer clinical and molecular data**

110 The normalized gene-level RNA-seq data (TPM, transcripts per million) for 31
111 TCGA cohorts were downloaded from Xena (<https://xenabrowser.net/>). Pre-
112 compiled curated somatic mutations for TCGA cohorts were downloaded from
113 Xena, and missense variants are selected for downstream analysis.
114 ABSOLUTE-annotated MAF file which contains cancer cell fraction (CCF)
115 information of mutations was downloaded from GDC PanCanAtlas publications,
116 and then we used liftover function from R package “rtracklayer” to convert the
117 hg37 genome coordinates to hg38. Clinical data was obtained from GDC
118 PanCanAtlas publications (<https://gdc.cancer.gov/about-data/publications/pancanatlas>). HLA typing data was downloaded from

120 Thorsson et.al study (15). The downloaded mutation data, HLA typing data and
121 CCF values for TCGA samples have also been validated with in house
122 algorithm. Immune cell infiltration data for all TCGA tumors was downloaded
123 from ImmuneCellAI study (16), which estimates the abundance of 24 immune
124 cells comprised of 18 T-cell subtypes and 6 other immune cells. For TCGA
125 tumors which do not have HLA typing data in the mutation data set (2404
126 samples), we downloaded raw bam files, and performed HLA typing as
127 described below. Driver mutation data was downloaded from Bailey et al study
128 (17). This study used three different categories of tools to identify driver
129 mutations: (1) tools distinguishing benign versus pathogenic mutations using
130 sequence (CTAT population); (2) tools distinguishing driver versus passenger
131 mutations using sequence (CTAT cancer); and (3) tools discovering statistically
132 significant three-dimensional clusters of missense mutations (structure based).
133 We keep mutations identified by ≥ 2 approaches as the final high confident
134 driver mutations, including 2165 unique mutations.

135

136 **Somatic mutation calling**

137 Sequences were aligned to the reference human genome (hg38) using
138 Burrows–Wheeler Alignment (BWA) tool (18). Preprocessing followed the
139 GATK4 best practices workflow, including duplicate removal, base quality score
140 recalibration. Somatic mutations were identified on processed data using
141 Mutect2 (19). BCFtools was used to filter genome variants that passed all
142 quality control filters (20). The resulting VCF files were annotated by VEP and
143 further converted to MAF file by vcf2maf.pl (<https://github.com/mskcc/vcf2maf>).
144 The MAF file was loaded into R, analyzed and visualized by Maftools (21).

145

146 **Gene expression analysis**

147 Paired-end RNA-seq data were processed using hisat2 (22) aligner on the
148 basis of the hg38 human genome assembly with default parameters. Then the
149 aligned SAM files were transformed to BAM files using samtools (23).

150 Normalized RNA expression values (TPM) were calculated by TPMCalculator
151 (24).

152

153 **Cancer cell fraction (CCF) calculation**

154 We followed the GATK4 copy number analysis pipeline to get copy number
155 segment files. CCF information for each mutation was calculated based on
156 segment files and somatic mutation MAF files using ABSOLUTE software (25).
157 Briefly, read counts for each of the exome targets were collected from all
158 samples and calculated the coverage by count reads that overlap intervals
159 which were formed by padding the target regions. Each of the tumor samples
160 was compared to a panel of normal (PoN) controls for normalization and
161 denoising. The tool standardizes counts by the PoN median counts. The
162 normalization process includes log2 transformation and normalizing the counts
163 data to center around one. Then, the tool denoised the standardized copy ratios
164 using the principal components of the PoN. These normalized coverage profiles
165 were then segmented using Gaussian-kernel binary-segmentation algorithm,
166 which were fed into ABSOLUTE algorithm to determine CCF.

167

168 **HLA typing and neoantigen prediction**

169 HLA genotyping was performed with Optitype (26), using default parameters.
170 Mutect2 mutation files were first transformed into VCF format by maf2vcf tools,
171 and we used NeoPredPipe to predict neoantigen (27). Single-nucleotide
172 variants leading to a single amino acid change are the focus of this study. From
173 the output results, if the IC50 of a novel peptide is less than 50, the bind level
174 is SB (strong binder, rank is less than 0.5%), and the expression level (TPM) is
175 greater than 1, then this peptide is labeled as neoantigen. A mutation is
176 considered neoantigenic if there is at least one peptide derived from the
177 mutated site is predicted as neoantigen.

178

179 **Enrichment Score calculation**

180 The Kolmogorov–Smirnov (K-S) statistic can be used to quantify the distance
181 between two cumulative distributions. We constructed a K-S like statistic to
182 quantify the difference between the distribution of the CCF (or mRNA
183 expression) of immunogenic mutations and non-immunogenic mutations in
184 each individual sample.

185

186 **ES_{CCF} quantification in individual cancer patient**

187 We equally divided the whole CCF range (0-1) into 100 intervals (in descending
188 order) and assigned each interval a rank value (from 100 to 1). To make the
189 heavier weights on two tails of the rank distribution, we further normalized the
190 ranks:

$$191 \quad R_i = \left| \frac{L}{2} - r \right| + 1$$

192 Where R_i is the normalized rank value, L is the total number of intervals (here
193 L=100), and r is the original rank value.

194

195 Then we counted the number of mutations lied in each interval ($m(i)$) and
196 assigned a value $a(i)$ to each interval depending on mutation counts ($m(i)$) and
197 interval rank ($R(i)$):

$$198 \quad a_i = \frac{m_i \times R_i}{\sum m_i \times R_i}$$

199 We can calculate the empirical cumulative distribution of random variable $a(i)$
200 by walking from top to bottom (Fig 2):

$$201 \quad F(n) = \sum_i^n a_i, \quad n = 1, 2, \dots, 100$$

202 n is the index of intervals.

203

204 We then constructed two distributions for neoantigenic mutations ($F_N(n)$) and
205 non-neoantigenic mutations ($F_M(n)$) of individual sample, respectively. Then
206 K-S like statistics can be obtained by taking distance of two distributions:

207 $D = F_N(n) - F_M(n)$

208 Similar to GSVA (28), the enrichment score (ES) was defined as:

209 $ES = |D(n)^+| - |D(n)^-| = \max(0, D(n)) - |\min(0, D(n))|$

210 Where $D(n)^+$ and $D(n)^-$ are the largest positive and negative random walk
211 deviations from zero, respectively.

212

213 **ES_{RNA} quantification in individual cancer patient**

214 For a sample, using z to denote the mRNA expression. To reduce the influence
215 of potential outliers, we first convert z to rank z' , and normalize further
216 to $r = |P/2 - z'| + 1$, making the ranks symmetric around 1 (P is number of
217 mutations in a sample), making the heavier weights on two tails of the rank
218 distribution. Then we got two cumulative distributions, for mutations which are
219 neoantigens:

220
$$D_{neo}(S, i) = \sum_{r_j \in S, j \leq i} \frac{|r_j|}{\sum_{r_j \in S} |r_j|}$$

221 For mutations which are not neoantigens:

222
$$D_{noneo}(S, i) = \sum_{r_j \notin S, j \leq i} \frac{|r_j|}{\sum_{r_j \notin S} |r_j|}$$

223

224 Where S is the set of mutations which are neoantigens, the size of the set is P_s ,
225 P is the number of mutations in a sample, r is normalized rank of mutations.

226 Then we constructed a K-S like statistics:

227

228 $T = D_{neo}(S, i) - D_{noneo}(S, i)$

229

230 We transform the K-S like statistic into neoantigen enrichment score (ES) as
231 difference between the largest positive and negative distribution deviations from
zero:

232 $ES = \max(0, T) - |\min(0, T)|$

233

234 **Estimation of significance level of ES.**

235 We employed a permutation method to derive a null distribution to calculate p
236 value of the ES (ES_{CCF} or ES_{RNA}). For each sample, the same number of
237 mutations as neoantigens are randomly selected from mutation list and the
238 corresponding ES is calculated. This process is repeated 1000 times to get the
239 ES null distribution. The p value is calculated from the positive or negative
240 region of the empirical null distribution:

241
$$P = \begin{cases} \frac{1}{1000} \sum_{n=1}^{1000} I(ES_n \geq ES), ES \geq 0 \\ \frac{1}{1000} \sum_{n=1}^{1000} I(ES_n < ES), ES < 0 \end{cases}$$

242 Where I is an indicator function.

243

244 **Neutral simulation**

245 For each sample, we permute the neoantigen labeling (ie. randomly select the
246 same number of mutations as the actual neoantigen number in the selected
247 sample and label them as neoantigenic mutations) and calculate ES value. For
248 pan-cancer or cancer type dataset, we can obtain the same number of
249 simulated samples and corresponding ES values, then calculate the median
250 ES of these simulated samples. This process was repeated many times (usually
251 2000 times) to get the simulated distribution of median ES. The actual pan-
252 cancer or cancer type median ES values are compared with this simulated ES
253 distribution, and p values are then reported.

254

255 **Cancer immunotherapy datasets analysis**

256 To investigate the predictive performance of the quantified immunoediting-
257 elimination signal in immune checkpoint inhibitor (ICI) therapy clinical response
258 prediction for individual patient, we searched for public ICI datasets with
259 available raw WES data and RNA-seq data. Three melanoma ICI datasets have

260 been identified for this study. The Hugo et al dataset was related to anti-PD-1
261 therapy in metastatic melanoma (29). This dataset has 37 samples with WES
262 data, 26 were also analyzed by RNA sequencing (RNA-seq). The Riaz et al
263 dataset was related to anti-PD-1therapy in metastatic melanoma, and it has 56
264 samples with WES data, 40 with RNA-seq (30). The Liu et al cohort includes
265 melanoma patients treated with anti-PD1 antibody, it has 119 samples with
266 WES data and 112 samples with RNA-seq (31). All three melanoma studies
267 used very similar definition for clinical endpoints. Clinical response for patients
268 was defined by RECIST v1.1, responding tumors were derived from patients
269 who have complete or partial responses (CR/PR) in response to anti-PD-1
270 therapy; non-responding tumors were derived from patients who had
271 progressive disease or stable disease (PD/SD). We only chose pre-
272 immunotherapy treatment sample for analysis. Mutation calling, neoantigen
273 prediction, expression quantified, CCF calculation and ES calculation were
274 performed as described above.

275

276 **Stochastic branching process model for cancer evolution and power 277 analysis**

278 Tumor evolution model constructed by Lakatos et al has been applied in this
279 study (32). In this model, tumor evolution was initiated by a single transformed
280 cell. At any simulation step, a cell is randomly selected and has three events
281 could happen: birth (divide to produce two offspring), death and waiting. For a
282 birth event, new cells could acquire some new mutations (counts is sampled
283 from Poisson distribution) and each mutation can become neoantigen as a
284 specific probability. Under negative selection on neoantigen, the death rate of
285 cells could increase from d_0 to d_i with neoantigen accumulation. Selection
286 strength (s) of neoantigen mediated negative selection can be calculated as:

287

$$1 + s \times n = \frac{b - d_i}{b - d_0}$$

288
289 n is the number of neoantigens in a cell, b is birth rate (for simplicity, set b=1)
290 In addition, every mutation has a probability (p_{esc}) to escape. Once a mutation
291 is escaped, death rate of the cell which contains this mutation decrease to basal
292 death rate d_0 . This simulation step continues until the population reach pre-
293 defined size. Similar to the original study (32), the following parameters were
294 applied: neoantigen probability $p=0.1$, birth rate $b=0.1$, basal death rate $d_0=0.1$,
295 Poisson distribution parameter (mutation rate) $\mu=1$, escape probability $p_{esc}=10^{-6}$,
296 selection strength $-0.25 \leq s \leq 0$, final population size $popSize=10^5$. At each
297 selection strength, we run simulation 100 times. The model was implemented
298 with Julia (v1.3.1, revised from original Julia code provided by Lakatos et.al).
299
300 Mutations harbored in at least 5 cells out of 10^5 were collected at the end of
301 each simulation and the CCF was calculated. To account for imperfect
302 sequencing measurements, CCF values were computed via a simulated
303 sequencing step introducing noise to these frequencies with the indicated read
304 depth. For a given read depth D, each frequency value f was substituted by a
305 new frequency f' sampled from a binomial distribution with parameters D and f:
306 $f' \sim Binom(D,f)/D$. We filtered for mutations with f' above 0 to discard mutations
307 that are not picked up due to limited detection power. In addition to sequencing
308 limitations, we also added different proportions of false positive neoantigen
309 when evaluating the power of detecting negative selection: we randomly
310 sampled nonantigenic mutations of simulated tumors (varied between 5 and
311 500% of the number of true neoantigen) that were falsely labeled as neoantigen.
312 To calculate the power of derivation from neutral VAF distribution method (32),
313 we used two side K-S test to detect the difference between the VAF distribution
314 of all mutations and neoantigenic mutations and reported K-S statistic and
315 corresponding p value.
316
317 **Statistical analysis**

318 All statistical tests were performed using R statistical language. In all boxplots,
319 the center lines represent the median, low and upper box limits are the first and
320 third quartiles, respectively, and whiskers represent the values up to 1.5 times
321 of the interquartile range. P values for comparisons between boxplots were
322 calculated by Wilcoxon rank sum test. Correlation and corresponding p values
323 were calculated by Pearson method using R function cor.test. Kaplan-Meier
324 survival analysis was performed using the R package “survival” with log-rank
325 test, and Cox-proportional hazard analysis was performed using the R package
326 “ezcox”. The cutoff value of E_{SCCF} in Kaplan-Meier overall survival analysis was
327 determined by surv_cutpoint function of “survminer” package. R function ks.test
328 was used to perform two-sided K-S test.

329

330 **Software and data availability**

331 Custom code for quantifying immunoediting-elimination and immunoediting-
332 escape are available in Supplementary Source Data 1. All code required to
333 reproduce the analysis outlined in this manuscript, and R markdown analysis
334 report are available in Supplementary Source Data 2.

335

336

337

338

339 **Results**

340 **Conceptual framework for the elimination and escape phases of cancer 341 immunoediting**

342 The interactions between cancer cells and immune cells are manifested as
343 immunoediting, which consists of three sequential phases: elimination,
344 equilibrium, and escape (1,2). In the elimination phase, tumor cells with genome
345 alterations encoding high immunogenicity are partially or completely eliminated
346 by immune cells, and this leads to the down-regulation of the cancer cell fraction
347 (CCF) of genome alterations encoding high immunogenicity (Fig. 1A, B). In the

348 escape phase, tumor cells escape the surveillance of immune system through
349 multiple mechanisms, including the following: 1. Suppress the transcription or
350 expression of genome alterations encoding high antigenicity; 2. Antigen
351 presentation pathway down-regulation; 3. Up-regulate the expression of
352 immune suppressive molecules, including PD-L1, CTLA-4, etc. (Fig. 1a, b).

353

354 The elimination phase of immunoediting will lead to the down-regulation of CCF
355 of neoantigenic mutations, and consequently this CCF down-regulation status
356 of neoantigenic mutations can reflect the strength of the elimination phase of
357 immunoediting. The mRNA down-regulation status of neoantigenic mutations is
358 a partial reflection of the strength of immunediting-escape phase. Here we use
359 TCGA pan-cancer dataset to investigate this immunoediting signal. TCGA
360 dataset includes 31 cancer types and 9511 samples with available WGS or
361 WES data and mRNA expression profiling (RNA-seq) data, and neoantigenic
362 genome alterations can be found in 9166 samples (Supplementary Fig. S1) (33).
363 In the following section we build a distribution pattern based method to quantify
364 the selection strength acting on the CCF or mRNA expression of neoantigenic
365 mutation.

366

367 **Method for quantifying neoantigen mediated negative selection**

368 For each genome mutation, we have CCF and normalized mRNA expression
369 (transcripts per million, TPM) information. The immunogenicity value of genome
370 mutation can be calculated as the possibility of the mutated peptide to be
371 presented by HLA type I, and mutated peptides with predicted HLA I binding
372 affinity (IC50) less than 50nM are labeled as neoantigens. A mutation was
373 considered neoantigenic if there was at least one peptide derived from the
374 mutated sequence is predicted as neoantigen. The consequence of
375 immunoediting-elimination phase will leads to an unbalanced distribution of the
376 CCF of neoantigenic mutations, and this CCF distribution pattern of genome
377 alterations encoding immunogenicity can reflect the selection strength of

378 immunoediting-elimination phase. This distribution enrichment status of CCF
379 was calculated following a similar principle of gene set variation analysis (GSVA)
380 or gene set enrichment analysis (GSEA), which was originally developed in the
381 estimation of the variation of pathway activity over a sample population in an
382 unsupervised manner (28,34).

383

384 The mutations in an individual sample or a cancer type as a whole, are ordered
385 by CCF or mRNA expression (TPM) as a ranked list L. The mutations with
386 immunogenicity are defined as a set S. The goal of this analysis is to determine
387 whether the members of S are randomly distributed throughout L or primarily
388 found at the top or bottom. There are two key steps of this method (Fig. 2A, B):
389

390 1. Enrichment score (ES) calculation based on the distribution of neoantigen.
391 We calculate an ES that reflects the degree to which a set S is overrepresented
392 at the extremes (top or bottom) of the entire ranked list L. The score is
393 calculated by walking down the list L, increasing a running-sum statistic when
394 we encounter a mutation in S and decreasing it when we encounter mutations
395 not in S. The ES is calculated based on the maximum deviations from zero
396 during the random walk, it corresponds to a weighted Kolmogorov–Smirnov (K-
397 S) like statistic (see details in the Methods). The CCF values of mutations are
398 in the range of 0-1, and in TCGA dataset, the CCF values of mutations do not
399 show normal distribution, and many mutations have CCF values equal to 1
400 (Supplementary Fig. S2). A fixed CCF rank from 1 to 100 has been constructed
401 in the quantification of CCF distribution enrichment status of neoantigenic
402 mutation (ES_{CCF}).

403

404 2. Estimation of the significance level of ES.

405 We estimate the statistical significance (nominal P value) of the ES by using a
406 permutation test procedure, this procedure permute the neoantigen labels and
407 recomputed the ES of each patient, and this generates a null distribution for the

408 ES. The p value of the observed ES is then calculated according to this null
409 distribution. For ES significance analysis in TCGA pan-cancer or individual
410 cancer type cohort level, the observed median ES value of the test cohort is
411 compared with the distribution of median ES values from 2000 simulations (Fig.
412 2C).

413

414 Tumor cells can evolve multiples strategies to escape the surveillance of
415 immune system, and down-regulating the mRNA expression of neoantigenic
416 mutation is one of these strategies (Fig. 1A, B). Similar to CCF values, the
417 mRNA expression values of mutations are independent variables from
418 immunogenicity IC50 values. Similar strategy can be applied to quantify this
419 mRNA expression down-regulation mediated immunoediting, and the resulting
420 ES_{RNA} is a partial reflection of the strength of immunoediting-escape signal (Fig.
421 1 and Fig. 2B).

422

423 **The existence of significant immunoediting signal**

424 Previous studies have debated the existence of neoantigen depletion signals in
425 cancer evolution. Van den Eynden J. *et al.* reported that neoantigen depletion
426 signal is undetectable in TCGA pan-cancer dataset (5). However as pointed out
427 in a preprint, their method for neoantigen depletion signal detection is
428 problematic, as the actual neoantigens with immunogenicity are not located in
429 their defined “HLA-binding regions” (10). We investigate the status of this
430 immunoediting signal with the new method developed in this study using TCGA
431 pan-cancer dataset. The immunogenicity IC50 value is calculated based on the
432 mutated DNA sequence and HLA status, and the CCF information is
433 independently obtained from high-throughput sequencing. Mutation types do
434 not influence the CCF values, and the distribution of immunogenicity are not
435 influenced by mutation types either. The immunogenicity IC50 values are thus
436 independent variables from CCF values, and this is different from the
437 calculation of dN/dS, where the two variables dN, dS are interconnected and

438 are both significantly influenced by mutation types (3,35).

439

440 Since the variables (CCF and HLA binding IC50 status) are independent, we
441 use random simulation to generate a null distribution of ES_{CCF} . For TCGA pan-
442 cancer or individual cancer type cohort, the median ES_{CCF} values are recorded
443 after each simulation. The observed median ES_{CCF} values are compared with
444 the simulated ES_{CCF} values. In TCGA pan-cancer cohort with at least 1
445 neoantigenic and 1 subclonal mutation ($CCF < 0.6$) ($n=5900$), the observed
446 ES_{CCF} is -0.017 ($p=0.051$) (Fig. 3A, D). In PAAD and LUAD, the observed ES_{CCF}
447 values are significant lower compared with the random simulations, suggesting
448 the existence of immunoediting-elimination signal (Fig. 3A and Supplementary
449 Fig. S3). Since some neoantigenic mutations can be cancer drivers, which are
450 known to undergo positive selection during the evolution of cancer.
451 Neoantigens that happens to be cancer drivers are not undergoing immune
452 based negative selection (Supplementary Fig. S4), probably because the
453 driving force could override the negative selection, or additional immune escape
454 mechanism could evolve. In TCGA pan-cancer cohort, when samples with
455 neoantigenic and driver mutations lied on same genes are not included, the
456 observed median ES_{CCF} is -0.023 ($n=5420$, $P=0.006$) (Fig. 3B, D). Several
457 cancer types including ACC, CHOL, UCEC, LUAD show significant low ES_{CCF}
458 values (Fig. 3B and Supplementary Fig. S5). This data demonstrates the
459 existence of immunoediting-elimination signal in TCGA dataset.

460

461 Similarly random simulations were performed to evaluate the significance of the
462 observed ES_{RNA} values. Compared with ES_{CCF} , the observed ES_{RNA} show much
463 strongly significant difference when compared with the random simulated
464 values. In TCGA pan-cancer dataset with at least 1 neoantigenic mutation and
465 accompanied mRNA expression information ($n=7151$), the observed $ES_{RNA} =$
466 0.046 ($p<0.0005$) (Fig. 3C, D). In majority of cancer types (including KICH,
467 DLBC, CHOL, SARC, LUAD, PAAD, UCS, STAD, LGG, LUSC, HNSC, OV,

468 UCEC, BRCA, LIHC, COAD), a significant low ES_{RNA} values are observed (Fig.
469 3C and Supplementary Fig. S6). This study demonstrates that the
470 immunoediting escape through down-regulating the expression of neoantigenic
471 alteration is prevalent in human cancer (Fig. 3C). Furthermore, the
472 immunoediting-escape signal is more prevalent than the immunoediting-
473 elimination signal (Fig. 3A-C). This is in line with the fact that tumors need to
474 escape the immune surveillance to provide specimen for analysis, and it is not
475 practical to get the specimen of tumors that have already been eliminated by
476 immune system.

477

478 Interesting we observed that in cancer types with strong immunoediting-
479 elimination signal, a weak immunoediting-escape signal exist, and vice versa
480 (Fig. 3E, F). In pan-cancer or individual cancer type level, the immunoediting
481 elimination and escape signal exist, however in majority of cancer patients, both
482 of immunoediting-elimination and escape signals are weak or undetectable
483 (Supplementary Fig. S7 and Supplementary Fig. S8). Sufficient sequencing
484 depth is required for the detection of this immunoediting signal, and the required
485 sequencing depth is not reached in many TCGA samples.

486

487 **Neoantigen enrichment score and immune negative selection strength 488 quantification**

489 Recently, immune based negative selection has been simulated using a
490 stochastic branching process model (32). The neoantigen mediated negative
491 selection strength (s) is an inherent feature of each patient. However method
492 for accurately quantifying this immune based negative selection strength is still
493 lacking. Here we investigate the connections between neoantigen enrichment
494 score (ES_{CCF}) and immune negative selection strength s using a stochastic
495 branching cancer evolution model as previously described (32). For each fixed
496 selection strength s , the resulting ES_{CCF} was calculated. ES_{CCF} show near linear
497 correlation with s values (Fig. 4A). This analysis suggests that the quantified

498 ES_{CCF} can be used to infer the immune selection strength in patient. The
499 median ES_{CCF} in TCGA datasets is -0.023, suggesting a median immune
500 negative selection strength $s=-0.08$ (Fig. 4A).

501

502 Proportional neoantigen burden measures the percentage of neoantigenic
503 mutations in individual sample or cancer types. Proportional neoantigen burden
504 was originally designed to compare the immune negative selection strength
505 between two or more samples (32). The value of control proportional
506 neoantigen burden in the condition of neutral selection cannot be obtained, and
507 the proportional neoantigen burden method could not be applied in the
508 quantification of immune negative selection strength in individual cancer patient,
509 individual cancer type, or all cancer types combined. Derivation from neutral
510 VAF distribution ($1/f$ dependence of the cumulative VAF distribution) has been
511 suggested to reflect the selection status (11,32). However neutral VAF
512 distribution method is not suitable in negative selection quantification due to
513 strict requirement in sequencing depth and neoantigen prediction accuracy
514 (Supplementary Fig. S9).

515

516 **Pan-cancer features and correlations of immunoediting signal**

517 Human cancer evolve over a long time interval, usually in decades. The
518 immunoediting-elimination signal quantified in this study suggest the existence
519 of an already happened neoantigen mediated tumor elimination process. While
520 the quantified immune cell infiltration level represent the current immune
521 response status. We calculated the immunoediting status in TCGA pan-cancer
522 datasets (Fig. 3A). The unbalanced distribution of CCF in neoantigenic vs non-
523 neoantigenic mutations quantified as ES_{CCF} could reflect the status of
524 neoantigen mediated tumor elimination. In tumors with detectable
525 immunoediting-elimination signal ($ES_{CCF}<0$, $p<0.05$), a slightly increased CD8⁺
526 T plus natural killer (NK) cell infiltration status compared with the remaining
527 samples were observed, and the difference does not reach statistical

528 significance ($p=0.2$) (Fig. 5A, B). This data suggests that historically happened
529 immunoediting-elimination process may not be reflected in the current immune
530 cell infiltration status.

531

532 The down-regulation of immunogenic mutation encoded mRNA is a partial
533 reflection of immunoediting-escape phase (Fig. 1). Pan-cancer status of this
534 ES_{RNA} is shown, and different cancer types have different median ES_{RNA} score
535 (Fig. 3B). The immunoediting-escape signal quantified as ES_{RNA} also do not
536 show statistically significant difference between samples with detectable
537 immunoediting-escape signal ($ES_{RNA}<0$, $p<0.05$) and the remaining samples in
538 CD8⁺ T plus NK cell infiltration (Fig. 5C). However we observe a significant up-
539 regulated regulatory T cell (Treg) percentage in samples with detectable
540 immunoediting-escape signal (Fig. 5D), and up-regulated Treg has been
541 reported to stimulate tumor immune escape (36).

542

543 **Quantified immunoediting-elimination signal predicts the clinical
544 response of cancer immunotherapy**

545 Immunotherapy, represented by immune checkpoint inhibitors (ICI), is
546 transforming the treatment of cancer. However, only a small percentage of
547 patients show response to ICI, and effective biomarkers for ICI clinical response
548 prediction is still urgent needed (37). To investigate the predictive performance
549 of the quantified immunoediting-elimination signal (ES_{CCF}) in ICI response
550 prediction for individual patient, we searched for public ICI datasets with raw
551 WES data and RNA-seq data available, and three melanoma ICI datasets have
552 been identified (29-31) (Supplementary Fig. S10).

553

554 We calculate the immunoediting-elimination signal (ES_{CCF}) for each patient. In
555 univariate Cox proportional hazards regression analysis, quantified ES_{CCF}
556 value is significantly associated with cancer patients' survival ($p=0.03$), and low
557 ES_{CCF} value (suggest the presence of high immunoediting-elimination signal) is

558 associated with improved ICI clinical response (Hazard ratio (HR)=3.74,
559 95%CI=1.11-12.6) (Fig. 6A). Patients are divided into three groups based on
560 ES_{CCF} value, patients with lowest ES_{CCF} values (indicate the presence of
561 immunoediting-elimination signal) show the best survival after ICI (Fig. 6B).

562

563 Logistic regression is the appropriate regression analysis to conduct when the
564 dependent variable is dichotomous (binary). Here we use logistic regression to
565 compare the efficiency of ES_{CCF} , tumor mutational burden (TMB) and
566 neoantigenic mutation count in predicting immunotherapy clinical response.
567 Relationship between prognosis (patients with clinical response or without
568 clinical response) and ES_{CCF} , TMB and neoantigenic mutation count was
569 analyzed. The goodness of fit was performed by Hosmer–Lemeshow test (H-L
570 test). The H-L test P-value of TMB is 0.051 (Fig. 6C, middle), close to 0.05,
571 implicate the difference between prediction and expectation is close to
572 significant. The H-L test P-value of ES_{CCF} is 0.771 (Fig. 6C, right), higher than
573 the H-L test P-value of TMB and neoantigen count, suggesting ES_{CCF} is more
574 suitable for predicting prognosis of patients than TMB and neoantigen count.
575 This study suggests that the quantified immunoediting-elimination signal can
576 be biomarker for ICI clinical response prediction.

577

578

579

580 **Discussion**

581 The existence of neoantigen depletion signal has been debated. Here we
582 provide reliable evidence to demonstrate the pan-cancer existence of
583 immunoediting signal. Importantly, the elimination and escape phases of
584 immunoediting can be separately quantified with our method. Cancer types with
585 strong immunoediting elimination signal usually have low immunoediting
586 escape signal, and vice versa. Furthermore, the quantified immunoediting
587 elimination signal predict cancer immunotherapy clinical response.

588

589 This study provides an initial method to reliably quantify immunoediting signal
590 in individual cancer patient. To quantify the immunoediting signal for individual
591 patient, at least one neoantigenic mutation is required. The mechanisms
592 employed by tumor cells to escape immune surveillance is very complex, and
593 the shutdown of the expression of neoantigen mutation is only one of the
594 mechanisms. In addition, the mRNA expression is the combination of both wild
595 type and mutated alleles. Lack of ES_{RNA} signal does not mean that the immune
596 escape signal does not exist in the specific cancer or cancer types.

597

598 The existence of neoantigen mediated negative selection status in untreated
599 cancer has been debated (4-7). Existing method for negative selection study
600 including dN/dS and population genetics method. Rooney *et al.* use TCGA pan-
601 cancer CDS as the control sequence to calculate the expected neoantigen
602 number per non-silent mutation (Bpred/Npred), then the actual observed
603 neoantigen number per non-silent mutation (Bobs/Nobs) are compared with
604 Bpred/Npred (38). Since the pan-cancer CDS sequence has already been
605 immune edited, the neoantigen depletion signal reported in this study is
606 systematically underestimated. In addition, as pointed out by Van den Eynden
607 *et al.* There is HLA typing mistake in this study (5). Van den Eynden J. *et al.*
608 reported that neoantigen depletion signal is undetectable in untreated cancer
609 (5). This study select the “non HLA-binding regions” as the control, and
610 compare the nonsynonymous vs synonymous mutation ratio (n/s) in “HLA-
611 binding regions” vs “non HLA-binding regions”. They did not identify any
612 difference in these two regions in regard to n/s using TCGA pan-cancer dataset.
613 However their method is problematic, as the actual neoantigen with
614 immunogenicity are not located in their defined “HLA-binding regions” (10).
615 Martincorena *et al.* performed a comprehensive gene level evolution selection
616 study with dN/dS method, and reported significant neutral and positive selection,
617 but not negative selection in cancer genome (3). Since immunogenic mutations

618 occupy less than 5% of total mutations. In gene level, the selection on
619 neoantigen mutations are overshadowed by other driving or neutral mutations.
620 Neoantigen mediated negative selection cannot be observed in gene-level
621 does not mean the absence of immune based neoantigen depletion. Zapata *et*
622 *al.* investigated immune based negative selection with dN/dS method (4).
623 However same problem exist in the selection of control DNA sequence. Similar
624 to Van den Eynden J. *et al.* CDS was divided into epitope region and non-
625 epitope region, dN/dS was compared in these two regions. Since neoantigen
626 with immunogenicity are not necessarily located in “epitope region”, the results
627 reported in this study is also questionable. Population genetics model for
628 neutral evolution has been proposed to detect neutral evolution based on
629 cumulative VAF distribution. For instance, a recent study test the neutrality of
630 cancer based on the VAF of mutations in a limited subclonal frequency range
631 (11). However, that test of neutrality has been questioned, the frequency
632 distribution of mutated alleles in a limited frequency range is not an accurate
633 statistic for detecting selection in cancer (13,14). Furthermore the application
634 of this population genetics method in neoantigen mediated negative selection
635 quantification in individual cancer patient has not been established (32)
636 (Supplementary Fig. S9).

637

638 The quantification of negative selection in cancer evolution has been a major
639 scientific challenge, the method developed here for neoantigen mediated
640 negative selection quantification could be instructive for the future design of
641 strategies for studying negative selection in cancer evolution. The existence of
642 neoantigen mediated negative selection has been demonstrated with our new
643 method. Importantly we observed a strong immunoediting-escape signal
644 reflected as the down-regulation of mRNA encoded by neoantigenic mutations.
645 The quantification of immunoediting provides an evolutionary perspective for
646 the design of neoantigen vaccine for cancer therapy. The immune based
647 negative selection is an inherent feature of cancer patient, the quantified

648 immunoediting signal can be used in cancer precision stratification, including
649 the clinical response prediction for cancer immunotherapy.

650

651

652

653

654 **Acknowledgments**

655 We thank ShanghaiTech University high performance computing public service
656 platform for computing services. This work was supported by the national
657 natural science foundation of China (31771373), Shanghai science and
658 technology commission (21ZR1442400), and startup funding from
659 ShanghaiTech University.

660

661 **Contributions**

662 TW collected the data, developed the immunoediting quantification analysis
663 method and performed the computational analysis. GW participated in data
664 collection and preprocessing. XW participated in neoantigen prediction. SW
665 help to build the method for immunoediting quantification. XZ, CW, WN, ZT, FC
666 participated in critical project discussion. XSL conceptualized the idea,
667 designed, supervised the study and wrote the manuscript.

668

669 **Conflict of interest**

670 The authors declare no competing interests.

671

672

673

674

675

676

677

678 **Reference:**

- 679 1. Schreiber RD, Old LJ, Smyth MJ. Cancer Immunoediting: Integrating Immunity's Roles in Cancer
680 Suppression and Promotion. *Science* **2011**;331:1565-70
- 681 2. O'Donnell JS, Teng MWL, Smyth MJ. Cancer immunoediting and resistance to T cell-based
682 immunotherapy. *Nat Rev Clin Oncol* **2019**;16:151-67
- 683 3. Martincorena I, Raine KM, Gerstung M, Dawson KJ, Haase K, Van Loo P, *et al.* Universal Patterns
684 of Selection in Cancer and Somatic Tissues. *Cell* **2017**;171:1029-+
- 685 4. Zapata L, Pich O, Serrano L, Kondrashov FA, Ossowski S, Schaefer MH. Negative selection in
686 tumor genome evolution acts on essential cellular functions and the immunopeptidome.
687 *Genome Biol* **2018**;19
- 688 5. Van den Eynden J, Jimenez-Sanchez A, Miller ML, Larsson E. Lack of detectable neoantigen
689 depletion signals in the untreated cancer genome. *Nat Genet* **2019**;51:1741-+
- 690 6. Marty R, Kaabinejadian S, Rossell D, Slifker MJ, van de Haar J, Engin HB, *et al.* MHC-I Genotype
691 Restricts the Oncogenic Mutational Landscape. *Cell* **2017**;171:1272-+
- 692 7. Claeys A, Luijts T, Marchal K, Van den Eynden J. Low immunogenicity of common cancer hot
693 spot mutations resulting in false immunogenic selection signals. *Plos Genet* **2021**;17
- 694 8. Goldman N, Yang ZH. Codon-Based Model of Nucleotide Substitution for Protein-Coding DNA-
695 Sequences. *Mol Biol Evol* **1994**;11:725-36
- 696 9. Yang ZH, Bielawski JP. Statistical methods for detecting molecular adaptation. *Trends Ecol Evol*
697 **2000**;15:496-503
- 698 10. Wang S, Wang X, Wu T, He Z, Li H, Sun X, *et al.* Revisiting neoantigen depletion signal in the
699 untreated cancer genome. *bioRxiv* **2020**
- 700 11. Williams MJ, Werner B, Barnes CP, Graham TA, Sottoriva A. Identification of neutral tumor
701 evolution across cancer types. *Nat Genet* **2016**;48:238-44
- 702 12. Williams MJ, Werner B, Heide T, Curtis C, Barnes CP, Sottoriva A, *et al.* Quantification of
703 subclonal selection in cancer from bulk sequencing data. *Nat Genet* **2018**;50:895-+
- 704 13. Tarabichi M, Martincorena I, Gerstung M, Leroi AM, Markowitz F, Spellman PT, *et al.* Neutral
705 tumor evolution? *Nat Genet* **2018**;50:1630-3
- 706 14. McDonald TO, Chakrabarti S, Michor F. Currently available bulk sequencing data do not
707 necessarily support a model of neutral tumor evolution. *Nat Genet* **2018**;50:1620-3
- 708 15. Thorsson V, Gibbs DL, Brown SD, Wolf D, Bortone DS, Yang THO, *et al.* The Immune Landscape
709 of Cancer. *Immunity* **2018**;48:812-+
- 710 16. Miao YR, Zhang Q, Lei Q, Luo M, Xie GY, Wang HX, *et al.* ImmuCellAI: A Unique Method for
711 Comprehensive T-Cell Subsets Abundance Prediction and its Application in Cancer
712 Immunotherapy. *Adv Sci* **2020**;7
- 713 17. Bailey MH, Tokheim C, Porta-Pardo E, Sengupta S, Bertrand D, Weerasinghe A, *et al.*
714 Comprehensive Characterization of Cancer Driver Genes and Mutations. *Cell* **2018**;173:371-+
- 715 18. Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler transform.
716 *Bioinformatics* **2010**;26:589-95
- 717 19. Benjamin D, Sato T, Cibulskis K, Getz G, Stewart C, Lichtenstein L. Calling Somatic SNVs and
718 Indels with Mutect2. *bioRxiv* **2019**
- 719 20. Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, *et al.* Twelve years of SAMtools
720 and BCFtools. *Gigascience* **2021**;10
- 721 21. Mayakonda A, Lin DC, Assenov Y, Plass C, Koeffler HP. Maftools: efficient and comprehensive

- 722 analysis of somatic variants in cancer. *Genome Res* **2018**;28:1747-56
- 723 22. Kim D, Paggi JM, Park C, Bennett C, Salzberg SL. Graph-based genome alignment and
724 genotyping with HISAT2 and HISAT-genotype. *Nat Biotechnol* **2019**;37:907-+
- 725 23. Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, et al. The Sequence Alignment/Map
726 format and SAMtools. *Bioinformatics* **2009**;25:2078-9
- 727 24. Alvarez RV, Pongor LS, Marino-Ramirez L, Landsman D. TPMCalculator: one-step software to
728 quantify mRNA abundance of genomic features. *Bioinformatics* **2019**;35:1960-2
- 729 25. Carter SL, Cibulskis K, Helman E, McKenna A, Shen H, Zack T, et al. Absolute quantification of
730 somatic DNA alterations in human cancer. *Nat Biotechnol* **2012**;30:413-+
- 731 26. Szolek A, Schubert B, Mohr C, Sturm M, Feldhahn M, Kohlbacher O. OptiType: precision HLA
732 typing from next-generation sequencing data. *Bioinformatics* **2014**;30:3310-6
- 733 27. Schenck RO, Lakatos E, Gatenbee C, Graham TA, Anderson ARA. NeoPredPipe: high-throughput
734 neoantigen prediction and recognition potential pipeline. *Bmc Bioinformatics* **2019**;20
- 735 28. Hanzelmann S, Castelo R, Guinney J. GSVA: gene set variation analysis for microarray and RNA-
736 Seq data. *Bmc Bioinformatics* **2013**;14
- 737 29. Hugo W, Zaretsky JM, Sun L, Song CY, Moreno BH, Hu-Lieskovian S, et al. Genomic and
738 Transcriptomic Features of Response to Anti-PD-1 Therapy in Metastatic Melanoma. *Cell*
739 **2016**;165:35-44
- 740 30. Riaz N, Havel JJ, Makarov V, Desrichard A, Urba WJ, Sims JS, et al. Tumor and Microenvironment
741 Evolution during Immunotherapy with Nivolumab. *Cell* **2017**;171:934-+
- 742 31. Liu D, Schilling B, Liu D, Sucker A, Livingstone E, Jerby-Amon L, et al. Integrative molecular and
743 clinical modeling of clinical outcomes to PD1 blockade in patients with metastatic melanoma.
744 *Nat Med* **2019**;25:1916-+
- 745 32. Lakatos E, Williams MJ, Schenck RO, Cross WCH, Househam J, Zapata L, et al. Evolutionary
746 dynamics of neoantigens in growing tumors. *Nat Genet* **2020**;52:1057-+
- 747 33. Hoadley KA, Yau C, Hinoue T, Wolf DM, Lazar AJ, Drill E, et al. Cell-of-Origin Patterns Dominate
748 the Molecular Classification of 10,000 Tumors from 33 Types of Cancer. *Cell* **2018**;173:291-+
- 749 34. Subramanian A, Tamayo P, Mootha VK, Mukherjee S, Ebert BL, Gillette MA, et al. Gene set
750 enrichment analysis: A knowledge-based approach for interpreting genome-wide expression
751 profiles. *P Natl Acad Sci USA* **2005**;102:15545-50
- 752 35. Van den Eynden J, Larsson E. Mutational Signatures Are Critical for Proper Estimation of
753 Purifying Selection Pressures in Cancer Somatic Mutation Data When Using the dN/dS Metric.
754 *Front Genet* **2017**;8
- 755 36. Facciabene A, Motz GT, Coukos G. T-Regulatory Cells: Key Players in Tumor Immune Escape and
756 Angiogenesis. *Cancer Res* **2012**;72:2162-71
- 757 37. Nishino M, Ramaiya NH, Hatabu H, Hodi FS. Monitoring immune-checkpoint blockade:
758 response evaluation and biomarker development. *Nat Rev Clin Oncol* **2017**;14:655-68
- 759 38. Rooney MS, Shukla SA, Wu CJ, Getz G, Hacohen N. Molecular and Genetic Properties of Tumors
760 Associated with Local Immune Cytolytic Activity. *Cell* **2015**;160:48-61
- 761
- 762
- 763
- 764
- 765

766 **Figure legend**

767 **Figure 1. Conceptual framework for the quantification of elimination and**
768 **escape phases of immunoediting.**

769 **A**, Phases of immunnoediting and the manifestations of the elimination and
770 escape phases of immunoediting.

771 **B**, Diagram showing the manifestations of immunoediting phases. Red star
772 indicates neoantigenic mutation. The elimination phase will decrease the CCF
773 of neoantigenic mutations, and cancer cell can eventually evolve multiple
774 mechanisms to escape the surveillance of immune system, including shutdown
775 the transcription of mRNA encoded by neoantigenic mutations.

776

777 **Figure 2. Distribution pattern based method for the quantification of**
778 **neoantigen mediated negative selection in cancer evolution.**

779 **A**, Detailed steps for CCF down-regulation based immunoediting-elimination
780 (ES_{CCF}) quantification. 1. Equally divide the whole CCF range (0-1) into 100
781 intervals and calculate the distribution of CCF of neoantigenic mutations and
782 non-neoantigenic mutations in these intervals; 2. Construct the Kolmogorov–
783 Smirnov (K-S) statistics based on difference between the two distributions; 3.
784 Calculate the enrichment score (ES_{CCF}).

785 **B**, Detailed steps for mRNA down-regulation based immunoediting-escape
786 (ES_{RNA}) quantification. 1. Rank mutations by corresponding mRNA expression
787 and calculate the distribution of mRNA expression of neoantigenic mutations
788 and non-neoantigenic mutations; 2. Construct the K-S statistics based on
789 difference between the two distributions; 3. Calculate the enrichment score
790 (ES_{RNA}).

791 **C**, Random simulation to obtain the null distribution of ES. For each sample, we
792 permute the mutation labeling (ie. randomly select the same number of
793 mutations as the observed number in the sample, and label them as
794 neoantigenic mutations) and calculate ES value, the processes are repeated
795 for 2000 times, and the actual ES values are compared with the simulated

796 values.

797

798 **Figure 3. Pan-cancer distributions and features of the quantified
799 immunoediting signals (ES_{CCF} and ES_{RNA}).**

800 **A**, Distribution of ES_{CCF} in pan-cancer (left) and in specific cancer type (right).

801 The p values are calculated from simulated median ES distributions. ns: p >

802 0.05, *: p <= 0.05, **: p <= 0.01, ***: p <= 0.001, ****: p <= 0.0001.

803 **B**, Distribution of ES_{CCF} in pan-cancer (left) and in specific cancer type (right),
804 after removing samples with neoantigenic and driver mutations located in the
805 same gene. The p values are calculated from simulated median ES distributions.

806 **C**, Distribution of ES_{RNA} in pan-cancer (left) and in specific cancer type (right).

807 The p values are calculated from simulated median ES distributions.

808 **D**, From left to right, simulated median ES distribution and the observed median
809 ES for Fig 3a, 3b and 3c respectively.

810 **E**, Correlation between median ES_{RNA} and ES_{CCF} of TCGA cancer types.

811 Pearson correlation coefficient and p value are shown.

812 **F**, Correlation between the percent of escape samples ($ES_{RNA} < 0$ and p < 0.05)
813 and median ES_{CCF} in TCGA cancer types. Pearson correlation coefficient and
814 p value are shown.

815

816 **Figure 4. Immunoediting-elimination signal (ES_{CCF}) and neoantigen-
817 mediated negative selection strength quantification.**

818 **A**, ES_{CCF} as a function of neoantigen-mediated negative selection strength s,
819 computed from n=100 tumors, with simulated read depth of 200× for each
820 indicated selection strength s. The observed median ES_{CCF} of TCGA samples
821 is indicated with a horizontal dashed line.

822 **B**, The proportion of 100 simulated tumors with significant ES_{CCF} (FDR
823 corrected p value less than 0.1) in each selection strength s.

824

825 **Figure 5. Immunoediting-elimination and escape signals and tumor**

826 **immune cell infiltration status.**

827 **A, B,** Comparisons between TCGA cancer patients with detectable
828 Immunoediting-elimination signal ($ES_{CCF} < 0$, $p < 0.05$) and the remaining patients
829 in CD8⁺ T plus natural killer (NK) cell (**a**) and Treg cell (**b**) infiltration status.

830 **C, D,** Comparisons between TCGA cancer patients with detectable
831 Immunoediting-elimination signal ($ES_{RNA} < 0$, $p < 0.05$) and the remaining patients
832 in CD8⁺ T plus NK cell (**c**) and Treg cell (**d**) infiltration status. Wilcoxon rank
833 sum test p value is shown.

834

835 **Figure 6. Quantified immunoediting-elimination signal (ES_{CCF}) predicts**
836 **cancer immunotherapy clinical response.**

837 **A,** Univariate Cox regression analysis was performed to estimate the hazard
838 ratio (HR) associated with ES_{CCF} values. The length of horizontal line
839 represents the 95% confidence interval (CI) and the vertical dashed line
840 represents HR = 1.

841 **B,** Kaplan-Meier overall survival curves show the comparison between different
842 groups stratified by ES_{CCF} value. Samples with ES_{CCF} values higher than the
843 cutoff (-0.222, determined by `surv_cutpoint` function of "survminer" package)
844 were classified as "high" group, and samples with ES_{CCF} value less than the
845 cutoff were classified as "low" group. The remaining samples (without
846 neoantigen or minimum CCF is higher than 0.6) were classified as "other" group.

847 **C,** The goodness-of-fit is performed by Hosmer-Lemeshow test, which shows
848 that ES_{CCF} is more suitable for predicting the prognosis of patients after ICI than
849 TMB or neoantigen burden.

850

851

852

853

854

855

856 **Supplementary Figure Legends**

857

858 **Supplementary Fig. S1. Overview of the TCGA pan-cancer dataset**
859 **included in this study.**

860 **A**, Samples used for calculating ES_{CCF}. Number of the patient with CCF
861 information and number of patient with at least one neoantigenic mutation and
862 subclonal mutation (CCF<0.6) are shown for each cancer type.

863 **B**, Samples used for calculating ES_{mRNA}. Number of the patient with mRNA
864 expression information and number of patient with at least one neoantigenic
865 mutation and accompanied mRNA expression information are shown for each
866 cancer type.

867

868 **Supplementary Fig. S2. CCF distribution of all mutations in TCGA dataset.**

869

870 **Supplementary Fig. S3. Simulated median ES distribution and**
871 **corresponding p value (blue area) of TCGA cancer types.** Related to
872 Figure3a.

873

874 **Supplementary Fig. S4. Comparison of ES_{CCF} between samples with and**
875 **without neoantigenic and driver mutations located in the same gene.**
876 Wilcoxon sum test p value is shown.

877

878 **Supplementary Fig. S5. Simulated median ES distribution and**
879 **corresponding p value (blue area) of TCGA cancer types.** Related to
880 Figure3b.

881

882 **Supplementary Fig. S6. Simulated median ES distribution and**
883 **corresponding p value (blue area) of TCGA cancer types.** Related to
884 Figure3c.

885

886 **Supplementary Fig. S7. Distribution of patients with significant**
887 **immunoediting-elimination and expression based escape signal (ES_{CCF} or**
888 **$ES_{RNA} < 0$, $p < 0.05$) in TCGA pan-cancer dataset.**

889

890 **Supplementary Fig. S8. Proportion of TCGA cancer patients with**
891 **detectable immunoediting signal.**

892 A, Proportion of TCGA cancer patients with $ES_{CCF} < 0$ (and) $p < 0.05$.

893 B, Proportion of TCGA cancer patients with $ES_{RNA} < 0$ (and) $p < 0.05$.

894

895 **Supplementary Fig. S9. Comparison ES_{CCF} with other methods in negative**
896 **selection strength quantification.**

897 **A**, Derivation from neutral VAF distribution quantification as a function of
898 negative selection strength s , computed from individual tumor of a simulation
899 cohort ($n=100$), with a simulated read depth of $200\times$.

900 **B**, Proportion of 100 simulated tumors with significant signal (FDR corrected p
901 value less than 0.1) quantified using derivation from neutral VAF distribution
902 method under each negative selection strength s . Of note, no tumors show
903 significant signal under the same simulated conditions as the data show in Fig.
904 4b.

905 **C**, Power to detect negative selection as a function of sequencing read depth
906 (x axis) and false neoantigen rate (y axis) using the enrichment score method
907 developed in this study. Power is the proportion of 100 simulated tumors with
908 significant negative ES value (FDR corrected p value less than 0.1).

909 **D**, Power to detect negative selection as a function of sequencing read depth
910 (x axis) and false neoantigen rate (y axis) using derivation from neutral VAF
911 distribution method. Power is the proportion of 100 simulated tumors with
912 significant difference (two-sided K-S test, FDR corrected p value less than 0.1)
913 between the distribution of all mutations and neoantigenic mutations.

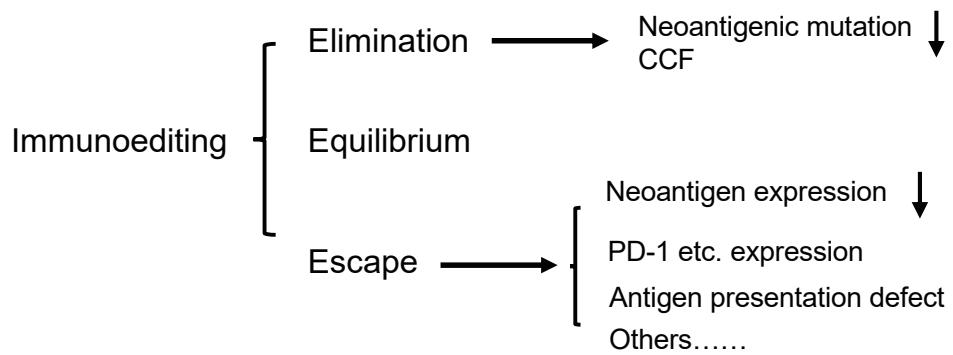
914

915 **Supplementary Fig. S10. Clinical parameters of the three cancer ICI**

916 **datasets used in this study.**

Figure 1

A



B

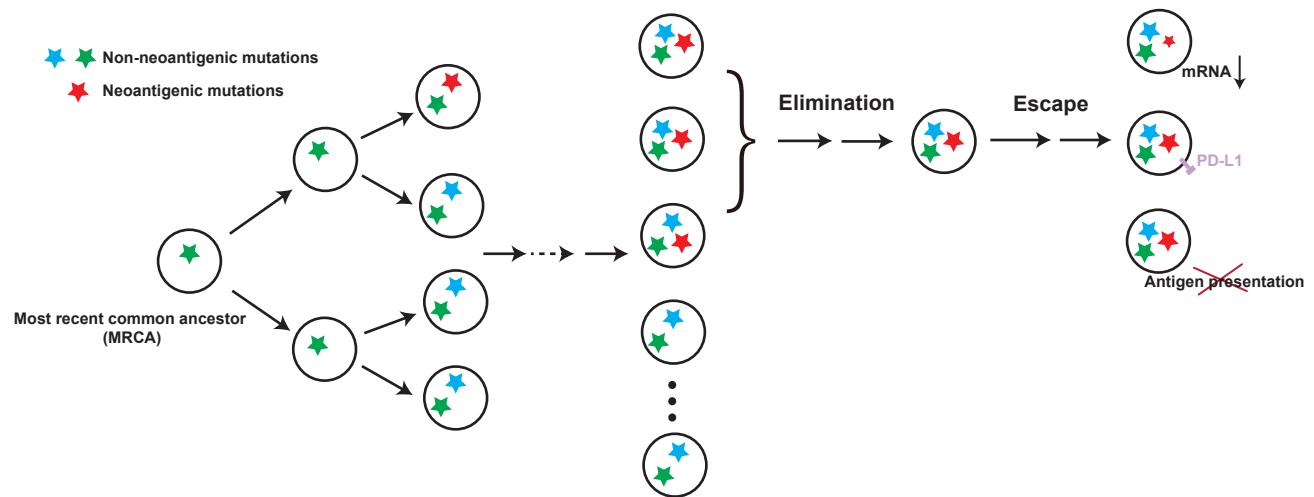
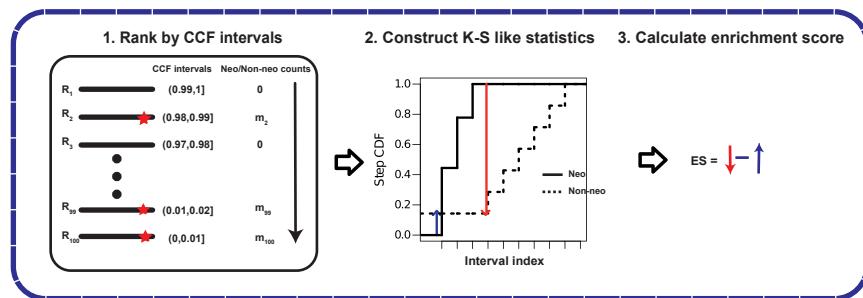


Figure 2

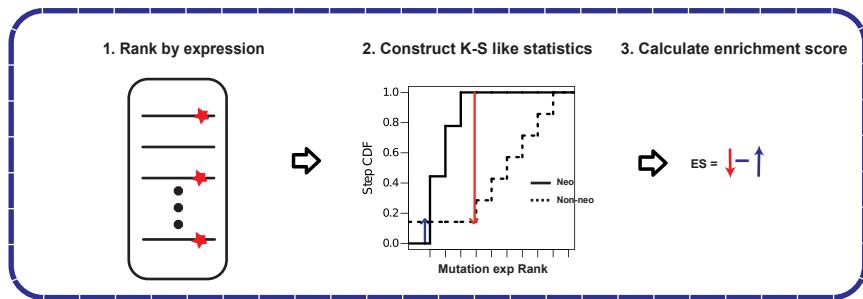
A

CCF down-regulation based immunoediting-elimination quantification

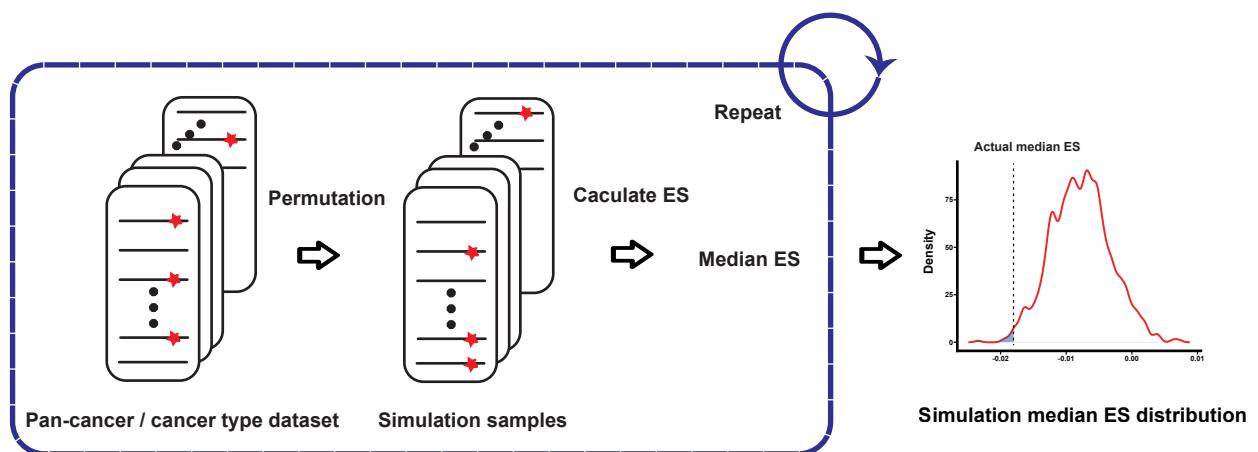


B

mRNA down-regulation based immunoediting-escape quantification



C



Permute neoantigen labeling, and calculate ES score

Figure 3

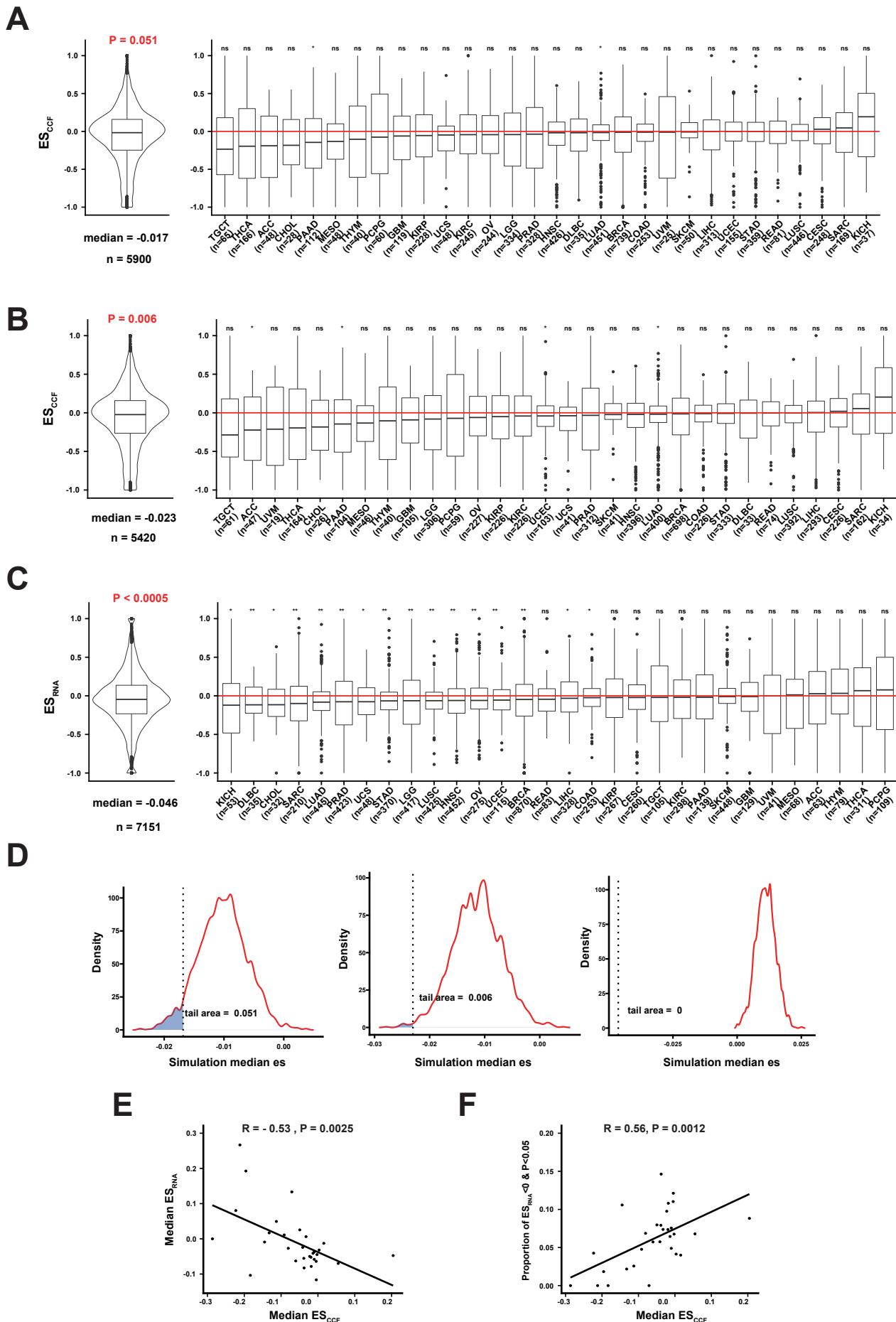
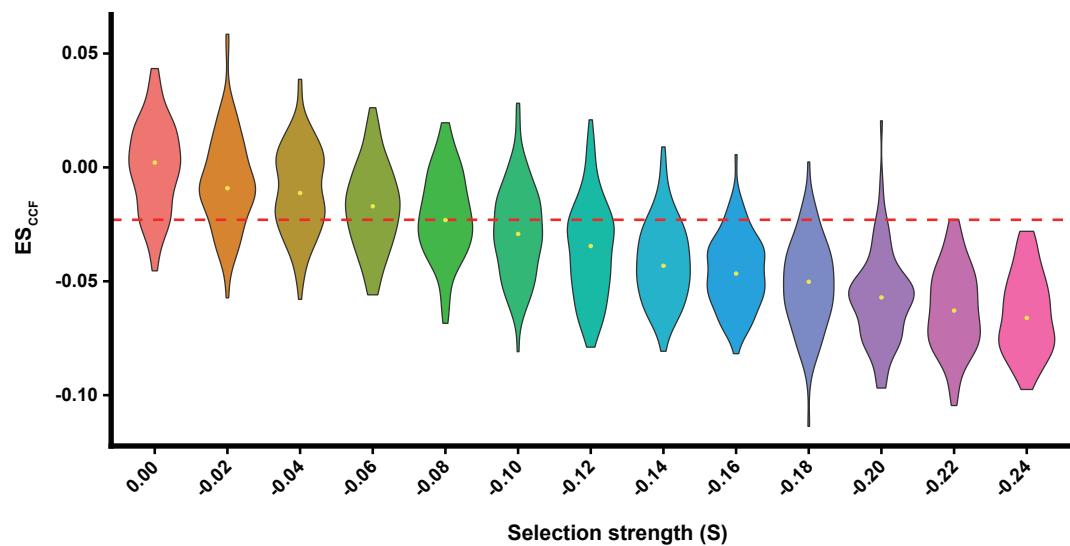


Figure 4

A



B

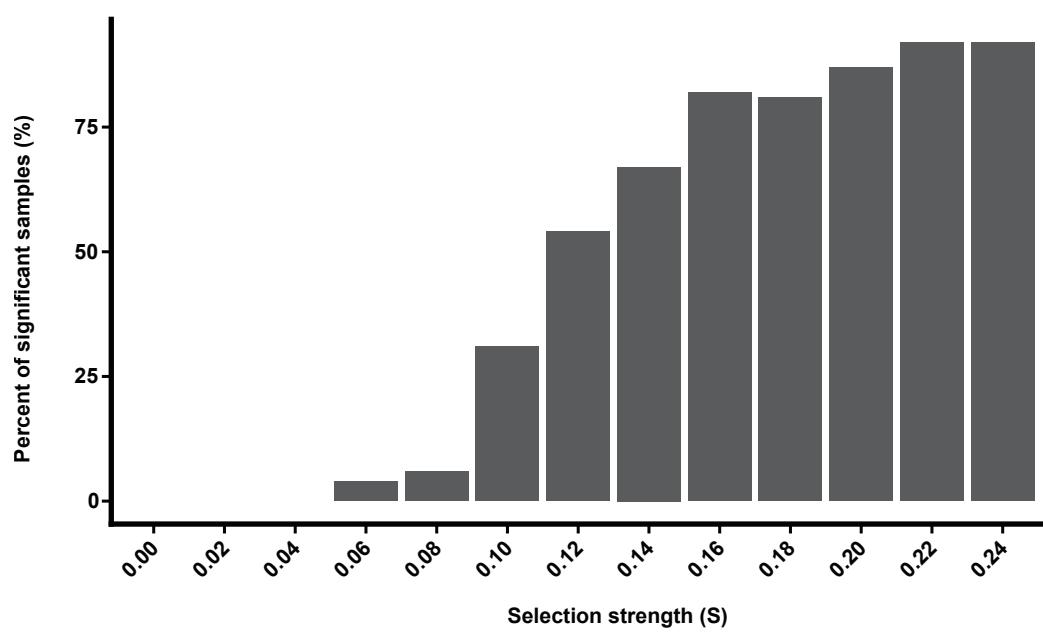


Figure 5

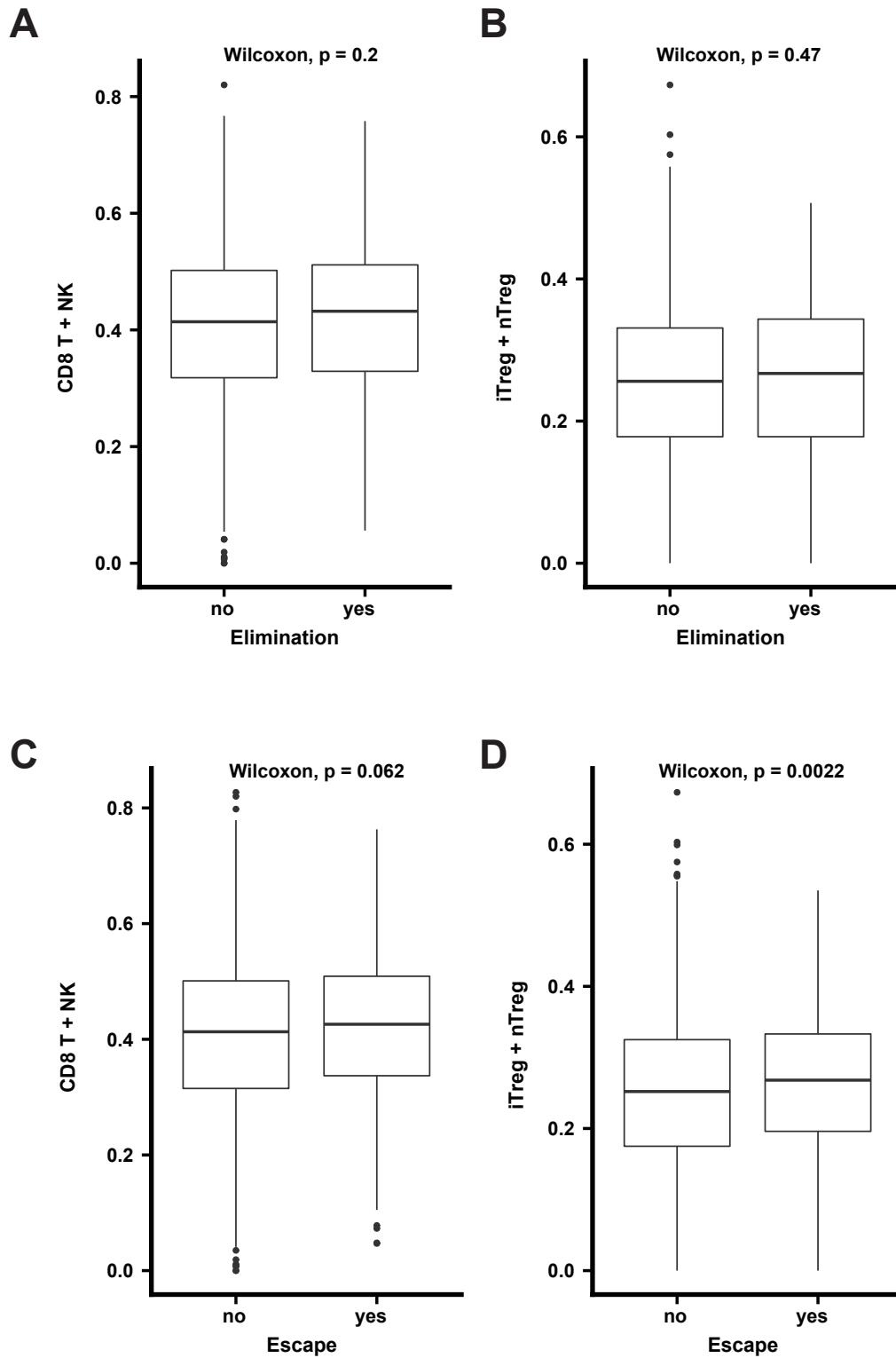


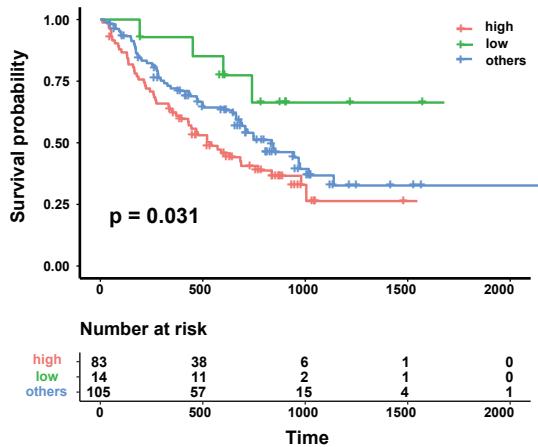
Figure 6

A

Variable	N	Hazard ratio	P
ES_{CCF}	97	3.74 (1.11, 12.60)	0.03

Cox analysis for variable ES_{CCF}

B



C

