# Lack of detectable neoantigen depletion signals in the untreated cancer genome

Jimmy Van den Eynden[1,2,3*], Alejandro Jiménez-Sánchez[3,4], Martin L. Miller[3] and Erik Larsson[1,3]

[1]Department of Medical Biochemistry and Cell Biology, Institute of Biomedicine, The Sahlgrenska Academy, University of Gothenburg, Gothenburg, Sweden.

[2]Department of Human Structure and Repair, Anatomy and Embryology Unit, Ghent University, Ghent, Belgium.

[3]Cancer Research UK Cambridge Institute, University of Cambridge, Li Ka Shing Centre, Robinson Way, Cambridge, UK.

[4]Program for Computational and Systems Biology, Sloan Kettering Institute, Memorial Sloan Kettering Cancer Center, New York, NY, USA.

*Email: jimmy.vandeneynden@ugent.be

16   **Somatic mutations can result in the formation of neoantigens, immunogenic peptides**

17   **that are presented on the tumor cell surface via HLA molecules. These mutations are**

18   **expected to be under negative selection pressure, but the extent of the resulting**

19   **neoantigen depletion remains unclear. Based on HLA affinity predictions, we**

20   **annotated the human genome for its translatability to HLA binding peptides and**

21   **screened for reduced single nucleotide substitution rates in large genomic datasets**

22   **from untreated cancers. Apparent neoantigen depletion signals became negligible**

23   **when considering trinucleotide-based mutational signatures, either due to lack of**

24   **power or efficient immune evasion mechanisms active early during tumor evolution.**

25   Cancer is caused by somatic mutations in driver genes. These genomic alterations result in

26   a selective growth advantage and positive selection of the affected cells[1]. With the rise of

27   next-generation sequencing technologies, increasing insights into the cancer genome have

28   led to a comprehensive characterization of the frequencies and patterns of somatic

29   mutations across different cancers[2,3]. For a tumor to evolve, it also needs to develop ways to

30   avoid immune destruction, a process referred to as immunoediting and one of the more

31   recent hallmarks of cancer[4,5]. Mouse studies have shown that T lymphocyte recognition of

32   tumor-specific antigens is crucial for immunoediting to occur[6]. The accumulation of somatic

33   mutations in the tumor genome results in the formation of neoantigens, small peptides

34   presented on the cell surface that can stimulate cytotoxic (CD8+) T lymphocytes (CTLs). To

35   attenuate these CTL responses, a cancer cell can upregulate ligands for checkpoint

36   receptors[7]. Therapeutically blocking these checkpoint pathways has been shown effective in

37   several cancers such as metastatic melanoma and non-small cell lung cancer[7–9]. However,

38   responses to immune checkpoint blockade (ICB) therapy are still largely unpredictable, and

39   it is not completely understood why some tumors do not respond or develop resistance to

40   therapy.

41         Several genomic alterations (e.g. *CASP8* mutations, *B2M* mutations, *HLA* loss) have

42   been discovered that can partially explain this ICB therapy unresponsiveness[10–16].

43   Furthermore, as stimulation of CTLs is critically dependent on the formation and presentation

44   of neoantigens, it is not surprising that one of the main determinants of therapy

45   responsiveness is mutation burden[17–19]. Indeed, the higher the mutation burden, the higher

46   the number of potential neoantigens and hence ways to stimulate the immune system. On

47   the other hand, negative (or purifying) selection is expected to act on neoantigen-forming

48   mutations. This should result in a depletion of such mutations and escape from immune-

49   induced cancer cell death. The presence of neoantigen depletion has been suggested in

50   several cancers such as colorectal cancer, metastatic melanoma, esophageal, bladder,

51   cervical and lung cancer[10,13,20,21]. As the main determinant of CTL immunogenicity is a

52   peptide's capacity to bind to the cell's human leukocyte antigens (HLA) from the type I major

3

53    histocompatibility complex (MHC-I), the conclusions of these studies are mostly based on

54    lower-than-expected numbers of non-synonymous somatic mutations in predicted HLA-

55    binding peptides, using the number of synonymous mutations as a reference.

56          Somatic mutations are caused by different mutational processes that are active

57    during tumor evolution. A widely used method for characterizing the properties of mutational

58    processes are trinucleotide-based mutational signatures, which describe frequencies for all

59    single nucleotide substitutions in all possible sequence contexts in terms of adjacent

60    upstream and downstream nucleotides, resulting in a total of 96 substitution types[3]. This

61    implies that the mutation probability at any genomic position is dependent on the immediate

62    sequence context in combination with the active mutational processes. It has now been

63    clearly demonstrated that mutational signatures need to be accounted for in any model

64    aiming at finding signals of selection in cancer[22–24]. However, it is currently not clear whether

65    and how mutational signatures and their sequence context preferences influence signals of

66    neoantigen depletion.

67          Here we show that, when mutational signatures are considered, putative signals of

68    neoantigen depletion become weak to absent in cancer genomics data from treatment-naïve

69    tumor samples. Our results are in line with the overall weak signals of negative selection in

70    cancer and challenge the idea that neoantigen depletion signals are detectable based on

71    HLA affinity predictions in large-scale cancer mutation datasets.


72    **RESULTS**


73    **Annotation of HLA-binding regions in the human genome.**

74    Somatic mutations are expected to result in neoantigen formation when (i) the resulting

75    peptides are presented via MHC-I and (ii) they are recognized by CTLs through specific T-

76    cell receptor (TCR) binding, which only occurs when there is no immune tolerance, i.e. when

77    presented peptides are new to the immune system. Given sufficient co-stimulatory signals,

78    this will result in CTL-mediated killing of neoantigen-presenting cancer cells, enforcing a

4

79  negative selection pressure during tumor evolution (Fig. 1a). We hypothesized that this

80  specific form of negative selection and hence neoantigen depletion should be detectable as

81  reduced mutation rates in genomic regions that can be translated to HLA-binding peptides.

82  Therefore, our first aim was to define these regions, thereby generating an HLA-binding

83  genomic annotation.

84      HLA-binding affinities are determined by both the amino acid sequence and by

85  patient-specific HLA genotypes, composed of a combination of two HLA-A, two HLA-B and

86  two HLA-C alleles. We initially considered a single prototypical HLA genotype consisting of

87  the two most common HLA-alleles (HLA-A01:01, HLA-A02:01 HLA-B07:02, HLA-B08:01,

88  HLA-C07:01 and HLA-C07:02; Supplementary Fig. 1), enabling us to define a single HLA-

89  binding genome annotation to use throughout the analyses. For these six HLA alleles, the

90  affinities were predicted for all possible nonapeptides (9-mers) translated from the coding

91  genome and were aggregated in a single affinity, a similar approach to what has been

92  described recently[25] (see Methods and Supplementary Fig. 1). By considering a nonapeptide

93  HLA-binding when the aggregated $K_d$ was lower than 500 nM[26], we found that the complete

94  pool of HLA-binding nonapeptides mapped to 22.1% of the exome (Fig. 1b).

95  **Apparent negative selection signals in HLA-binding regions.**

96  Having annotated the human exome for the HLA-binding properties of its translated peptides,

97  we next aimed to search for signals of immune-induced negative selection in the cancer

98  genome. All available synonymous and non-synonymous (i.e. missense) somatic mutation

99  data were downloaded from The Cancer Genome Atlas (TCGA), encompassing 1,836,369

100 mutations from 8,683 different samples and spanning 32 cancer types (Supplementary Table

101 1). As only non-synonymous mutations in HLA-binding regions are expected to be under

102 immunogenic selection pressure, we used the number of synonymous mutations as a

103 background reference and determined the ratio between the observed numbers of non-

104 synonymous and synonymous mutations (n/s) in HLA-binding as well as non-binding regions.

105 We found that n/s was lower in HLA-binding regions on a pan-cancer level (n/s 2.23 in HLA-

5

106     binding vs. 2.58 in non-binding regions, $P = 3.24 \times 10^{-298}$, Fisher's exact test; Fig. 2a,b). To

107     quantify the extent of this putative neoantigen depletion signal, we defined an HLA-binding

108     mutation ratio (HBMR) as the ratio of n/s in HLA-binding to non-binding peptides. This way,

109     negative immunogenic selection of somatic mutations is expected to result in HBMR values

110     lower than 1 (or higher than 1 if these mutations have been influenced by positive selection).

111     For the pan-cancer analysis this implied an HBMR of 0.87, suggesting the overall loss of 13%

112     of non-synonymous mutations due to negative selection (Fig. 2a,b).

113         We next aimed to determine how these signals differed between cancer types

114     focusing on the 19 cancer types with at least 10,000 mutations in the TCGA dataset. Given

115     the observed mutation burdens, we estimate sufficient power (0.8 at $P < 0.05$) to detect

116     negative selection operating on between 2% (UCEC) and 13% (KIRP) of the predicted

117     neoantigens (Supplementary Fig. 2). We observed HBMR values that were significantly

118     below 1 for 12 out of 19 analyzed cancer types, including bladder cancer (BLCA, HBMR =

119     0.66, $P = 1.5 \times 10^{-127}$), metastatic melanoma (SKCM, HBMR = 0.69, $P = 0$), cervical cancer

120     (CESC, HBMR = 0.72, $P = 1.3 \times 10^{-51}$), lung adenocarcinoma (LUAD, HBMR = 0.77, $P = 2.3$

121     $\times 10^{-60}$), head and neck cancer (HNSC, HBMR = 0.78, $P = 6.6 \times 10^{-36}$) and squamous cell

122     lung cancer (LUSC, HBMR = 0.80, $P = 1.4 \times 10^{-34}$) (Fig. 2c and Supplementary Table 2).

123     **Reduced non-synonymous mutations in HLA-binding regions are not caused by**

124     **selection processes.**

125     To be able to determine whether and to what extent selection processes and hence

126     neoantigen depletion are indeed responsible for the observed reduction in non-synonymous

127     mutations in HLA-binding regions, we determined the expected mutation rates in the

128     absence of any selection pressure. For every observed somatic mutation, we simulated one

129     mutation by randomly sampling from all possible point mutations with the same trinucleotide

130     substitution type (e.g. TCC>TTC), resulting in a simulated mutation dataset with a similar

131     size as the observed data. As expected, all signals of positive selection in driver genes

132     disappeared in the simulated mutation data (Supplementary Fig. 3). Using this simulated

6

133   mutation database, we recalculated the mutation rates and HBMR values. Strikingly, a

134   strong signal of apparent negative selection and hence neoantigen depletion, similar to the

135   real mutation data, was still present (HBMR = 0.83, $P$ = 0; Fig. 2b). This similarity was also

136   present for the individual cancer types (Pearson's $r$ = 0.91, $P$ = 7.5 × $10^{-8}$), with the strongest

137   signals again observed for bladder cancer and metastatic melanoma (Fig. 2c). The fact that

138   a set of randomly generated mutations, upon which selection cannot have acted, gave

139   results that closely mimicked those from actual mutation data casts doubt on the apparent

140   neoantigen depletion signals. As the simulated and real mutations were only matched with

141   respect to trinucleotide substitution types, this analysis suggests that sequence differences

142   between HLA-binding and non-binding regions, combined with specific sequence

143   preferences of relevant mutagenic exposures, introduce biases in n/s ratios, leading to

144   apparent signals of neoantigen depletion.

145        We noted that these findings were robust to the way HLA-binding capacity was

146   determined. Determining HLA affinities using patient-specific HLA genotypes (rather than the

147   six most frequent alleles), focusing on the best binding allele only and using a more stringent

148   $K_d$ cut-off of 50 nM or a percentile-based cut-off of 1%, did not substantially alter the

149   observed reduction in HBMR values (Supplementary Fig. 4). Similar results were also

150   obtained when the analysis was restricted to genomic regions encoding known epitopes

151   from IEDB (immune epitope database; Supplementary Fig. 5).

152        While the exclusion of non-expressed or cancer driver genes did not change the

153   observed differences between tumor types either (Supplementary Fig. 4), we observed a

154   lower overall percentage of somatic mutations in HLA-binding regions for expressed

155   compared to non-expressed genes (21.7% vs. 28.3% respectively for the pan-cancer

156   dataset, Fisher's exact test $P$ = 0), and an opposite effect for driver compared to non-driver

157   genes (18.8% vs. 22.8% respectively, $P$ = 7.6 × $10^{-238}$; Supplementary Fig. 6). Similar

158   differences were observed for both non-synonymous and synonymous mutations, again

159   raising doubts about a putative interpretation as immunogenic selection signals. These

160   findings also imply that mutations in non-expressed transcripts should not be used as

7

161   background reference when studying immunogenic selection pressures in cancer genomics

162   data.

**Different trinucleotide substitution probabilities explain lower non-synonymous**

**mutation rates in HLA-binding regions.**

165   To better understand the association between trinucleotide substitution types and HLA-

166   binding regions, we simulated all possible point mutations in 17,992 genes (21,203,704

167   synonymous and 67,766,542 non-synonymous mutations; Fig. 3a) and used the HBMR

168   metric to quantify the difference between expected mutation rates in HLA-binding and non-

169   binding regions for each trinucleotide substitution type. There was a notable variability

170   between the trinucleotide substitution types, with HBMR values ranging from 0.35 for

171   TCT>TGT substitutions to 2.07 for ATG>ACG substitutions (Fig. 3b). The trinucleotide

172   substitution types with the lowest HBMR values were the most abundant in the cancer types

173   with low overall HBMRs (e.g. 23.9% of all malignant melanoma mutations are TCC>TTC, the

174   trinucleotide substitution type with the second to lowest HBMR; Supplementary Fig. 7).

175   Remarkably, many of the substitution types with the lowest HBMR values were TCN>TNN

176   (Fig. 3b), and a strong negative correlation was indeed observed between a cancer type's

177   HBMR value and its proportion of TCN>TNN mutations (Pearson's $r$ = -0.81, $P$ = 2.4 × $10^{-5}$;

178   Supplementary Fig. 7). Mutational signatures 2 and 3 (APOBEC-related) and the UV-

179   induced signature 7, which are both related to these patterns, consequently had the lowest

180   HBMR values (Supplementary Fig. 7).

**High synonymous mutation probabilities in hydrophobic amino acid codons correlate**

**to lower perceived mutation rates in HLA-binding regions.**

183   We next aimed to explain the association between trinucleotide substitution types and HLA-

184   binding properties. Because different sequence contexts imply different amino acid codon

185   probabilities on the one hand, while different physicochemical properties of amino acids

186   influence binding to HLA on the other hand, we investigated the relationships between

8

187  trinucleotide substitution types, the amino acid content of peptides, and their expected

188  HMBR values.

189      We first focused on the correlation between HBMR values and amino acid classes

190  (hydrophobic, polar or charged) in our annotated genome. For synonymous mutations, a

191  strong negative correlation was observed between a trinucleotide substitution type's HBMR

192  value and the frequency of hydrophobic amino acid codons (Spearman's $r$ = -0.61, $P$ = 8.1 ×

193  $10^{-11}$; Fig. 3b), while an opposite, weaker and positive correlation was noted for non-

194  synonymous mutations (Spearman's $r$ = 0.30, $P$ = 4.2 × $10^{-3}$; Fig. 3b). This effect was mainly

195  related to Leu, Val and Iso (Supplementary Fig. 8); hydrophobic amino acids encoded by

196  codons with a thymine on the second codon position (Supplementary Fig. 9). Combined with

197  the observation that most of the corresponding trinucleotide substitution types conform to the

198  pattern TCN>TNN, this association can be explained by the upstream T of the substitution

199  type matching with the T at the second codon position and the substituted nucleotide

200  matching with the third codon position (Fig. 3c). Indeed, when a codon with a T at the

201  second position is hydrophobic, any mutation involving the third position of a Leu or Val

202  codon always results in a synonymous mutation. This is also the case for most mutations

203  that affect the same position in Ile and for some mutations at the Phe codon as exemplified

204  in Figure 3c.

205      Secondly, as hydrophobic amino acids are known to influence HLA-binding

206  affinities[27], we determined the correlation between the number of amino acids from a certain

207  class in a nonapeptide and its HLA-binding capacity. By randomly sampling from 1 million

208  coding regions and determining the translated peptides' HLA-binding affinity, we observed a

209  positive association between the number of hydrophobic amino acids in a peptide and its

210  HLA-binding capacity (logistic regression coefficient $\beta$ = 0.48, Fig. 3d and Supplementary

211  Fig. 10).

212      These results demonstrate that certain trinucleotide substitution types, like

213  TCN>TNN, which occur frequently in metastatic melanoma, bladder cancer and cervical

214  cancer, are likely to lead to synonymous mutations in Leu, Val and Ile codons. Because

9

215 these amino acids are more frequent in HLA-binding peptides, this leads to lower perceived

216 non-synonymous mutation rates when synonymous mutations are used as a background

217 reference. The earlier described difference in apparent neoantigen depletion in expressed vs.

218 non-expressed genes is also related to hydrophobic amino acid content, as a gene

219 enrichment analysis of non-expressed genes showed a strong membrane protein

220 enrichment (e.g. olfactory receptors; Supplementary Fig. 6).

221 **Weak to absent neoantigen depletion signals after correcting for trinucleotide**

222 **substitution effects.**

223 Our study shows that differential mutation rates between HLA binding and non-binding

224 peptides mainly result from differences in trinucleotide substitution probabilities. We next

225 aimed to determine whether any remaining signal of neoantigen depletion would be

226 detectable after correcting for these trinucleotide substitution effects.

227 As a first approach, we normalized the observed HBMR value to its expected value

228 for each cancer, under a trinucleotide substitution model and considering the HLA-binding

229 annotation developed in this study (see Methods). We reanalyzed all cancers and observed

230 a disappearance of neoantigen depletion signals, except for a limited signal in lung cancer

231 (Fig. 4a and Supplementary Table 2). In line with our earlier findings (Supplementary Fig. 4),

232 results did not substantially change when different criteria were used to calculate HLA

233 binding capacity or when mutations were called using the more recent MC3 mutation caller[28]

234 (Supplementary Fig. 11). Similarly, dN/dS values did not suggest any signal of negative

235 selection after correcting for differing trinucleotide sequence contexts in HLA binding vs.

236 non-binding regions (Supplementary Fig. 12).

237 Notably, correcting using mutation probabilities derived from the SSB7 or other

238 models that do not consider the complete adjacent sequence context resulted in corrected

239 signals falsely suggestive of neoantigen depletion in e.g. melanoma and bladder cancer (Fig.

240 4a and Supplementary Fig. 12). Conversely, normalization using an extended sequence

10

241 context (pentanucleotide substitution model) further decreased the apparent selection

242 signals, with loss of significance in lung squamous cell carcinoma (Fig. 4a,b).

243       The previous results were all derived for a prototypical HLA genotype and for the

244 reference genome (i.e. wild-type peptides). While this approach was useful in gathering new

245 insights into associations between substitution types and HLA affinities, there is a risk of

246 missing selection signals that are HLA genotype-specific and/or only act on mutations that

247 result in new HLA binders (i.e. hit the HLA-binding residues of a nonapeptide, rather than the

248 CTL contact residues). We thus searched for neoantigen depletion signals in mutated HLA-

249 binding peptides, where binding affinities were predicted for sample-specific genotypes. We

250 noted that only 1.88% of all non-synonymous mutations resulted in a non-binding peptide

251 gaining HLA-binding properties (Supplementary Fig. 13). Similar numbers (1.92%) were

252 found using our simulated mutation database, thus again providing no convincing support for

253 selection acting on these specific mutations.

254       Finally, given that we have shown that synonymous mutation counts are particularly

255 vulnerable to the effects of mutational signatures, we considered a selection metric

256 ($dN_{HLA}/dN_{nonHLA}$) that was independent of synonymous mutations. This metric compares the

257 observed ratio between the number of non-synonymous mutations in HLA-binding and non-

258 binding peptides with the corresponding expected ratio. The latter was determined for each

259 HLA genotype from all TCGA samples, using mutated peptides from 960,000 randomly

260 simulated mutations (10,000 for each trinucleotide substitution type) and considering the

261 aggregated mutational signature from each cancer type (Fig. 5a,b). By normalizing the

262 observed to the expected ratios for each sample, all tumor types were reanalyzed for

263 putative selection signals. This analysis generally confirmed the absence of detectable

264 neoantigen depletion, except for a signal in cervical cancer (median $dN_{HLA}/dN_{nonHLA}$ = 0.91,

265 one-sample Wilcoxon signed-rank test $P = 2.4 \times 10^{-4}$; Fig. 5c and Supplementary Table 2).

266 Further, $dN_{HLA}/dN_{nonHLA}$ did not correlate with immune cytolytic activity (Supplementary Fig.

267 14). Notably, 3 out of 19 tumor types had values significantly above 1. These signals were

268 comparable in effect size to cervical cancer, most pronounced in melanoma (median

11

269     $dN_{HLA}/dN_{nonHLA} = 1.08$, $P = 1.2 \times 10^{-10}$), and remained when using a pentanucleotide rather

270     than trinucleotide model (Supplementary Fig. 14). As these positive signals are unlikely to

271     indicate true positive selection, they may rather reflect limitations of the $dN_{HLA}/dN_{nonHLA}$ model,

272     which does not consider synonymous mutation rates. Finally, neoantigen depletion signals

273     were absent when the number of non-synonymous mutations in HLA-binding peptides was

274     normalized to an expected number that was estimated directly from the pan-cancer dataset,

275     as suggested previously[10] (Supplementary Fig. 14). Notably, we observed that the

276     neoantigen depletion signals in colorectal and kidney cancer, as reported by Rooney et al.[10] ,

277     disappeared after excluding samples with miscalled HLA genotypes from the original dataset

278     (results obtained using authors' source code; Supplementary Fig. 14).

279        Taken together, these results point to a general absence of detectable neoantigen

280     depletion signals in large-scale mutation data from untreated tumors and emphasize the

281     importance of using accurate background mutation models to correct for sequence biases

282     introduced by relevant mutational processes.

283     **DISCUSSION**

284     In this study, we initially observed an apparent reduction of somatic point mutations in

285     genomic regions encoding HLA-binding nonapeptides. Rather than being an effect of

286     negative selection acting on immunogenic mutations, we demonstrated correlative

287     relationships between the probability of mutagenesis in different nucleotide sequences and

288     predicted HLA affinities for corresponding peptides. In particular, the number of hydrophobic

289     amino acids are a major determinant of HLA binding capacity for a peptide while

290     simultaneously being a strong determinant of mutation rate, depending on the mutational

291     processes at play. When correcting for these correlations, detectable negative selection

292     signals were weak to absent. Our results demonstrate that mutation rate differences

293     between peptides with variable HLA affinities should be interpreted with care and have broad

294     relevance for other studies that derive selection signals from HLA affinity predictions.

295     To detect immunogenic selection signals, we initially annotated the human exome

296     with respect to HLA-binding capacity by determining which segments are translatable to

297     HLA-binding peptides, for simplicity assuming a single prototypical HLA genotype for all

298     samples. This implies a focus on wild-type peptides under the hypothesis that mutations in

299     CTL contact residues are subject to negative selection pressures. Using this annotation, the

300     approach can be easily reproduced on any mutation dataset, without the need for complex

301     and time-consuming HLA-typing or HLA affinity predictions. The theoretical drawback is that

302     this does not capture neoantigenic mutations that lead to new HLA-binding peptides (i.e.

303     increase the HLA affinities) and/or effects that are HLA genotype-specific. However,

304     additional analyses addressing patient-specific HLA genotypes as well as *de novo* HLA-

305     binding peptides likewise failed to produce strong support for neoantigen depletion signals.

306     Synonymous mutations are often used as a background mutation reference when

307     analyzing non-synonymous substitutions with respect to selection, resulting in metrics such

308     as dN/dS. Recent studies have shown that these metrics get confounded when not

309     considering the adjacent sequence context[22,23]. A key finding of our study is that simplistic

310     substitution models will lead to biased immunogenic selection signals, due to HLA affinity

311     predictions also being sequence dependent. An important advantage of any metric that

312     considers synonymous mutations as a background reference (like HBMR) is that any

313     unexpected property that equally effects synonymous and non-synonymous mutation rates

314     will be cancelled out (such as differential mutation burdens in expressed and non-expressed

315     genes). However, given that we observed strong dependencies specifically between

316     synonymous mutation probabilities and HLA-binding properties of corresponding encoded

317     peptides, leaving synonymous mutations out of the equation may also have advantages. We

318     did this by considering the ratio between the observed number of non-synonymous

319     mutations in HLA-binding and non-binding regions and normalizing this ratio to an expected

320     ratio, estimated under a trinucleotide substitution model for each individual HLA genotype.

321     Calculation of the resulting $dN_{HLA}/dN_{nonHLA}$ metric for each sample did not provide clear

322     evidence of neoantigen depletion, similar to our initial analysis taking synonymous mutations

13

323  into account. We could only detect a weak signal in cervical cancer and demonstrated that

324  the previously reported neoantigen depletion signal in colorectal adenocarcinoma[10] was due

325  to HLA genotyping problems in samples that were later removed from TCGA. Notably, the

326  $dN_{HLA}/dN_{nonHLA}$ approach also indicated positive signals in some cancers, at effect sizes

327  comparable to the depletion in cervical cancer. Since positive selection in HLA-binding

328  regions is improbable, this likely reflects limitations in the accuracy of the expectation model,

329  casting doubt on the observed negative signal in cervical cancer as well. While this may

330  reflect exclusion of synonymous mutations in this metric, it can also be noted that mutational

331  signatures were here determined at the tumor type level, and it is possible that consideration

332  of patient-specific mutational signatures from whole genome sequencing datasets may

333  potentiate more refined analyses in the future.

334       In addition to point mutations, which have been the main focus of studies of

335  neoantigen depletion, future studies should also address frameshifting indels in this context.

336  This is a different challenge, as single indels may generate large numbers of unnatural

337  peptides through introduction of novel open reading frames, which may or may not be

338  subject to nonsense-mediated decay[29]. Consistently, indels have been described as more

339  strongly associated with response to immunotherapy[30], and it can be noted that

340  microsatellite unstable colon cancers, which harbor larger numbers of indels, appear

341  responsive to checkpoint inhibitors while normal colon carcinomas are not[31].

342       In summary, our results indicate that signals of neoantigen depletion, detected using

343  HLA affinity predictions, are overall weak to absent in the untreated cancer genome. While

344  we cannot exclude that this is related to poor accuracy to predict neoantigen formation

345  (Supplementary Fig. 2), it is noteworthy that signals of negative selection in general are

346  weak in cancer mutation data[22,23,32,33]. Therefore, either only a very small fraction of

347  predicted neoantigenic sites are immunogenic, or the lack of negative selection signals

348  suggests that developing tumors possess or evolve efficient immune evasion mechanisms

349  (e.g. *HLA* loss or *PDL1* amplification). If this is indeed the case, detectable signals of

14

350 neoantigen depletion are only expected in the absence of these escape mechanisms, such

351 as after ICB therapy[21].

## ACKNOWLEDGEMENTS

## AUTHOR CONTRIBUTIONS

365 J.V.d.E., E.L. and M.L.M. designed and conceptualized the study. A.J.-S. provided input on

366 study design and contributed to the interpretation of the results. J.V.d.E. was responsible for

367 data analysis and drafted the manuscript. All authors discussed the results, edited and

368 finalized the manuscript.

## COMPETING INTERESTS STATEMENT

370 The authors declare they have no conflicts of interest.

371

## REFERENCES

1. Vogelstein, B. *et al.* Cancer genome landscapes. *Science* **339**, 1546–1558 (2013).

2. Cancer Genome Atlas Research, N. *et al.* The Cancer Genome Atlas Pan-Cancer analysis project. *Nat. Genet.* **45**, 1113–1120 (2013).

3. Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature* **500**, 415–421 (2013).

4. Hanahan, D. & Weinberg, R. A. Hallmarks of cancer: The next generation. *Cell* **144**, 646–674 (2011).

5. Dunn, G. P., Old, L. J. & Schreiber, R. D. The three Es of cancer immunoediting. *Annu. Rev. Immunol.* **22**, 329–360 (2004).

6. DuPage, M., Mazumdar, C., Schmidt, L. M., Cheung, A. F. & Jacks, T. Expression of tumour-specific antigens underlies cancer immunoediting. *Nature* **482**, 405–409 (2012).

7. Pardoll, D. M. The blockade of immune checkpoints in cancer immunotherapy. *Nat. Rev. Cancer* **12**, 252–264 (2012).

8. Hodi, F. S. *et al.* Improved survival with ipilimumab in patients with metastatic melanoma. *N. Engl. J. Med.* **363**, 711–723 (2010).

9. Sharma, P. & Allison, J. P. Immune checkpoint targeting in cancer therapy: toward combination strategies with curative potential. *Cell* **161**, 205–214 (2015).

10. Rooney, M. S., Shukla, S. A., Wu, C. J., Getz, G. & Hacohen, N. Molecular and genetic properties of tumors associated with local immune cytolytic activity. *Cell* **160**, 48–61 (2015).

11. Shukla, S. A. *et al.* Comprehensive analysis of cancer-associated somatic mutations in class I HLA genes. *Nat. Biotechnol.* **33**, 1152–1158 (2015).

12. McGranahan, N. *et al.* Allele-specific HLA loss and immune escape in lung cancer evolution. *Cell* **171**, 1259-1271.e11 (2017).

16

398   13.   Davoli, T., Uno, H., Wooten, E. C. & Elledge, S. J. Tumor aneuploidy correlates with

399          markers of immune evasion and with reduced response to immunotherapy. *Science*

400          **355**, eaaf8399 (2017).

401   14.   Rutledge, W. C. *et al.* Tumor-infiltrating lymphocytes in glioblastoma are associated

402          with specific genomic alterations and related to transcriptional class. *Clin. Cancer Res.*

403          **19**, 4951–4960 (2013).

404   15.   Brown, S. D. *et al.* Neo-antigens predicted by tumor genome meta-analysis correlate

405          with increased patient survival. *Genome Res.* **24**, 743–750 (2014).

406   16.   Rosenthal, R. *et al.* Neoantigen-directed immune escape in lung cancer evolution.

407          *Nature* **567**, 479–485 (2019).

408   17.   Van Allen, E. M. *et al.* Genomic correlates of response to CTLA-4 blockade in

409          metastatic melanoma. *Science* **350**, 207–211 (2015).

410   18.   Snyder, A. *et al.* Genetic basis for clinical response to CTLA-4 blockade in melanoma.

411          *N. Engl. J. Med.* **371**, 2189–2199 (2014).

412   19.   Rizvi, N. A. *et al.* Cancer immunology. Mutational landscape determines sensitivity to

413          PD-1 blockade in non-small cell lung cancer. *Science* **348**, 124–128 (2015).

414   20.   Zapata, L. *et al.* Negative selection in tumor genome evolution acts on essential

415          cellular functions and the immunopeptidome. *Genome Biol.* **19**, 67 (2018).

416   21.   Riaz, N. *et al.* Tumor and microenvironment evolution during immunotherapy with

417          nivolumab. *Cell* **171**, 934-949.e15 (2017).

418   22.   Van den Eynden, J. & Larsson, E. Mutational signatures are critical for proper

419          estimation of purifying selection pressures in cancer somatic mutation data when

420          using the dN/dS metric. *Front. Genet.* **8**, 74 (2017).

421   23.   Martincorena, I. *et al.* Universal patterns of selection in cancer and somatic tissues.

422          *Cell* **171**, 1029-1041.e21 (2017).

423   24.   Lawrence, M. S. *et al.* Mutational heterogeneity in cancer and the search for new

424          cancer-associated genes. *Nature* **499**, 214–218 (2013).

425   25.   Marty, R. *et al.* MHC-I genotype restricts the oncogenic mutational landscape. *Cell* **17**,

426      1272–1283.e15 (2017).

427   26.   Paul, S. *et al.* HLA class I alleles are associated with peptide-binding repertoires of

428      different size, affinity, and immunogenicity. *J. Immunol.* **191**, 5831–5839 (2013).

429   27.   Chowell, D. *et al.* TCR contact residue hydrophobicity is a hallmark of immunogenic

430      CD8+ T cell epitopes. *Proc. Natl. Acad. Sci. USA* **112**, E1754-E1762 (2015).

431   28.   Ellrott, K. *et al.* Scalable open science approach for mutation calling of tumor exomes

432      using multiple genomic pipelines. *Cell Syst.* **6**, 271-281.e7 (2018).

433   29.   Turajlic, S. *et al.* Insertion-and-deletion-derived tumour-specific neoantigens and the

434      immunogenic phenotype: a pan-cancer analysis. *Lancet Oncol.* **18**, 1009–1021 (2017).

435   30.   Mandal, R. *et al.* Genetic diversity of tumors with mismatch repair deficiency

436      influences anti-PD-1 immunotherapy response. *Science* **364**, 485–491 (2019).

437   31.   Stein, A. & Folprecht, G. Immunotherapy of colon cancer. *Oncol. Res. Treat.* **41**, 282–

438      285 (2018).

439   32.   Van den Eynden, J., Basu, S. & Larsson, E. Somatic mutation patterns in hemizygous

440      genomic regions unveil purifying selection during tumor evolution. *PLoS Genet.* **12**,

441      e1006506 (2016).

442   33.   Weghorn, D. & Sunyaev, S. Bayesian inference of negative and positive selection in

443      human cancers. *Nat. Genet.* **49**, 1785–1788 (2017).

444

445

**FIGURE LEGENDS**

447 **Figure 1 | Development of an HLA-binding genomic annotation to detect somatic**

448 **mutations under immunogenic selective pressure.**

449 **a**, Neoantigen formation is expected when a non-synonymous mutation leads to a structural

450 change in the CTL (CD8+ cytotoxic T lymphocyte) contact residues of an HLA-binding

451 nonapeptide. This can result in CTL-mediated apoptotic cell death and hence negative

452 selection of the underlying somatic mutation. TCR, T cell receptor; MHC-I, type I major

453 histocompatibility complex. **b**, Binding affinities of all possible nonapeptides were determined

454 for the six most common HLA alleles as indicated. Peptides were considered HLA-binding

455 when their aggregated $K_d$ over the six alleles was below 500 nM (see Methods); HLA-

456 binding peptides mapped to 22.1% of the exome as indicated.

457 **Figure 2 | Analysis of somatic mutation rates in HLA-binding annotated genomic**

458 **regions.**

459 **a**, Contingency table showing the total number of synonymous (s) and non-synonymous (n)

460 mutations in the HLA-binding and non-binding exome. The HLA-binding mutation ratio

461 (HBMR) indicates the ratio of n/s in HLA-binding to non-binding regions. **b**, Bar plot

462 comparing the n/s ratios of observed and simulated mutations. **c**, HBMR calculated for

463 observed and simulated mutations from 19 cancer types containing at least 10,000 somatic

464 mutations per cancer type. Error bars indicate 95% confidence intervals, calculated using

465 two-sided Fisher's exact test. Pearson correlation coefficient (*r*) and *P* value indicated on top

466 left. See Supplementary Table 1 for cancer type abbreviations and sample sizes.

467 **Figure 3 | Association between trinucleotide substitution types and HLA-binding**

468 **properties.**

469 **a**, All possible synonymous and non-synonymous mutations were determined in 17,992

470 genes. Pie charts indicate the proportions of mutations that are located in HLA-binding

471 regions. **b**, Bar plot on top indicates the expected HBMR value for each trinucleotide

472     substitution type, determined from all possible mutations from a given type in the complete

473     exome (numbers shown in **a**). Main substitution types are colored as indicated by the legend

474     on top left. Note that HBMR values are not derivable for four trinucleotide substitution types

475     (ATT>AAT, ATT>AGT, ACT>AGT and ACT>AAT) due to the absence of synonymous

476     mutations resulting from these substitution types (e.g. an ATT>AAT substitution can never

477     be synonymous). TCN>TNN substitutions are indicated by red asterisks. Below the bar plot,

478     the frequency of synonymous and non-synonymous mutations hitting hydrophobic amino

479     acids is indicated for each substitution type (scale indicated on bottom right). Loess

480     regression line in red with Spearman correlation coefficient (*r*) and *P* value indicated on top

481     right (correlation between HBMR and mutation frequency for 92 different substitution types).

482     **c**, Illustration of TCN>TNN mutations mainly resulting in synonymous mutations in

483     hydrophobic amino acid codons. **d**, Logistic regression line indicating the correlation

484     between a nonapeptide's mean number of hydrophobic/charged/polar amino acids (0 to 9)

485     and the HLA-binding probability. Regression coefficients (β) are given for each amino acid

486     class. The mean number of amino acids for each class was determined for 1 million random

487     exome locations (9 nonapeptides per position) to make the analysis comparable to the other

488     analyses. A similar analysis on individual nonapeptides is shown in Supplementary Figure

489     10.


490     **Figure 4 | Weak to absent neoantigen depletion signals after correcting for**

491     **trinucleotide-based mutational signature effects.**

492     **a**, Bar plot showing normalized HBMR values for 19 different cancer types. HBMR values

493     were obtained by normalization of the observed HBMR values to the expected tumor-type

494     specific values. The latter were calculated using mutation probabilities derived from different

495     models as indicated on top left. Error bars indicate 95% confidence intervals, calculated

496     using two-sided Fisher's exact test. See Supplementary Table 1 for cancer type

497     abbreviations and sample sizes and Supplementary Table 2 for detailed results. **b**,

498     Comparison of HBMR deviations from 1 after normalization using different substitution

499    models as indicated. Each dot represents a cancer type. Median values are indicated by

500    horizontal lines.

501    **Figure 5 | An HLA genotype-specific analyses of mutated peptides confirms the**

502    **absence of neoantigen depletion signals in most tumor types.**

503    **a**, Methodological approach. For each trinucleotide substitution type (i), 10,000 mutations

504    were randomly simulated (960,000 mutations in total). The expected number of non-

505    synonymous mutations in HLA-binding and non-binding peptides were derived for each

506    substitution type considering the mutated peptides' HLA affinities for the sample-specific

507    HLA genotype (heatmap on bottom). From these numbers, the expected ratio between non-

508    synonymous mutations in HLA-binding and non-binding peptides was calculated using the

509    substitution probabilities of the corresponding cancer type (legoplot on top). **b**, Scatter plot

510    shows the correlation between observed and expected ratios, with Pearson correlation

511    coefficients (*r*) and *P* values indicated on top left. **c**, $dN_{HLA}/dN_{nonHLA}$ values were calculated

512    for each TCGA sample and grouped by tumor types. Boxplots indicate median values and

513    lower/upper quartiles with whiskers extending to 1.5x the interquartile range. Two-sided

514    Wilcoxon signed-rank test was used to test deviation from 1. *P* values are given for cancers

515    with *q* values below 0.1. Mutations in cancer driver genes or non-expressed genes were

516    excluded. See Supplementary Table 1 for cancer type abbreviations and sample sizes and

517    Supplementary Table 2 for detailed results.

518    **METHODS**

519    **TCGA mutation and expression data.**

520    MuTect2-called whole exome sequencing (WES) mutation annotation format (maf) files from

521    all 33 available cancer types from The Cancer Genome Atlas (TCGA) were downloaded

522    from the Genomic Data Commons (GDC) Data Portal (data release v7). Colon and rectal

523    adenocarcinoma were considered as a single cancer type for the analysis. All mutation data

524    were fused in a single mutation database and were converted from hg38 to hg19 using

525　UCSC's liftOver[34]. Variants were reannotated using ANNOVAR[35]. For each mutation, the

526　main substitution type (i.e. C>A, C>G, C>T, T>A, T>C and T>G) was derived by converting

527　each purine substitution to its complementary base substitution. To determine the

528　trinucleotide substitution type, additional information was added regarding the identity of the

529　upstream and downstream base. Sequence information was derived from UCSC hg19[34].

530　　　　TCGA Level 3 RNASeqV2 (RSEM normalized) mRNA expression data were

531　downloaded from the Broad Institute TCGA Genome Data Analysis Center (2016): Firehose

532　stddata__2016_01_28 run (Broad Institute of MIT and Harvard; doi:10.7908/C11G0KM9).

533　Expression data were fused in a single gene x sample matrix. Each mutation's gene

534　expression value was added to the mutation database.

535　**HLA typing.**

536　HLA typing of all TCGA samples was performed using Polysolver[11]. WES normal bam files

537　from all available TCGA samples were accessed using FireCloud[36], the HLA regions from

538　the main HLA-alleles (HLA-A, HLA-B and HLA-C) in chromosome 6 (coordinates

539　6:29909037-29913661; 6:31321649-31324964; 6:31236526-31239869) were extracted and

540　the resulting bam files were downloaded. Polysolver was run on these bam files using

541　default settings and without setting prior population probabilities, resulting in the successful

542　genotyping of 8,968 TCGA samples. The resulting output was converted in a sample x HLA

543　allele matrix. To validate this HLA typing, the derived frequencies for each HLA allele were

544　compared with the allele frequencies from a healthy US blood donor population, downloaded

545　from Allele frequency net[37] (Supplementary Fig. 1).

546　**HLA affinity predictions and annotation of the HLA-binding genome.**

547　Using the R *GenomicRanges* package[38] and UCSC hg19 genome sequence information, a

548　*GPos* object was created containing information about the complete exome. For each coding

549　DNA sequence (CDS) position, the amino acid sequences of the nine possible translated 9-

550　mers (nonapeptides) were determined using Ensembl 75. Genes with unavailable or

ambiguous protein information in Ensembl were discarded, resulting in a *GPos* object

containing nonapeptide information of 17,992 genes. HLA affinities of these nonapeptides

were predicted for the most frequent HLA alleles (A02:01, A01:01, B07:02, B08:01, C07:01,

C07:02; a combination referred to as the prototypical genotype) using netMHCPan3.0[39]. For

each CDS position, the best binding peptide (peptide with the lowest predicted $K_d$ value) was

determined for each of the six HLA alleles. Finally, one aggregated $K_d$ value was calculated

using the harmonic mean value of the $K_d$ values of the six different peptides (one from each

allele) and all genomic regions with aggregated $K_d$ values below 500 nM were considered as

HLA-binding regions. The same methodology was used to predict HLA affinities in TCGA

somatic mutation data. These TCGA predictions were done for both the prototypical and the

sample-specific HLA genotype (specific combination of two HLA-A, two HLA-B and two HLA-

C alleles) and for wild-type as well as mutated peptides.

**Simulation of somatic mutations.**

All possible point mutations were determined for 17,992 genes by considering for each CDS

position the three possible substitutions (any nucleotide can be substituted in three different

nucleotides). ANNOVAR[35] was used to annotate the variants and determine the reference

and alternative amino acids for each mutation. This information was added to the higher

described *GPos* object.

To determine the expected somatic mutation rates in the absence of any selection

pressure, a simulated mutation database was created, with a similar size as the TCGA

mutation database. To match this simulation database for differences in trinucleotide

substitution probabilities, we randomly sampled the observed number of mutations from

each corresponding substitution type from the *GPos* object. Like for the observed TCGA

mutations, HLA affinities were predicted for the wild-type and the mutated nonapeptides and

for both the prototypical and the sample-specific genotype. The later was determined by

scrambling the columns from the sample x HLA allele matrix. This way, completely random

23

577    HLA genotypes were generated, with the same allele frequency and mutation frequency per

578    type as in the real data.

579    **Amino acid analysis.**

580    To derive the probability of any substitution type to hit a certain amino acid or class of amino

581    acids, we used the *GPos* object containing all possible mutations and determined the amino

582    acid frequency for each substitution type and separately for synonymous and non-

583    synonymous mutations. Amino acids were grouped in three classes: hydrophobic (Gly, Ala,

584    Pro, Val, Leu, Iso, Met, Trp and Phe), polar (Ser, Thr, Tyr, Asn, Gln and Cys) and charged

585    (Lys, Arg, His, Asp and Glu).

586    **Calculation of the HLA-binding mutation ratio (HBMR) and related metrics.**

587    To quantify putative signals of immunogenic selection, we defined an HLA-binding mutation

588    ratio (HBMR):

589    $$HBMR = \frac{n_+/s_+}{n_-/s_-}$$

590    where n+ and n- are the total number of non-synonymous mutations located in HLA-binding

591    and non-binding regions, respectively. Similarly, s+ and s- are the number of synonymous

592    mutations in- and outside HLA-binding genomic regions. A similar metric, called the epitope

593    mutation ratio (EMR) was calculated for the analysis of the IEDB epitopes. Here, + and –

594    refer to the location inside and outside of epitope mapped regions. HBMR *P* values and 95%

595    confidence intervals were calculated using a two-sided Fisher's exact test.

596        dN/dS was calculated considering differences in specific trinucleotide substitution

597    probabilities between cancer types[22]:

598    $$\frac{dN}{dS} = \frac{n/\sum_i N_i P_i}{s/\sum_i S_i P_i}$$

       $$with\ i\ \in\ \{A[C > A]A, ..., T[T > G]T\}\ (96\ substition\ types)$$

599

       24

600     where $N_i$ and $S_i$ are the number of (non-)synonymous sites with class i substitutions and $P_i$ is

601     the probability of substitution class i.

602        The normalized HBMR was calculated as follows:

603
$$Normalized\ HBMR = \frac{HBMR_{obs}}{HBMR_{exp}}$$

604
$$with\ HBMR_{exp} = \frac{N_+/S_+}{N_-/S_-} = \frac{\sum_i N_{i+} P_i / \sum_i S_{i+} P_i}{\sum_i N_{i-} P_i / \sum_i S_{i-} P_i}$$

605     where $N_{i+}$ and $S_{i+}$ are the number of (non-)synonymous sites with class i substitutions in

606     HLA-binding regions, $N_{i-}$ and $S_{i-}$ are the number of (non-)synonymous sites with class i

607     substitutions in non-HLA-binding regions respectively and $P_i$ is the probability of substitution

608     class i.

609        The $dN_{HLA}/dN_{nonHLA}$ ratio was calculated for each TCGA sample as follows:

610
$$\frac{dN_{HLA}}{dN_{nonHLA}} = \frac{n_+/n_-}{N_+/N_-} = \frac{n_+/n_-}{\sum_i N_{i+} P_i / \sum_i N_{i-} P_i}$$

611     with variables as defined above, but with HLA affinities determined for mutated peptides

612     from individual genotypes. The number of HLA-binding and non-binding sites was

613     determined for each individual TCGA genotype, under a trinucleotide substitution model. To

614     achieve this, 960,000 substitutions were randomly sampled from the complete exome

615     (10,000 for each substitution type) and HLA affinities were predicted for all the mutations,

616     considering the cancer-type-specific mutational signature.

617        The ratio R of observed to expected neoantigens as described by Rooney *et al.*[10]

618     was calculated for each TCGA sample as follows:

619
$$R = \frac{n_+/n}{N_+/N} = \frac{n_+/n}{\sum_i S_i \frac{\overline{N_i}}{\overline{S_i}} \frac{\overline{N_{i+}}}{\overline{N_i}} / \sum_i S_i \frac{\overline{N_i}}{\overline{S_i}}}$$

620  where $\overline{N_i}/\overline{S_i}$ is the expected number of non-synonymous mutations per synonymous site and

621  $\overline{N_{i+}}/\overline{N_i}$ refers to the expected number of HLA-binders per non-synonymous site, both for

622  substitution type i and estimated empirically from the pan-cancer dataset. Note that these

623  variables are similar to the originally defined variables $\overline{N}_{s(m)}$ and $\overline{B}_{s(m)}$, respectively.

624  Similarly, n+ and N+ were originally called $B_{obs}$ and $B_{pred}$, while n and N were originally

625  referred to as $N_{obs}$ and $N_{pred}$. They were defined here as such to be consistent with the rest

626  of the methodology.

627  Calculation of these metrics was always based on a trinucleotide substitution model

628  as indicated (i index). The normalized HBMR, dN/dS and $dN_{HLA}/dN_{nonHLA}$ were also

629  calculated using alternative substitution models, either based on the six main substitution

630  classes, pentanucleotide substitution classes or using the SSB7 model. The latter is based

631  on the six main substitution classes but considers CpG mutations as a separate class[20].


632  **Neoantigen depletion simulation and power analysis.**

633  All metrics developed in this study were evaluated using an *in silico* analysis of neoantigen

634  depletion by removing increasing amounts of non-synonymous mutations hitting HLA-

635  binding regions from the mutation dataset.

636  Statistical power of the HBMR metric was evaluated using the R *exact2x2* package

637  (Fisher's exact test at significance level 0.05) for different amounts of neoantigen depletion,

638  numbers of mutations and neoantigen prediction accuracies. For this analysis, the non-

639  synonymous mutation proportion (71%) and HLA-binding proportion (22.1%) were fixed to

640  values derived from the pan-cancer dataset and the HLA-binding annotation respectively.

641  For the power analysis of the $dN_{HLA}/dN_{nonHLA}$ ratio, the ratios obtained from the

642  simulated mutation database (containing no selection signals) were log-transformed to

643  obtain a normal distribution. After resampling 1,000 times a predefined amount of values

644  from this normal distribution and adding an *in silico* amount of neoantigen depletion, power

645  was determined based on the number of significant deviations from 0 (corresponding to 1 in

646 non-logtransformed data) using Wilcoxon signed-rank test at $P < 0.05$. This analysis was

647 performed again for different amounts of neoantigen depletion, sample numbers and

648 neoantigen prediction accuracies.

649 **Human epitope mapping.**

650 Data from 66,698 known human IEDB (Immune Epitope Database) epitopes were

651 downloaded from synapse at https://www.synapse.org/ (id syn11935058)[20]. These epitopes

652 were mapped to the human genome (hg19) using the *proteinToGenome* function from the

653 *ensembldb* R package and the *EnsDb.Hsapiens.v75* R library. Mapping was successful for

654 66,536 (99.8%) epitopes.

655 **Statistical analysis.**

656 The R statistical package was used for all data processing and statistical analysis. Details on

657 statistical tests used are reported in the respective sections. Further information on research

658 design is available in the Life Sciences Reporting Summary.

659 **DATA AVAILABILITY**

660 This study is based on public data (open or controlled access) from The Cancer Genome

661 Atlas Network. Downstream data and source code are available at zenodo

662 (https://doi.org/10.5281/zenodo.2621365 and https://doi.org/10.5281/zenodo.3461642,

663 respectively).

664 **METHODS-ONLY REFERENCES**

665 34. Rosenbloom, K. R. *et al.* The UCSC Genome Browser database: 2015 update.

666 *Nucleic Acids Res.* **43**, D670-D681 (2014).

667 35. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic

668 variants from high-throughput sequencing data. *Nucleic Acids Res* **38**, e164 (2010).

669 36. Birger, C. *et al.* FireCloud, a scalable cloud-based platform for collaborative genome

670     analysis: Strategies for reducing and controlling costs. *bioRxiv* 209494 (2017).

671     doi:10.1101/209494

672  37.  González-Galarza, F. F. *et al.* Allele frequency net 2015 update: new features for HLA

673     epitopes, KIR and disease and HLA adverse drug reaction associations. *Nucleic*

674     *Acids Res.* **43**, D784–D788 (2015).

675  38.  Lawrence, M. *et al.* Software for computing and annotating genomic ranges. *PLoS*

676     *Comput. Biol.* **9**, e1003118 (2013).

677  39.  Nielsen, M. & Andreatta, M. NetMHCpan-3.0; improved prediction of binding to MHC

678     class I molecules integrating information from multiple receptor and peptide length

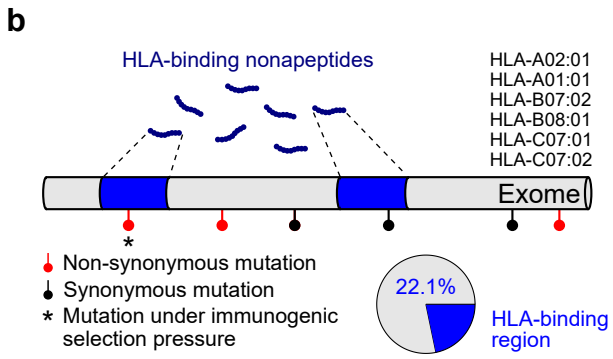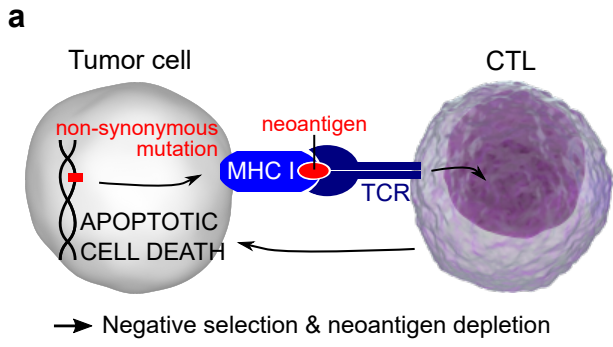679     datasets. *Genome Med.* **8**, 33 (2016).

**a**

Tumor cell

CTL

non-synonymous
mutation

neoantigen

MHC I

TCR

APOPTOTIC
CELL DEATH

→ Negative selection & neoantigen depletion

**b**

HLA-binding nonapeptides

HLA-A02:01
HLA-A01:01
HLA-B07:02
HLA-B08:01
HLA-C07:01
HLA-C07:02

Exome

*

↓ Non-synonymous mutation
↓ Synonymous mutation
* Mutation under immunogenic
  selection pressure

22.1%

HLA-binding
region

*Figure 1*

**a**

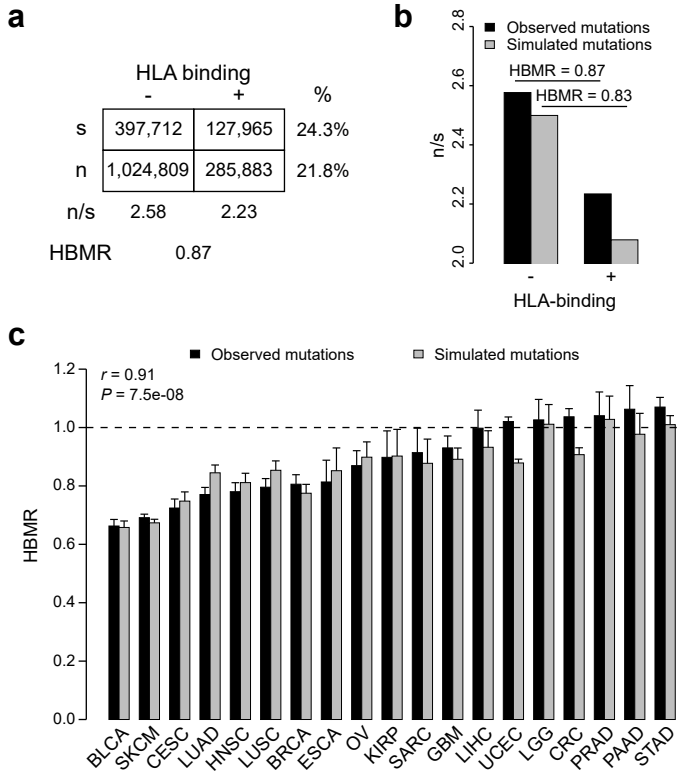|   | HLA binding | | % |
|---|---|---|---|
|   | **-** | **+** |   |
| s | 397,712 | 127,965 | 24.3% |
| n | 1,024,809 | 285,883 | 21.8% |
| n/s | 2.58 | 2.23 |   |
| HBMR | 0.87 | | |

**b**

**c**

*Figure 2*

Figure 3

*Figure 4*

**a** Trinucleotide substitution probability ($P_i$) per cancer

Simulated $N_i$ HLA / $N_i$ nonHLA per HLA genotype

Expected ratio = $\Sigma_i N_{i\,HLA} P_i / \Sigma_i N_{i\,nonHLA} P_i$

**b**

Expected ratio $N_{HLA}/N_{nonHLA}$

Observed ratio $n_{HLA}/n_{nonHLA}$

$P$=0
r=0.62

**c**

$dN_{HLA}/dN_{nonHLA}$ (obs/exp)

CESC   $P$ = 2.4e-04
CRC    $P$ = 2.0e-02
LUSC   $P$ = 5.0e-03
SKCM   $P$ = 1.2e-10

*Figure 5*