1  **A neoantigen fitness model predicts tumor response to checkpoint**
2  **blockade immunotherapy**
3
4  Marta Łuksza[1,*], Nadeem Riaz[2,3], Vladimir Makarov[3,4], Vinod P. Balachandran[5,6,7],
5  Alexander Solovyov[8], Naiyer A. Rizvi[9], Taha Merghoub[7,10,11], Arnold J. Levine[1], Timothy
6  A. Chan[2,3,4,7], Jedd D. Wolchok[7,10,11,12], Benjamin D. Greenbaum[8,*]
7
8  [1]The Simons Center for Systems Biology, Institute for Advanced Study, Princeton, NJ,
9  USA.
10  [2]Departments of Radiation Oncology, [5]Surgery and [12]Medicine, Memorial Sloan
11  Kettering Cancer Center, New York, NY, USA.
12  [3]Immunogenomics and Precision Oncology Platform, Memorial Sloan Kettering Cancer
13  Center, New York, NY, USA.
14  [4]Human Oncology and Pathogenesis Program, Memorial Sloan Kettering Cancer
15  Center, New York, NY, USA.
16  [6]David M. Rubenstein Center for Pancreatic Cancer Research, Memorial Sloan Kettering
17  Cancer Center, New York, NY, USA.
18  [7]Parker Institute for Cancer Immunotherapy, Memorial Sloan Kettering Cancer Center,
19  New York, NY, USA.
20  [8]Tisch Cancer Institute, Departments of Medicine, Oncological Sciences, and Pathology,
21  Icahn School of Medicine at Mount Sinai, New York, NY, USA.
22  [9]Department of Medicine, Columbia University Medical Center, New York, NY, USA
23  [10]Ludwig Collaborative and Swim Across America Laboratory, Memorial Sloan Kettering
24  Cancer Center, New York, NY, USA.
25  [11]Melanoma and Immunotherapeutics Service, Department of Medicine, Memorial Sloan
26  Kettering Cancer Center, New York, NY, USA; Weill Cornell Medical College, Cornell
27  University, New York, NY, USA.
28
29  [*]**Corresponding Authors:**
30  Marta Łuksza, PhD
31  The Simons Center for Systems Biology
32  School of Natural Sciences
33  Institute for Advanced Study
34  Princeton, NJ 08540
35  Tel: (609) 734-8387
36  Fax: (609) 951-4438
37  E-mail: mluksza@ias.edu
38
39  Benjamin D. Greenbaum, PhD
40  Tisch Cancer Institute
41  Icahn School of Medicine at Mount Sinai
42  New York, NY 10029
43  Tel: (212) 824-8434
44  E-mail: benjamin.greenbaum@mssm.edu
45

Checkpoint blockade immunotherapies enable the host immune system to recognize and destroy tumor cells[1]. Their clinical activity has been correlated with activated T-cell recognition of neoantigens, which are tumor-specific, mutated peptides presented on the surface of cancer cells[2,3]. How these underlying processes determine the success of immunotherapies has remained unclear. Here, we show that a fitness model for tumors based on immune interactions of neoantigens predicts response to immunotherapy. Two factors determine a neoantigen's fitness cost. First, the cost depends on its presentation by the major histocompatibility complex (MHC), estimated as a function of that neoantigen's relative MHC binding affinity. Second, it depends on T-cell recognition of a neoantigen, which is modeled as a non-linear function of its sequence similarity to known antigens. To describe the evolution of a heterogeneous tumor, we evaluate its fitness as a weighted average over dominant neoantigens in the tumor's subclones. Our model predicts survival in anti-CTLA-4 treated melanoma patients[4,5] and anti-PD-1 treated lung cancer patients[6]. Importantly, low-fitness neoantigens identified by our method may be leveraged for developing novel immunotherapies. By using an immune fitness model to study immunotherapy, we reveal broad evolutionary similarities between cancers and fast-evolving pathogens[7-9].

67 Recent clinical trials using immune checkpoint blocking antibodies, such as anti-
68 cytotoxic T-lymphocyte-associated protein 4 (anti-CTLA-4), or anti-programmed
69 cell death protein-1 (anti-PD-1), have improved overall survival in many
70 malignancies by disinhibiting the immune system[1]. However, only a minority of
71 patients achieves a durable clinical benefit, suggesting there may be genetic
72 determinants of response. De novo somatic mutations within coding regions can
73 create *neoantigens* – novel protein epitopes specific to tumors, which MHC
74 molecules present to the immune system and which may be recognized by T-
75 cells as non-self. An elevated number of mutations or neoantigens has been
76 linked to improved response to checkpoint blockade therapy in multiple
77 malignancies[4-6]. Hence, inferred neoantigen burden is a coarse-grained proxy for
78 whether a tumor is likely to respond to therapy. Other implicated biomarkers of
79 response include T-cell receptor (TCR) repertoire profiles[10], assays of checkpoint
80 status[11,12], immune based microenvironment signatures[4,13], and tumor
81 heterogeneity[14]. Despite high overall mutational load, a heterogeneous tumor
82 may have immunogenic neoantigens present only in certain subclones. As a
83 result, therapies targeting only a fraction of the tumor could disrupt clonal
84 competitive balance and inadvertently stimulate growth of untargeted clones[16,17].
85 Moreover, mass spectrometry-based validation of neoantigens, already limited by
86 sensitivity, does not sample all of the many relevant clones in heterogeneous
87 tumors nor account for clonal variations across metastases[15]. A mathematical
88 model using genomic data has the advantage of broad consideration of
89 neoantigen space. Worldwide efforts are being undertaken to model neoantigens
90 and quantify neoantigen features from genomic data, and a predictive
91 neoantigen-based model for immunotherapy response is therefore a highly
92 sought-after goal.
93
94 We propose a fitness model of immune interactions to describe the evolutionary
95 dynamics of cancer cell populations under checkpoint-blockade immunotherapy
96 (Fig. 1). Fitness models of this kind have been successful in capturing immune
97 interactions for human influenza[7], HIV[8] and chronic viral infections[9], and we aim
98 to introduce this approach to the study of immunotherapy. Checkpoint blockade
99 exposes cancer cells to strong immune pressure on their neoantigens and
100 thereby reduces their reproductive success. Our fitness model, which is detailed
101 below, predicts the evolution of a cancer cell population under such selection
102 pressure. Specifically, we compute $n(\tau)$, the predicted effective size of a cancer
103 cell population in a tumor relative to its effective size at the start of therapy. This
104 effective size is a weighted sum over tumor's genetic clones (Fig. 1a, Methods),
105

$$n(\tau) = \sum_{\alpha} X_{\alpha} \exp(F_{\alpha}\tau) \tag{1}$$

106
107 where $F_{\alpha}$ is the fitness and $X_{\alpha}$ is the initial frequency of clone $\alpha$ and $\tau$ is a
108 characteristic evolutionary time scale (Methods). Our effective size estimates the
109 number of cancer cells required to generate the observed population diversity
110 and is not an estimate of the physical tumor size. Patients with less

111    immunologically fit tumors will have more significant size reductions and, we
112    assume, improved survival prognosis, which is it what we aim to predict. To
113    reconstruct the clonal tree structure of a tumor from exome sequencing data, we
114    use a likelihood scheme based on the allele frequencies of its mutations[18]. Unlike
115    in previous approaches[14], here we learn the ancestral dependencies between
116    clones. These determine the mutations and neoantigens that are inherited by
117    clones from their ancestors (Fig. 1a). Our fitness model assigns to subclones the
118    same or lower fitness than their ancestral clones, depending on whether they
119    acquired new dominant neoantigens.
120

121    Our approach quantifies two essential factors that determine immunogenicity of
122    a neoantigen: an amplitude determined by MHC-presentation, $A$, and the
123    probability of TCR-recognition, $R$ (defined below). We call the product of these
124    two factors, $A \times R$, the *cross-reactivity load* of the neoantigen. Next, we quantify
125    total fitness for cancer cells in a tumor clone by aggregating over the fitness
126    effects due to its neoantigens (Fig 1b, Methods). Specifically, we model the
127    fitness of a given clone $\alpha$ by the cross-reactivity load of the immunodominant
128    neoantigen,
129

$$F_\alpha = - \max_{i \in \text{Clone } \alpha} (A_i \times R_i) \qquad (2)$$

130

131    where index $i$ runs over all neoantigens in clone $\alpha$ (we discuss other choices for
132    aggregating neoantigen fitness effects in Methods). We utilize nonamer
133    neoantigens inferred by a consistent identification pipeline with dissociation
134    constants for both mutant and wildtype peptides for a patient's HLA type[18] (SI).
135

136    We quantify the MHC-presentation factor for a neoantigen using the relative
137    MHC affinity between the wildtype and the mutant peptide. This ratio, which was
138    used to analyze computational neoantigen predictions[20], defines our amplitude $A$
139    (Methods). We show that, unlike considering the mutant or wildtype affinity value
140    alone, the ratio has consistent predictive value within our model (Extended Data
141    Table 1). The interpretation of this model component is consistent with the
142    competitive advantage gained by a neoantigen due to increased concentration,
143    and a neoantigen being less likely to have immune tolerance due to presentation
144    of its closest self-peptide (see discussion in Methods).
145

146    For TCR-recognition, we model cross-reactivity of neoantigens with positive,
147    class I restricted T-cell antigens from the Immune Epitope Database[21] (IEDB).
148    This approach does not assume preexisting immunity due to this set of epitopes.
149    Rather, we posit that neoantigens predicted to be more cross-reactive with
150    a member of this set are more "non-self" and, therefore, more likely to be
151    immunogenic. As cross-reactivity is caused by physical binding of a TCR and a
152    neoantigen, we use an established thermodynamic model to estimate this
153    binding probability from sequence[22]. For a neoantigen with peptide sequence $s$
154    and IEDB epitope with sequence $e$, the alignment score between $s$ and $e$ is used

155  as a proxy for the binding energy between this neoantigen and a TCR specific to
156  epitope **e**. Under this assumption, each mutation that changes a residue in **e** into
157  a corresponding residue in **s** in their alignment will increase the binding energy
158  between **s** and the TCR recognizing epitope **e**, proportionally to the alignment
159  mismatch cost. Importantly, the probability a TCR binds a neoantigen is given by
160  a nonlinear logistic dependence on sequence alignment score (Fig. 2). A similar
161  nonlinear dependence on sequence similarity was previously used to estimate
162  cross-immunity between influenza strains: strains with homologous epitope
163  regions are likely to be antigenically similar[7]. Our model does not require full 9-
164  amino acid identity of the neoantigen and epitope sequences for recognition.
165  The total TCR-recognition probability, $R$, is defined as the probability that
166  neoantigen **s** is recognized by at least one TCR corresponding to an IEDB
167  epitope (Methods).
168
169  We apply the model to three datasets: two melanoma patient cohorts treated with
170  anti-CTLA-4[4,5], and one lung tumor cohort treated with anti-PD-1[6]. We assess
171  our predictions with available patient survival data: total survival times of patients
172  in the melanoma cohorts and progression free survival data on the lung cohort.
173  Neoantigen amino-acid anchor positions, 2 and 9, are constrained due to their
174  molecular function and display a hydrophobic bias, which is also reflected by
175  non-informative MHC affinity amplitudes (Extended Data Fig. 1a). Hence,
176  neoantigens with mutations in these positions are excluded from predictions with
177  our model. Amino-acid diversity in remaining positions is unconstrained
178  (Extended Data Fig. 1b)[23]. Parameter $\tau$, a characteristic evolutionary time scale
179  for a patient cohort, is a finite value at which we expect cancer populations from
180  responding tumors to have been affected by the therapy. This is the time at
181  which, following equation (1), samples are predicted to have a resolved
182  heterogeneity, with their highest fitness clone dominating the evolutionary
183  dynamics. We show that we are able to choose a consistent value of $\tau$ across
184  datasets and that predictions are stable in a broad interval around it (Methods
185  and Extended Data 2). Two model parameters are optimized: the midpoint and
186  the steepness of the logistic binding function, which describes the probability of
187  binding between neoantigens and epitope-specific TCRs (Methods). We
188  maximize the survival log-rank test score to fit the binding curve parameters to
189  the data on the largest dataset, Van Allen et al.[5] (103 metastatic patients). The
190  parameter choice is confirmed to give high log-rank test scores also in the two
191  smaller datasets from Snyder et al., and Rizvi et al., (64 and 34 patients
192  respectively) (Fig. 2 and Extended Data Fig. 3). When using these logistic
193  function parameters in all three datasets, the binding probability of 0.5 is obtained
194  by alignments of average length of 6.55 amino acids; for almost certain binding of
195  probability above 0.95 the average alignment length is 6.98 amino acids.
196
197  The predicted evolutionary dynamics of tumors separates long- and short-term
198  survivors in our datasets (Fig. 3). Long-term survivors (patients with survival time
199  longer than two years in the Van Allen et al. and Snyder et al. datasets, and one
200  year of progression free survival in Rizvi et al. dataset) are predicted to have

201 faster decreasing relative population sizes $n(\tau)$. Moreover, our fitness model
202 results in highly significant separation of patients in survival analysis of all three
203 datasets (Fig. 4). We use the median value of $n(\tau)$ to separate patients into high
204 and low predicted response groups. Using the median as opposed to an
205 optimized threshold[4,5,14] prevents overfitting and allows for robust validation. Log-
206 rank test *p*-values are *p*=0.001 for the Van Allen et al., *p*=0.011 for Snyder et al.,
207 and *p*=7.8e-5 for Rizvi et al. For comparison, a model considering only total
208 neoantigen burden is significant only for Rizvi et al. (*p*=0.007), when also using
209 unsupervised median partitioning of patients. We also use an alternative
210 neoantigen load model that accounts for clonal structure (Methods). Again, only
211 the Rizvi et al. cohort has a significant patient survival separation (*p*=0.0009,
212 Extended Data Table 1).
213
214 The success of our model strongly depends on the joint contribution of two
215 fitness components, the MHC presentation amplitude and TCR recognition
216 probability in equation (2). We deconstruct the model by removing each of
217 the components one at a time (Fig 4, bottom panels and Extended Data Table 1).
218 The MHC-only model, achieved by fixing $R_i = 1$, results in consistently worse
219 segregation of patients (not significant in Snyder et al., decreased significance in
220 Van Allen, et. al, and Rizvi, et. al, *p*=0.027 and *p*=0.004 respectively). The TCR-
221 recognition-only model, achieved by fixing $A_i = 1$, does not result in a significant
222 segregation in any cohort. It is important to assess the clonal structure of a tumor
223 when trying to identify dominant neoantigens. We compare the performance of
224 the full model to one assuming homogenous, single-clone tumor structure, with
225 all neoantigens at tumor frequency = 1 (Methods). This model does not
226 segregate patients significantly in Van Allen et al., and performs worse in Rizvi et
227 al. (*p* =0.019). In Snyder et al., the homogenous structure model gives slightly
228 better separation than the full model (*p*=0.008); however, the score difference
229 between the two is within error bars of the original model.
230
231 In a broader context, our model suggests strong similarities in the evolution of
232 cancers and fast-evolving pathogens. In both systems, immune interactions
233 govern the dynamics of a genetically heterogeneous population; fitness models
234 can predict these dynamics over limited periods, as recently shown for seasonal
235 human influenza[7]. Yet there are important differences between the immune
236 interactions of these systems. Influenza evolution is determined by antigenic
237 similarity with previous strains in the same lineage. Cancer cells originate from
238 normal cells and acquire mutations in a large set of proteins. Hence, their
239 immune interactions are distributed in a larger antigenic space. The fitness
240 effects of these interactions have a specific interpretation: they capture
241 neoantigen "non-selfness"; that is, they formalize aspects of what makes a tumor
242 immunologically different from its host[24]. Thus, our fitness model quantifies
243 the presence of non-self peptides in cancers, which offers insight into adaptive
244 immunity analogous to that for innate recognition of non-self nucleic acids[25].
245
245 Our approach can be naturally extended to other fitness effects, such as positive
246 selection due to acquisition of driver mutations, the impact of other components

247   in the tumor microenvironment or the hypothesized role of the microbiome[26,27,28].
248   Modeling evolutionary dynamics of a cancer cell population can also be useful in
249   studies of acquired resistance to therapy, which is a more distant response
250   effect. The proposed fitness model is based on biophysical interactions
251   underlying the presentation of neoantigens and their immune cross-reactivity.
252   Therefore, besides its predictive function, it may also inform the choice of
253   therapeutic targets for tumor vaccine design. Moreover, this insight may be
254   crucial for understanding when cross-reactivity with self-peptides may result in
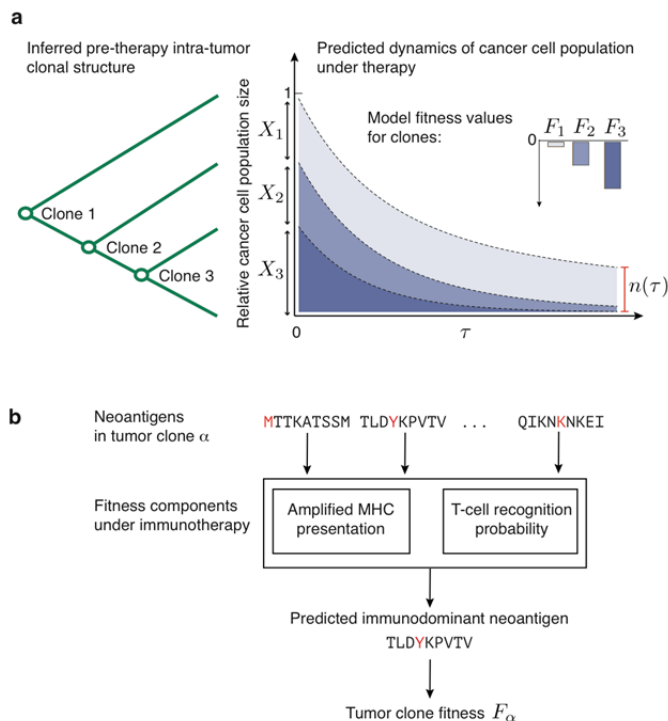255   side effects[29,30].

256

257  **References**
258  1.  Topalian, S.L. et al. Immune checkpoint blockade: a common denominator approach to cancer therapy. *Cancer Cell*
259      **27**, 450-61 (2015).
260  2.  Schumacher, T.N. & Schreiber, R.D. Neoantigens in cancer immunotherapy. *Science* **348**, 69-74 (2015).
261  3.  Gubin, M.M., Artyomov, M.N., Mardis, E.R. & Schreiber, R.D. Tumor neoantigens: building a framework for
262      personalized cancer immunotherapy. *J. Clin. Invest.* **125**, 3413-3421 (2015).
263  4.  Snyder, A. et al. Genetic Basis for Clinical Response to CTLA-4 Blockade in Melanoma. *N. Engl. J. Med.* **371**, 2189-
264      2199 (2014).
265  5.  Van Allen, E.M. et al. Genomic correlates of response to CTLA-4 blockade in metastatic melanoma. *Science* **350**,
266      207-211 (2015)
267  6.  Rizvi, N.A. et al. Mutational landscape determines sensitivity to PD-1 blockade in non–small cell lung cancer.
268      *Science* **348**, 124-128 (2015).
269  7.  Łuksza, M. & Lässig, M. Predictive fitness model for influenza. *Nature* **507**, 57-61 (2014).
270  8.  Wang, S. et al. (2015) Manipulating the selection forces during affinity maturation to generate cross-reactive HIV
271      antibodies. *Cell* **160**, 785–797 (2015).
272  9.  Nourmohammad, A., Otwinowski, J., & Plotkin, J.B. Host-pathogen coevolution and the emergence of broadly
273      neutralizing antibodies in chronic infections. *PLoS Genet* **12**, e1006171 (2016).
274  10. Tumeh, P.C., et al. PD-1 blockade induces responses by inhibiting adaptive immune resistance. *Nature* **515**, 568–
275      571 (2014).
276  11. Topalian, S.L., et al. Safety, activity, and immune correlates of anti-PD-1 antibody in cancer. *N. Engl. J. Med.*, **366**,
277      2443–2454 (2012).
278  12.  Herbst, R.S., et al. Predictive correlates of response to the anti-PD-L1 antibody MPDL3280A in cancer patients.
279      *Nature,* **515**, 563-567 (2014).
280  13. de Henau, O. et al. Overcoming resistance to checkpoint blockade therapy by targeting PI3Kγ in myeloid cells.
281      *Nature* **539**, 443–447 (2016).
282  14. McGranahan, N. et al. Clonal neoantigens elicit T cell immunoreactivity and sensitivity to immune checkpoint
283      blockade. *Science* **351**, 1463-1469 (2016).
284  15. Purcell, A.W., McCluskey, J., & Rossjohn, J. More than one reason to rethink the use of peptides in vaccine design.
285      *Nature Rev. Drug Discov.* **6**, 404-414.
286  16. Fisher, A., Vazquez-Garcia, I., Mustonen V. The value of monitoring to control evolving populations. *Proc. Natl. Acad.
287      Sci.* **112(4)**, 1007-1012 (2015)
288  17. Anagnostu, V. et al. Evolution of neoantigen landscape during immune checkpoint blockade in non-small cell lung
289      cancer. Cancer Discov. ***in press*** (2016).
290  18. Deshwar, A. G. *et al.* PhyloWGS: reconstructing subclonal composition and evolution from whole-genome
291      sequencing of tumors. *Genome Biol.* **16**, 35 (2015).
292  19. Andreatta, M. & Nielsen, M . Gapped sequence alignment using artificial neural networks: application to the MHC
293      class I system. *Bioinformatics* **32**, 511-517 (2016).
294  20. Hundal, J. et al. pVAC-Seq: A genome-guided in silico approach to identifying tumor neoantigens. *Genome* Med. **8**,
295      1-11 (2016).
296  21. Vita, R. et al. The immune epitope database (IEDB) 3.0. *Nucleic Acids Res.* **43**, D405-D412 (2014).
297  22. Berg, O.G. and von Hippel, P.H. Selection of DNA binding sites by regulatory proteins: Statistical-mechanical theory
298      and application to operators and promoters. *J. Mol. Biol.* **193**, 723-743 (1987).
299  23. Lehmann, J., Libchaber, A., & Greenbaum, B.D. Fundamental amino acid mass distributions and entropy costs in
300      proteomes. *J. Theor. Biol.* **410**, 119-124 (2016).
301  24. Old, L.J. & Boyse, E.A. Immunology of experimental tumors. *Ann. Rev. Med.* **15**, 167-186 (1964).
302  25. Tanne, A et al. Distinguishing the immunostimulatory properties of noncoding RNAs expressed in cancer cells. *Proc.
303      Natl. Acad. Sci. USA* **112**, 5154-15159 (2015).
304  26. Vétizou, M. *et al.* Anticancer immunotherapy by CTLA-4 blockade relies on the gut microbiota. *Science* **350**, 1079–
305      1084 (2015).
306  27. Zitvogel, L., Ayyoub, M., Routy, B. & Kroemer, G. Microbiome and Anticancer Immunosurveillance. *Cell* **165**, 276–
307      287 (2016).
308  28. Dubin, K., et al. Intestinal microbiome analyses identify melanoma patients at risk for checkpoint-blockade-induced
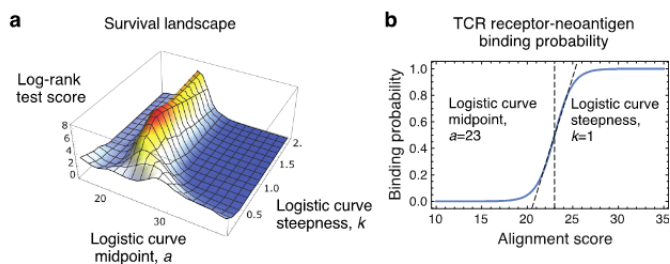309      colitis. *Nat. Commun.*, **7,** 10391 (2016)

310  29. Johnson D.B. et al. Fulminant myocarditis with combination immune checkpoint blockade. *New Engl. J. Med.* **375**,
311      1749-1755 (2016).
312  30. Hofmann, L. et al. Cutaneous, gastrointestinal, hepatic, endocrine, and renal side-effects of anti-PD-1 therapy.
313      *European J. Cancer* **60**, 190-209 (2016).
314

## Figures



**Figure 1 | Evolutionary tumor dynamics under strong immune selection and a neoantigen fitness model based on immune interactions. a**, Clones are inferred from a tumor's phylogentic tree. We predict $n(\tau)$, the future effective size of the cancer cell population, relative to its size at the start of therapy (equation (1)), by evolving clones forward under the fitness model over a fixed time-scale, $\tau$. Application of therapy can decrease fitness of tumor clones depending on their neoantigens. Tumors with strongly negative fitness have a greater loss of population size than more fit tumors. **b**, Our fitness model accounts for the presence of dominant neoantigens within a clone, $\alpha$, by modeling the presentation and recognition of inferred neoantigens and assigning a fitness to a clone, $F_\alpha$.
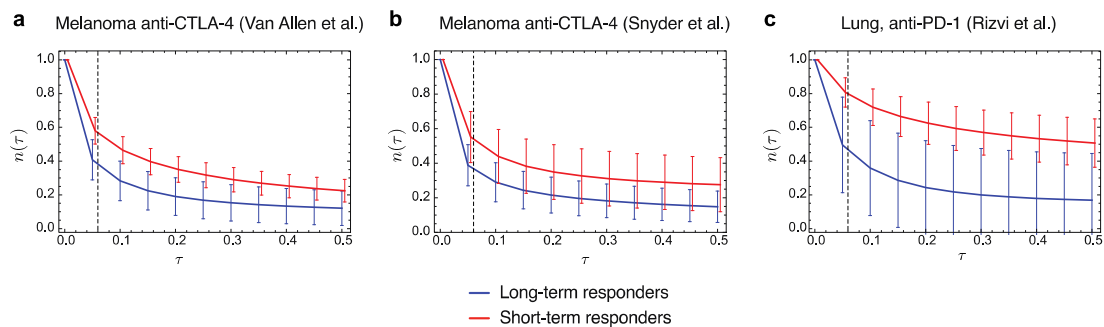
330



331
332 **Figure 2 | Survival landscape as a function of TCR binding model.**
333 **a**, The landscape is a contour plot of log-rank test scores in survival analysis with
334 patient data split by median relative population size (equation (1)). The locally
335 smoothed landscape is plotted for the Van Allen et al. dataset as a function of the
336 model parameters for the logistic curve midpoint ($a$) and steepness ($k$)
337 (Methods). **b**, Logistic binding curve at inferred midpoint and steepness para-
338 meters used across all three datasets from parameters in Van Allen et al.
339 The curve represents the binding probability of a neoantigen and a T-cell
340 receptor specific to a given IEDB epitope as a function of alignment score
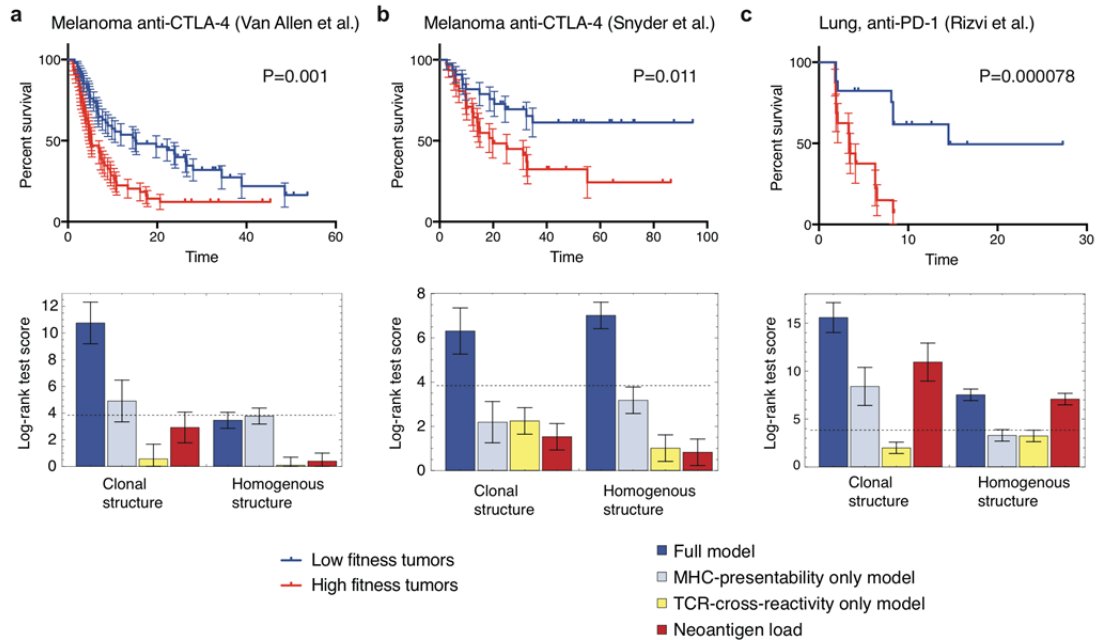341 between the neoantigen's peptide sequence and the epitope.
342

343



344
345

346 **Figure 3 | Evolutionary dynamics predictions in patient cohorts. a,** Relative
347 population size predictions for long-term and short tem survivors across the **a,**
348 Van Allen et al.; **b,** Snyder et al.; and **c,** Rizvi et al. cohorts. Long-term survivors
349 are defined in the text. Error bars are 95% confidence intervals around the
350 population average. The dashed line indicates the consistent choice of $\tau = 0.06$
351 used across all three datasets for patient survival predictions (Methods and
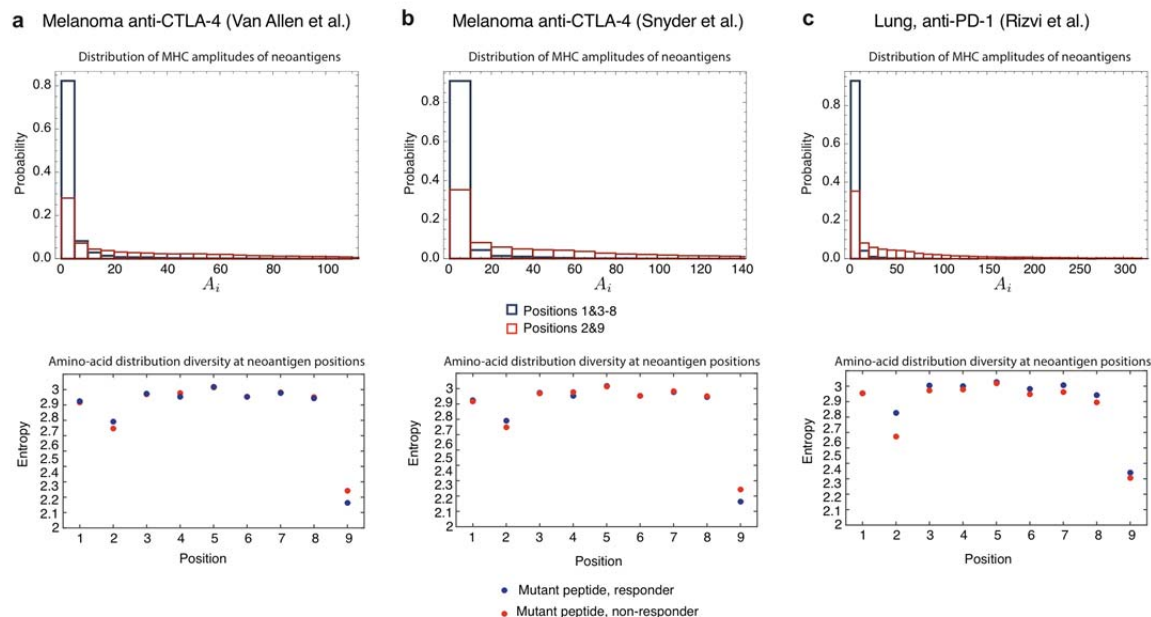352 Extended Data Figure 3).
353

354



355
356
357 **Figure 4 | Neoantigen fitness model is predictive of patient survival after**
358 **checkpoint blockade immunotherapy. a**,**b**, Kaplan-Meier survival curves are
359 calculated across two melanoma patient datasets with patient survival data,
360 which were treated with anti-CTLA-4 antibodies[4,5] and **c,** one dataset of lung
361 patients with progression free survival data, which were treated with anti-PD-1
362 antibodies[6]. The samples are split in an unsupervised manner by the median
363 value of their tumor's relative population size $n(\tau)$ defined in equation (1); the
364 error bars represent the standard error. For comparison we show the log-rank
365 test score for models, which account for removal of one feature of our model
366 (bottom panels, higher score values indicate better patient segregation): an
367 MHC-presentability only model (light blue) and a TCR-recognition only model
368 (yellow). We compare their values with a tumors' neoantigen burden (red). All
369 models are computed both over a tumor's clonal structure (clonal, left) and
370 without taking heterogeneity into account (homogenous, right). Dashed lines on
371 the bottom panels marks the score value corresponding to the significance
372 threshold of 5%, scores above that threshold raise significant patient
373 segregation. The error bars are the standard deviation of log-rank test score
374 acquired from the survival analysis with one sample removed from the cohort at a
375 time.
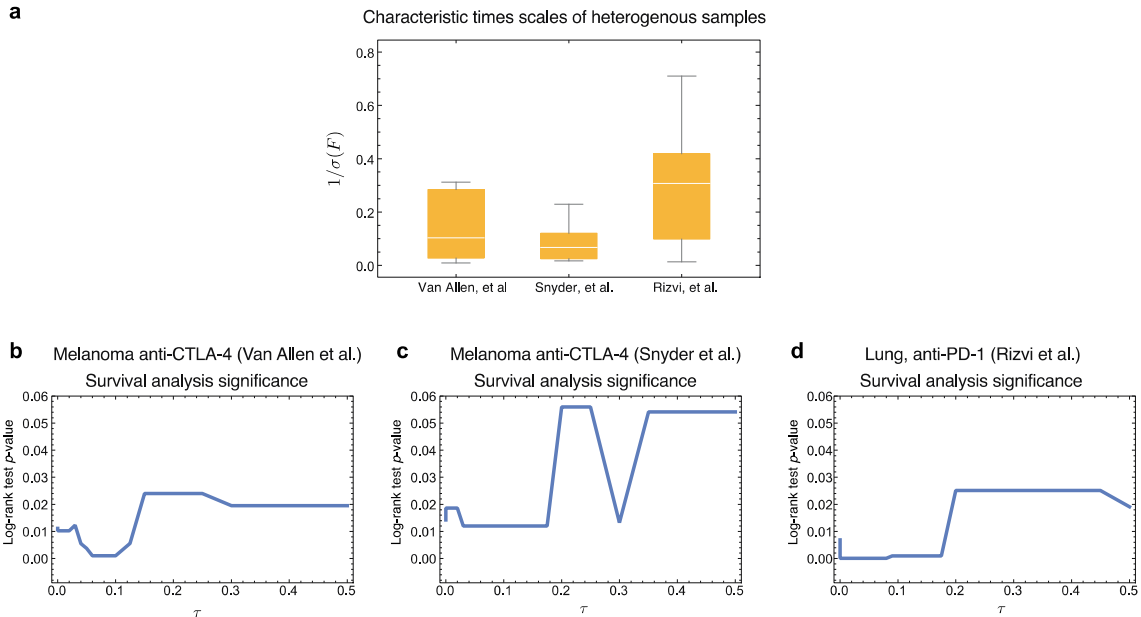376

377 **Extended Data Figures**



378
379 **Extended Data Figure 1 | Positions 2 and 9 in neoantigens are of less**
380 **predictive value. a,** Neoantigens with mutations at anchor residues at position 2
381 and 9 have highly diverging amplitude values and are of less overall predictive
382 value than neoantigens at other positions. **b,** Patients classified in studies as
383 responders are marked in blue and non-responders are marked in red. Positions
384 2 and 9 are highly constrained by a bias to be hydrophobic. Their Shannon
385 entropy is lower than that of other residues, across all three datasets regardless
386 of classification of their neoantigens in those datasets. Other residue sites have
387 the same entropy as the overall proteome[22] and are therefore unconstrained.
388

389



**a** Characteristic times scales of heterogenous samples

**b** Melanoma anti-CTLA-4 (Van Allen et al.)
Survival analysis significance

**c** Melanoma anti-CTLA-4 (Snyder et al.)
Survival analysis significance

**d** Lung, anti-PD-1 (Rizvi et al.)
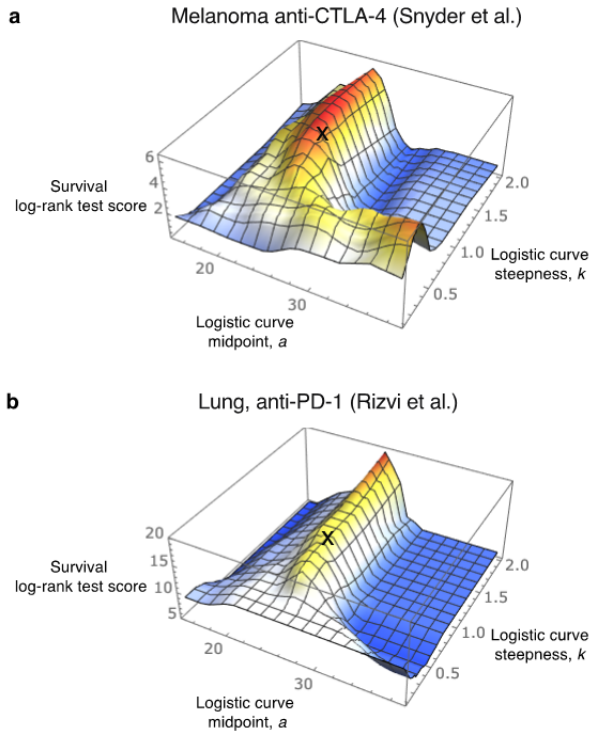Survival analysis significance

390

391 **Extended Data Figure 2 | Consistency of evolutionary time-scale across**
392 **datasets. a,** The distribution of characteristic times scales of samples with clonal
393 fitness heterogeneity (Methods) for the three patient cohorts. These distributions
394 consistently define the interval for relevant time scales of $\tau$, in all datasets we
395 subsequently investigate $\tau \in [0,0.5]$ . **b-d,** Significance of survival analysis
396 reported as the result of the log-rank test on the three datasets with sample split
397 at a median value $n(\tau)$ plotted as a function of $\tau$. The chosen value of parameter
398 $\tau =0.06$ and a broad surrounding interval gives highly significant sample
399 segregation in each of the datasets.
400

401



**a** Melanoma anti-CTLA-4 (Snyder et al.)

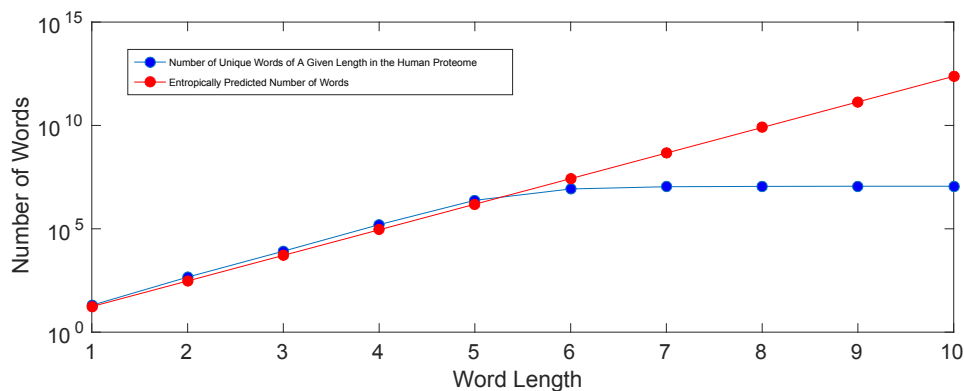**b** Lung, anti-PD-1 (Rizvi et al.)

402
403
404   **Extended Data Figure 3 | Survival landscape for Snyder et al., and Rizvi et
405   al., cohorts.** The survival landscape is defined by the log-rank test score as a
406   function of the model parameters for the logistic curve shape, i.e the midpoint ($a$)
407   and steepness ($k$) (Methods). The locally smoothed landscape is plotted for the
408   **a,** Snyder et al., and **b,** Rizvi et al., datasets. An X marks the optimal parameters
409   from Van Allen et al., $a$ =23 and $k$ =1 (cf. Fig. 2), which are used to derive
410   survival curves for these two datasets and are at high score regions of the
411   landscapes.
412

413



414
415
416
**Extended Data Figure 4 | Word usage in the proteome is exhausted between 5 and 6 letter words.** Given the entropy of the genome from Ref. 23, we calculate the expected number of words of a given length in the proteome as a function of word length. We compare that to the number of unique words in the proteome of a given length. Between 5 and 6 letters the two curves diverge due to the finite size of the genome. By the time one reaches 9 letter nonamers (the length of a neoantigen) this divergence is of several orders of magnitude.
424

425 **Extended Data Table 1 | Ranking of fitness models.** We compare survival
426 prediction of our full fitness model (Methods, equation (9)) with alternative
427 models described in Methods: **(1)** models that eliminate one of the features of the
428 full model, namely the MHC-presentability only model (Methods, equation (13)) a
429 nd a TCR-recognition only model (Methods, equation (14)); absolute MHC-
430 amplitude model and absolute wildtype MHC-amplitude model (Methods,
431 equations (15) and (16) respectively); simple neoantigen load model and
432 mutational load model (Methods, equations (17) and (18)); and finally an additive
433 neoantigen fitness model (Methods, equation (19)), which summates fitness
434 contributions of neoantigens in a clone as opposed to maximizing them as in our
435 original model. **(2)** Above models evaluated without accounting for clonal
436 structure structure of tumors. For each model we report the following parameters
437 (if applicable): the aggregating function for neoantigen effects within a clone
438 (MAX or SUM replacing $Ag$ in equation (11)), the value of parameter $\tau$ used in
439 predictions, the parameters of the logistic function $a$ and $k$ (Methods, equation
440 (7)). Finally, we report the predictive value of the models as the log-rank test *p*-
441 value and the corresponding log-rank test score. The comparison is shown on all
442 three immunotherapy datasets.
443

463

## Author Contributions
465 M.Ł. and B.D.G. designed the mathematical model and wrote the manuscript with
466 critical comments from all the authors. M.Ł., N.R., V.M., V.P.B., A.S., N.A.R.,
467 T.M., A.J.L., T.A.C., J.D.W., and B.D.G contributed to data acquisition and
468 analysis. M.Ł., T.A.C., J.D.W., and B.D.G. contributed to study conception and
469 design. M.Ł., N.R., V.M., V.P.B., A.S., N.A.R., T.M., A.J.L., T.A.C., J.D.W., and
470 B.D.G. interpreted the data and provided a critical reading of the manuscript.

471

## Methods

### 1. Evolutionary dynamics of a cancer cell population in a tumor

The fitness of a cancer cell in a genetic clone $\alpha$ is its expected replication rate, i.e.

$$\frac{dN_\alpha}{d\tau} = F_\alpha N_\alpha \tag{3}$$

where $N_\alpha$ is the population size of clone $\alpha$. Checkpoint-blockade immunotherapy introduces a strong selection challenge, which we anticipate overshadows pre-therapy fitness effects in a productive response. For a given clone $\alpha$ the dynamics of its absolute size are hence given by $N_\alpha(\tau) = N_\alpha(0)\exp(F_\alpha\tau)$, and the total cancer cell population size is computed as a sum over its clones

$$N(\tau) = \sum_\alpha N_\alpha(\tau) = \sum_\alpha N_\alpha(0)\exp(F_\alpha\tau). \tag{4}$$

The absolute size $N(\tau)$ is meant as the effective population size, the number of cells estimated to have generated the observed clonal diversity; it is not to be understood as the physical tumor size. As our diagnostic of survival we use the relative effective population size $n(\tau) = N(\tau)/N(0)$, which compares the predicted evolved population size after a characteristic time scale of evolution $\tau$ (discussed below) to the initial pretreatment effective size $N(0)$. We denote the initial clone $\alpha$ frequency $X_\alpha = N_\alpha(0)/N(0)$, these frequencies are inferred from bulk exome reads from a tumor sample[17]. Hence, to compute $n(\tau)$ we only require estimates of the initial frequencies and fitness values for each clone, as shown in equation (1); the absolute population size estimates are not needed.

**Clonal structure of a tumor and clone frequencies.** Tumor clones are reconstructed using the PhyloWGS software package[18] (SI). The trees estimate the nested clonal structure of the tumor and the frequency of each clone, $X_\alpha$. The differences between the high scoring trees are marginal on our data, concerning only peripheral clones and small differences in frequency estimates. We compute the predicted relative size of a cancer population $n(\tau)$ as an averaged prediction over 10 trees with the highest likelihood score.

### 2. Fitness model

**MHC-amplitude.** The amplitude due to the dissociation constant between a neoantigen and its wildtype peptide is defined as

$$A = K_D^{WT}/K_D^{MT}. \tag{5}$$

The dissociation constants are inferred for each peptide sequence and patient HLA type[19]; all mutant peptide sequences considered as neaontigens meet the

512  standard cutoff, $K_D{}^{MT} < 500$ nM (SI). The amplitude in this form has a high
513  predictive value for patient survival predictions (discussed in section 4.,
514  demonstrated in Fig. 4 and Extended Table 1), consistently over the three patient
515  cohorts, which is not the case of either the mutant or wildtype dissociation
516  constants on their own.

517

518  We offer two interpretations of why this amplitude is relevant, which are not
519  mutually exclusive of one another. The first is that the amplitude can be thought
520  of as an approximate form derived with the use of simple equilibrium kinetics,
521  where the concentration of peptide bound to MHC is given by their individual
522  concentrations and inferred binding constant $K_D$, derived from NetMHC[19]. The
523  underlying dependencies are

524

$$A = \frac{[\text{MHC: neoantigen}]^{MT}}{[\text{MHC: neoantigen}]^{WT}} = \frac{K_D^{WT}[\text{MHC}]^{MT}[\text{neoantigen}]^{MT}}{K_D^{MT}[\text{MHC}]^{WT}[\text{neoantigen}]^{WT}}, \qquad (6)$$

525

526  where $[\text{MHC: neoantigen}]^{MT}$ is the concentration of the mutant form of the
527  neoantigen to MHC, with the $WT$ superscript representing the same quantity for
528  the wild-type peptide. This interpretation assumes the above quantity is
529  dominated by the ratio of dissociation constants, which derives the formula for $A$
530  in equation (5). In this sense, the amplitude reflects the relative concentration of
531  mutant to wildtype peptide and therefore the likelihood that the mutant peptide
532  would be presented versus its wildtype peptide. As such it may reflect the
533  competitive advantage a neoantigen has acquired in terms of presentation
534  through mutation, as posited in other in silico analyses[31].

535

536  The second interpretation is that the amplitude reflects the likelihood a
537  neoantigen is similar to a peptide that has undergone immune tolerance. As we
538  exclude neoantigens with mutations on positions 2 and 9, a high value of
539  amplitude means the wildtype peptide is also likely to have hydrophobic residues
540  at the anchor position and hence can be presented by the MHC. Since
541  neoantigens differ from their wildtype peptides by a single mutation, and given
542  the uniqueness of nonamer sequences in the proteome (Extended Data Fig. 4),
543  the self-nonamer in the genome with the greatest similarity to a neoantigen is
544  likely to be its wildtype peptide. We verified that this is the case for 92% of all
545  neoantigens, with the remainder largely emanating from gene families with many
546  paralogs (SI). Therefore a high amplitude usually stands for the self peptide most
547  similar to a neoantigen not being likely to have been abundantly presented by the
548  MHC. Following this reasoning, the mutant peptide with high affinity is likely to be
549  novel to T-Cells as its immunogenicity is not mitigated by a homologous self-
550  peptide.

551

552  **TCR-recognition.** We model $R$, the cross-reactivity of a neoantigen with a TCR-
553  pool defined as the probability that a neoantigen cross-reacts with at least one
554  TCR corresponding to a known immunogenic epitope. We profile *in silico* the
555  cross-reactivity of neoantigen with a set of epitopes given by the Immune Epitope

556     Database and Analysis Resource[21] (IEDB). We restrict ourselves to IEDB
557     epitopes that are positively recognized by T-cells after class I MHC presentation.
558     We hypothesize that a neoantigen that is predicted to cross-react with a TCR
559     from this pool of immunogenic epitopes is a neoantigen more likely to be
560     immunogenic itself.
561
562     The probability that a TCR for a given epitope binds a given neoantigen is
563     defined by a simple two-state thermodynamic model with logistic shape. In this
564     model we use sequence alignment as a proxy for binding energy[22]. To assess
565     sequence similarity between a neoantigen with peptide sequence $\mathbf{s}$ and an IEDB
566     epitope $\mathbf{e}$, we compute a gapless alignment between the two sequences with a
567     BLOSUM62 amino-acid similarity matrix[32]. For an alignment score, $|\mathbf{s}, \mathbf{e}|$, we
568     compute the binding probability as
569

$$\mathrm{Pr}_{\mathrm{binding}}(\mathbf{s}, \mathbf{e}) = \frac{1}{1 + e^{-k(|\mathbf{s},\mathbf{e}|-a)}}, \tag{7}$$

570
571     where $a$ represents the horizontal displacement of the binding curve and $k$ sets
572     the steepness of the curve at $a$. These are two free parameters to be fit in our
573     model (see below). The parameters that we use in predictions are $a$=23 and $k$=1;
574     these parameters give binding probability $\mathrm{Pr}_{\mathrm{binding}}(\mathbf{s}, \mathbf{e})$ =0.5 at alignment score
575     $|\mathbf{s}, \mathbf{e}|$=23; the probability drops to below 0.05 at $|\mathbf{s}, \mathbf{e}|$=20 and reaches value of
576     above 0.95 at $|\mathbf{s}, \mathbf{e}|$=26 (Fig. 2b). The corresponding alignment score span of 6 is
577     close to the average identity match score in the BLOSUM62 matrix (5.64). The
578     average alignment length corresponding to score 26 is 6.98 amino acids in our
579     datasets and it is 6.55 for binding probability 0.5. The logistic function is therefore
580     a strongly nonlinear function of the alignment score, where a mismatch on 1-2
581     positions can decide about lack of binding between the neoantigen and the
582     epitope specific TCR.
583
584     For a given neoantigen $\mathbf{s}$ we calculate the probability it is recognized by a TCR
585     within a repertoire as the probability it cross-reacts with at least one IEDB
586     epitope:
587

$$R = 1 - \prod_{\mathbf{e} \in \mathrm{IEDB}} [1 - \mathrm{Pr}_{\mathrm{binding}}(\mathbf{s}, \mathbf{e})] \tag{8}$$

588
589     **Neoantigen-based fitness cost for a tumor clone.** Our model associates each
590     neoantigen with a fitness cost, the *cross-reactivity load*, defined as the product of
591     the MHC-amplitude in equation (5) and TCR-recognition probability in equation
592     (8), $A \times R$.
593     To assess the total fitness effect for a clone $\alpha$ with multiple neoantigens, we
594     aggregate the individual neoantigen fitness effects as $F_{\alpha} = -\max_{i \in \mathrm{Clone}\ \alpha}(A_i \times$
595     $R_i)$, where $i$ is an index running over neoantigens in the clone. Therefore, the full
596     form of the predicted relative cancer cell population size is given by
597

$$n(\tau) = \sum_{\alpha} X_{\alpha} \exp[- \max_{i \in \text{Clone } \alpha} (A_i \times R_i) \, \tau]. \tag{9}$$

598
599    One could use a more general model for fitness model of a clone,
600

$$F_{\alpha} = - \underset{i \in \text{Clone } \alpha}{\text{Ag}} (A_i \times R_i) \tag{10}$$

601    and use different function $\text{Ag}$ to aggregate over cross-reactivity fitness effects of
602    neoantigens within a clone, such as a summation over all neoantigens (Extended
603    Data Table 1), summation over a fixed set, or other nonlinear dependency.
604    Taking the best score within a clone is consistent with the notions of
605    heterologous immunity and immunodominance – that a small set of antigens
606    drive the immune response, whereas summing over neoantigens would imply a
607    more uniform distribution of contributions.
608
609    **3. Model parameters**
610
611    **Logistic binding function parameter optimization.**    To choose model
612    parameters $a$ and $k$ in equation (7) we investigate the log-rank-test scores of
613    patient segregation as a function of these parameters. The survival analysis is
614    performed by splitting patient cohort into *high* and *low fitness* groups by the
615    median cohort value of $n(\tau)$, the predicted relative cancer cell population size at
616    a characteristic time $\tau$ (we discuss the choice of $\tau$ below). The survival score
617    landscapes (Fig. 2 and Extended Data Fig. 3) appear to be consistent between
618    the datasets, with an optimal value of parameter $a$ around 23 and parameter $k$
619    living on a trivial axis above value 1, suggesting strong nonlinear fitness
620    dependence on the sequence alignment score. We choose parameters that
621    optimize the log-rank-test score in the largest dataset in our study, the melanoma
622    anti-CTLA4 cohort from Van Allen, et al[5].
623
624    **Characteristic time scale parameter estimation.** In the survival analysis the
625    samples are split by the median cohort value $n(\tau)$ at a specified time scale $\tau$.
626    Intuitively, this time should be set to a finite value at which the tumors are
627    expected to have responded to therapy. At this value of $\tau$ the clonal
628    heterogeneity of tumors is supposed to have decreased, with the highest fitness
629    clone dominating in the population. For one tumor this time scale is inversely
630    proportional to the standard deviation of intra-tumor fitness (i.e. of the order of
631    $1/\sigma(F)$), where
632

$$\sigma^2(F) = \sum_{\alpha} X_{\alpha} F_{\alpha}^2 - \left( \sum_{\alpha} X_{\alpha} F_{\alpha} \right)^2. \tag{11}$$

633
634    In each cohort we determined the interval of characteristic times of heterogenous
635    samples (Extended Data Fig. 2a) and we tested the dependence of prediction
636    power on $\tau$ by performing log-rank test (Extended Fig. 2b-d). The optimal values

637 of $\tau$ in each cohort belong to a relatively wide interval. The consistent broadness
638 of these intervals suggests low sensitivity of predictive power on $\tau$. Moreover, the
639 parameter intervals giving highly significant patient segregation are also
640 consistent between the cohorts. We choose $\tau = 0.06$ for our predictions in all
641 datasets. As $\tau$ is an inverse fitness it also defines a typical maximum cross-
642 reactivity load in a clone beyond which one would expect to have a clone that
643 responded to therapy. For instance, at $\tau = 0.06$ this typical fitness value would be
644 about 16.67. This would indicate that a neoantigen with a TCR recognition
645 probability $R = 1$ would on average lead to a productive response if the ratio of its
646 dissociation constants would be greater 16.67. Well beyond that value
647 amplitudes would essentially carry the same predictive value.
648
649 Heterogenous samples were selected with criterion $e^{H_F} \geq 2$, where $H_F$ is clonal
650 fitness entropy defined as
651

$$H_F = -\sum_{\beta} Y_\beta \log Y_\beta, \tag{12}$$

652
653 where the frequencies of clones with the same fitness are added together and
654 denoted as $Y_\beta$. The index $\beta$ then refers to all clones with a given fitness.
655
656 **4. Alternative fitness models**
657
658 We compare our full model in equation (9) to the following alternative models
659 (Extended
660  Data Table 1):
661
662    1. **Heterogenous structure models**
663
664       a. **MHC-presentability only model:**
665          In this model the recognition factor is ignored and fitness is
666          assumed to be determined only by MHC-amplitude of neoantigens.
667          The defining equation is given by

$$n(\tau) = \sum_{\alpha} X_\alpha \exp[-\max_{i \in \text{Clone } \alpha} A_i \tau]. \tag{13}$$

668
669       b. **TCR-recognition only model:**
670          Conversely, in this model the MHC-presentation factor is ignored
671          and fitness is assumed to be determined only by TCR-recognition
672          of neoantigens. The defining equation is given by
673

$$n(\tau) = \sum_{\alpha} X_\alpha \exp[-\max_{i \in \text{Clone } \alpha} R_i \tau]. \tag{14}$$

674

675 **c. Absolute MHC-amplitude model**
676 In this model the likelihood of MHC presentation for a neoantigen is
677 inversely correlated with its inferred dissociation constant,
678 $A^{\mathrm{abs}} = 1/K_D^{MT}$ (cf. equation (5)). The model is defined as
679

$$n(\tau) = \sum_\alpha X_\alpha \exp[-\max_{i \in \mathrm{Clone}\,\alpha} \left(A_i^{\mathrm{abs}} \times R_i\right)\tau]. \tag{15}$$

680 **d. Absolute MHC-amplitude model**
681 In this model the likelihood of MHC presentation for a neoantigen is
682 inversely correlated with its inferred dissociation constant,
683 $A^{\mathrm{abs,WT}} = K_D^{WT}$ (cf. equation (5)). The model is defined as
684

$$n(\tau) = \sum_\alpha X_\alpha \exp[-\max_{i \in \mathrm{Clone}\,\alpha} \left(A_i^{\mathrm{abs,WT}} \times R_i\right)\tau]. \tag{16}$$

685 **e. Neoantigen load model**
686 This model assigns uniform fitness cost to each neoantigen. For $L_\alpha$,
687 the number of neoantigens in clone α, this model is defined by
688

$$n(\tau) = \sum_\alpha X_\alpha \exp[-L_\alpha\tau]. \tag{17}$$

689 We do not exclude neoantigens with mutations on positions 2 and 9
690 in the neoantigen load model.
691
692 **f. Mutational load model**
693 This model assigns uniform fitness cost to each somatic mutations.
694 For, $L_\alpha^M$, the number of somatic mutations (with respect to a normal
695 cell) in clone α, this model is defined by
696

$$n(\tau) = \sum_\alpha X_\alpha \exp[-L_\alpha^M\tau]. \tag{18}$$

697
698 **g. Additive neoantigen fitness model**
699 This model implements an additive neoantigen aggregating
700 function, namely
701

$$n(\tau) = \sum_\alpha X_\alpha \exp\left[-\left(\sum_{i \in \mathrm{Clone}\,\alpha} A_i \times R_i\right)\tau\right]. \tag{19}$$

702
703 2. **Homogenous structure models**

24

704         For each model defined in point (1) we can define its homogenous
705         structure equivalent, which assumes tumor is strictly clonal with all
706         neoantigens in the same clone at frequency 1.
707
708 We assess the predictive power of these models with a survival analysis, by
709 separating patients by the median value of $n(\tau)$ in each patient cohort and
710 computing the log-rank test for such segregation. For stringency of comparisons,
711 we adjust the value of parameter $\tau$ in a supervised manner to optimize the
712 performance of each alternative model (Extended Data Table 1).
713
714 **5. Data availability**
715
716 Mutation data and inferred neoantigen peptide data for each dataset are
717 submitted as supplementary data.
718
719 **References**
720

721 31.     Khalili, J.S., Hanson, R.W., & Szallasi, Z. In silico prediction of tumor antigens derived from functional missense
722           mutations of the cancer gene census. *Oncoimmunology* **1**,1281-1289 (2012).
723 32.     Henikoff, S. & Henikoff, J.G. Amino acid substitution matrices from protein blocks. *Proc. Natl. Acad. Sci. USA* **89**,
724           10915–10919 (1992).
725

726 **Supplementary Information**
727
728 **Computational identification of neoantigens**
729
730 Neoantigens from the three datasets were inferred using a consistent pipeline
731 established at Memorial Sloan Kettering Cancer Center. Raw sequence data
732 reads were aligned to the reference human genome (hg19) using the Burrows-
733 Wheeler Alignment tool. Base-quality score recalibration, and duplicate-read
734 removal were performed, with exclusion of germline variants, annotation of
735 mutations, and indels as previously described[4]. Local realignment and quality
736 score recalibration were conducted using the Genome Analysis Toolkit (GATK)
737 according to GATK best practices[33,34]. For sequence alignment and mutation
738 identification, the FASTQ files were processed to remove any adapter sequences
739 at the end of the reads using cutadapt (v1.6)[35]. The files were then mapped using
740 the BWA mapper (bwa mem v0.7.12[36], the SAM files sorted, and read group tags
741 added using the PICARD tools. After sorting in coordinate order, the BAM's were
742 processed with PICARD MarkDuplicates. First realignment was carried out using
743 the InDel realigner followed by base quality value recalibration with the Base-
744 QRecalibrator.
745
746 A combination of 4 different mutation callers (Mutect 1.1.4, Somatic Sniper 1.0.4,
747 Varscan 2.3.7, and Strelka 1.013) were used to identify single nucleotide variants
748 (SNVs)[37-39]. As previously described, SNVs with an allele read count of less than
749 4 or with corresponding normal coverage of less than 7 reads were filtered out[40].
750
751 The assignment of a somatic mutation to a neoantigen was estimated using a
752 previously described bioinformatics tool called NASeek[4]. Briefly, NASeek is a
753 computational algorithm that first translates all mutations in exomes to strings of
754 17 amino acids, for both the wild type and mutated sequences, with the amino
755 acid resulting from the mu-tation centrally situated. Secondly, it evaluates
756 putative MHC Class I binding for both wild type and mutant nonamers using a
757 sliding window method using NetMHC3.4[19]
758 (http://www.cbs.dtu.dk/services/NetMHC-3.4/) for patient-specific HLA types, to
759 gene-rate predicted binding affinities for both peptides. NASeek finally assesses
760 for similarity between nonamers that predicted to be presented by patient-specific
761 MHC Class I. All nonamers with binding scores (i.e. the inferred dissociation
762 constants $K_D^{MT}$) below 500 nM are defined as neoantigens.
763
764 **Clonal tree reconstruction with PhylowWGS**
765
766 Tumor clones are reconstructed using the PhyloWGS software package[18]. The
767 input data for the algorithm is extracted from exome sequencing data: (1)
768 mutation reads obtained with the pipeline described above, and (2) allele-specific
769 copy-number variant data, obtained with FACETS v0.5.0[41]. Briefly, the package
770 clusters mutations into clones by the frequency of their reads and it infers
771 possible nesting of clones (ancestral relations) between pairs of clones.

26

772 Intuitively, an ancestral clone needs to have higher frequency then its derived
773 clone. From this information PhyloWGS reconstructs high likelihood tumor
774 geneological trees.
775
776 **Amino acid diversity**
777
778 We define the amino acid diversity at $i$-th position in a neoantigen as $e^{H_i}$, where
779 $H_i$ is Shannon entropy[42] of amino acid usage at this position, i.e.
780
781 $$H_i = - \sum_{j=1}^{20} f(a_{ij}) \log(f(a_{ij})),$$
782
783 where $f(a_{ij})$ is frequency of the $i$-th position in all neoantigens in a group.
784 Inferred neoantigens are nonamers, so $i$ ranges in value from 1 to 9. The
785 diversity of neoantigens at a given site were compared to the values found in the
786 human proteome in Lehman, et al.[23].
787
788 To calculate the expected number of words in the proteome we utilize the
789 frequency of amino acids from Lehman, et al. We compute the entropy
790 associated with the frequency of amino acids in the human genome:
791
792 $$H(a) = - \sum_{j=1}^{20} f(a_j) \log(f(a_j)),$$
793
794 where $f(a_j)$ is the frequency of the $j$-th amino acid in the human genome. The
795 expected number of words of length $n$ is therefore $e^{nH(a)}$. This value is compared
796 to the observed number of words of length $n$ in the reference proteome for
797 GRCh38.p7
798
799 **Identification of closest nonamers in human proteome to neoantigens**
800
801 We have mapped the WT and MT 9-mer peptides to all proteins in the current
802 human reference genome (GRCh38.p7) with at least 8 out of 9 matches and no
803 gaps (allowing only mismatches). For this we used LAST[43] (version 819) with the
804 following parameters:
805 lastal -f BlastTab -j1 -r2 -q1 -e15 -y2 -m100000000 -l4 -L4 -P0
806 (9-mer mapping with at most one mismatch is guaranteed to have a matching 4-
807 mer word).
808
809 One expects the mutated peptide to only map to the same location as the WT
810 peptide, WT mapping exactly (9 matches) and MT mapping with one mismatch (8
811 matches). The expected case is that the WT peptide maps to the proteome
812 exactly and the MT peptide maps to the proteome with one mismatch and only to
813 the loci WT peptide maps to.
814
815 This rule can be violated in the following cases, sorted from the most to the least
816 severe:

817 1. WT peptide does not map to the proteome exactly. Some possible reasons
818 are:        a difference in the reference assemblies used for mutation calling and
819 peptide mapping, a germline mutation mistakenly identified as somatic, or a
820 difference between the pa-tient genome and the reference genome used for
821 alignments.
822 2. WT peptide maps to the proteome exactly (9 matches), MT peptide maps to
823 the pro-teome exactly (9 matches) but to a different locus.
824 3. WT peptide maps to the proteome exactly, MT peptide maps to the proteome
825 with one mismatch; however, MT peptide maps with one mismatch to the
826 subjects WT does not map exactly.
827 4. WT peptide maps to the proteome exactly, MT peptide maps to the proteome
828 with one mismatch; however, MT peptide maps with one mismatch to a different
829 locus on    the gene WT maps to.
830
831 We have examined each peptide for the worst possible scenario. We have gone
832 from category 1 to 4 in the list. Category 1 indicates a difference in the reference
833 genome. Categories 2-4 typically are due to mutations that occur in repetitive
834 gene families with many paralogs. Once we identified that a peptide belongs to
835 any category, we excluded it from further considerations. This way the numbers
836 of peptides in each category add up to the total number of peptides. Below is a
837 summary for the different datasets utilized in this study:
838
839 Van Allen, et al.[5]:
840 39373 total peptides, (1) 42 WT unmapped, leaving 39331
841 36783 expected peptides (93.42%), (2) 387 have 9 matches in MT, (3) 2076
842 have other alignments, (4) 85 have other alignments to the same subject.
843
844 Snyder, et al[4].:
845 29781 total peptides, (1) 35 WT unmapped, leaving 29746
846 27674 expected peptides (92.93%), (2) 361 have 9 matches in MT, (3) 1644
847 have other alignments, (4) 67 have other alignments to the same subject.
848
849 Rizvi, et al.[6]:
850 5581 total peptides, (1) 6 WT unmapped, leaving 5575
851 5125 expected peptides (91.83%), (2) 105 have 9 matches in MT, (3) 323 have
852 other alignments, (4) 22 have other alignments to the same subject.
853
854 Additional supplementary files for each dataset are included as Supplementary
855 Data:
856
857 mt-with-9.tsv – list of peptides from category 2 and the subjects each one aligns
858 to .
859
860 peptides-with-extra-aln.tsv – peptides from group 3 and the subjects each one
861 aligns to.
862

863 peptides-multimapping-same-subj.tsv – peptides from group 4 and their
864 alignments including the start and end coordinates
865
866 **Additional References**
867

868 33.   DePristo, M., et al. A framework for variation discovery and genotyping using next-generation DNA sequencing
869       data. *Nature Genet.* **43**, 491-498 (2011).
870 34.   Van der Auwera, G.A., et al. From FastQ Data to High-Confidence Variant Calls: The Genome Analysis Toolkit
871       Best Practices Pipeline. *Curr. Prot. in Bioinformatics* **43**, 11.10.1-11.10.33 (2013).
872 35.   Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet.journal* **17**, 10-
873       12 (2011).
874 36.   Li, H., & Durbin R. Fast and accurate short read alignment with Burrows-Wheeler Transform. *Bioinformatics* **25**,
875       1754-1760 (2009).
876 37.   Wei, L. et al. MAC: identifying and correcting annotation for multi-nucleotide variations. *BMC Genomics* **16**, 569
877       (2015).
878 38.   Snyder, A. & Chan, T.A. Immunogenic peptide discovery in cancer genomes. *Curr Opin Genet Dev* **30**, 7-16
879       (2015).
880 39.   Nielsen, M. et al. Reliable prediction of T-cell epitopes using neural networks with novel sequence
881       representations. *Protein Sci* **12**, 1007-1017 (2003).
882 40.   Riaz, N. et al. Recurrent SERPINB3 and SERPINB4 mutations in patients who respond to anti-CTLA4
883       immunotherapy. *Nat. Genet.* **48**, 1327-1329 (2016).
884 41.   Shen, R. & Seshan, V.E. FACETS: allele-specific copy number and clonal heterogeneity analysis tool for high-
885       throughput DNA sequencing. *Nucleic Acids Res.* **44**, e131 (2016).
886 42.   Shannon, C.E. A mathematical theory of communication. *Bell Sys. Tech. J.* **27**, 379-423 (1948).
887 43.   Kiełbasa, S.M., Wan, R., Sato, K., Horton, P., & Frith, M.C. Adaptive seeds tame genomic sequence comparison.
888       *Genome Res.* **21**, 487-493 (2011).
889
890