

1 Evolutionary dynamics of neoantigens in growing tumors

2 Eszter Lakatos¹, Marc J. Williams¹, Ryan O. Schenck^{2,3}, William C. H. Cross¹, Jacob
3 Househam¹, Luis Zapata⁴, Benjamin Werner⁴, Chandler Gatenbee², Mark Robertson-Tessi²,
4 Chris P. Barnes⁵, Alexander R. A. Anderson², Andrea Sottoriva^{4,*}, Trevor A. Graham^{1,*}

5

6 ¹ Evolution and Cancer Laboratory, Centre for Genomics and Computational Biology, Barts Cancer Institute, School
7 of Medicine and Dentistry, Queen Mary University of London, London, UK.

8 ² Integrated Mathematical Oncology, Moffitt Cancer Center, Tampa, FL, USA.

9 ³ Wellcome Centre for Human Genetics, University of Oxford, Oxford, UK.

10 ⁴ Evolutionary Genomics & Modelling Lab, Centre for Evolution and Cancer, Institute of Cancer Research, London,
11 UK.

12 ⁵ Department of Cell and Developmental Biology, University College London, London, UK.

13 * Correspondence: andrea.sottoriva@icr.ac.uk; t.graham@qmul.ac.uk

14

15 ABSTRACT

16 **Cancers accumulate mutations that lead to neoantigens, novel peptides that elicit an**
17 **immune response, and consequently undergo evolutionary selection. Here we establish**
18 **how negative selection shapes the clonality of neoantigens in a growing cancer, by**
19 **constructing a mathematical model of neoantigen evolution. The model predicts that,**
20 **without immune escape, tumor neoantigens are either clonal or at low frequency, and**
21 **hyper-mutated tumors can only establish following the evolution of immune escape.**
22 **Moreover, the site frequency spectrum of somatic variants under negative selection**
23 **appears more neutral as the strength of negative selection increases, consistent with**
24 **classical neutral theory. These predictions are corroborated by the analysis of**
25 **neoantigen frequencies and immune escape in exome and RNA sequencing data from**
26 **879 colon, stomach and endometrial cancers.**

27

28 INTRODUCTION

29 Mutations accrue throughout tumor development and provide ‘fuel for the fire’ of cancer
30 evolution. However, mutations can also hinder tumor evolution if they lead to an anti-tumor
31 immune response, via the generation of *neoantigens*, novel peptides presented on the cell’s
32 surface and recognized as ‘non-self’ by cells of the adaptive immune system^{1,2}. The immune
33 system is a major determinant of tumor evolution, most starkly demonstrated by the prognostic
34 value of immune-infiltration³ and the success of immunotherapy^{4,5}.

35 The landscape of neoantigenic mutations is shaped by ecological and evolutionary interactions
36 between a tumor and its microenvironment^{1,6,7}. In the absence of an immune system,
37 neoantigens accumulate as a ‘side-effect’ of mutation acquisition⁸, and are expected to follow
38 neutral evolutionary dynamics⁹. *Immuno-editing* refers to immune-cell killing of antigenic cells¹
39 and so represents a negative selective pressure⁶. Tumor cells can also experience positive
40 selection upon the evolution of mechanisms to inhibit the immune system’s ability to recognize
41 or react to cancer-associated antigens. These are termed *immune escape* mechanisms^{7,8,10}.
42 Cancer evolution in response to immune control is a ‘hallmark of cancer’¹¹ and it is well-
43 recognized that the tumor-specific immune microenvironment shapes the neoantigenic
44 repertoire found in tumors^{12–14}.

45
46 Therapies that (re)activate the immune response following escape have achieved exceptional
47 success (reviewed in ref¹⁵), especially in cancers of high mutational load^{16–18}. Neoantigen-
48 profiling is predictive of treatment response¹⁹ and long-term survival²⁰. However, a significant
49 number of patients do not respond to immunotherapy regardless of a high mutational load and
50 the presence of molecular markers of immune escape²¹, and there is a need to better predict
51 the likelihood of treatment response.

52

53 The evolutionary dynamics of tumor development can be partially decoded from the pattern of
54 intra-tumor genetic heterogeneity²². Positive and negative selection, respectively, cause the
55 expansion and contraction of subclones. Consequently, the site frequency spectrum of
56 mutations, as measured by variant allele frequencies (VAF)^{9,23} from genome sequencing data,
57 and cohort-wide mutation frequencies (e.g. dN/dS analysis) can be used to infer the
58 evolutionary dynamics that shaped the mutational landscape²⁴⁻²⁸. Population genetics has long
59 been concerned with the dynamics of negative selection in *constant population sizes*²⁹⁻³³, which
60 has been extended for expanding populations with *rare mutations*^{34,35}. However, cancer
61 evolution represents a distinct evolutionary regime because neoantigens are common, making
62 negative selection pervasive, immune escape can diminish selection; and tumors are growing
63 populations. Therefore, the dynamics resulting from negative selection acting on neoantigens in
64 a growing tumor remain to be determined.

65 Here we use stochastic modelling to study how the clonal structure and immunological
66 phenotype of growing tumors is shaped by negative selection in response to neoantigenic
67 mutations. We establish the dynamics expected under different selective environments and
68 tumor mutator phenotypes. We characterize the emerging VAF distribution under pervasive
69 negative selection, and determine the power to identify negative selection from genomics data.
70 We compare our modelling predictions with whole-exome sequencing and RNA sequencing
71 data from human cancers of the colon, stomach and endometrium.

72 RESULTS

73 Modelling predicts antigen-hot and antigen-cold tumors

74 We created a mathematical model of neoantigen evolution during tumor growth, based on a
75 stochastic branching process (Fig. 1a and Methods). At each step, tumor cells of lineage i
76 produced two surviving offspring at birth rate $b=1$ per unit time and offspring accumulated
77 mutations at rate μ , which had antigenicity drawn from a pre-specified distribution. Cells died
78 with death rate determined by the strength of negative selection, s , against the cumulative
79 antigenicity of neoantigens in the lineage. s can be interpreted as the effectiveness of immune
80 predation against an antigen: $s=0$ indicates no selection pressure (neutral evolution), and $s \ll 0$
81 strong negative selection (following ref³⁴). Tumor growth was simulated until the tumor reached
82 a predefined population size (approximating a clinically detectable size) or until a sufficiently
83 long time elapsed without the tumor reaching detectable size.

84

85 We first examined the temporal neoantigen burden in simulated tumors. We defined the ‘antigen
86 score’ of a tumor as the proportion of tumor cells carrying cumulative antigenicity $\geq T_c$. Tumors
87 simulated with identical parameters separated into two distinct groups due to the stochasticity of
88 neoantigen accrual: ‘antigen-hot’ and ‘antigen-cold’. Antigen-hot tumors had an antigen score
89 close to 1, corresponding to every tumor cell in the population being highly antigenic, whereas in
90 antigen-cold tumors the majority of cells lacked immunogenic mutations (Fig. 1b-c). The
91 proportion of antigen-hot tumors depended on the negative selection strength (Extended Data
92 Fig. 1a): increased negative selection for neoantigens decreased the probability of observing
93 antigen-hot tumors. In antigen-cold tumors, the proportion of neoantigen-carrying cells also
94 decreased inversely with the strength of negative selection.

95

96 In the simulations, the antigenicity of newly accrued neoantigens was sampled from a ‘prior’ pre-

97 specified distribution. Regardless of the shape of the prior distribution, surviving lineages always
98 showed enrichment for low-antigenicity alterations with an exponential-like distribution of final
99 antigenicity values (Fig. 1d and Extended Data Fig. 1b).

100 We next simulated hyper-mutated tumors that generated a high number of mutations per cell
101 division, causing lineages to rapidly accrue antigenicity. Consequently, most lineages rapidly
102 became neoantigen-hot and were eradicated by negative selection (Fig. 1e). In rare tumors that
103 survived to detectable size, high-frequency neoantigens were absent (Extended Data Fig. 2a-b
104 and Supplementary Note).

105 Overall, we observed that negative selection prevented subclonal neoantigens rising to high
106 frequency in a tumor, and this effect was exacerbated at higher mutation rates.

107 We compared the dynamics observed in our model to the dynamics of neoantigen accrual in a
108 constant population size (Supplementary Note). Models of negative selection with constant
109 population size²⁹⁻³² can lead to a broad range of evolutionary dynamics as the mutation rate and
110 strength of negative selection are varied. In contrast, here we observed that allowing the
111 population size to vary led to broadly consistent dynamics across the parameters space
112 (Extended Data Fig. 2). We considered three scenarios: (i) High s , low μ . When negative
113 selection was strong and mutations rare, selection operated efficiently in a constant population
114 rendering it devoid of neoantigenic mutations, but was attenuated in a variable-sized population
115 due to population expansion decreasing the efficiency of selection, as previously reported for
116 positive selection²³. (ii) Low s , high μ . Due to weak selection, only lineages with multiple
117 mutations experienced non-negligible selection. As in the previous case, population growth
118 attenuated the influence of selection relative to the constant-sized population model. (iii) High s ,
119 high μ . In constant size populations, the population could not go extinct, and dynamics were
120 determined by the relative strength of negative selection between lineages all accruing

121 neoantigenic mutations. The additive effect of any single mutation on fitness was proportionally
122 diminished as mutation burden increased due to a Muller's Ratchet-like effect³³, leading to
123 weakly selected dynamics. In a variable-sized population the dynamics were markedly different:
124 populations where all lineages were strongly negatively selected went extinct, and surviving
125 populations consisted of the 'lucky' lineages that had not accrued neoantigens (Extended Data
126 Fig. 2a,d). These extinction-driven dynamics persisted in the growing population even in the
127 special case of extremely high μ and low s , while the constant population became effectively
128 neutral.

129 **Immune escape leads to antigen-hot and antigen-warm tumors**

130 We next simulated *immune escape alterations* acquired by one cell that renders descendants
131 less susceptible to immune predation^{36,37}. Specifically, we set the death rate of immune escaped
132 cells to the baseline non-immunogenic death rate irrespective of the cell's burden of antigenic
133 mutations.

134 If the founder cell of the tumor contained an escape mutation (*clonal escape*), tumors with a
135 continuum of antigenicity scores emerged (Fig. 1f). We termed these tumors 'antigen-warm' as
136 they contained strong high-frequency and/or several subclonal neoantigens.

137 We then simulated tumors which could acquire immune escape at a random time (*probabilistic*
138 *escape*) and evaluated the detectable neoantigen load in the emerging tumors (Methods). When
139 the mutation rate was low, tumors that reached detectable size had rarely evolved immune
140 escape, and the strength of negative selection imposed on growth was inversely correlated with
141 the subclonal neoantigen burden observed in the final tumor (Fig. 1g). When the mutation rate
142 was high, lineages rapidly accrued neoantigens and were driven to extinction by negative
143 selection (Fig. 1e). Tumors only grew to detectable frequency if the founder lineage

144 stochastically acquired immune escape to ‘rescue’ them. Consequently, at high mutation rates,
145 detectable tumors were exclusively immune escaped and had a high burden of high-frequency
146 neoantigens (Fig. 1h).

147 Taken together, these results suggest that there is a non-linear relationship between the levels
148 of immune surveillance in the microenvironment and the magnitude of immuno-editing seen in
149 tumors of detectable size. Moving from low to moderate negative selection, the dynamics
150 increasingly depart from strictly neutral dynamics as expected, and correspondingly the clonal
151 and subclonal neoantigen burden is progressively decreased. At strong negative selection,
152 detectable tumors are those that have stochastically accrued immune escape, and
153 consequently show a high proportion of neoantigen-warm and –hot cases and evolve effectively
154 neutrally. We also note that the mutation rate is a determinant of the strength of negative
155 selection experienced by a lineage: at high mutation rates a lineage is likely to accrue multiple
156 negatively selected variants and so experience stronger negative selection.

157 **Immune-infiltrated cancers are antigen-hot and escaped**

158 To compare model predictions to experimentally measured neoantigen landscapes, we
159 analyzed neoantigens in 363 colorectal, 146 stomach and 370 endometrial cancers (CRC,
160 STAD and UCEC, respectively) from The Cancer Genome Atlas (TCGA) (Fig. 2a). We focused
161 on these cancer types because of the prevalence of mutator phenotypes, namely cancers with:
162 polymerase- ϵ mutation (POLE – very high mutation rate), mismatch repair deficiency (MMR –
163 high mutation rate, often responding well to immunotherapy^{18,38}), and microsatellite stable
164 tumors (MSS – lower mutation rate). Therefore, they provide a good model to explore the effect
165 of different tumor-immune dynamics. TCGA samples filtered for high sequencing depth and
166 purity were first HLA-typed *in silico*³⁹, and their neoantigens called and filtered¹⁹ using the
167 NeoPredPipe pipeline⁴⁰ (see Methods). We also evaluated T-cell infiltration from paired RNA-

168 seq data⁴¹ as a measure analogous to negative selection strength s experienced by
169 neoantigens.

170 The vast majority of tumors (90%) had clonal neoantigens (Supplementary Table 1), and so
171 were defined as 'antigen-hot'. We observed that the mutation-antigenicity distribution of tumors
172 (see Methods) was enriched for low binding neoantigens irrespective of the level of T-cell
173 infiltrate, but still contained a tail of high-scoring neoantigens (Fig. 2b). Subclonal neoantigen
174 burden varied significantly between cancers: cancers with low or medium T-cell infiltration
175 (putative small or moderate s) had proportionally fewer subclonal neoantigens than high T-cell
176 infiltrate tumors (high s) (Fig. 2c), suggesting a critical role of immune escape in early evolution.
177 Interestingly, this trend was absent in STAD tumors, suggesting a more homogeneous evolution
178 due to either widespread or rare immune escape.

179 We therefore sought evidence of immune escape in the cancers: alterations in antigen
180 presentation and over-expression of immune checkpoint genes (Methods). Overall, 57% of all
181 cancers showed evidence of at least one escape mechanism, with increased prevalence of
182 escape in MMR (71%) and POLE (98%) cases and significantly different patterns of immune
183 escape (Fig. 2d and Extended Data Fig. 3a), in agreement with previous studies^{18,41,42}. STAD
184 cancers in particular had a high proportion of immune escaped cancers – potentially a result of
185 strong early immune predation. Further work is needed to confirm that these differences
186 between mutational subtypes arose from differential selective pressures on immune escape.

187 Consistent with the predictions and previous studies⁴³, tumors with immune escape had a
188 higher neoantigen burden, and the majority of highly antigenic tumors (neoantigen burden >100)
189 were immune-escaped (Fig. 2e). Increased immune infiltration level was strongly associated
190 with immune escape, even in non-hyper-mutated (MSS) samples (Fig. 2f). We expected
191 neoantigen-associated mutations to be most under-represented amongst high-cancer cell

192 fraction (CCF) subclonal mutations, as selection had the longest time to act on these mutations.
193 Therefore, we compared the number of neoantigens at high CCF (present in 30%-60% of cells)
194 between MMR cases with and without immune escape, and found greater depletion in non-
195 escaped cancers (Fig. 2g), consistent with immuno-editing shaping the clonal structure of hyper-
196 mutated tumors without immune escape. The above phenomena were also observed in a meta-
197 cohort that combined the three cancer types (Extended Data Fig. 3).

198 Together, these data suggest that these cancer types usually evolve in the face of stringent
199 immune-selective pressures (analogous to the moderate/high s regime in simulated tumors) and
200 consequently immune-escape is frequently selected for at the onset of tumor growth, permitting
201 the development of tumors with high and clonal neoantigen load.

202 **Subclonal immune escape shapes local neoantigen evolution**

203 Next, we explored the evidence for subclonal immune escape in a previously published multi-
204 region sequenced colorectal tumor dataset⁴⁴. Overall, loss of heterozygosity (LOH) at HLA loci,
205 called with the LOHHLA tool³⁷, was found in 5/10 (50%) carcinomas and 1/6 (17%) adenomas,
206 and some of these events were present subclonally, in spatially distinct region(s) of the tumor
207 (Fig. 3a-b).

208 Simulations of subclonal immune escape in our model predicted that subclones should become
209 proportionally enriched for neoantigens following escape (Fig. 3c), consistent with previous
210 observations³⁷. In our primary tumor data, a significantly higher proportion of detected
211 neoantigens were associated with the lost allele in escaped clones than in clones without LOH
212 (Fig. 3d). These results confirm that locally different immune-mediated negative selection
213 pressures shape individual subclones inside a tumor.

214 To study how subclonal immune escape mechanisms can influence the efficiency of therapy, we

215 extended our simulations to model immunotherapy. We introduced two different types of escape
216 stochastically during tumor growth, *active* and *passive*, that notionally represented reversible
217 escape mechanisms affecting interactions with the microenvironment (e.g. expression of PDL1)
218 and irreversible cell-intrinsic escape (e.g. genomic loss of an HLA allele) respectively (Methods).
219 After the tumor population grew up to detectable size, we simulated immunotherapy by
220 cancelling the effect of active immune escape, and also increasing the negative selection
221 pressure s against neoantigens. The clonal population(s) with active escape rapidly shrank, but
222 clones with passive-type escape continued growing (Fig. 3e). Neoantigens were progressively
223 pruned from the expanding clone, leading eventually to an immune-cold tumor. Thus, the
224 immune landscape of a tumor post-immunotherapy is predicted to be distinct from the original
225 tumor (consistent with observations^{45,46}), with potential implications for the choice of the next
226 line of therapy.

227 **Negative selection leads to neutral VAF distribution**

228 We sought to explore how negative selection shapes the distribution of subclonal mutation
229 frequencies within an individual cancer. We considered the VAF distribution in simulated tumors
230 with moderate and high negative selection. Evidence for positive selection in the VAF
231 distribution is provided by an over-abundance of passenger mutations at high-frequency that are
232 within the expanding clone²³, whereas under pervasive negative selection, antigenic clones are
233 continually depleted and so rarely grow to a large size (rarely reach high VAF). Thus, the vast
234 majority of higher-VAF mutations are neutral passengers, that evolve according to neutral
235 dynamics and so exhibit a characteristic $1/f^2$ dependence (leading to a $1/f$ dependence of the
236 cumulative VAF distribution, Fig. 4a)⁹. As negative selection strength increases, the
237 phenomenon is exacerbated: antigenic subclones are more rapidly depleted and so more
238 neutral-like VAF distributions are observed (Fig. 4b). We note that pervasive negative selection

239 was part of the original neutral theory⁴⁷, and our observations are consistent with the classical
240 theory.

241 The VAF distribution computed of solely neoantigens shows depletion relative to the neutral
242 expectation (red lines in Fig. 4a-b), consistent with population genetics theory of constant-sized
243 models^{29,34} (Extended Data Fig. 2c,f). The magnitude of deviation from the neutral curve
244 depends on the strength of negative selection, which means that, in theory, negative selection
245 could be detected from neoantigen-VAF distributions (Extended Data Fig. 4). However, in
246 practice, the few persisting neoantigens are at *very low* VAFs and so are problematic to
247 measure accurately⁴⁸, severely hindering the power to quantify negative selection strength
248 directly from neoantigen VAF distributions.

249

250 **Negative selection is elusive in VAF distribution**

251 We performed *in silico* sequencing on simulated tumors, and explored the effect of read depth
252 and false-positive neoantigen identification⁴⁹ on the identifiability of negative selection in
253 individual tumors (see Methods). The simulations predicted that very high depth sequencing
254 was required to robustly call negative selection from VAF distributions, and the efficacy strongly
255 depended on the strength of selection against neoantigens (Fig. 4c-d). Erroneously labelling
256 neoantigens also had substantial impact on the power, but could be mitigated by very high-
257 depth sequencing. Detection was mostly limited by the tumors retaining too few neoantigens to
258 reliably evaluate their VAF distribution, a phenomenon further exacerbated when concentrating
259 on strongly immunogenic mutations alone (Extended Data Fig. 5a-d).

260 In order to overcome the technical issues of limited sequence depth and low antigen numbers,
261 we pooled mutations from groups of identically simulated and comparable TCGA tumors

262 (Methods) and considered their combined VAF distribution (Fig. 4e), in a similar manner to how
263 cohort-wide positive selection by dN/dS analysis is evaluated^{24,25}. In the pooled TCGA cohort,
264 we investigated essential genes⁵⁰ that are expected to be constitutively expressed and under
265 selection^{25,51}. In cancers with medium T-cell score and no evidence of immune escape, there
266 was a depletion of all neoantigens and neoantigens in essential genes compared to the neutral
267 expectation (Fig. 4f). In contrast, there was no neoantigen depletion in cancers with low T-cell
268 score. Neoantigens in CRC and UCEC cancers individually, as well as frameshift and nonsense
269 mutations in essential genes, showed similar trends (Extended Data Fig. 5e-f), suggesting a
270 more stringent selection in moderately infiltrated tumors and on essential genes.

271 **Proportional burden can measure negative selection**

272 Depletion of neoantigens relative to the overall non-synonymous mutation is a well-established
273 signature of immuno-editing^{52,53}. We investigated the relationship between the degree of
274 neoantigen depletion and strength of negative selection experienced by neoantigens.

275 First, we simulated tumors with a known neoantigen production rate ($p_a=0.075$ per non-
276 synonymous mutation, Supplementary Note) to evaluate how the proportion of immunogenic to
277 non-synonymous mutations changed with negative selection strength. As expected, stronger
278 negative selection led to proportionally fewer observed neoantigens in the final tumor (Fig. 5a).
279 We also measured the effective mutation rate (the ratio of the per cell division mutation and
280 survival rate), derived from the linear slope of the neutral VAF curve⁹, as a function of increasing
281 negative selection for neoantigens (Supplementary Note). Stronger negative selection caused
282 higher effective mutation rates in antigenic tumors (Fig. 5b), as a consequence of increased
283 death rate. We suggest that the higher cell death rate inferred in hyper-mutated tumors⁵⁴ is
284 likely to be, at least in part, a direct consequence of immuno-editing.

285 Next, we examined the proportional neoantigen burden in TCGA cancers stratified by cancer
286 type and predicted immune escape status. We observed no difference in overall proportional
287 neoantigen burden according to cancer type (Extended Data Fig. 6a), and so combined all data
288 into a single meta-cohort. We detected no significant difference in overall proportional burden
289 between MSS and MMR, and immune escaped or non-escaped cancers (Extended Data Fig.
290 6b-c). The observed uniformity in overall proportional burden across the cohort is consistent
291 with the lack of neoantigen depletion signal reported in ref⁵². The majority of mutations
292 considered in these analyses were clonal, and so were likely accrued prior to tumor expansion
293 and acquisition of immune escape. To better delineate the decrease in negative selection
294 expected following immune escape, we computed subclonal proportional neoantigen burden for
295 mutations with CCF<0.6. Comparing total and subclonal proportional burden (considering all
296 tumors with >30 subclonal mutations) showed a lower subclonal proportional burden in non-
297 escaped cancers, but no shift was detected in cancers with immune escape (Fig. 5c), consistent
298 with stronger negative selection in non-escaped cancers. When cancer types were considered
299 independently, UCEC and CRC cancers showed a similar pattern, but no subclonal depletion
300 was evident in STAD cancers (Extended Data Fig. 6d).

301 To examine the potential confounding effect of different mutational processes, we generated
302 synthetic cohorts analogous to real tumors (Methods). Comparing the synthetic cohorts
303 matching the overall mutation composition of CRCs showed no significant difference in
304 proportional burden, suggesting that MMR-specific mutational processes (e.g. Signature 6 from
305 ref⁵⁵) are not strongly biased for neoantigen generation (Extended Data Fig. 6e). A synthetic
306 matched cohort of Fig. 5c confirmed that the observed difference in subclonal proportional
307 neoantigen burden was also independent of mutational processes (Extended Data Fig. 6f).
308 Burden normalized to this synthetic cohort showed a trend for lower than random subclonal
309 neoantigen burden (Extended Data Fig. 6g). These observations imply the presence of active

310 immune surveillance when escape has not occurred, and highlight the high inter-patient
311 variability in evolutionary dynamics.

312

313 DISCUSSION

314 Here we have investigated the evolutionary dynamics of neoantigens and immune escape in
315 growing tumors using a mathematical model of tumor evolution. Our analysis shows how
316 negative selection by the immune system (immuno-editing) sculpts the clonal architecture of the
317 tumor: the hallmark of negative selection is the lack of neoantigens at intermediate subclonal
318 frequency within a tumor, and conversely, the presence of numerous neoantigens at
319 intermediate frequency is a hallmark of immune escape. Moreover, strong negative selection
320 for neoantigens inevitably provides a *strong selective pressure* for the evolution of immune
321 escape. Consequently, the observation that many cancers are both (neo)antigenic and have
322 immune escape points to a critical role for immune escape in the genesis of malignancy. Further
323 work directly measuring the immune repertoire at the time invasion first occurs is required.

324

325 Our simulations show that under negative selection, the *overall VAF distribution* of a tumor will
326 be effectively-neutral, as it will be dominated by the neutral passenger mutations that are able to
327 spread through the tumor unimpeded by immune predation. In constant size models, neutral
328 mutations linked to disadvantageous alterations show a pattern of background selection^{30–33}, but
329 in growing populations selection can only be observed on the selected mutations directly. The
330 VAF distribution observable in cancer genome sequencing data becomes more neutral-like as
331 the strength of negative selection increases, as negatively selected clones are pushed to
332 harder-to-detect frequencies leaving only neutrally evolving lineages at high VAF. Furthermore,
333 our analysis suggests that the majority of tumors with high mutational burden – where in theory
334 VAF distributions and so evolutionary dynamics should be easier to resolve – are most likely to
335 be immune escaped and so only exhibit effectively neutral dynamics. Consequently, we suggest
336 that the lack of immune-related selection signal (e.g. as identified by ref⁵²) could be due to
337 unclassified immune escape or false-positive neoantigen calls that together mean the mutations

338 studied are likely to be overall only very weakly negatively selected. Pooling data across
339 cancers increases power to resolve clone size distributions and detect negative selection, and
340 could be combined with dN/dS methods to evaluate selection of gene sets, such as natural
341 HLA-binders^{52,56} and MHC-II presented peptides⁵⁷.

342 Our modelling offers insight into the challenges of predicting immunotherapy response using
343 tumour mutation burden (TMB) alone. Strong negative selection (effective immune surveillance)
344 leads to a high rate of cell death, a corresponding increase in the effective mutation rate of
345 tumors, and the net result of high TMB with severe neoantigen depletion. Thus, despite having
346 high TMB, such tumors would be unlikely to respond to immune checkpoint blockade.
347 Assessment of neoantigens should be more predictive: tumors with clonal or numerous
348 subclonal neoantigens are very likely to have evolved immune escape – particularly if the
349 patient's immune system is highly predatory – and to respond to therapies reactivating immune
350 predation. This is consistent with previous studies suggesting that clonal antigens predict
351 sensitivity to immune checkpoint blockade⁴³. We illustrate that immune therapies targeted
352 against a specific neoantigen or immune mechanism are vulnerable to intra-tumor
353 heterogeneity, as subclones in which this target is altered or lost (e.g. neoantigen depleted or
354 HLA haplotype mutated) will experience net positive selection when the therapy is applied⁵⁸⁻⁶⁰.
355 Relatedly, a subclone that escapes immune blockade therapy and reforms the tumor is
356 predicted to have a different immune landscape due to the action of immune predation during
357 clone emergence, with potential implications for additional lines of therapy.

358 In summary, our mathematical framework provides insights into the evolutionary dynamics of
359 negatively selected neoantigens in growing tumors and the detectability of these dynamics in
360 genomic data.

361 **ACKNOWLEDGEMENTS**

362 This work was supported by the Wellcome Trust (202778/B/16/Z to A.S.; 202778/Z/16/Z to T.A.G.;
363 105104/Z/14/Z to the Centre for Evolution and Cancer, Institute of Cancer Research; 108861/7/15/7 to
364 R.O.S.; 097319/Z/11/Z to C.P.B.) and Cancer Research UK (A22909 to A.S.; A19771 to T.A.G.
365 supporting E.L.). A.R.A.A. and C.G, and A.S and T.A.G. received support from the US National Institutes
366 of Health National Cancer Institute (grant no. U54CA143970) and (U54 CA217376) respectively. R.O.S.
367 was also supported by the Wellcome Centre for Human Genetics (grant no. 203141/7/16/7). B.W. was
368 supported by the Geoffrey W. Lewis Postdoctoral Training fellowship. L.Z. is supported by the European
369 Union's Horizon 2020 research and innovation programme under the Marie Skłodowska-Curie Research
370 Fellowship scheme (846614).

371 **AUTHOR CONTRIBUTIONS**

372 E.L., A.R.A.A., A.S. and T.A.G. conceptualized the study. A.R.A.A., A.S and T.A.G. acquired funding for
373 the project. E.L., A.S. and T.A.G. led the investigation, analysed data, and wrote the original manuscript.
374 E.L., M.J.W., W.C.H.C., B.W., R.O.S., C.G., J.H., L.Z., M.R.T., and C.P.B. contributed to the
375 mathematical model, computational framework and bioinformatics analysis. All authors reviewed and
376 approved the final manuscript.

377 **COMPETING INTERESTS**

378 The authors declare no competing interests.

379

380

381 **REFERENCES**

382

- 383 1. Schumacher, T. N. & Schreiber, R. D. Neoantigens in cancer immunotherapy. *Science* (80-.).
384 **348**, 69 (2015).
- 385 2. Lu, Y.-C. & Robbins, P. F. Cancer immunotherapy targeting neoantigens. *Semin Immunol* **28**, 22–
386 27 (2016).
- 387 3. Galon, J. *et al.* Towards the introduction of the 'Immunoscore' in the classification of malignant
388 tumours. *J. Pathol.* **232**, 199–209 (2014).
- 389 4. Sharma, P. & Allison, J. P. The future of immune checkpoint therapy. *Science* (80-.). **348**, 56–61
390 (2015).
- 391 5. Larkin, J. *et al.* Combined Nivolumab and Ipilimumab or Monotherapy in Untreated Melanoma. *N.*
392 *Engl. J. Med.* **373**, 23–34 (2015).
- 393 6. Milo, I. *et al.* The immune system profoundly restricts intratumor genetic heterogeneity. *Sci.*
394 *Immunol.* **3**, eaat1435 (2018).
- 395 7. Dunn, G. P., Bruce, A. T., Ikeda, H., Old, L. J. & Schreiber, R. D. Cancer immunoediting: from

- 396 immunosurveillance to tumor escape. *Nat. Immunol.* **3**, 991–998 (2002).
- 397 8. DuPage, M., Mazumdar, C., Schmidt, L. M., Cheung, A. F. & Jacks, T. Expression of tumour-
- 398 specific antigens underlies cancer immunoediting. *Nature* **482**, 405–409 (2012).
- 399 9. Williams, M. J., Werner, B., Barnes, C. P., Graham, T. A. & Sottoriva, A. Identification of neutral
- 400 tumor evolution across cancer types. *Nat. Genet.* **48**, 238–244 (2016).
- 401 10. Koebel, C. M. *et al.* Adaptive immunity maintains occult cancer in an equilibrium state. *Nature* **450**,
- 402 903–907 (2007).
- 403 11. Hanahan, D. & Weinberg, R. A. Hallmarks of Cancer: The Next Generation. *Cell* **144**, 646–674
- 404 (2011).
- 405 12. Koebel, C. M. *et al.* Adaptive immunity maintains occult cancer in an equilibrium state. *Nature* **450**,
- 406 903 EP- (2007).
- 407 13. Marty, R. *et al.* MHC-I Genotype Restricts the Oncogenic Mutational Landscape. *Cell* **171**, 1272-
- 408 1283.e15 (2017).
- 409 14. Rosenthal, R. *et al.* Neoantigen-directed immune escape in lung cancer evolution. *Nature* **567**,
- 410 479–485 (2019).
- 411 15. Yarchoan, M., Johnson III, B. A., Lutz, E. R., Laheru, D. A. & Jaffee, E. M. Targeting neoantigens
- 412 to augment antitumour immunity. *Nat. Rev. Cancer* **17**, 209–222 (2017).
- 413 16. Rizvi, N. A. *et al.* Mutational landscape determines sensitivity to PD-1 blockade in non-small cell
- 414 lung cancer. *Science (80-.)*. **348**, 124–128 (2015).
- 415 17. Lennerz, V. *et al.* The response of autologous T cells to a human melanoma is dominated by
- 416 mutated neoantigens. *Proc Natl Acad Sci U S A* **102**, 16013–16018 (2005).
- 417 18. Le, D. T. *et al.* Mismatch repair deficiency predicts response of solid tumors to PD-1 blockade.
- 418 *Science (80-.)*. **357**, 409–413 (2017).
- 419 19. Łuksza, M. *et al.* A neoantigen fitness model predicts tumour response to checkpoint blockade
- 420 immunotherapy. *Nature* **551**, 517–520 (2017).
- 421 20. Balachandran, V. P. *et al.* Identification of unique neoantigen qualities in long-term survivors of
- 422 pancreatic cancer. *Nature* **551**, 512 (2017).
- 423 21. Gibney, G. T., Weiner, L. M. & Atkins, M. B. Predictive biomarkers for checkpoint inhibitor-based
- 424 immunotherapy. *Lancet. Oncol.* **17**, e542–e551 (2016).
- 425 22. Turajlic, S., Sottoriva, A., Graham, T. & Swanton, C. Resolving genetic heterogeneity in cancer.
- 426 *Nat. Rev. Genet.* **20**, 404–416 (2019).
- 427 23. Williams, M. J. *et al.* Quantification of subclonal selection in cancer from bulk sequencing data.
- 428 *Nat. Genet.* **50**, 895–903 (2018).
- 429 24. Ostrow, S. L., Barshir, R., DeGregori, J., Yeger-Lotem, E. & Hershberg, R. Cancer Evolution Is
- 430 Associated with Pervasive Positive Selection on Globally Expressed Genes. *PLOS Genet.* **10**,
- 431 e1004239 (2014).
- 432 25. Martincorena, I. *et al.* Universal Patterns of Selection in Cancer and Somatic Tissues. *Cell* **171**,
- 433 1029-1041.e21 (2017).
- 434 26. Temko, D., Tomlinson, I. P. M., Severini, S., Schuster-Böckler, B. & Graham, T. A. The effects of
- 435 mutational processes and selection on driver mutations across cancer types. *Nat. Commun.* **9**,
- 436 1857 (2018).
- 437 27. Cannataro, V. L., Gaffney, S. G. & Townsend, J. P. Effect Sizes of Somatic Mutations in Cancer.
- 438 *JNCI J. Natl. Cancer Inst.* **110**, 1171–1177 (2018).
- 439 28. Williams, M. J. *et al.* Measuring the distribution of fitness effects in somatic evolution by combining
- 440 clonal dynamics with dN/dS ratios. *Elife* **9**, e48714 (2020).
- 441 29. Cvijović, I., Good, B. H. & Desai, M. M. The Effect of Strong Purifying Selection on Genetic
- 442 Diversity. *Genetics* **209**, 1235 (2018).
- 443 30. Good, B. H., Walczak, A. M., Neher, R. A. & Desai, M. M. Genetic Diversity in the Interference
- 444 Selection Limit. *PLOS Genet.* **10**, e1004222 (2014).
- 445 31. Neher, R. A. & Hallatschek, O. Genealogies of rapidly adapting populations. *Proc. Natl. Acad. Sci.*
- 446 *U. S. A.* **110**, 437–442 (2013).
- 447 32. Charlesworth, B., Morgan, M. T. & Charlesworth, D. The effect of deleterious mutations on neutral
- 448 molecular variation. *Genetics* **134**, 1289–1303 (1993).
- 449 33. Haigh, J. The accumulation of deleterious genes in a population—Muller’s Ratchet. *Theor. Popul.*
- 450 *Biol.* **14**, 251–267 (1978).
- 451 34. Kessler, D. A. & Levine, H. Scaling solution in the large population limit of the general asymmetric

- 452 stochastic Luria-Delbrück evolution process. *J. Stat. Phys.* **158**, 783–805 (2015).
- 453 35. Antal, T. & Krapivsky, P. L. Exact solution of a two-type branching process: models of tumor
454 progression. *J. Stat. Mech. Theory Exp.* **2011**, P08018 (2011).
- 455 36. Vinay, D. S. *et al.* Immune evasion in cancer: Mechanistic basis and therapeutic strategies. *Semin.*
456 *Cancer Biol.* **35**, S185–S198 (2015).
- 457 37. McGranahan, N. *et al.* Allele-Specific HLA Loss and Immune Escape in Lung Cancer Evolution.
458 *Cell* **171**, 1259–1271.e11 (2017).
- 459 38. Kather, J. N., Halama, N. & Jaeger, D. Genomics and emerging biomarkers for immunotherapy of
460 colorectal cancer. *Semin. Cancer Biol.* **52**, 189–197 (2018).
- 461 39. Shukla, S. A. *et al.* Comprehensive analysis of cancer-associated somatic mutations in class I
462 HLA genes. *Nat. Biotechnol.* **33**, 1152–1158 (2015).
- 463 40. Schenck, R. O., Lakatos, E., Gatenbee, C., Graham, T. A. & Anderson, A. R. A. NeoPredPipe:
464 high-throughput neoantigen prediction and recognition potential pipeline. *BMC Bioinformatics* **20**,
465 264 (2019).
- 466 41. Grasso, C. S. *et al.* Genetic mechanisms of immune evasion in colorectal cancer. *Cancer Discov.*
467 **8**, 730–749 (2018).
- 468 42. Xie, T. *et al.* A Comprehensive Characterization of Genome-Wide Copy Number Aberrations in
469 Colorectal Cancer Reveals Novel Oncogenes and Patterns of Alterations. *PLoS One* **7**, e42001
470 (2012).
- 471 43. McGranahan, N. *et al.* Clonal neoantigens elicit T cell immunoreactivity and sensitivity to immune
472 checkpoint blockade. *Science (80-.)*. **351**, 1463–1469 (2016).
- 473 44. Cross, W. *et al.* The evolutionary landscape of colorectal tumorigenesis. *Nat. Ecol. Evol.* **2**, 1661–
474 1672 (2018).
- 475 45. Riaz, N. *et al.* Tumor and Microenvironment Evolution during Immunotherapy with Nivolumab. *Cell*
476 **171**, 934–949.e16 (2017).
- 477 46. Anagnostou, V. *et al.* Evolution of Neoantigen Landscape during Immune Checkpoint Blockade in
478 Non-Small Cell Lung Cancer. *Cancer Discov.* **7**, 264–276 (2017).
- 479 47. Kimura, M. *The Neutral Theory of Molecular Evolution*. (Cambridge University Press, 1983).
480 doi:DOI: 10.1017/CBO9780511623486
- 481 48. Stead, L. F., Sutton, K. M., Taylor, G. R., Quirke, P. & Rabbitts, P. Accurately Identifying Low-
482 Allelic Fraction Variants in Single Samples with Next-Generation Sequencing: Applications in
483 Tumor Subclone Resolution. *Hum. Mutat.* **34**, 1432–1438 (2013).
- 484 49. Yadav, M. *et al.* Predicting immunogenic tumour mutations by combining mass spectrometry and
485 exome sequencing. *Nature* **515**, 572–576 (2014).
- 486 50. Blomen, V. A. *et al.* Gene essentiality and synthetic lethality in haploid human cells. *Science (80-.*
487 *)*. **350**, 1092–1096 (2015).
- 488 51. Van den Eynden, J., Basu, S. & Larsson, E. Somatic Mutation Patterns in Hemizygous Genomic
489 Regions Unveil Purifying Selection during Tumor Evolution. *PLOS Genet.* **12**, e1006506 (2016).
- 490 52. Van den Eynden, J., Jiménez-Sánchez, A., Miller, M. L. & Larsson, E. Lack of detectable
491 neoantigen depletion signals in the untreated cancer genome. *Nat. Genet.* **51**, 1741–1748 (2019).
- 492 53. Rooney, M. S., Shukla, S. A., Wu, C. J., Getz, G. & Hacohen, N. Molecular and genetic properties
493 of tumors associated with local immune cytolytic activity. *Cell* **160**, 48–61 (2015).
- 494 54. Werner, B. *et al.* Measuring single cell divisions in human tissues from multi-region sequencing
495 data. *Nat. Commun.* **11**, 1035 (2020).
- 496 55. Alexandrov, L. B. *et al.* Signatures of mutational processes in human cancer. *Nature* **500**, 415–421
497 (2013).
- 498 56. Zapata, L. *et al.* Negative selection in tumor genome evolution acts on essential cellular functions
499 and the immunopeptidome. *Genome Biol.* **19**, 67 (2018).
- 500 57. Marty Pyke, R. *et al.* Evolutionary Pressure against MHC Class II Binding Cancer Mutations. *Cell*
501 **175**, 416–428.e13 (2018).
- 502 58. Kim, J. M. & Chen, D. S. Immune escape to PD-L1/PD-1 blockade: seven steps to success (or
503 failure). *Ann. Oncol.* **27**, 1492–1504 (2016).
- 504 59. Sharma, P., Hu-Lieskovan, S., Wargo, J. A. & Ribas, A. Primary, Adaptive, and Acquired
505 Resistance to Cancer Immunotherapy. *Cell* **168**, 707–723 (2017).
- 506 60. Iorgulescu, J. B., Braun, D., Oliveira, G., Keskin, D. B. & Wu, C. J. Acquired mechanisms of
507 immune escape in cancer following immunotherapy. *Genome Med.* **10**, 87 (2018).

508
509 **FIGURE LEGENDS**
510
511 **Figure 1: Tumor growth model predicts two distinct types of immune phenotypes and the necessity**
512 **of immune escape. (a)** Schematic representation of the model. Left panel: tumor growth for four
513 generations. Filled circles represent cells, colored by immunogenicity. Related cells are connected with lines.
514 Middle panel: cell division/mutation process. Right panel: prior distribution of newly generated
515 neoantigenicities. For details, see Methods. **(b)** Growth curve of six simulated tumors at $s=-0.8$. Line color
516 shows the antigen score of the tumor population over time. **(c)** Cancer cell fraction (CCF) of the most
517 common antigenic mutation of $n=100$ tumors at the final time-point. **(d)** Distribution of antigenicity values of
518 all neoantigens generated (grey) and only neoantigens present in >10 surviving cells (blue). Thin lines:
519 individual simulations; thick dashed line: ensemble mean. Inset: Mann-Whitney two-sided test. **(e)**
520 Distribution of maximum tumor size reached by hyper-mutated tumors at $s=-0.8$. Inset: growth curve of a
521 single tumor colored by antigenic score as in (b), blue line: number of non-immunogenic cells. **(f)** Neoantigen
522 scores in $n=100$ tumors at $s=-0.8$, without (left) and with (right) clonal immune escape. **(g-h)** Number of
523 detectable neoantigens (read depth $\sim 50x$) in $n=50$ simulated tumors as a function of negative selection
524 strength. Middle panel: mean clonal neoantigen burden. Bottom panel: clonality of immune escape. Only
525 non-hyper-mutated (g) and hyper-mutated (h) tumors that reached a detectable size are shown. Violin widths
526 represent raw data density.

527

528 **Figure 2: Colorectal, stomach and endometrial tumors from TCGA are antigen-hot and enriched for**
529 **immune escape. (a)** Cancer type and mutator subtype of the TCGA cancers analyzed. The size and shade of
530 each circle represent the number of tumors (also shown) in that sub-category. **(b)** Distribution of normalized
531 binding strength of neoantigens in TCGA cancers with low, medium and high immune infiltration. The thick line
532 shows the mean density of all distributions from tumors in each category, the shaded regions represent ± 1
533 standard deviation around this mean. **(c)** Distribution of the number of subclonally detected (in $<60\%$ of the
534 tumor) neoantigen-associated mutations in cancers according to immune infiltration (T-cell average) score. Two-
535 sided Mann-Whitney tests are reported on each plot. **(d)** Prevalence of immune escape in MSS, MMR and
536 POLE samples. Two-sided chi-squared test is indicated on top of each panel. **(e)** Distribution of the number of
537 subclonal antigenic mutations in cancers with and without immune escape (magenta and grey, respectively)
538 Two-sided Mann-Whitney test is reported on each panel. **(f)** Prevalence of immune escape in MSS cancers
539 according to their immune infiltration level. Two-sided chi-squared test is indicated on top of each panel. **(g)**
540 Number of antigenic mutations present in large subclones ($>30\%$ and $<60\%$ of cells) in MMR samples with and
541 without immune escape. One-sided Mann-Whitney test is reported above each plot. Violin widths in (c), (e) & (g)
542 represent raw data density with binned individual data points overlaid on top.

543

544 **Figure 3: Subclonal immune escape shapes neoantigen landscape and tumor growth after therapy. (a)**
545 Immune escape through loss of heterozygosity (LOH) at an HLA locus in the multi-region sequenced colorectal
546 cohort. LOH events are divided up according to whether the alteration is detected in all (clonal) or not all
547 (subclonal) biopsies. **(b)** HLA LOH in individual biopsies in tumors with at least one subclonal or clonal loss
548 event. Unfilled boxes represent homozygous HLA alleles. **(c)** The number of antigenic mutations detected in two
549 distinct (with and without immune escape) subclones of $n=25$ simulated tumor. Antigenic mutations are detected
550 at simulated read depth of 100x. Visual elements of the boxplot correspond to the following summary statistics:
551 centre line, median; box limits, upper and lower quartiles; whiskers, 1.5x inter-quartile range. **(d)** The proportion
552 of all neoantigens binding to the HLA allele lost in the LOH event in the colorectal tumors that show subclonal
553 HLA LOH ($n=6$). One-sided Wilcoxon signed-rank tests are reported on (c) and (d). **(e)** Growth curve of
554 simulated tumors following anti-PD-L1-type immunotherapy. The tumors have previously developed active
555 immune escape, but also harbor a small subclone with different escape mechanism. Black dashed lines show
556 the number of cells in this subclone over time. The inset shows growth around the point when the subclone
557 takes over, on a logarithmic scale.

558

559 **Figure 4: Negative selection leads to characteristic depletion of neoantigens and effectively-neutral**
560 **overall VAF distributions. (a-b)** Cumulative number of mutations as a function of the inverse of the frequency
561 for all mutations (grey, left axis) and neoantigen-associated mutations (red, right axis) harbored in at least 30
562 cells in (a) a tumor with $s=-0.8$; (b) a tumor with $s=-1.2$. **(c)** Power to detect negative selection from the VAF
563 distribution as a function of sequencing read depth (x axis) and false neoantigen rate (y axis). Power is the
564 proportion of 100 simulated tumors with significant difference (two-sided Kolmogorov-Smirnov test, $\alpha=0.1$)
565 between the distribution of all mutations and neoantigen-associated mutations. **(d)** Power (in $n=100$ tumors) to
566 identify negative selection as a function of selection strength (x axis) and the stringency of the two-sided
567 Kolmogorov-Smirnov test used for detection ($\alpha=0.1$, $\alpha=0.05$, and $\alpha=0.01$, shown in black, maroon and red,
568 respectively). **(e)** Cumulative VAF distribution as a function of the inverse of the frequency for all (in grey) and
569 neoantigen-associated mutations (in red) detected with a sequencing depth of 800x in antigen cold tumors from
570 a simulated set of $n=100$. The y axis shows proportion of mutations. The mutation-antigenicity threshold 0.2 is
571 used in all cases in (a)-(e). **(f)** Cumulative VAF distribution of mutations detected in any low- and medium-
572 immune infiltrated TCGA MSS cancers without immune escape. The distribution is shown for all mutations
573 (grey), exonic mutations (blue), exonic mutations in essential genes (purple), antigenic mutations (pink) and
574 neoantigen-associated mutations in essential genes (red).

575

576

577 **Figure 5: Proportional neoantigen burden as a measure of selection.** (a) The proportion of neoantigen-
578 associated mutations (the percentage of all mutations) as a function of negative selection pressure, computed
579 from n=100 tumors each, with a simulated read depth of 200x. The expected value of antigens per mutation is
580 indicated with a horizontal dashed line. The mutation-antigenicity threshold of 0.2 is used. (b) Effective mutation
581 rate (per cell division mutation rate divided by per cell division death rate) computed from the VAF distribution of
582 mutations in antigen-hot tumors as a function of negative selection pressure. Read depth = 200x. Colors in (a) &
583 (b) indicate selection strength also shown on x axis (c) Proportional neoantigen burden of escaped and non-
584 escaped TCGA samples, computed from all mutations (red) and only subclonal mutations (CCF<0.6, colored
585 salmon). Lines connect total and subclonal proportional burdens measured in the same sample. Paired two-
586 sided Wilcoxon test is reported above the violin plots. Violin widths represent raw data density with individual
587 data points in (c) also indicated by end-points of connecting lines.

588

589

590

591 METHODS

592 *Mathematical model of tumor growth and mutation accumulation*

593 We created a minimal stochastic branching process model to represent tumor growth and
594 accumulation of mutations under selection pressure from the environment⁶¹. The model
595 described the proliferation, death and mutation accumulation of tumor cells, and environmental
596 factors (e.g. the level of T-cell infiltration) were described implicitly through parameters that
597 quantified the strength of selection against tumor cells.

598 We made use of a rejection-kinetic Monte Carlo algorithm⁶² to permit efficient simulation of large
599 populations of cells. Tumor evolution was initiated by a single transformed cell that produced
600 two surviving offspring at birth rate b per unit time. Cells in clone i died at rate d_i per unit time,
601 where the death rate increased with the neoantigen burden of the clone. Each time a cell
602 divided, it acquired new unique mutations at overall rate μ (Poisson distribution), which were
603 assigned as neoantigens at rate p_a , or as passengers (evolutionary neutral) at rate $1-p_a$. Each
604 antigenic mutation was assigned an antigenicity value (denoted A_j for the j^{th} antigen in a given
605 cell) sampled from an exponential distribution with the rate parameter set to 5 to produce a
606 skewed distribution wherein >99% of antigenicity values fall between 0 and 1, and most
607 neoantigens are only negligibly immunogenic (Fig. 1a). Neoantigens caused the death rate d_i of
608 the lineage to increase from a basal rate of $d_b=0.1$ to a higher value determined by the strength
609 of negative selection against each new neoantigen, controlled by the parameter s . The overall
610 effect on the birth/death rate of cells was determined by the cumulative antigenicity of
611 neoantigens harbored in the lineage, $\sum A_j$. The death rate of a subclone was computed as:

$$612 \quad d_i = (1 + s * \sum_{j=1}^{n_i} A_j^i)(d_b - 1) + 1. \quad (1)$$

613 And we defined the selective (dis)advantage of a subclone by its effective proliferation rate (the
614 difference of its birth and death rate), as compared to a non-immunogenic clone:

615
$$1 + s * \sum_{j=1}^{n_i} A_j^i = fitness = \frac{b-d_i}{b-d_b}, \quad (2)$$

616 where A_j^i denotes the j th neoantigen in lineage i ; $s=0$ stands for neutral evolution with no
617 neoantigen-associated selection and negative selection is represented by $s<0$.

618

619 This antigenicity-dependent increase in the clone death rate represented an aggregate of the
620 many stochastic factors that lead to the negative selection of neoantigens, including; (i)
621 sufficient presentation of neoantigens on the cell surface; (ii) recognition of neoantigens by T-
622 cell; (iii) antigen-mediated recruitment of further T-cells; and (iv) T-cell killing efficiency. We
623 chose to integrate all variability into a single probabilistic rate to be able to observe general
624 qualities of the tumor-immune interaction without the need for precise parametrisation. For
625 details on the steps of *in silico* simulations, see Supplementary Note and code at
626 <https://zenodo.org/record/3601322#.XvKCGJJKi4>.

627 We also modelled the acquisition of immune escape during tumor growth. Known immune
628 escape mechanisms include mutations affecting the antigen presenting machinery and
629 expression of immune checkpoint molecules^{36,37}. Immune escape was modelled as a heritable
630 property of a cell (representing e.g. copy number alteration of the PD-L1 or HLA gene). Immune
631 escape occurred as a result of a mutation with probability p_e per nonsynonymous mutation; or
632 through manual introduction of the escape alteration at a pre-determined clone size to achieve
633 clonal or subclonal immune escape. We considered two different types of escape mechanism:
634 (i) *active escape*, which shields the clone from negative selection (decreasing the clone death
635 probability to d_b) but does not decrease the neoantigen burden of the cell (corresponding to
636 escape mechanisms such as PD-L1 overexpression); and (ii) *passive escape*, which renders a
637 portion of neoantigenic mutations neutral (by rendering their antigenicity, A_j to 0; representing,
638 for example, loss of a HLA allele that predicts a subset of neoantigenic peptides being
639 presented).

640 We also incorporated therapeutic intervention in our model by time-dependently changing model
641 parameters. The most commonly used agents in immunotherapy target and inhibit immune
642 checkpoint pathways, helping the immune system to overcome immune escape achieved by
643 checkpoint over-expression and re-activate immune predation of neoantigenic cancer cells. We
644 simulated this effect by rendering active type immune escape ineffective (death rate of escaped
645 cells is increased by antigenic load) and simultaneously increasing the negative selection
646 strength s experienced by each neoantigen.

647 We chose model parameters to represent a wide range of possible tumor-immune
648 environments, and correspond to phenotypic properties of real cancers (Extended Data Fig. 6).
649 The following parameters were used in all simulations: $b = 1$; $d_b = 0.1$; $\mu = 1$ (not hyper-mutated)
650 and $\mu = 10$ (hyper-mutated); $-2 \leq s \leq 0$ (as indicated in figures or in caption); $p_a = 0.075$ and $p_e =$
651 10^{-6} (where applicable). For analyses where cells and mutations were classified as antigenic or
652 not, the cell- and mutation-antigenicity thresholds $T_c = 0.5$ and $T_m = 0.2$ were used, unless
653 stated otherwise. For further discussion on the simplifications applied in the model, and the
654 choices of simulation parameters and how they influence results, see the Supplementary Note
655 and Extended Data Figs. 7-9.

656 ***Simulation of VAF/CCF distributions and power calculation***

657 To evaluate the mutation spectrum of simulated tumors, mutations harbored in at least 10 cells
658 out of 10^5 (0.01%) were collected at the end of each simulation and the number of carrier cells
659 reported. Real sequencing data naturally introduces uncertainty about mutated allele frequency
660 due to limited sequencing depth and several sources of sampling bias²². To account for
661 imperfect measurements, CCF values were either computed by taking the raw frequency values
662 or via a simulated sequencing step introducing noise to these frequencies with indicated read
663 depth. For a given read depth, D , each frequency value, f , was substituted by a new frequency

664 sampled from a binomial distribution with parameters D and f : $\bar{f} \sim \text{Binom}(D, f)/D$. We filtered for
665 mutations with \bar{f} above 0, to discard mutations that are not picked up due to limited detection
666 power.

667 In addition to sequencing limitations, neoantigen identification from DNA sequencing alone has
668 a high rate of false-positive calls⁴⁹, and therefore the VAF distribution of neoantigens is
669 expected to be ‘contaminated’ with a large proportion of neutrally-evolving passenger mutations.
670 To simulate this effect when evaluating the power of detecting selection, we randomly sampled
671 non-antigenic mutations of simulated tumors (varied between 5% to 500% of the number of true
672 neoantigens, Fig. 4c) that were falsely labelled as neoantigens and included in the neoantigen-
673 based VAF distribution.

674 We computed the power to detect selection by comparing the distribution of all detected
675 mutations to that of the neoantigen-labelled subset using a two-sample Kolmogorov-Smirnov
676 test, and identified any samples as under selection in which the p-value of the test was below
677 0.1 (Fig. 4c) or a pre-defined value (Fig. 4d).

678 ***TCGA sample acquisition and processing***

679 All samples from the TCGA COAD and READ (merged together as CRC), STAD and UCEC
680 domains were retrieved through the NCI Genomics Data Commons (GDC) portal⁶³ between
681 15/06/2018 and 13/11/2019. Only patients with matched germline (from blood samples) and
682 primary tumor information available were considered. For each sample, purity (fraction of tumor
683 cells in the sample) and overall ploidy were evaluated using ASCAT⁶⁴ on Affymetrix SNP array
684 data. Samples with purity below 0.4 and ploidy above 3.6 were excluded from the analysis,
685 leaving 363 CRC, 146 STAD and 370 UCEC samples for which HLA typing and neoantigen
686 calls were performed (Supplementary Table 1 and Fig. 2a).

687 For analyzing immune escape, the cohort was narrowed down to patients for whom gene
688 expression data was available in GDC; and at least one pair of their HLA A/B/C alleles were
689 heterozygous and distinct enough to allow for loss of heterozygosity calls (n(CRC) = 341,
690 n(STAD)=118, n(UCEC)=362).

691
692 For each patient considered, the following information was downloaded: blood derived normal
693 bam files; primary tumor bam files; unfiltered variant call (vcf) files processed with Mutect2; SNP
694 array files; gene expression HTSeq counts (where available); and clinical information. We used
695 the unfiltered controlled-access variant call format (vcf) files to avoid over-filtering and missing
696 antigenic variants. The variants were filtered to only include variants that passed all filters of the
697 vcf files and not present (allelic depth of 0 or 1 for bases covered with over 30 reads) in normal
698 samples.

699
700 Samples were divided into MSS, MMR and POLE subtypes using data integrated from (i)
701 clinical TCGA annotation⁶⁵; (ii) calls retrieved from ref⁶⁶ that used the computational tool
702 MANTIS to analyze repetitions in tumor-normal sample pairs over microsatellite loci; (iii) and
703 mutational signature activities computed using non-negative least squares regression^{26,55}.
704 Samples with a MANTIS score ≥ 0.5 and TCGA annotation of 'MSI-H' ('microsatellite instability',
705 where available) were considered MMR, and those with MANTIS < 0.5 and 'MSI-L'/'MSS' were
706 labelled MSS. In case the two sources of information contradicted each other, neither of the
707 categories was assigned. Samples with at least 1,000 mutations inferred to originate from the
708 characteristic POLE signature (signature 10 in ref⁵⁵) were labelled as POLE tumors regardless
709 of their MMR status.

710 ***Multi-region sequenced dataset processing***

711 The multi-region sequenced colorectal dataset was accessed from Cross *et al.*⁴⁴ (raw data

712 available from the European Genome-Phenome Archive (<https://ega-archive.org/>) at accession
713 code: EGAS00001003066). Bam files with marked duplicates were used for HLA calling and
714 HLA variant detection. As in the original work, variants were called using Platypus⁶⁷, annotated
715 by ANNOVAR⁶⁸, and filtered to only contain somatic single nucleotide variations that were
716 present in at least 1 tumor sample and in either 0 reads in the normal sample (for normal
717 coverage ≤ 30 reads) or in at most 1 read (for normal coverage above 30 reads).

718 ***HLA haplotyping and calling immune escape***

719 HLA-A, -B and -C haplotyping was performed on blood derived normal bam files using
720 POLYSOLVER³⁹. As POLYSOLVER takes into account the individual's race to compute the
721 likelihood of each allele haplotype, we supplied ethnicity data, where available from clinical
722 TCGA information, and ran haplotyping with race 'Unknown' otherwise.

723 Using exome and RNAseq data, we tested for the presence of three types of immune escape
724 mechanisms: (i) somatic mutations in either one of the HLA alleles or in the B2M gene^{39,41}; (ii)
725 loss of an HLA haplotype through loss of heterozygosity (LOH) in the corresponding genomic
726 locus³⁷; and (iii) PD-L1 or CTLA-4 over-expression⁶⁹.

727 Mutations in HLA alleles were called using the previously called HLA haplotypes and the
728 corresponding functionality of POLYSOLVER³⁹. Variant calling was run using default settings
729 and HLA was considered mutated if at least one allele had a nonsynonymous somatic mutation
730 located in an exon or at a splice-site. Mutations in B2M were called if the sample contained a
731 nonsynonymous somatic mutation located inside one of the exons of the B2M gene, as
732 annotated by ANNOVAR⁶⁸ and confirmed using Variant Effect Predictor⁷⁰. Loss of
733 heterozygosity at the HLA locus was assessed using the software LOHHLA³⁷, using blood
734 derived normal, and tumor bam files were used. Tumor purity and ploidy estimates were derived
735 from ASCAT (for TCGA data) and from Sequenza⁷¹ (for the multi-region sequenced colorectal

736 tumors). A sample was considered to have Allelic Imbalance at an HLA locus if the
737 corresponding p-value was below 0.01 and LOH if, in addition, the copy number prediction of
738 that allele was below 0.5, with the confidence interval strictly below 0.7. Immune checkpoint
739 over-expression was assessed using RNA-seq data. Normal expression values (in transcripts
740 per million (TPM)) of PD-L1 and CTLA-4 were established for each cohort from TCGA based on
741 RNA-seq counts of the two proteins in 'solid tissue normal' samples. Checkpoint over-
742 expression was called if either PD-L1 or CTLA-4 expression in the tumor was higher than the
743 mean plus two standard deviations of normal expression. Immune checkpoint over-expression
744 could not be inferred for the multi-region sequenced dataset as only genomic data were
745 available.

746 We note that the extent of the impact of these escape alterations is not always known –
747 especially for mutations altering antigen presenting proteins – but we argued that nonetheless
748 they represent a level of impairment in the tumor-immune interaction.

749 Immune infiltration levels were computed from RNA-seq data based on the method of Grasso et
750 al.⁴¹: a signature of 12 genes (*CCL2*, *CCL3*, *CCL4*, *CXCL9*, *CXCL10*, *CD8A*, *HLA-DOB*, *HLA-*
751 *DMB*, *HLA-DOA*, *GZMK*, *ICOS*, and *IRF1*) was extracted, and a continuous T-cell score derived
752 as their log(TPM) average. The continuous score was then divided into three equal sized
753 intervals (based on all cancers) to provide low, medium and high T-cell score levels.

754

755 ***Neoantigen prediction***

756 Neoantigens were predicted from variant call tables and HLA types using NeoPredPipe⁴⁰, a
757 neoantigen prediction and evaluation pipeline designed for parallel analysis of single- and multi-
758 region samples. We only evaluated single nucleotide variants leading to a single amino acid
759 change, and novel peptides of 9 and 10 amino acids were considered. The pipeline was run

760 with default analysis settings and preserving intermediate files (-p flag), using hg38 and hg19
761 ANNOVAR⁶⁸ reference files for annotation of the TCGA and multi-region CRC samples,
762 respectively. The analysis outputted a table of novel peptides binding the patient's MHC-I
763 molecules (considering all six alleles independently) and their respective recognition potential
764 calculated from their MHC-binding affinity and similarity to pathogenic peptides, as described in
765 ref¹⁹. For evaluating the recognizability (R) part of the recognition potential, we used the
766 parameter values derived in ref¹⁹. Unless stated otherwise, we labelled a peptide as neoantigen
767 if its recognition potential was $\geq 10^{-1}$ (with respect to any of the patient's HLA types) to focus
768 on antigens with the highest predicted probability of eliciting an immune response: both similar
769 to known pathogens and similar or stronger MHC-binders than their wild-type counterpart. A
770 mutation was considered (neo)antigenic if there was at least a single peptide produced from the
771 mutated base that got labelled as neoantigen.

772 To evaluate the antigenicity distribution of tumors, we used the predicted percentile rank of
773 neoantigens that ranks a putative antigen against a large set of random substrates to the same
774 HLA molecule, and thus eliminates bias introduced by structural properties of HLA alleles⁷², that
775 might be present in plain binding affinity values (considered in the recognition potential pipeline).
776 We inverted this value to obtain a normalized binding score that correlates with the importance
777 ranking of peptides, where values above ~ 1.3 represented strong putative antigens.

778 ***Computation of VAF and CCF values***

779 For each mutation, we calculated the VAF as the number of mutant reads spanning the position,
780 divided by the number of total reads of the position. The proportion of cancer cells carrying a
781 particular mutation (CCF) was calculated from the VAF of the mutation, sample purity (tumor
782 content), and copy number (CN) of the mutation's genomic locus as: $(VAF * CN) / purity$. CCF
783 values above 1 (arising from sequencing noise and copy-neutral loss-of-heterozygosity events)

784 were assumed to be 1. We only considered a mutation as subclonal if it had CCF<0.6, to
785 account for the possibility of ‘bleeding’ of clonal mutations into the subclonal frequency range
786 because of the limited sequence depth of TCGA samples.

787 For pooling together VAF distributions of a cohort of samples (Fig. 4f), we first filtered the set of
788 TCGA cancers: cancers with any evidence of immune escape (including allelic imbalance of
789 HLA locus), MMR or POLE cancers and cancers with purity <50% were discarded. The
790 remaining cancers were divided into low and medium immune infiltration groups (all highly T-cell
791 score cancers were immune escaped and previously discarded). Total and neoantigen-
792 associated cumulative VAF distributions were computed from all mutations detected at
793 subclonal frequencies in the two groups. In a similar manner, TCGA MSS cancers with purity
794 >70% (to ensure more accurate VAF and ploidy calls) were combined into a cohort to study
795 mutations in essential genes (Extended Data Fig. 5f). Essential genes, and antigenic mutations
796 located in essential genes were identified using the list of shared genes in ref⁵⁰.

797 ***Synthetic cohorts***

798 In order to evaluate the antigen-producing capacity of different mutational processes, we
799 generated synthetic tumor cohorts matching the mutation number and tri-nucleotide composition
800 of real cancers. We measured the average composition (as measured by 96-channel-
801 composition⁵⁵) of the real cohort (e.g. TCGA CRCs, Extended Data Fig. 6d), and randomly
802 sampled a matching number of exonic mutations at probability specified by the respective
803 channel intensities. Six HLA haplotypes were also randomly sampled from the complete list of
804 alleles in the real cohort. Sampling was repeated independently 100 times to generate a
805 synthetic cohort.

806 ***Statistical analysis***

807 Details of statistical analysis performed are summarized in the Life Science Reporting
808 Summary. All data processing and statistical tests were performed in R (version 3.5.0) using
809 built-in functions. The tests and functions used were as follows: Figs. 1d, 2c,e, Extended Data
810 Figs. 3c, 6a,b,c,e: Mann-Whitney U-test/ Wilcoxon sum-rank test (*wilcox.test*, default settings).
811 Figs. 2d,f and Extended Data Fig. 3a,b,d: Chi-squared test (*chisq.test*). Fig. 2g and Extended
812 Data Fig. 3e: One-sided Mann-Whitney U-test (*wilcox.test* with option *alternative='greater'*). Fig.
813 3c-d: One-sided paired Wilcoxon signed-rank test (*wilcox.test* with options *paired=TRUE* and
814 *alternative='greater'*). Fig. 4c-d and S5b-c: Kolmogorov-Smirnov test (*ks.test*) between the raw
815 VAF distribution of neoantigens and all mutations. The two distributions were deemed with
816 significance level $p < 0.1$ or as indicated in Fig. 4c-d and Extended Data Fig. 5b-c. Fig. 5c and
817 Extended Data Fig. 6d,f: Paired Wilcoxon signed-rank test (*wilcox.test*, option *paired=TRUE*).
818 Extended Data Fig. 6g: Students t-test against mean of 1 (*t.test*, *mu=1*).

819 All violin plots were generated with automatic smoothing bandwidth value of *geom_violin*.
820 Individual observations for TCGA samples are shown on top of violins, generated with
821 *geom_dotplot*.

822 DATA AVAILABILITY

823 The datasets analyzed during the current study are available from the NCI Genomics Data
824 Commons Portal (<https://portal.gdc.cancer.gov>) COAD, READ, STAD and UCEC domains, and
825 from the European Genome-Phenome Archive (<https://ega-archive.org/>) at accession code:
826 EGAS00001003066.

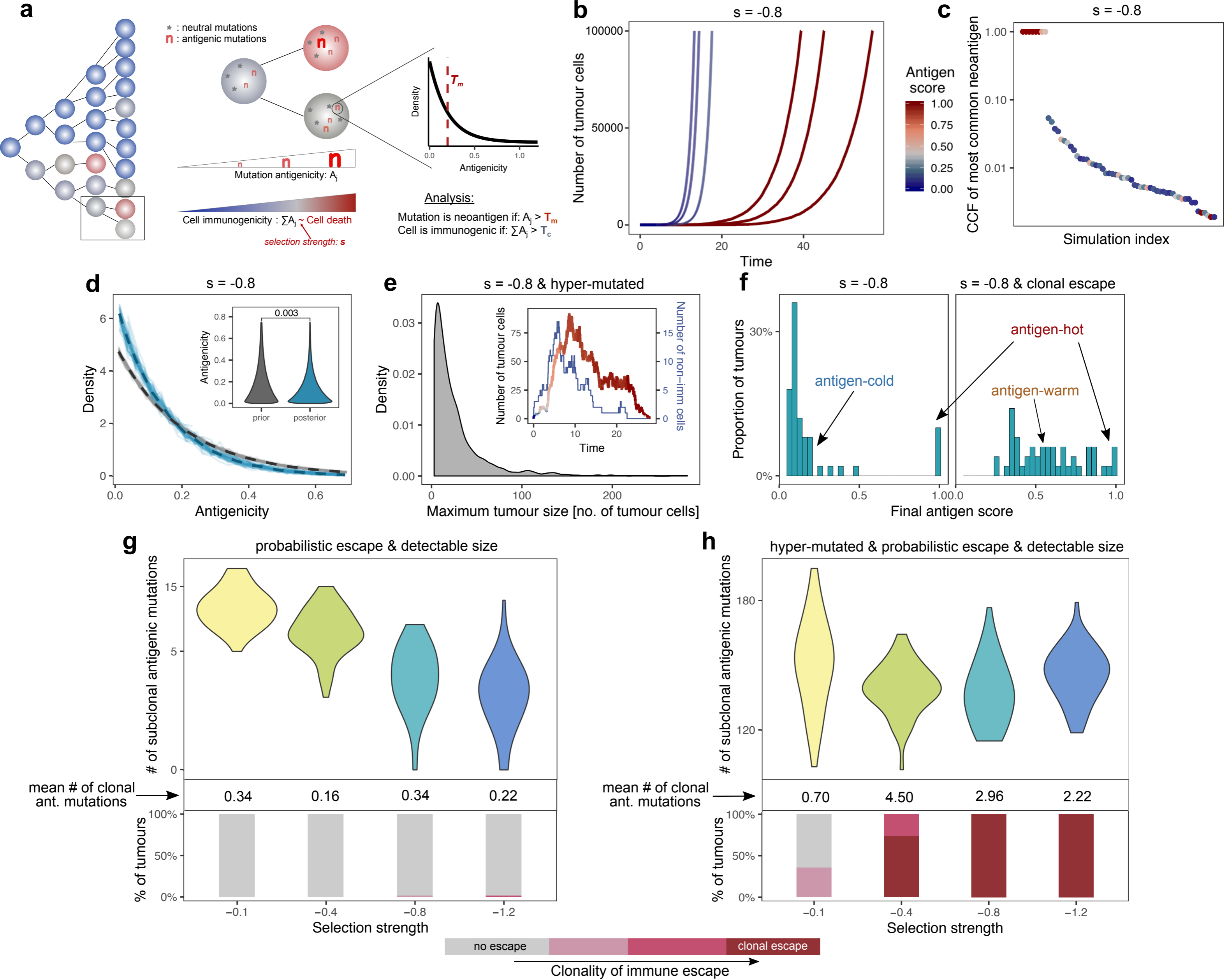
827 CODE AVAILABILITY

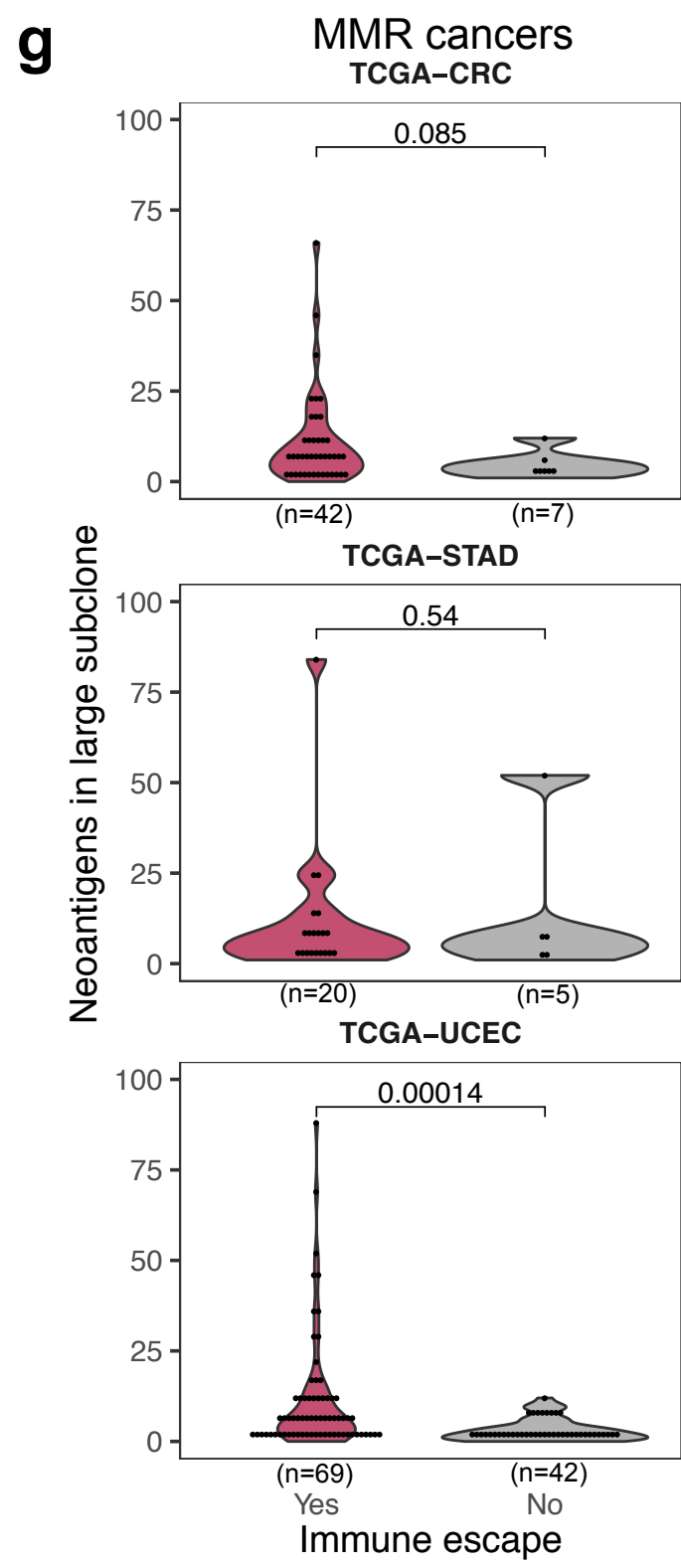
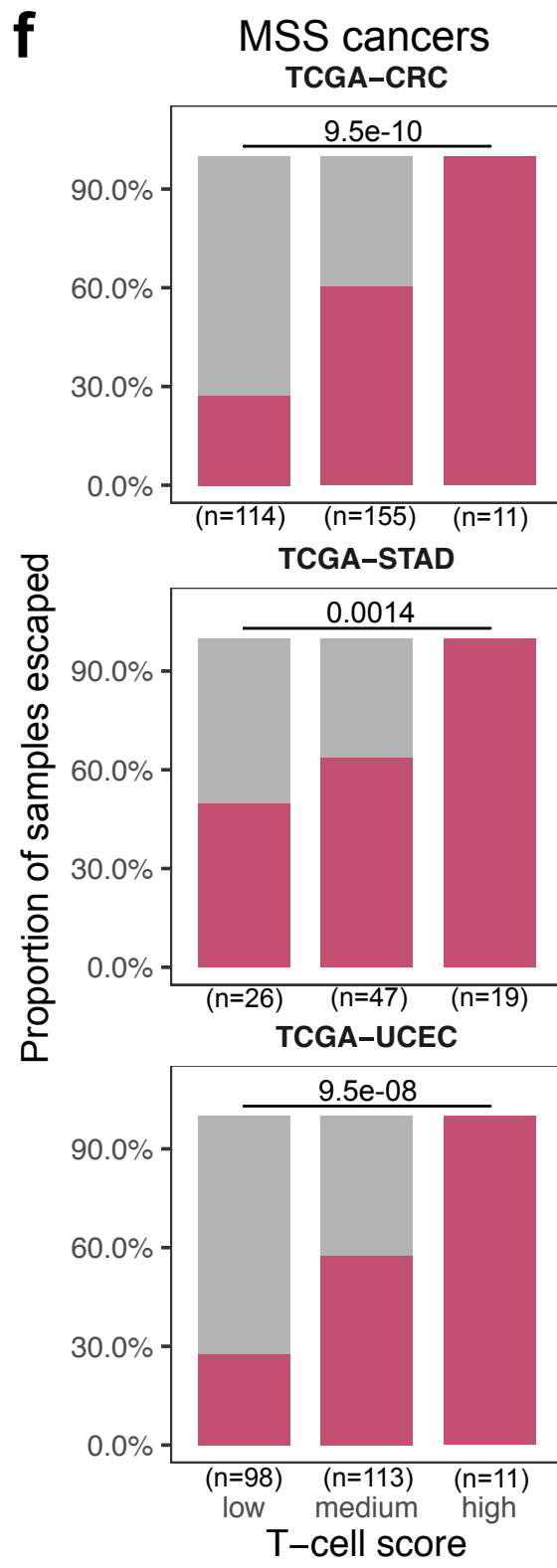
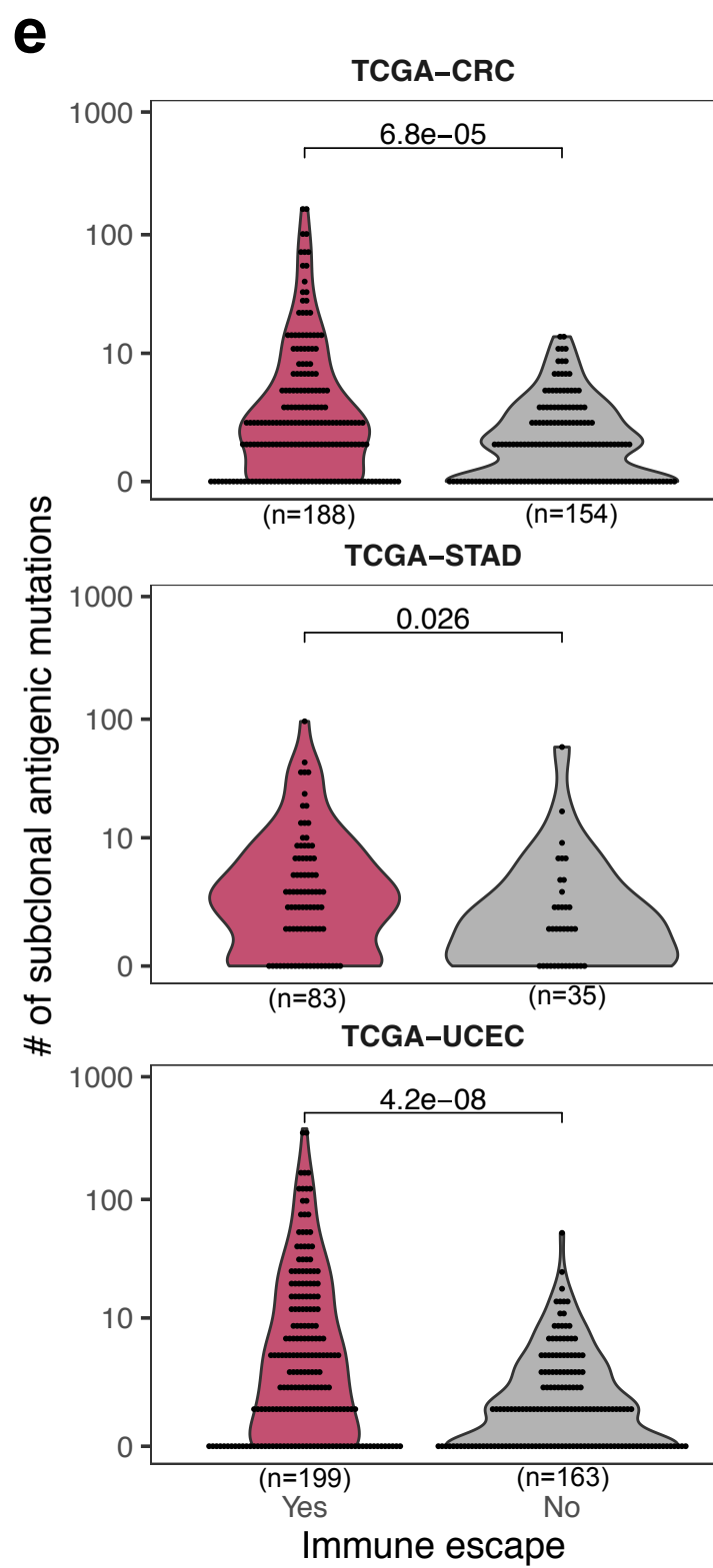
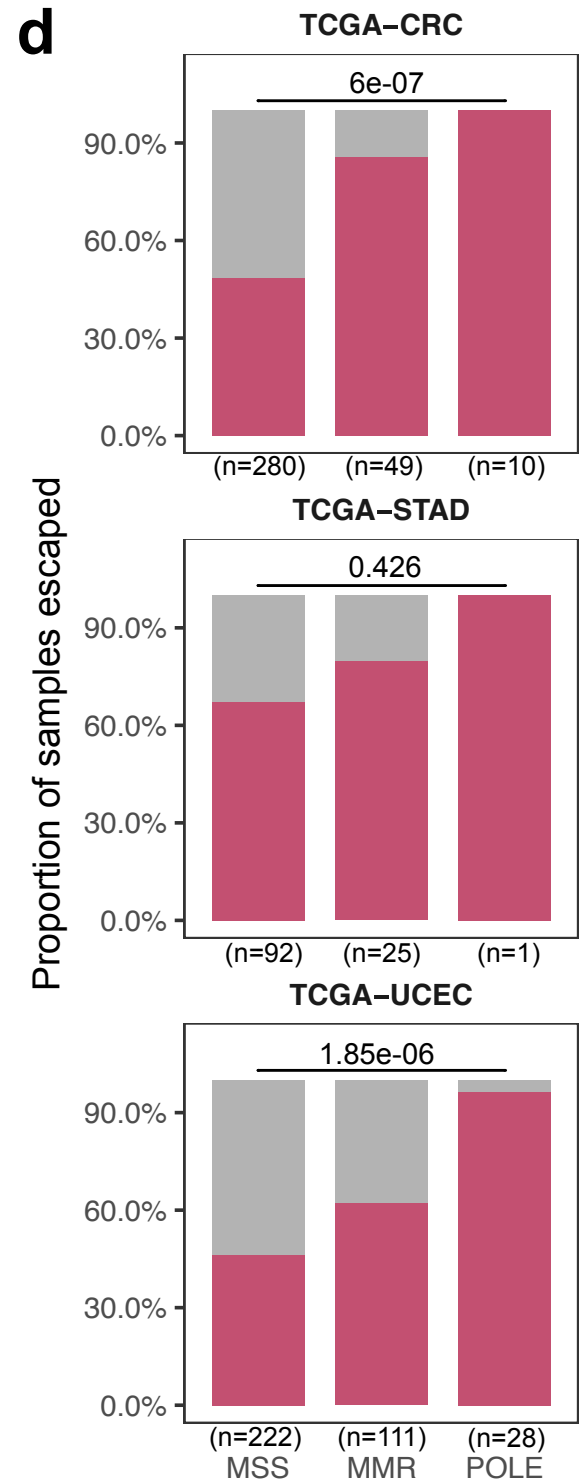
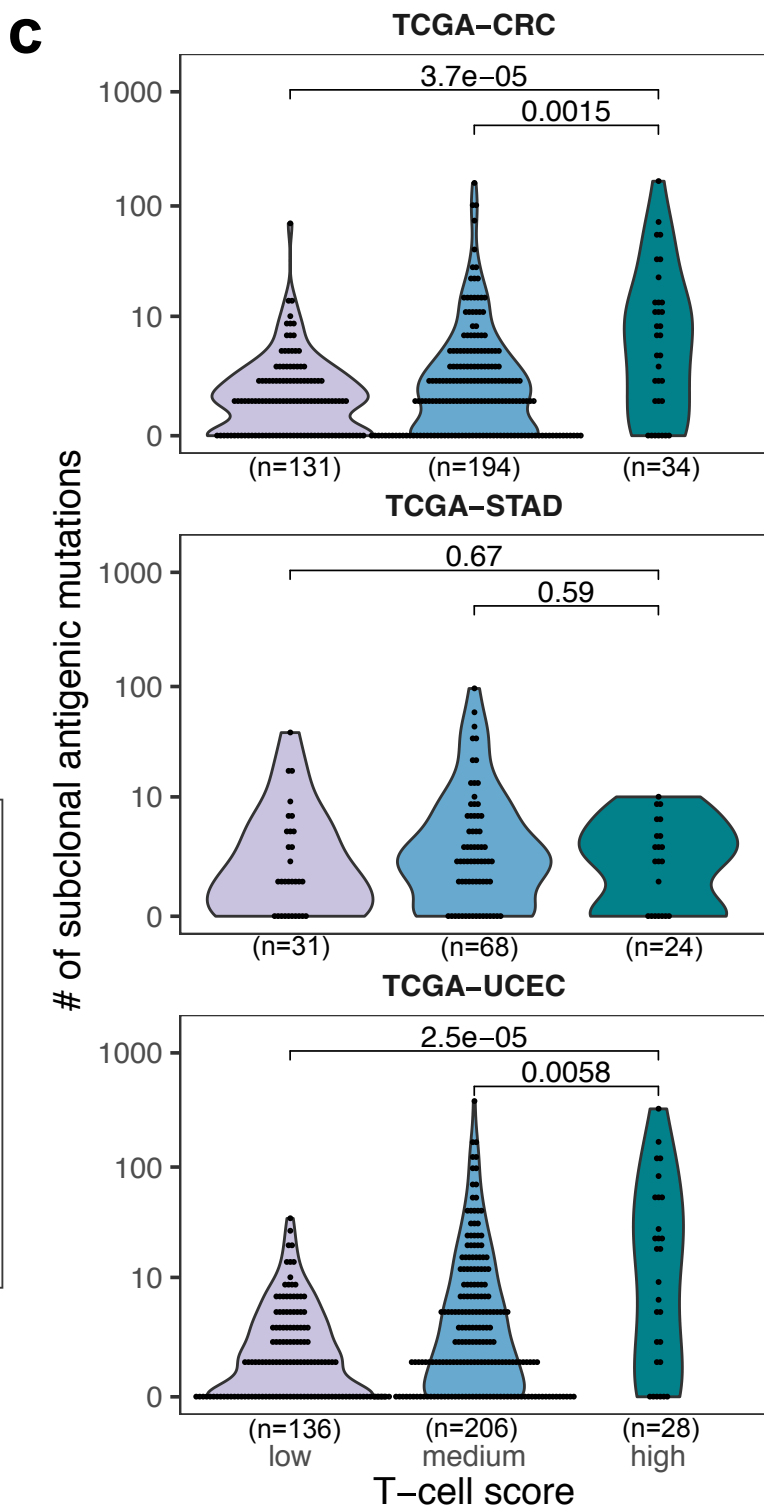
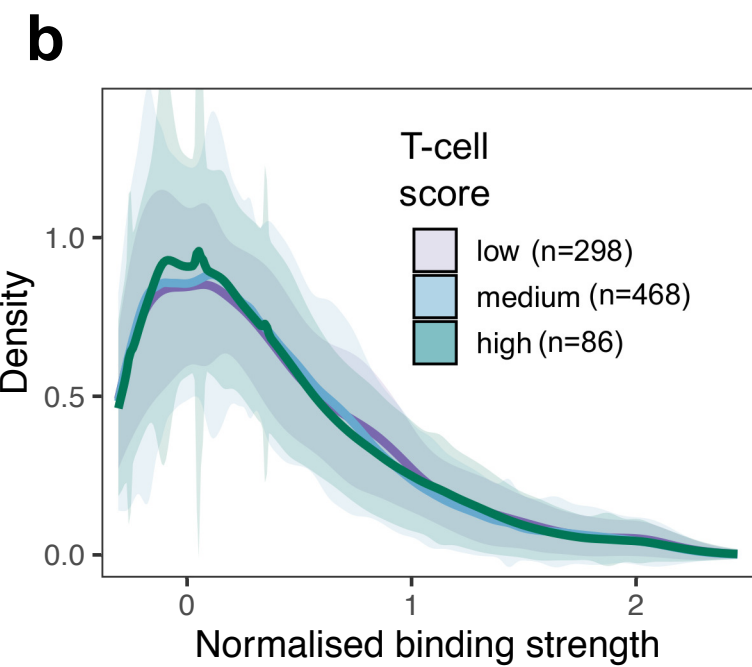
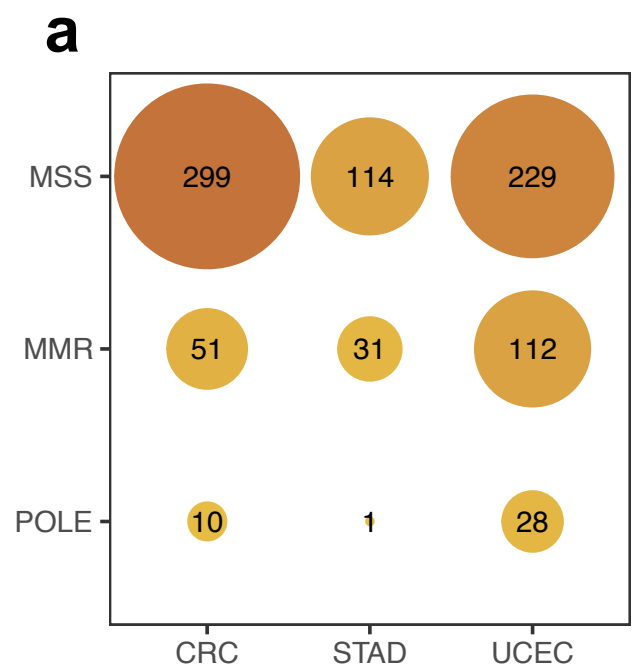
828 Julia (<https://julialang.org/>, version 0.5+) code implementing simulations of the tumor growth
829 model is available from Zenodo (doi: 10.5281/zenodo.3601322)⁶¹.

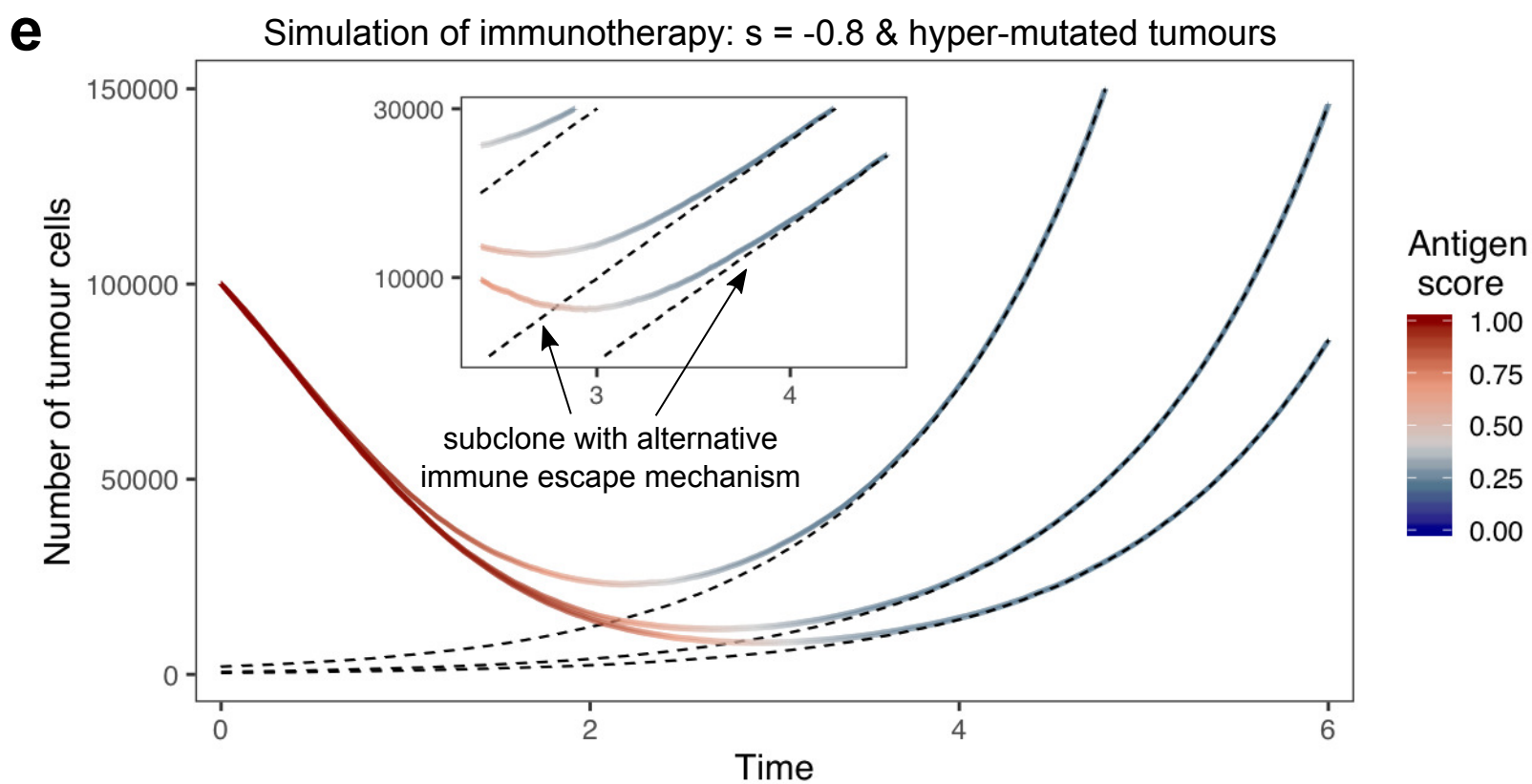
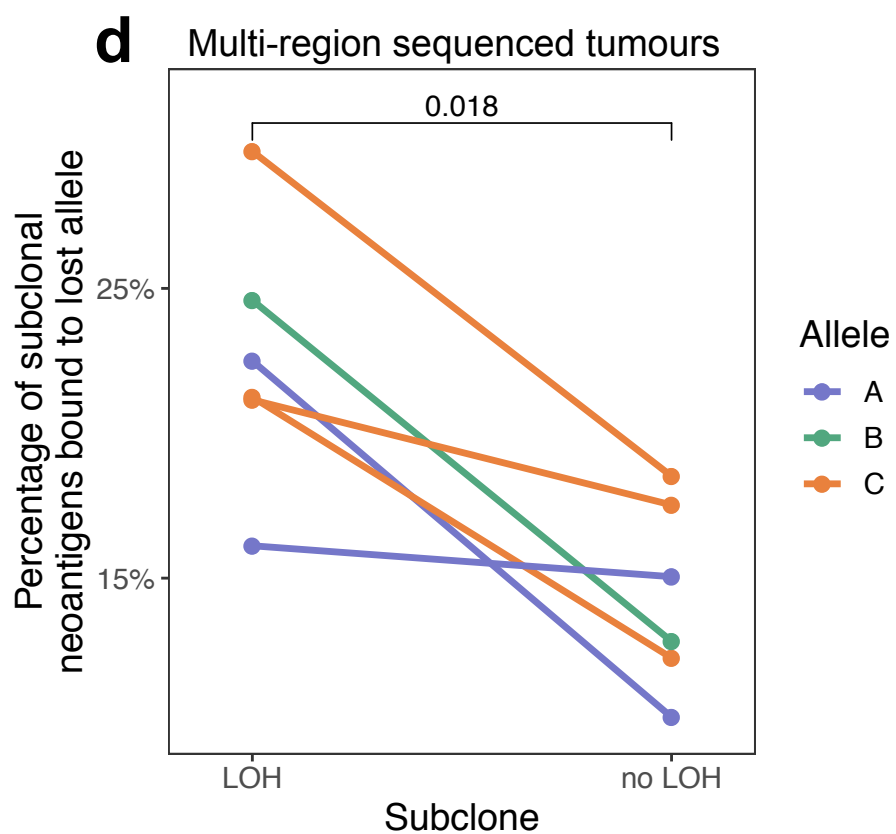
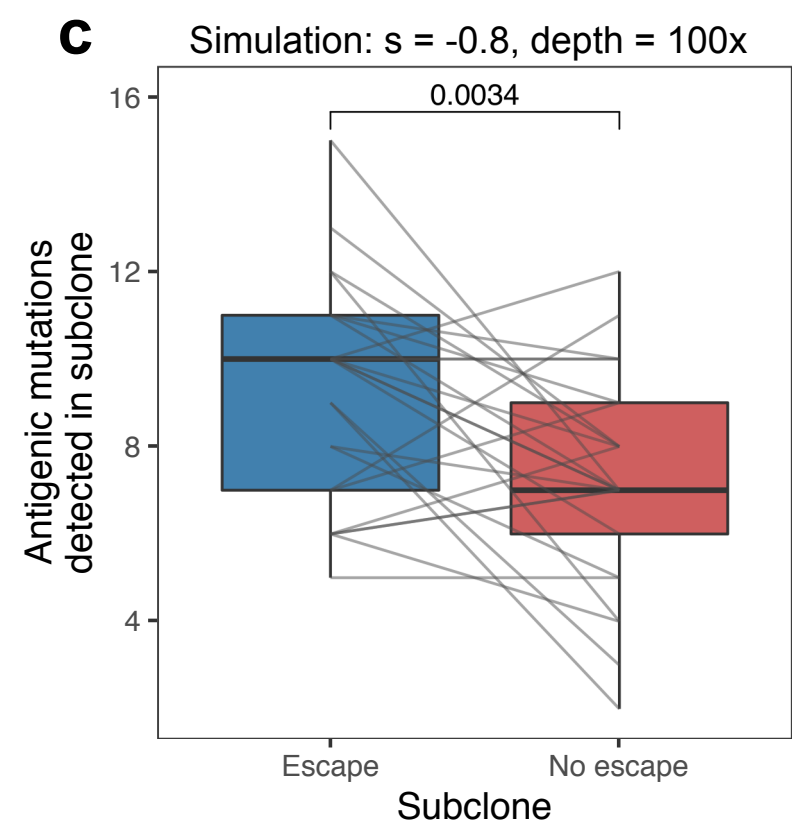
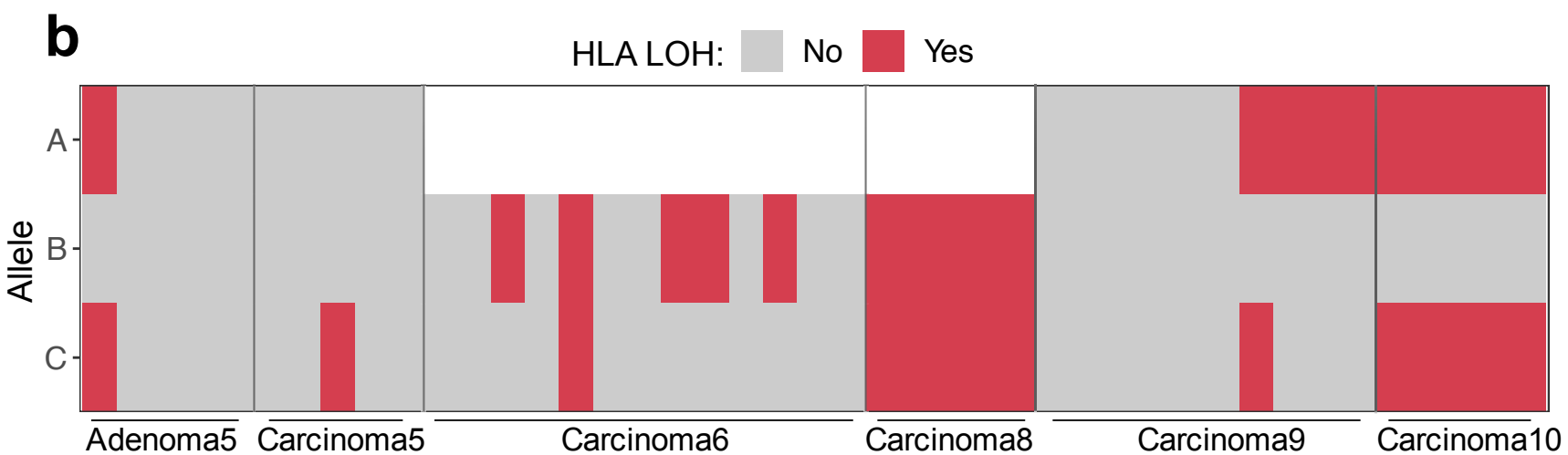
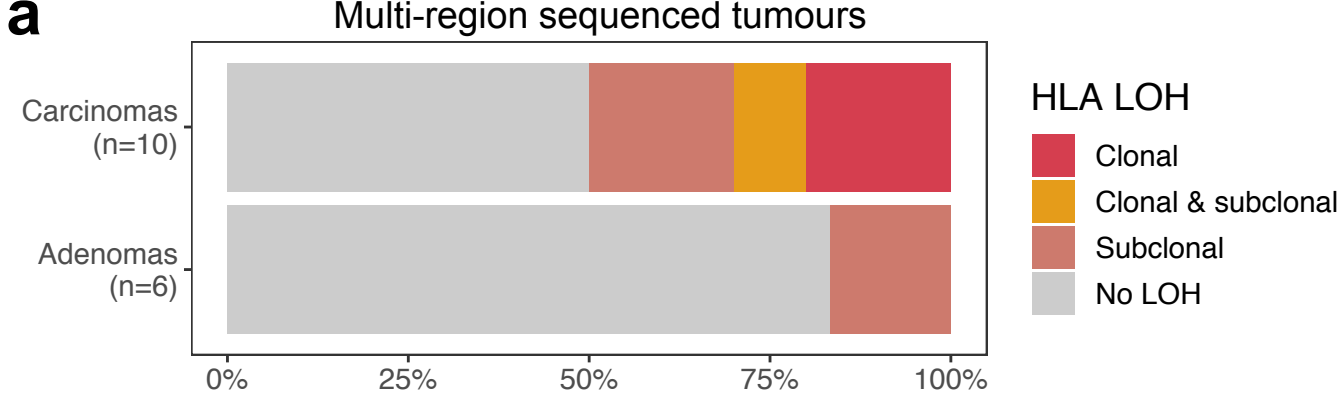
830

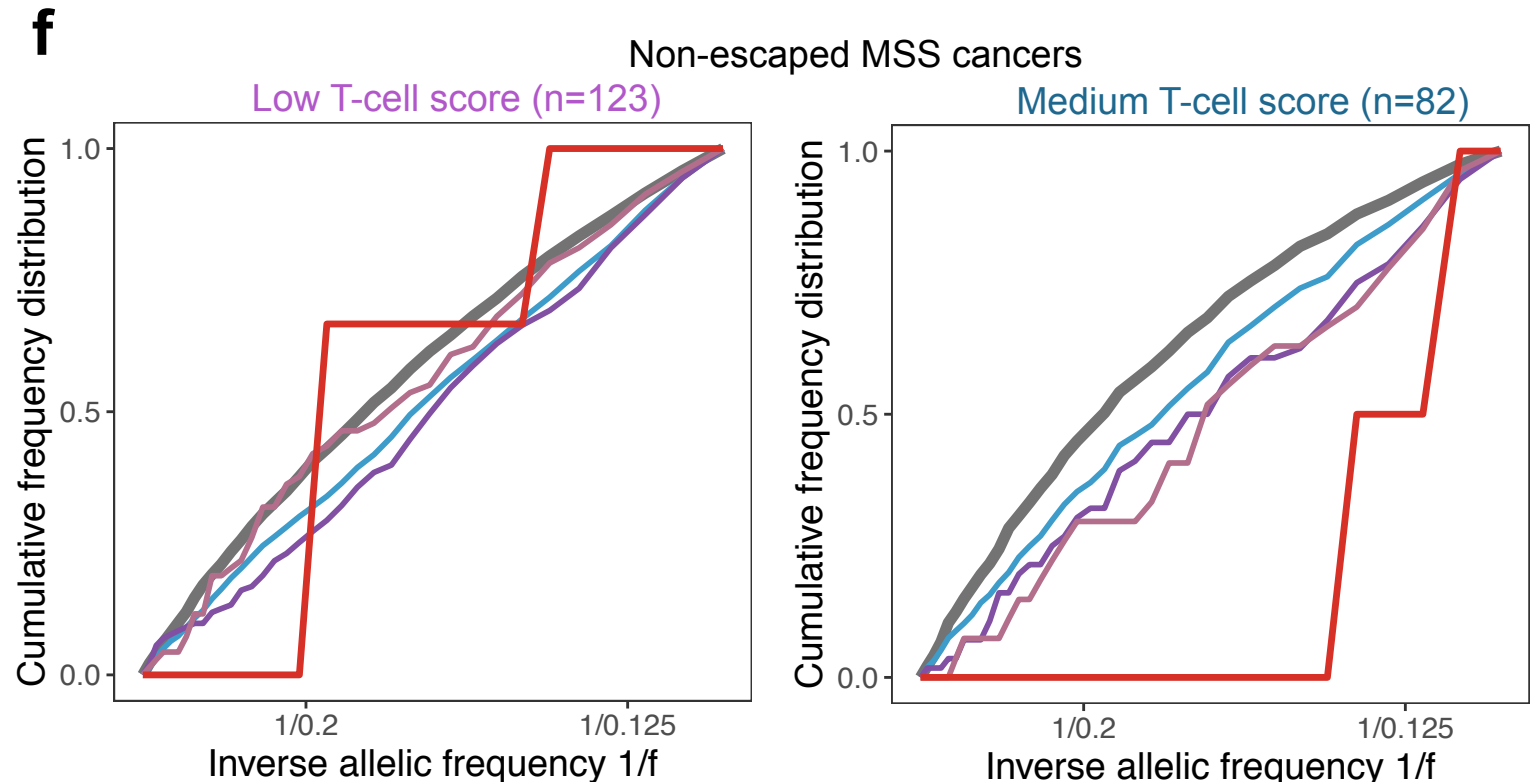
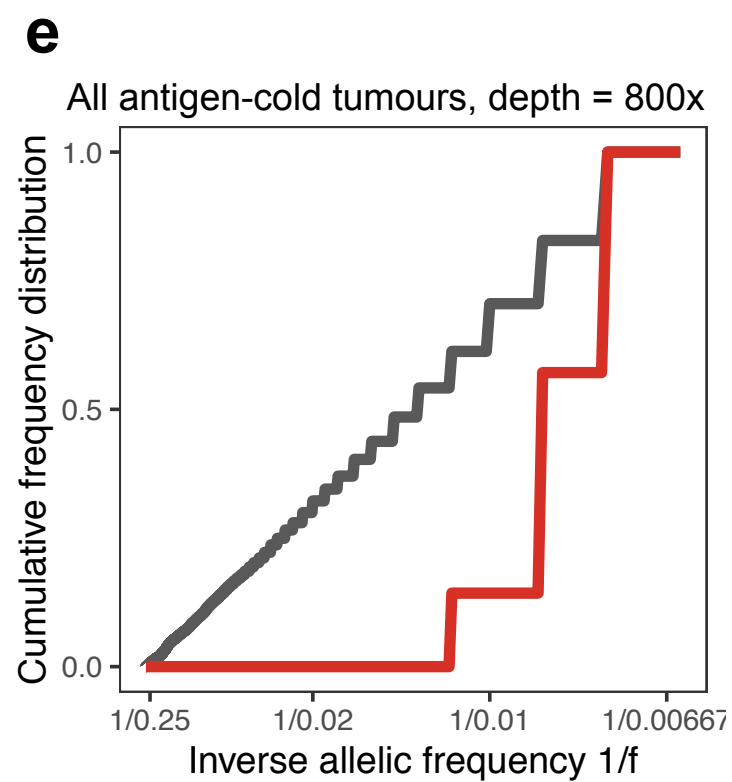
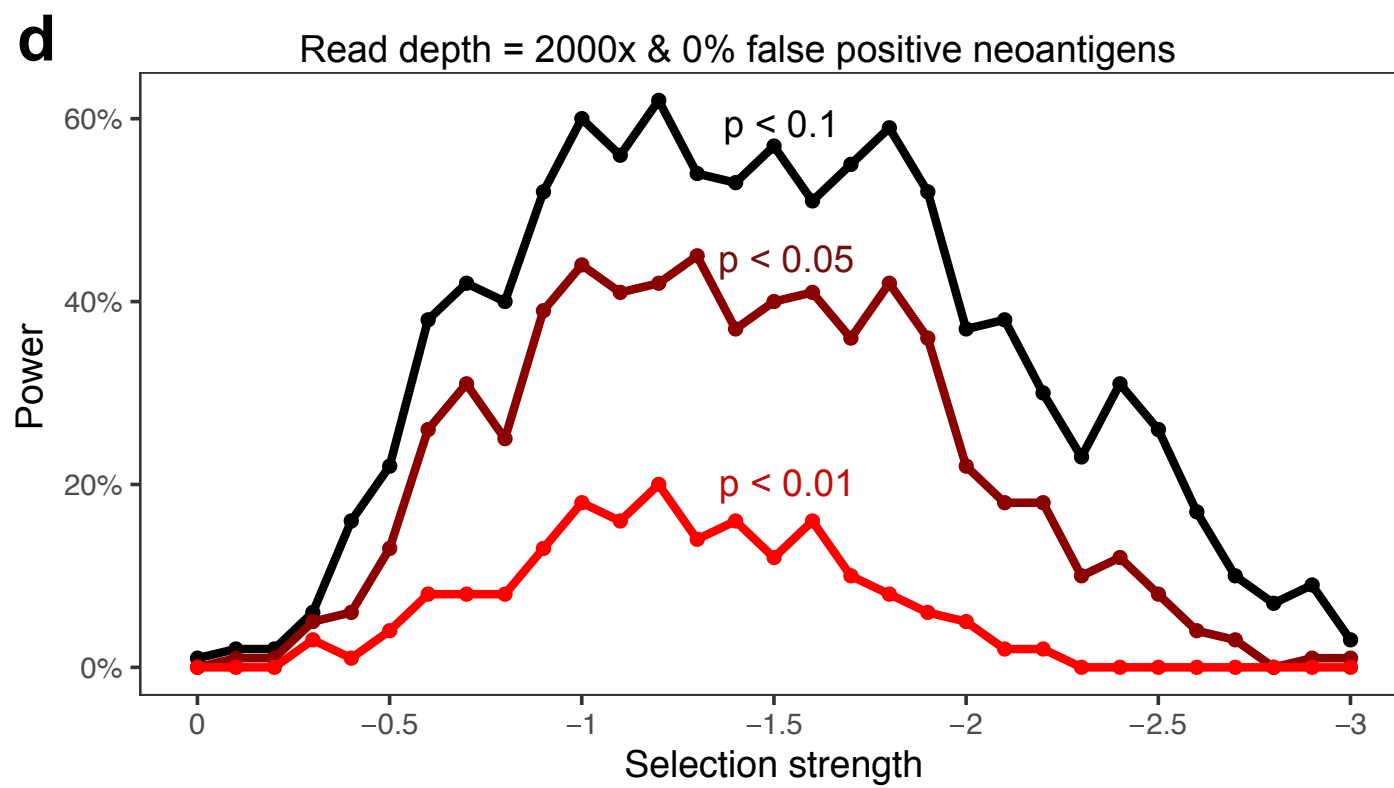
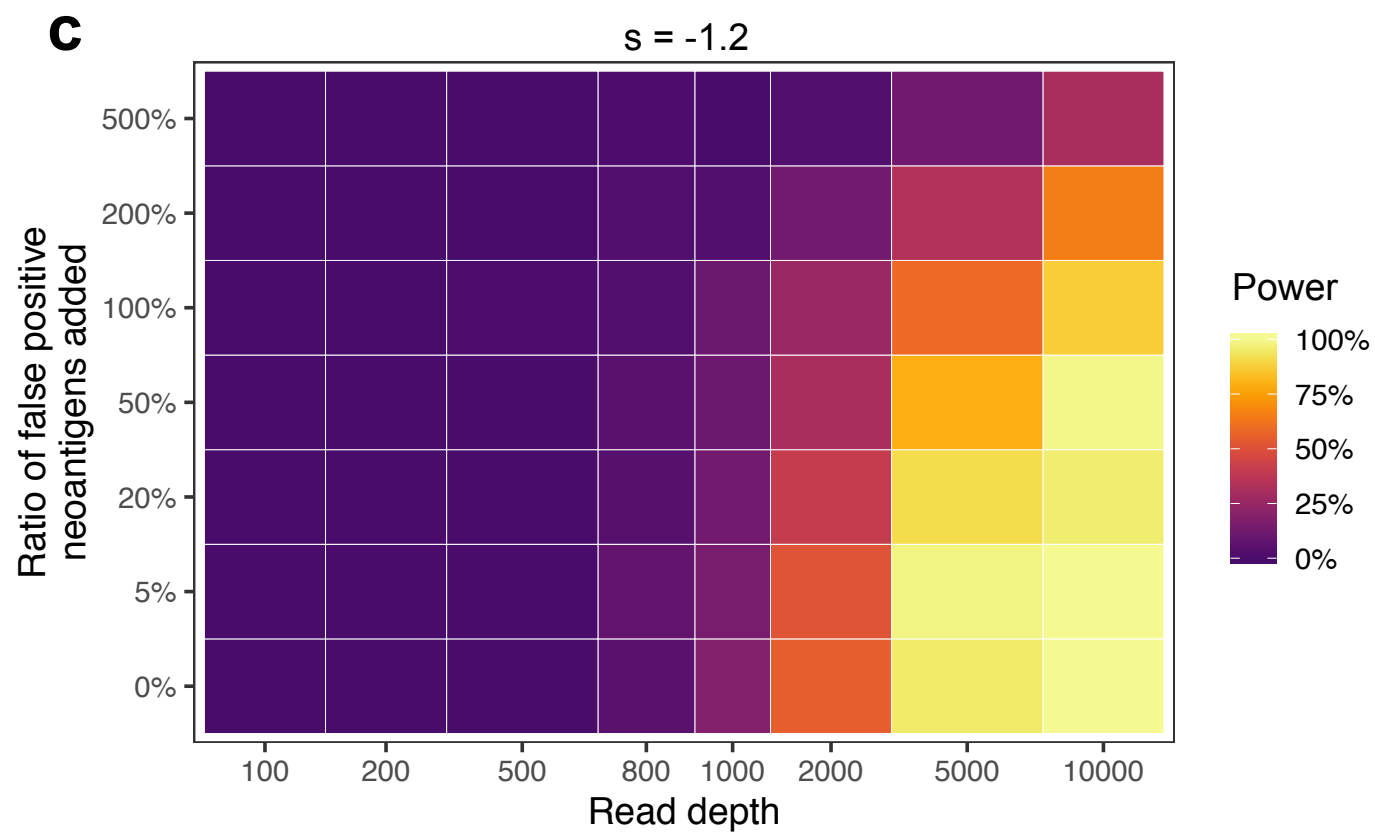
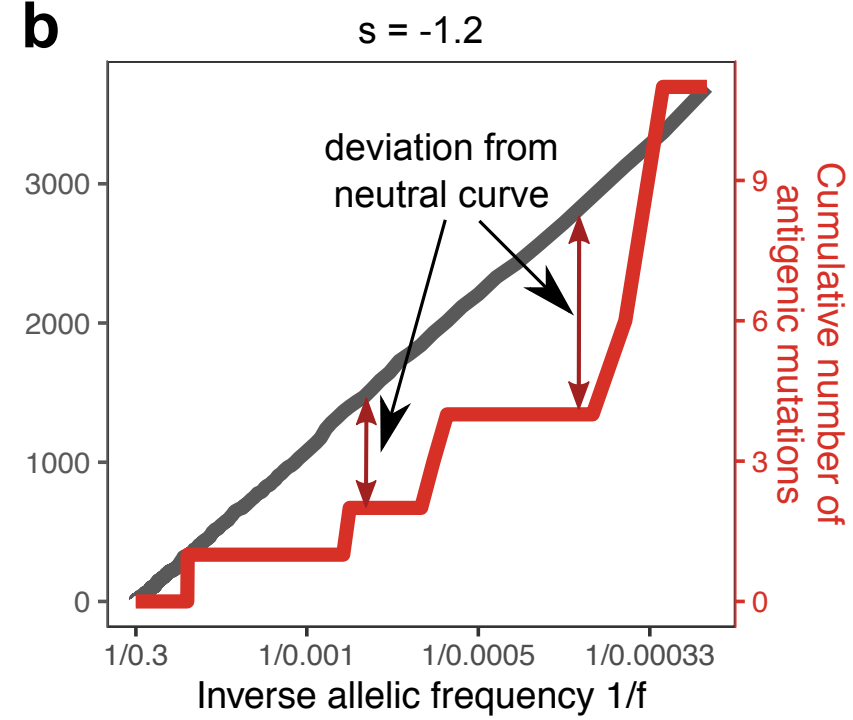
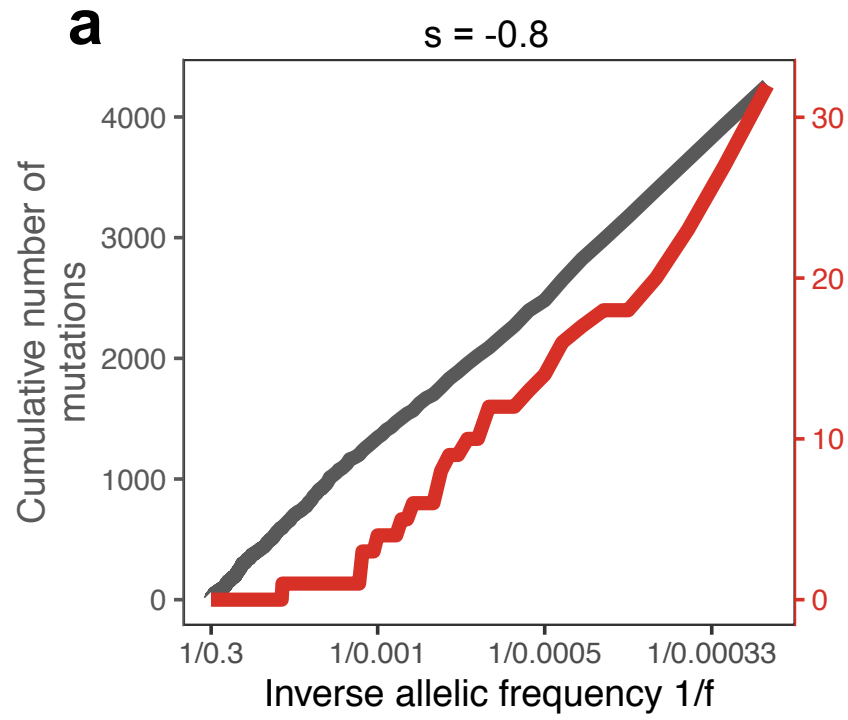
831 **REFERENCES (CONTINUED)**

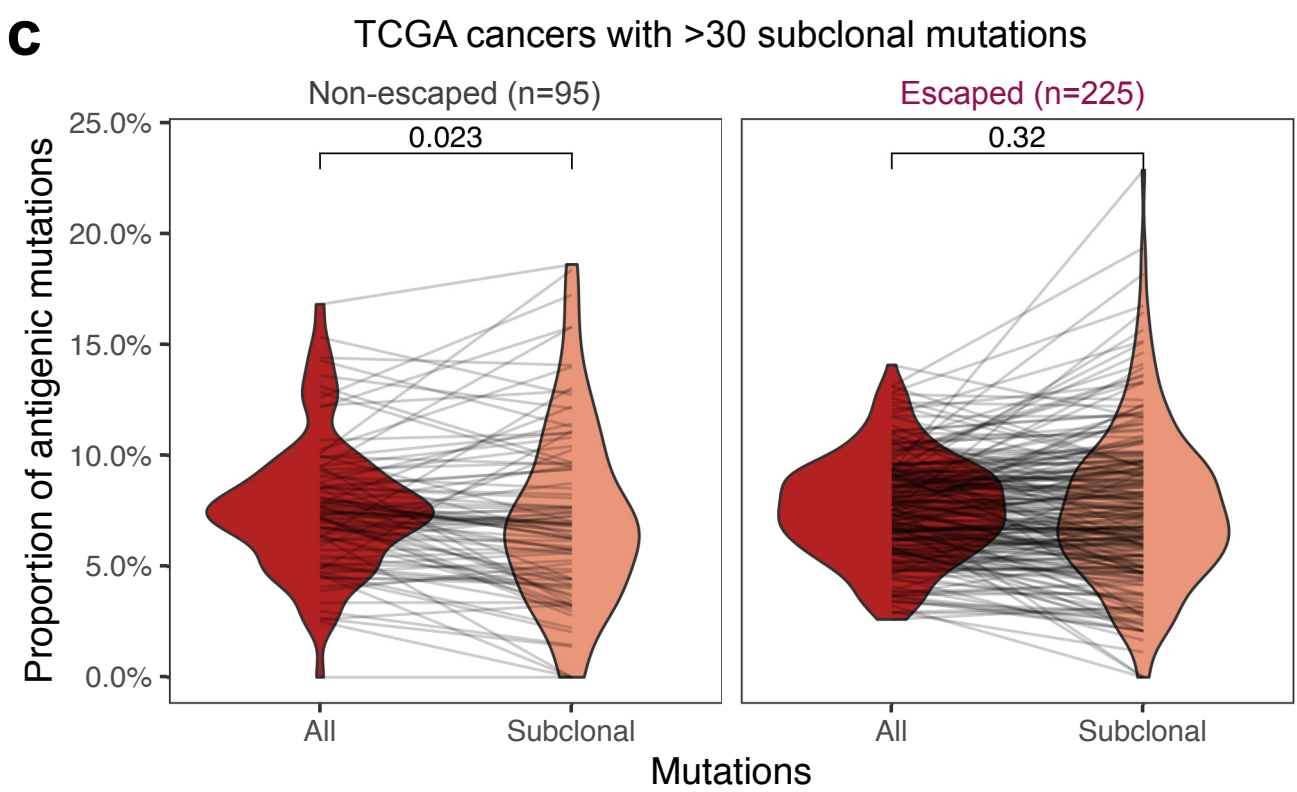
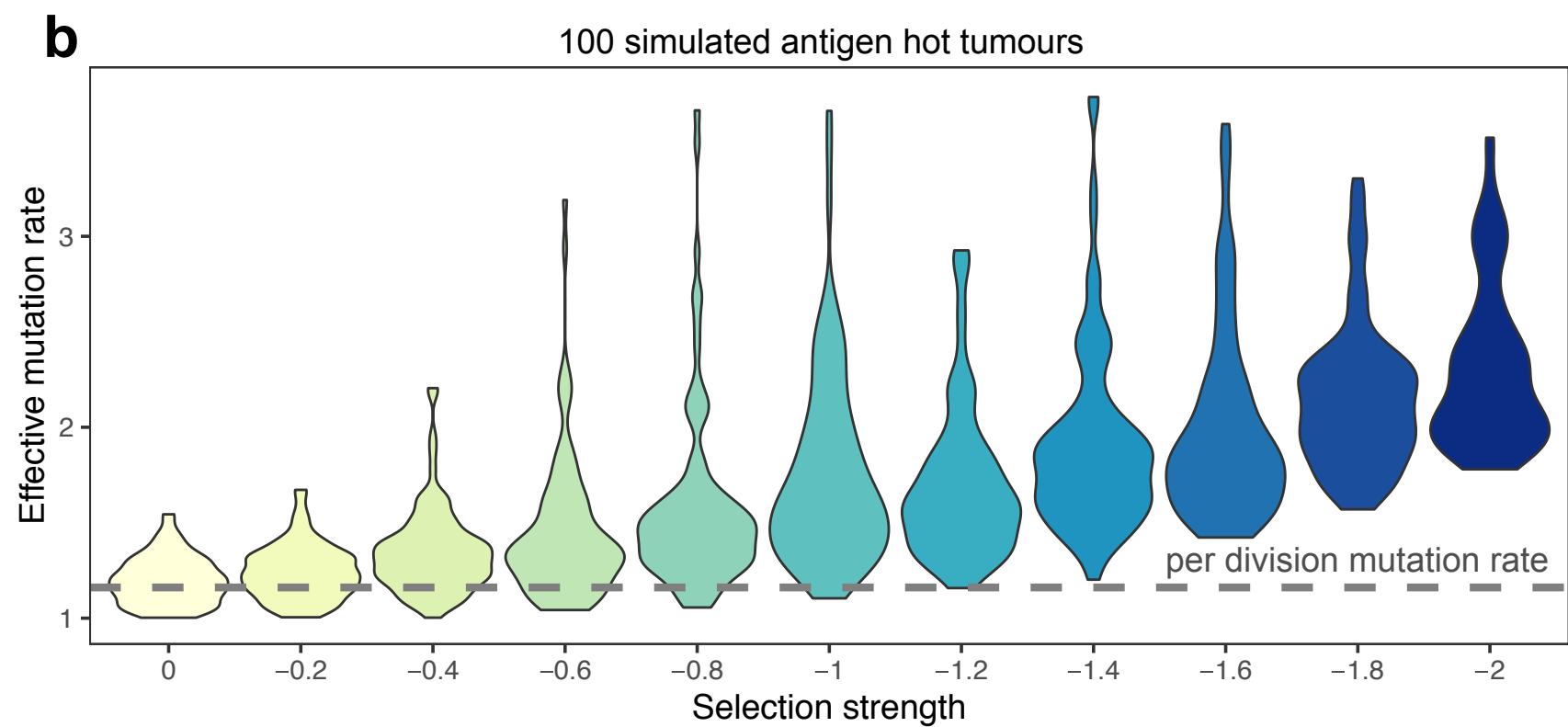
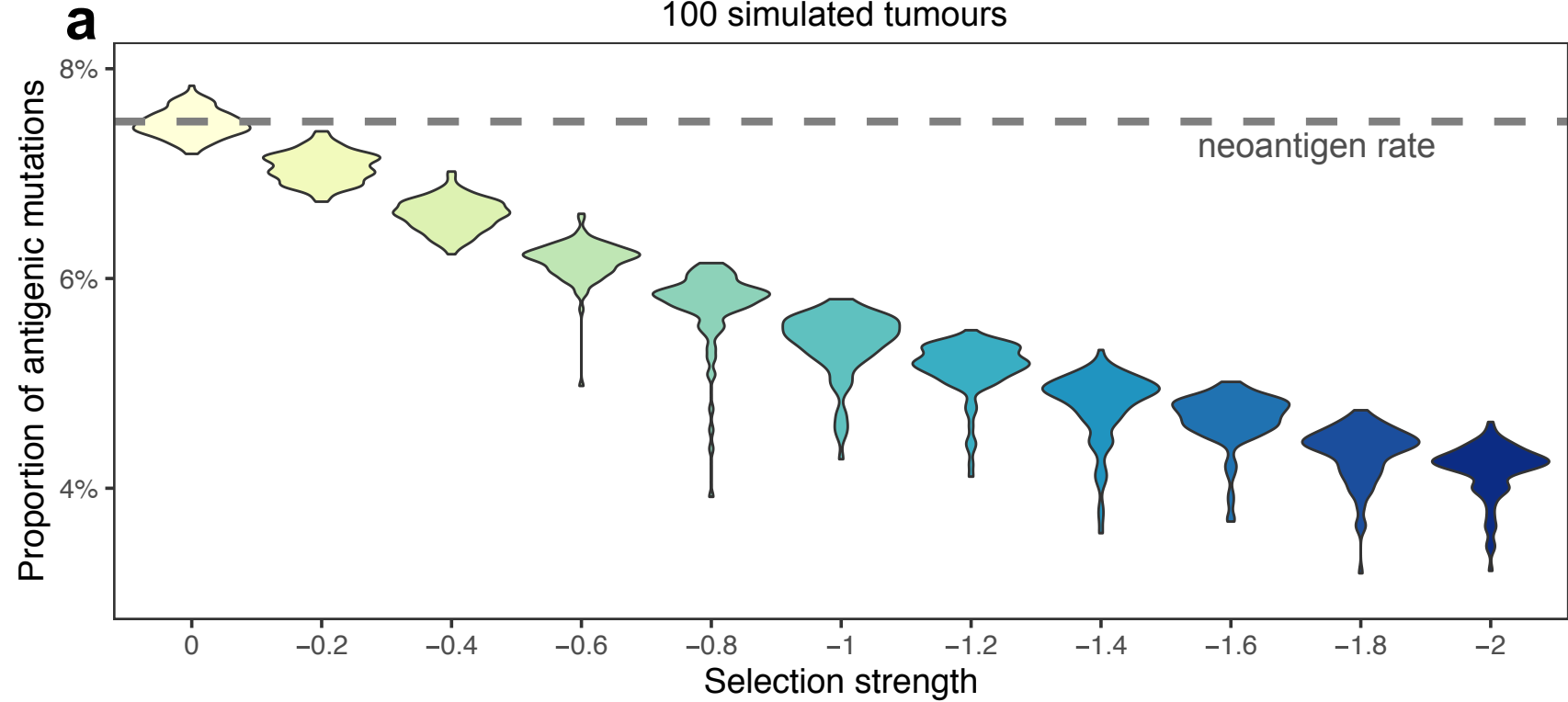
- 832 61. Lakatos, E. CloneGrowthSimulation. (2020). doi:10.5281/zenodo.3601322
- 833 62. Gillespie, D. T. A general method for numerically simulating the stochastic time evolution of
- 834 coupled chemical reactions. *J. Comput. Phys.* **22**, 403–434 (1976).
- 835 63. Grossman, R. L. *et al.* Toward a Shared Vision for Cancer Genomic Data. *N. Engl. J. Med.* **375**,
- 836 1109–1112 (2016).
- 837 64. Van Loo, P. *et al.* Allele-specific copy number analysis of tumors. *Proc. Natl. Acad. Sci.* **107**,
- 838 16910–16915 (2010).
- 839 65. The Cancer Genome Atlas Network. Comprehensive molecular characterization of human colon
- 840 and rectal cancer. *Nature* **487**, 330–337 (2012).
- 841 66. Kautto, E. A. *et al.* Performance evaluation for rapid detection of pan-cancer microsatellite
- 842 instability with MANTIS. *Oncotarget* **8**, 7452–7463 (2016).
- 843 67. Rimmer, A. *et al.* Integrating mapping-, assembly- and haplotype-based approaches for calling
- 844 variants in clinical sequencing applications. *Nat. Genet.* **46**, 912–918 (2014).
- 845 68. Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic variants from high-
- 846 throughput sequencing data. *Nucleic Acids Res.* **38**, e164 (2010).
- 847 69. Alsaab, H. O. *et al.* PD-1 and PD-L1 Checkpoint Signaling Inhibition for Cancer Immunotherapy:
- 848 Mechanism, Combinations, and Clinical Outcome. *Front. Pharmacol.* **8**, 561 (2017).
- 849 70. McLaren, W. *et al.* The Ensembl Variant Effect Predictor. *Genome Biol.* **17**, 122 (2016).
- 850 71. Favero, F. *et al.* Sequenza: allele-specific copy number and mutation profiles from tumor
- 851 sequencing data. *Ann. Oncol. Off. J. Eur. Soc. Med. Oncol.* **26**, 64–70 (2015).
- 852 72. Jurtz, V. *et al.* NetMHCpan-4.0: Improved Peptide–MHC Class I Interaction Predictions Integrating
- 853 Eluted Ligand and Peptide Binding Affinity Data. *J. Immunol.* **199**, 3360–3368 (2017).
- 854
- 855
- 856

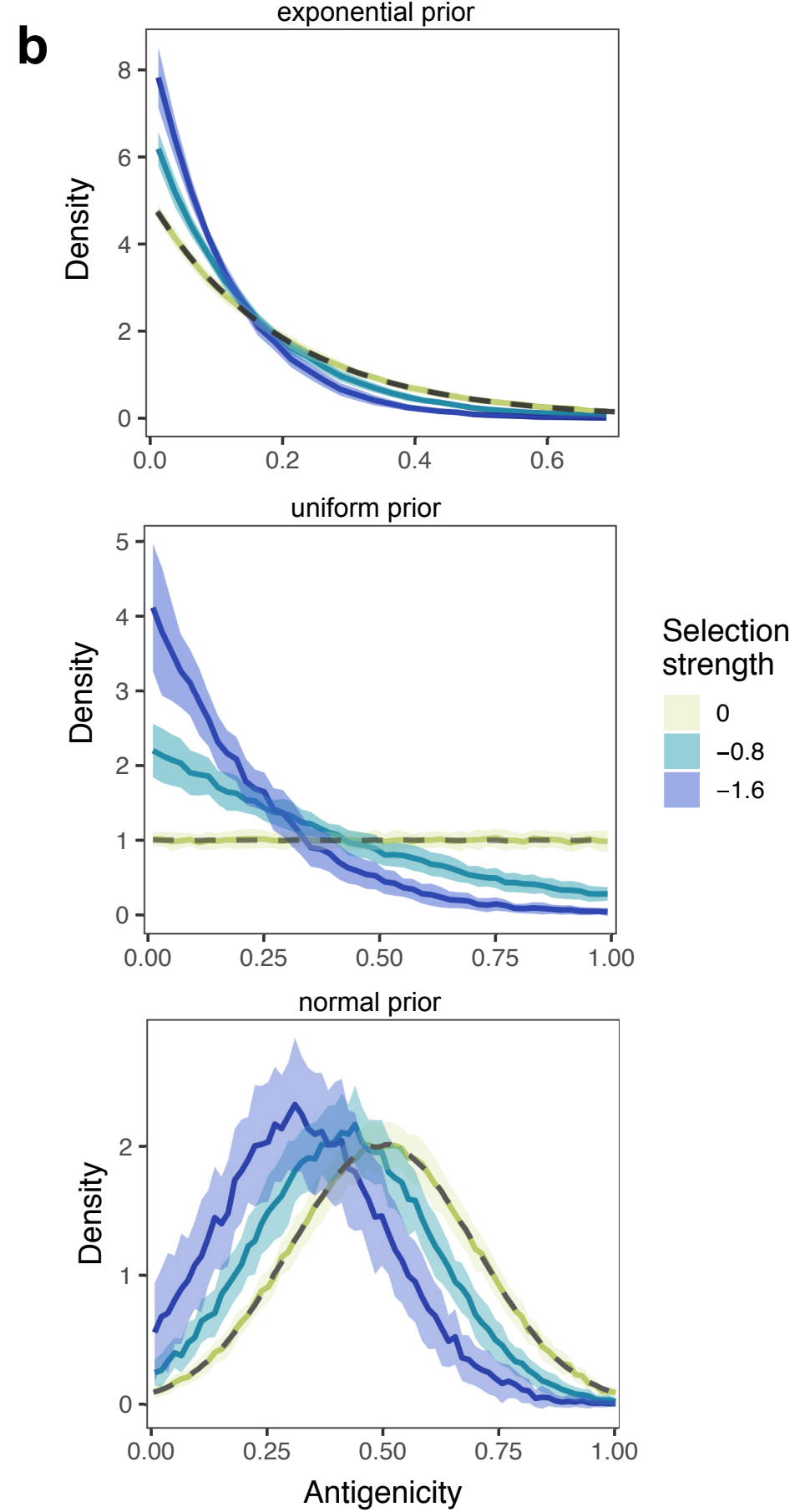
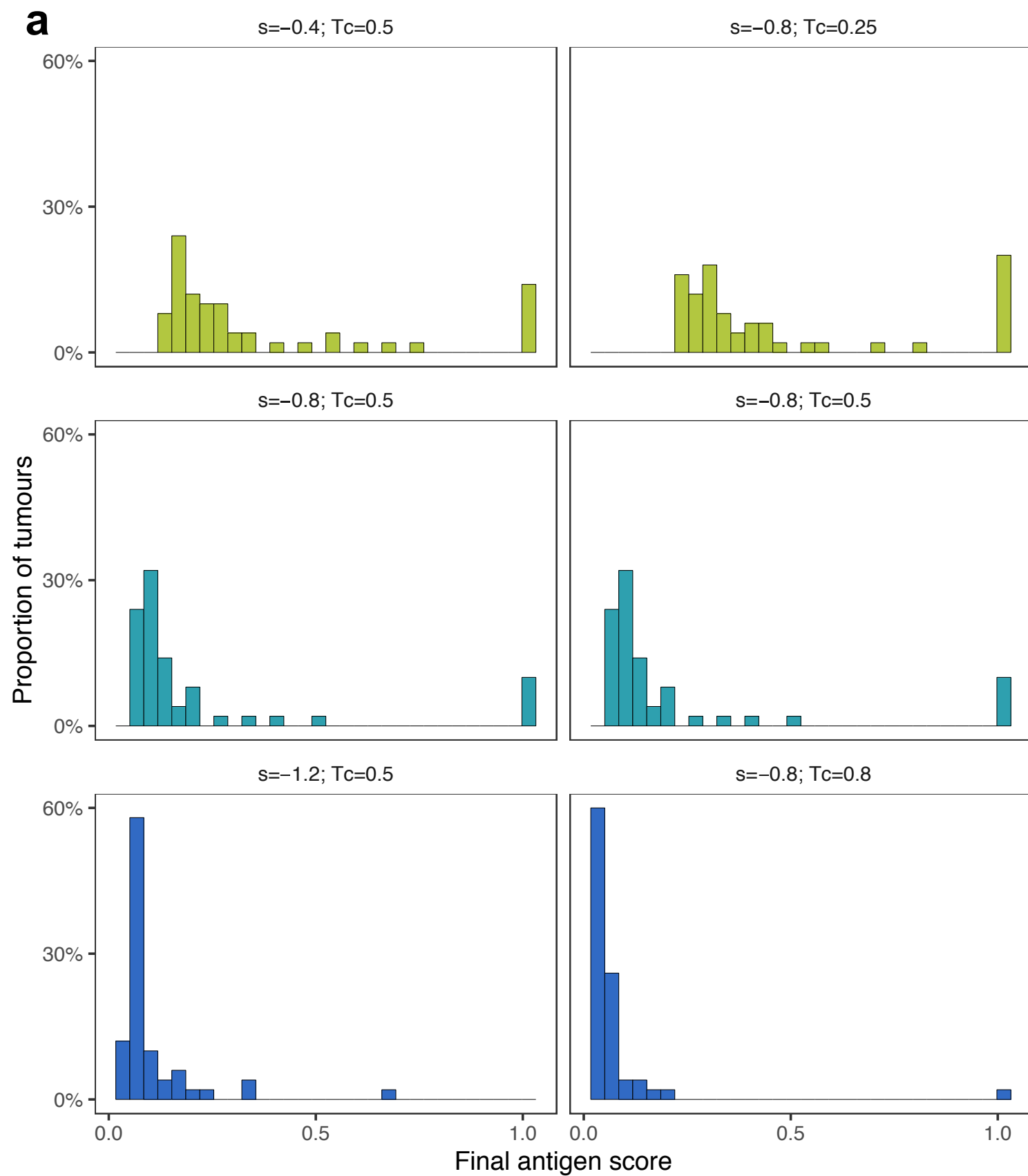












growing population

constant size population

