1  **Measuring the distribution of fitness effects in somatic evolution by combining clonal**
2  **dynamics with dN/dS ratios**
3
4  Marc J Williams[1], Luiz Zapata[2], Benjamin Werner[2], Chris Barnes[3], Andrea Sottoriva[2]*, Trevor
5  A Graham[1]*
6
7  [1] Barts Cancer Institute, Barts and the London School of Medicine and Dentistry, Queen Mary
8  University of London, UK
9  [2] Centre for Evolution and Cancer, Institute of Cancer Research, London, UK
10  [3] Department of Cell and Developmental Biology, University College London, UK
11
12  * Correspondence to: andrea.sottoriva@icr.ac.uk and t.graham@qmul.ac.uk
13
14
15  **Abstract**
16  The distribution of fitness effects (DFE) defines how new mutations spread through an
17  evolving population. The ratio of non-synonymous to synonymous mutations (dN/dS) has
18  become a popular method to detect selection in somatic cells, however the link, in somatic
19  evolution, between dN/dS values and fitness coefficients is missing. Here we present a
20  quantitative model of somatic evolutionary dynamics that yields the selective coefficients
21  from individual driver mutations from dN/dS estimates, and then measure the DFE for
22  somatic mutant clones in ostensibly normal oesophagus and skin. We reveal a broad
23  distribution of fitness effects, with the largest fitness increases found for TP53 and NOTCH1
24  mutants (proliferative bias 1-5%). Accurate measurement of the per-gene DFE in cancer
25  evolution is precluded by the quality of currently available sequencing data. This study
26  provides the theoretical link between dN/dS values and selective coefficients in somatic
27  evolution, and reveals the DFE for mutations in human tissues.
28
29
30  **Introduction**
31
32  One of the principal goals of large-scale somatic genome sequencing is to uncover genetic loci
33  under positive selection, so-called "driver" genes, that lead to clonal expansions.
34  Enumeration of the selective advantage of each driver mutation enables prediction of future
35  evolutionary dynamics[1]. In evolutionary biology, the distribution of fitness effects (DFE) is a
36  fundamental entity that describes the selective consequences of a (large) number of
37  individual mutations of an ancestral genome[2]. In somatic evolution, particularly cancer
38  genomes, we have an extensive knowledge of the catalogue of recurrent, and likely positively
39  selected, somatic mutations[3], but the fitness changes associated with each mutation remain
40  largely unquantified.
41
42  Extensive experimental effort is ongoing to determine the fitness effects of mutations. Most
43  prominently is lineage tracing of mutations in mouse models[4,5], but these methods are not
44  sufficiently high-throughput to produce the DFE for all somatic mutations. Other studies have
45  estimated the selective coefficient of somatic mutations by measuring the frequency of such
46  mutations over time in the same individual using longitudinal sampling[6,7] however this
47  method is broadly limited to somatic evolution in the blood (where it is feasible to take

48  samples from healthy individuals over time) and in rare cases of patients under active
49  surveillance.
50
51  An alternative approach is to infer selective coefficients directly from somatic genome
52  sequencing data. Methods to identify positively-selected (driver) mutations rely on finding
53  genes that have significantly more mutational 'hits' (typically hits are non-synonymous
54  mutations) than would be expected by chance, after correction for factors known to influence
55  the mutation rate across the genome[8]. Conversely, negatively selected genes are expected to
56  show a paucity of mutations[9,10]. This idea is formalised in the calculation of the dN/dS ratio
57  – a method originally developed in molecular species evolution – that has recently been
58  adapted for use to study somatic evolution (both cancer and normal tissue)[3,9-15]. The intuitive
59  idea behind dN/dS is to measure the rate of non-synonymous (dN) mutations (possibly under
60  selection) and compare that to the rate of synonymous (dS) mutations (presumed neutral).
61  The ratio of these two numbers, each normalised for the local sequence-specific biases in the
62  mutation rate, putatively identifies a signature of selection: dN/dS > 1 indicating positive
63  selection, dN/dS = 1 indicating neutral evolution and dN/dS < 1 indicating negative selection.
64
65  Transforming dN/dS values to selective coefficients in somatic evolution is an unaddressed
66  problem.  dN/dS was originally developed in the context of species evolution using the
67  Wright-Fisher process, a classical population genetics model that assumes that evolution
68  occurs over very long timescales, which permits new mutations to fix within lineages, and also
69  that the population size is constant, with all individuals having equal potency and non-
70  overlapping generations. Under the Wright-Fisher model, the dN/dS of a locus is related to its
71  selective coefficient by the relation[16]:

$$\frac{dN}{dS} = \frac{2Ns}{1 - e^{-2Ns}}$$

72

73
74  Where $N$ is the effective population size and $s$ the selection coefficient.
75
76  However, in somatic evolution the assumptions of the Fisher-Wright model are violated.
77  Somatic evolution is rapid and new mutations are infrequently fixed in the population[17],
78  clonal dynamics are complex and population sizes unlikely to be constant[18]. Further, the lack
79  of recombination in somatic evolution can result in strong hitchhiking effects. In addition,
80  since in somatic evolution the ancestral genome is known it circumvents the need to measure
81  dN/dS across a phylogeny (a necessary step for dN/dS analysis in species evolution).
82  Violations of some of these assumptions was previously recognised to make the
83  interpretation of dN/dS problematic[19,20], and consequently the relationship between
84  selective coefficients and dN/dS values is uncertain.
85
86  The size distribution of clones (called the site frequency spectrum in population genetics
87  nomenclature) also contains information on the selective coefficients of newly arising
88  mutations. Mathematical descriptions of the dynamics of populations of cells can make
89  predictions on the shape of the clone size distribution under different demographic and
90  evolutionary models[21,22], and this approach has been used to quantify the dynamics and cell
91  fate properties of stem cells across many tissues[23-25]. We and others have also used similar
92  approaches to infer the evolutionary dynamics of tumours in deep sequencing data[26-29].

93   To date, dN/dS analysis and the analysis of the clone size distribution have been performed
94   independently, with conflictual results[30,31]. Here we develop the mathematical population
95   genetics theory necessary to combine these approaches and explore how the inter-
96   individual measure of selection at a locus as provided by dN/dS values is related to the
97   underlying cell population dynamics that generate intra-individual clone size distributions.
98   This approach naturally accounts for the nuances in somatic evolution that can make the
99   interpretation of dN/dS difficult. We show how this unified approach allows for greater
100  insight into patterns of selection than either method in isolation, and importantly reveal the
101  precise mathematical relationship between dN/dS values and selective coefficients in
102  somatic evolution. We use this approach to infer the selective advantage of mutations in
103  normal tissue and examine the evolutionary dynamics of cancer subclones.
104
105  **Results**
106
107  **A general approach to integrate dN/dS and clone size distributions**
108  We present a general mathematical framework for the interpretation of frequency-
109  dependent dN/dS values in somatic evolution. First, we construct null models of the
110  evolutionary dynamics in the absence of selection, and then augment these models to
111  incorporate the consequences of selection. Evolutionary dynamics differ between normal
112  tissues and cancer cells: in normal tissues maintained by stem cells, the long-term
113  population dynamics is controlled by an approximately fixed-size set of equipotent stem
114  cells undergoing a process of neutral competition[32], whereas in tumour growth the overall
115  population increases over time. In each scenario, we develop a null model to predict the
116  expected genetic diversity in the population in the absence of selection. Positive selection
117  causes selected variants to rise to higher frequency than expected under neutral evolution
118  (Figure 1a), and negative selection has the opposite effect. This insight guides how we
119  model the effects of selection (i.e diversity of non-synonymous mutations).
120
121  Specifically, we defined the function $g(\theta, \mu, s, f)$ as the expected distribution of mutations
122  with selective (dis)advantage $s$ found at a frequency $f$, for a given evolutionary dynamics
123  scenario, where mutations accumulate at a rate $\mu$. For the remainder of the paper we use
124  passenger mutations to refer to those mutations that have no functional effect (s=0) and
125  driver mutations those that have s>0 . When comparing to data, driver mutations are taken
126  as equivalent to non-synonymous mutations and passengers equivalent to synonymous
127  mutations.
128
129  The functional form of $g(\theta, \mu, s, f)$ encapsulates the population dynamics of the system
130  with parameter vector $\theta$, which may, for example, include the growth rate of a tumour, or
131  loss replacement rate of stem cells in normal tissue. The direct interpretation of $s$ depends
132  on the system under question. Following the logic of the effect of selection above, for $s' >$
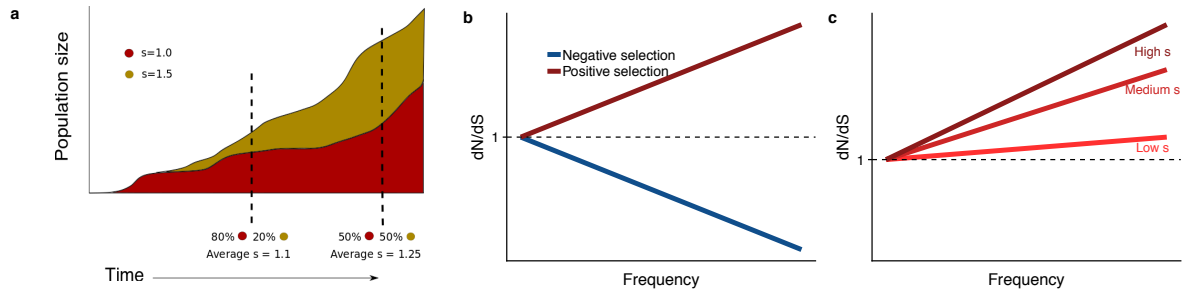133  $s$ we have that:
134
135  $$g(\theta, \mu, s', f) > g(\theta, \mu, s, f).$$
136
137
138
139

**Figure 1**

**A** Variants under positive selection are enriched at high frequency, this means dN/dS estimates are dependent on the frequency of mutation, **b.** The strength of selection influences the degree to which positively selected variants are enriched at high frequencies **c.**

Since dN/dS measures the excess or deficiency of mutations due to selection, taking the ratio of $g(\theta, s, m)$ when $s \neq 0$ to $s = 0$ and normalizing for the mutation rates, which may differ for passenger ($\mu_p$) and driver ($\mu_d$) mutations respectively, informs how dN/dS is expected to change as a function of the frequency $f$ of mutations in the population (equation 1).

$$\frac{dN}{dS} = \frac{\mu_p}{\mu_d} \frac{g(\theta, \mu_d, s, f)}{g(\theta, \mu_p, s=0, f)} \qquad [1]$$

We discuss the general properties of this model. Firstly, when $s = 0$ (neutral evolution), the numerator and denominator are equal resulting in $\frac{dN}{dS} = 1$, as expected. Secondly, dN/dS increases as a function of frequency $f$ (clone size) for positive selection, and decreases as a function of $f$ for negative selection (Figure 1b), for all $g(\theta, \mu, s, f)$ that we consider. Thirdly, the shape of the curves predicted by the underlying population model encodes the value of the selection coefficient; for example the steepness of the increase is proportional to the selection coefficient $s$ (Figure 1C). These observations are a natural consequence of positive selection driving selected mutations to higher frequency (Figure 1a).

Unfortunately, directly using equation [1] to measure selective coefficients from the slope of the dN/dS curve as function of frequency is often impractical. Real sequencing data often suffers from a limited number of mutations detected at any particular frequency and measurement uncertainties in these frequencies. To circumvent these issues, we introduce "interval dN/dS" (i-dN/dS) that aggregates over a frequency range to reduce the influence of these sources of noise. Interval dN/dS is defined as:

$$i\text{-}\frac{dN}{dS} = \frac{\mu_p}{\mu_d} \frac{\int_{f_{min}}^{f_{max}} g(\theta, \mu_d, s, f) df}{\int_{f_{min}}^{f_{max}} g(\theta, \mu_p, s=0, f) df} \qquad [2]$$

Fixing the integration range $[f_{min}, f_{max}]$ allows for robust inference of $s$ in potentially sparse and noisy sequencing data using maximum likelihood methods (see Methods).

**Frequency-dependent dN/dS values in stem cell populations**

In healthy tissue, only mutations that are acquired in the stem cells will persist over long times, and so we restrict our attention to these cells. Quantitative analysis of lineage tracing

4

178     data has shown that the stem cell dynamics of many tissues conform to a process of
179     population asymmetry[32]. In this paradigm, under homeostasis, the loss of stem cells through
180     differentiation is compensated by the replication of a neighbouring stem cell, thus
181     maintaining an approximately constant number of stem cells. These dynamics are
182     represented by the rate equations:
183

184

$$SC \quad \overset{r\lambda}{\to} \quad \begin{matrix} SC + SC \\ D + D \end{matrix} \quad \begin{cases} p = (1+\Delta)/2 \\ p = (1-\Delta)/2 \end{cases} \qquad [3]$$

185

186     where $SC$ refers to a single stem cell which divides symmetrically to produce either two
187     stem cells or two differentiated cells (denoted as $D$ above), $\lambda$ is the rate of cell division per
188     unit time, and $r$ is the probability of a symmetric divisions. The product $r\lambda$ is referred to as
189     the loss/replacement rate. Differentiated cells will ultimately be lost from the population
190     over long time scales. Under homeostasis, these processes should be exactly balanced with
191     $\Delta = 0$. With $\Delta \neq 0$, the fate of a stem cell is 'biased', introducing positive or negative
192     selection into the model. Previous mathematical analysis shows that this model is a good
193     description of the clonal dynamics in the oesophagus and skin[23,33,34]. Using the previous
194     analytical results describing the temporal evolution of the clone distribution (see
195     supplementary methods for detailed discussion) we derive the frequency distribution
196     $g(\theta, \mu, s, f)$ for oesophagus and skin as [21,23,35]:
197

198

$$g(\theta, \mu_x, s, f) = \frac{n_0 \mu_x}{f} e^{-\frac{f}{N(t)}} \qquad [4]$$
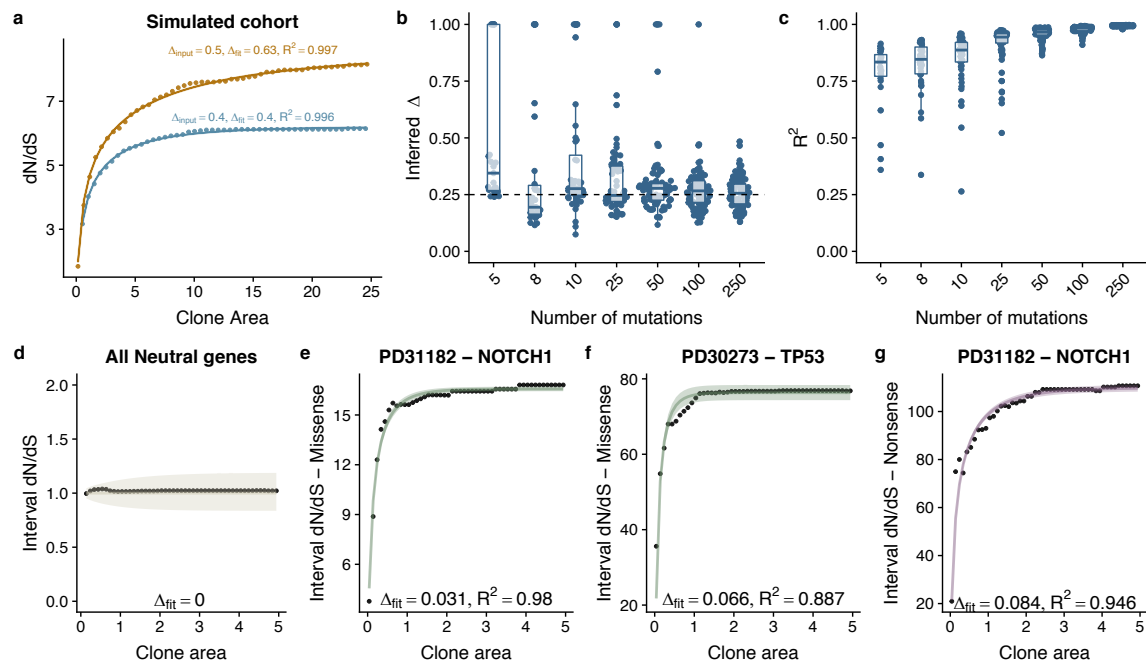
199

200     Where $n_0$ is the starting population size and $\mu_x$ the mutation rate, which may be different
201     for drivers ($s \neq 0$) and passenger mutations ($s = 0$). $N(t)$ is a scaling factor that depends
202     on $\Delta$, the bias toward self-renewal, which we interpret as our selection coefficient in this
203     system. Specifically:
204

205

$$N_{\Delta=0}(t) = 1 + r\lambda t \qquad [5]$$

206

$$N_\Delta(t) = \frac{(1+\Delta)e^{2r\lambda\Delta t} - (1-\Delta)}{2\Delta} \qquad [6]$$

207

208     We note that at long times (large $N(t)$) equation [4] converges to a $1/f$ distribution for the
209     site frequency spectrum of a fixed size population[36]. $N(t)$ can be interpreted as the average
210     size of a labelled clone after time $t$, which even under homeostasis grows over time and
211     compensates for some clones being lost due to drift. From these expressions, we can then
212     write down a closed-form expression for i-dN/dS as a function of clone frequency (see
213     methods) that allows for maximum likelihood estimation of parameter values ($\Delta, r\lambda$). We
214     confirmed the accuracy of our derivation using simulations (Figure 2a), and performed
215     power calculations to determine the minimum number of mutations required to correctly
216     infer the underlying population dynamics. We determined that 8 mutations per gene was
217     sufficient to accurately recover $\Delta$ (Figure 2b) with accuracy increasing for higher mutation
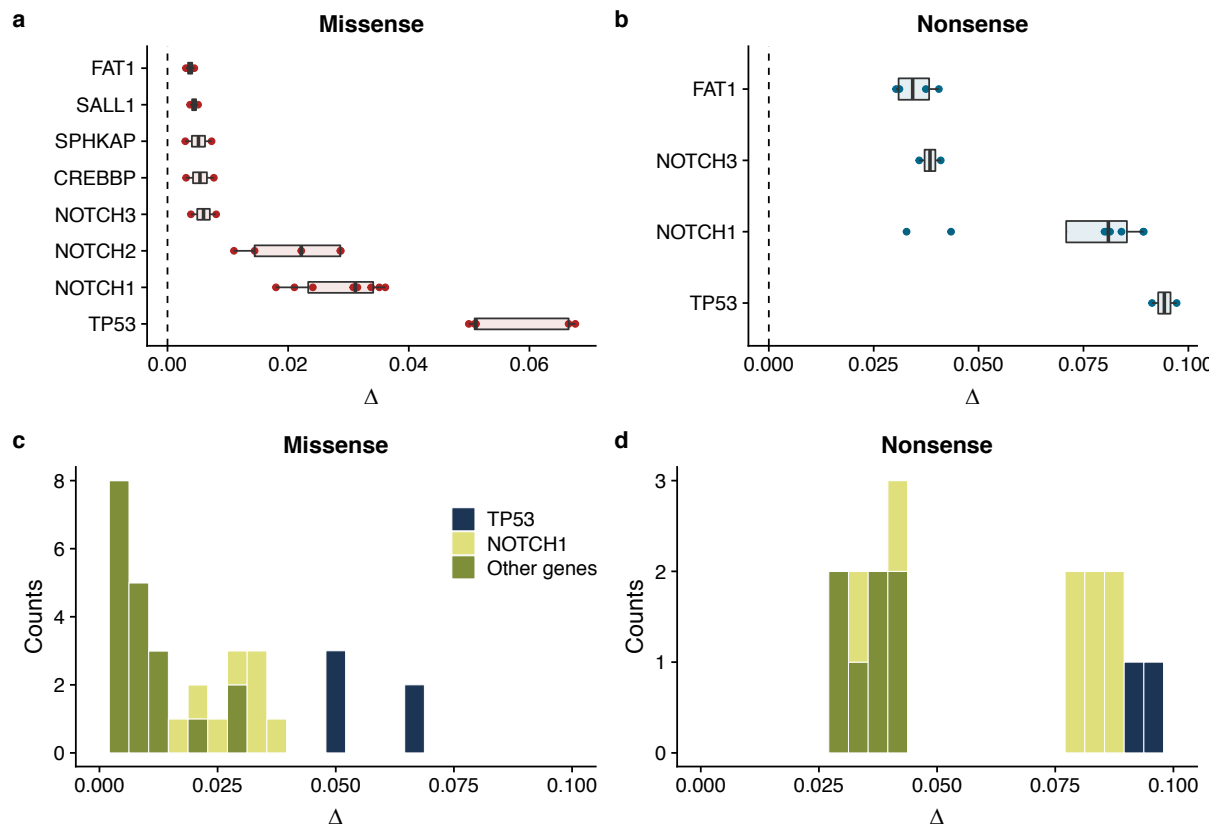218     burdens (Figure 2c).
219

220

221

**Figure 2**

**a** Interval dN/dS as a function of clone area for 2 simulated cohorts where driver mutations induce different biases, theoretical model captures the dynamics well and enables us to recover the bias Δ, accurately. As the number of mutations increases ability to recover the correct Δ and the model fit (measured using $R^2$) improves **b** and **c**. **d** Data and model fit for all neutral genes, shows i-dN/dS = 1 across the frequency range and inferred bias of 0. Data and model fit for **e** NOTCH1 missense mutations in patient PD31182, **f** missense TP53 mutations in PD30273 and NOTCH1 nonsense mutations in PD31182. Data are black points and model fits are solid lines with shaded areas denoting 95% CI.

## Selection advantages in histopathologically normal human oesophagus

We inferred the selective advantage of driver mutations in human oesophagus using published deep sequencing data from Martincorena and colleagues[14,37] that documents the clonal expansion of a panel of putative driver mutations in histopathologically-normal oesophageal biopsies.

We used the dndscv bioinformatics tool[3] to calculate frequency-dependent dN/dS values from these data (clone size measured in fraction of mutant reads multiplied by 2mm$^2$ – the area of the biopsy – and assuming 5,000 stem cells per mm$^2$ tissue). dN/dS values varied considerably as a function of mutation frequency (Figure S1).

We considered the average frequency-dependent dN/dS values across all genes in the panel, on a patient-by-patient basis. Our theoretical model of i-dN/dS calculated from these data fitted strikingly well (Figure S2). Estimates of the loss/replacement rate $r\lambda$ of the stem cell population were in the range 1.2-5.0 per year (Figure S2&S3). Inference of the selective advantage $s$ (measured in terms of the bias towards self renewal Δ) revealed an average bias of 0.004 (0.002 – 0.005 95% CI) per missense mutation (Figure S2). Nonsense mutations caused a five-fold greater bias towards self-renewal of 0.021 (0.008 – 0.032 95% CI) (Figure S3). After removal of all genes that are strongly selected, global dN/dS values on the remaining 48 genes show dN/dS of approximately 1 across the frequency range (Figure 2d), and i-dN/dS analysis revealed somatic mutation does not associate with a proliferative bias (Δ=0).

6

**Figure 3**
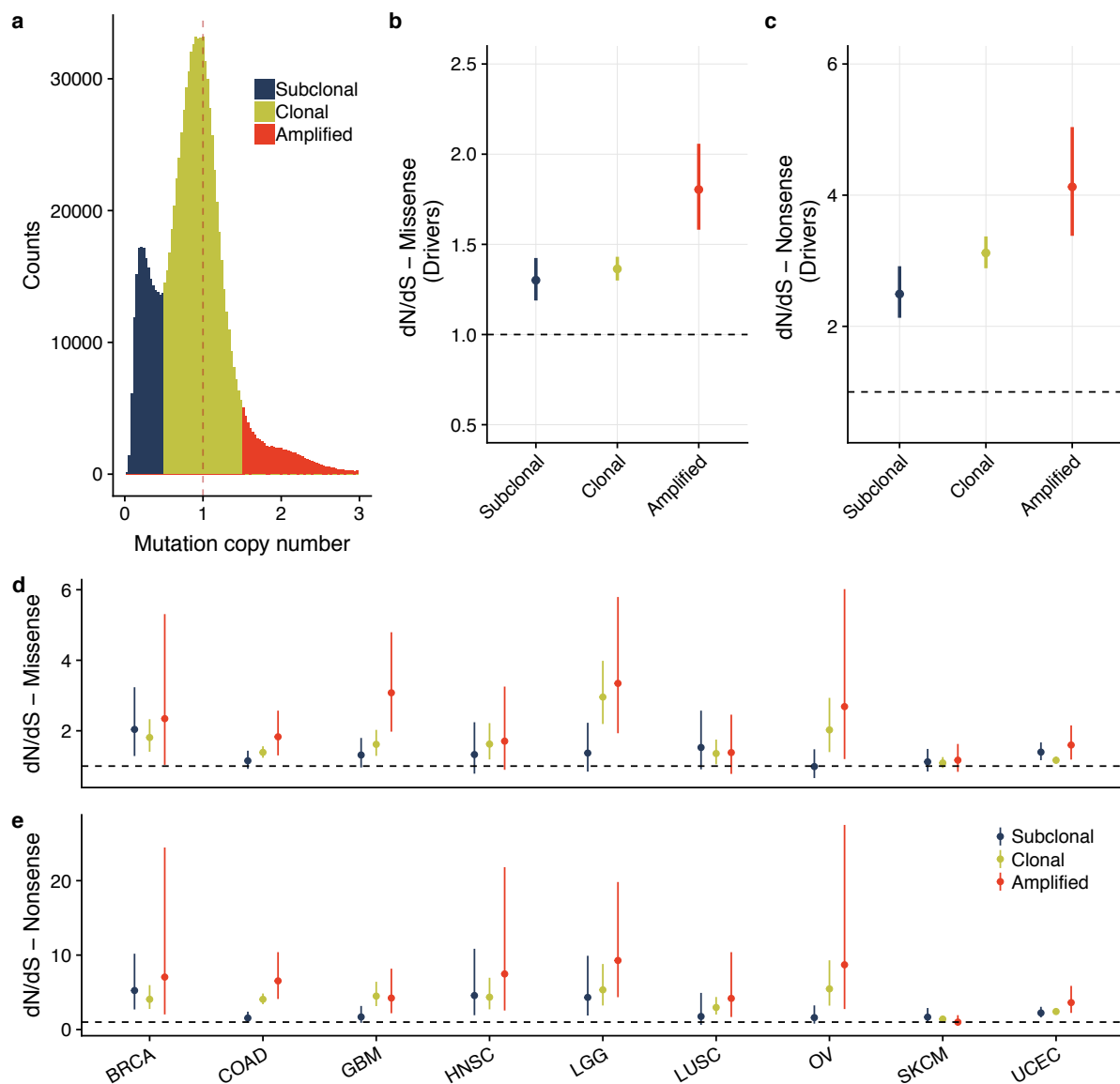
Summary of model fits across all patients for normal oesophagus data. Inferred biases Δ for genes where at least 2 patients had good model fits (R2 > 0.6 & >7 mutations) for missense mutations **a**, and nonsense mutations **b**. Inferred distribution of fitness effects for all genes across all patients for missense mutations **c**, and nonsense mutations **d.**

We then fitted the data on a gene-by-gene and patient-by-patient basis for cases where sufficient mutations were available to perform the fit (Figure 2e-g; Figure S4). A broad range of selective advantages were inferred (Figure S4&S5). Mutations in *TP53* showed large biases across all patients for both missense, Δ=0.057 (0.05-0.068 95% CI) and nonsense mutations, Δ=0.094 (0.091-0.097 95% CI) (Figure 3a-b). This was also true for mutations in NOTCH1 with Δ=0.029 (0.019-0.036 95% CI) for missense and Δ=0.072 (0.034-0.089 95% CI) for nonsense mutations. *NOTCH2, PIK3CA, CREBBP* and *FAT1* also showed a bias toward self-proliferation in multiple patients (Figures 3a-b), though most had a small effect on fitness (range 0.003 – 0.029 for missense mutations and 0.030 – 0.041 for nonsense mutations) . Together these data suggest a distribution of fitness effects (DFE) characterized by many small effect mutations with few large effect mutations (Figures 3c-d), as in seen in organismal evolution[2].

**Driver mutation selective advantage in normal skin**

Martincorena and colleagues have also published data on the expansion of driver mutations in ostensibly normal human skin[18]. Analyses of these data with interval dN/dS revealed a per-patient average selective advantage per mutation (again measured in terms of the bias towards self renewal Δ) of Δ=0.001 for missense mutations and four-fold higher for Δ=0.004 for nonsense mutations (Figures S6a-c). Performing the analysis on a gene-by-gene

279     basis was limited by the low detected number of mutations, and the limited frequency
280     range (clone size range). Good fits to the data were obtainable for *NOTCH1* missense
281     mutations in patient PD18003 with fitness estimated to be Δ=0.0149 (0.0148-0.0150 95%
282     CI), and TP53 missense mutations also in patient PD18003, Δ=0.0054 (0.0051-0.0058 95%
283     CI) Figure S6.  These fitness coefficients were similar to the oesophagus data. For missense
284     mutations we were also able to produce the distribution of fitness effects across the skin
285     cohort, which showed similar characteristics to the oesophagus data of a small number of
286     high effect mutations and a larger number of smaller effect mutations, Figure S6f.
287



**Figure 4**
Mutation copy number histogram across 2,619 TCGA samples coloured by mutation clonality, **a**. dN/dS by mutation clonality for missense, **b** and nonsense **c** mutations in a panel of 192 high confidence driver genes. The same analysis done per cancer type for missense **d** and nonsense **e**.

**Clonal mutations have greater dN/dS than subclonal mutations in cancers**

We next investigated the selective advantage of driver mutations in cancer. We first investigated whether or not differences existed between dN/dS values for clonal mutations (ie truncal, present in all cells in a cancer) and subclonal mutations (present in a subset of cells in a cancer) were apparent. Using sequencing data from 2,619 cancers from TCGA that had sufficient cellularity and depth (see Methods) we calculated the mutation copy number (MCN) for each mutation and grouped mutations into subclonal, clonal and amplified across the cohort, where mutations with MCN < 1 were subclonal, MCN == 1 were clonal and MCN > 1 were amplified (Figure 4a). We than calculated global dN/dS ratios for a panel of 198 high confidence driver genes (Methods).

Across all cancers, the signal of positive selection was more pronounced for clonal mutations (Figures 4b-e), with the highest dN/dS values found in amplified mutations[38]. Subclonal mutations on the other hand demonstrated much lower dN/dS values. The same pattern was also evident in individual cancer types (Figure 4e,d & S7). In many cancer types (colorectal, ovarian, glioblastoma) subclonal mutations showed no evidence of subclonal selection (neutral evolution; dN/dS = 1), Figure 4e,d & Figure S7.

**Interval dN/dS for cancer**

We applied our mathematical approach above to calculate i-dN/dS in cancer evolution. In cancer evolution $g(\theta, \mu, s, f)$ must account for tumour growth dynamics and subclonal mutations which may rise and fall in frequency due to selection and drift. The well-studied Luria-Delbrück distribution and its extensions describes these dynamics[39]. Specifically, the Luria-Delbrück distribution describes the expected number of mutational lineages at a particular frequency assuming an underlying birth-death process for individuals in the population. For neutral mutations the site frequency spectrum has a characteristic $\frac{1}{f^2}$ dependence, where $f$ is the frequency of the mutations [35,40]. Hence:

$$g(\theta, \mu_p, s = 0, f) = \frac{\mu_p}{\beta_p} \frac{1}{f^2} \qquad [7]$$

where $\mu_p$ is the passenger mutation rate and $\beta_p$ is the survival probability of a lineage at division. We previously showed that in many cancers across types (approx. 30% of cases), subclonal mutations closely follow the prediction of this neutral model[26].

Extensions to the classic Luria-Delbruck distribution describe the differential fitness of mutants. We defined the relative fitness advantage $s$ as the ratio of net growth rates between wildtype 'passenger' mutations ($\lambda_p$) and driver mutations ($\lambda_d$) :

$$s = \frac{\lambda_d}{\lambda_p} - 1 \qquad [8]$$

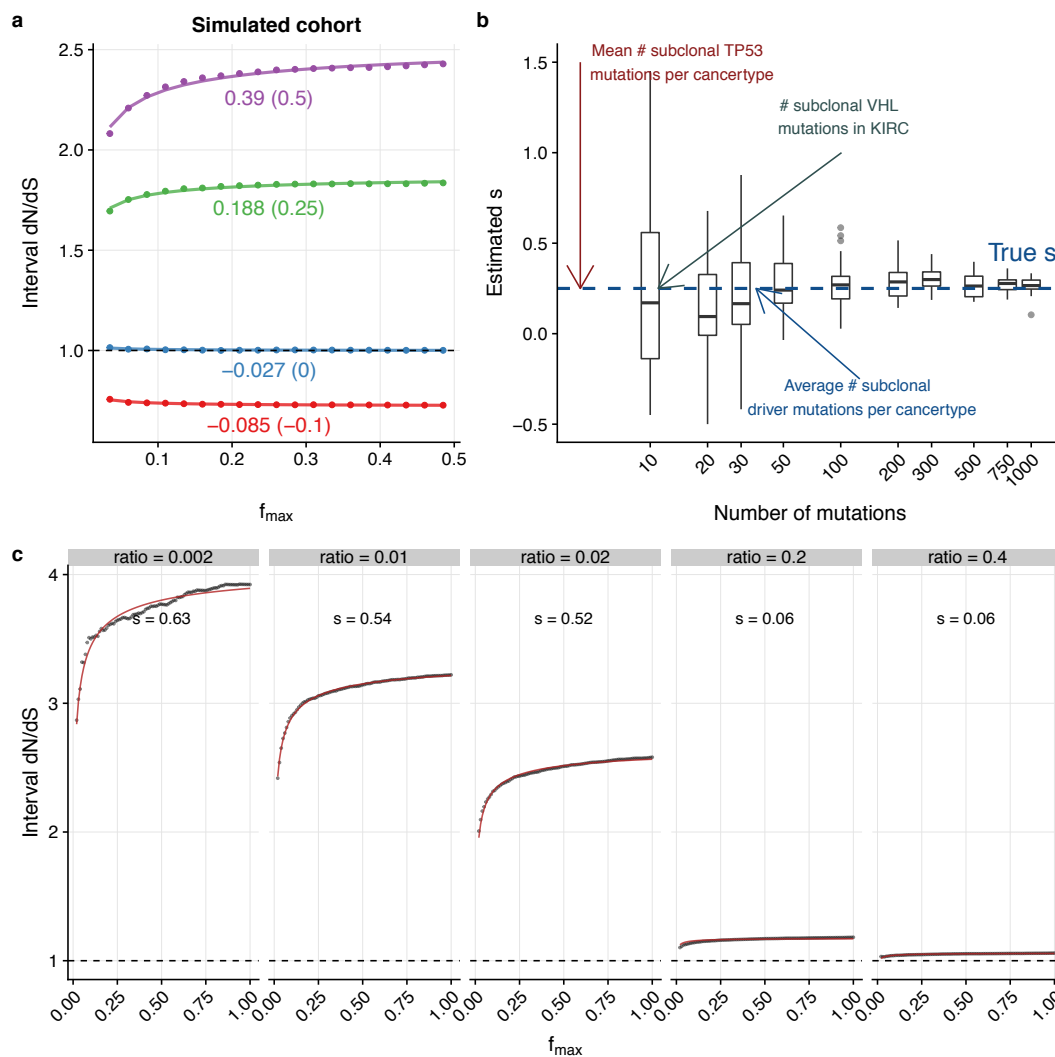s > 0 indicated positive selection while s < 0 indicated negative selection. We also defined the birth and death rates of the respective wildtype (passengers) and mutants (drivers) as $b_p$, $d_p$, $b_d$ and $d_d$. Here, the site-frequency distribution again follows a power law but with exponent dependent on the relative fitness advantage of the mutant [35,40]:

336
$$g(\theta, \mu_d, s \neq 0, f) = \frac{N\mu_d}{\beta_d^{\frac{1}{1+s}}} \frac{b_p}{b_d} \frac{\Gamma\left(\frac{2+s}{1+s}\right)}{N^{\frac{2+s}{1+s}}} \frac{1}{f^{\frac{2+s}{1+s}}}$$ [9]

337

338   Here, N is the tumour population size at the time of sampling. Using these expressions
339   (equations 7&9), we derive i-dN/dS (see Methods). The equation exhibits the same
340   qualitative behaviour as for the stem cell model, in that dN/dS increases as a function of
341   frequency for positive selection and decreases for negative selection (Figure 5a). Using a
342   simulation-based model to generate synthetic data, we confirmed the accuracy of the
343   model by accurately recovering the inputted selection coefficient by application of the
344   theoretical model and maximum likelihood inference (Figure 5a).

345



346

**Figure 5**

Interval dN/dS as a function of frequency for 4 simulated cohorts where driver mutations induce different selective advantages, **a**. Points are simulated data and lines are model fits, under each line is the inferred selective advantage and the true selective advantage in brackets. Power to correctly infer the selection coefficient depends on the number of mutations in the cohort, **b**. We generated a cohort of 1000 tumours and then subsamplesd the mutations (50 times) and inferred the selection coefficient. For TCGA we are limited by a small number of subclonal drivers to accurately perform the inference. The ratio of the driver mutation rate to passenger mutation rate has a strong influence on dN/dS, **c**. Here we generated synthetic cohorts where the strength of selection of driver mutations was 0.5, and different ratio of driver mutation rate to passenger mutation rate. When drivers are rare, dN/dS > 1 and we can accurately apply our model. When drivers are frequent compared to passengers we observe strong hitchhiking effects which results in dN/dS~1.

347 Subclonal dN/dS is strongly influenced by the ability to resolve low frequency variants. We
348 generated synthetic tumour cohorts that modelled subclonal selection, and simulated
349 'perfect sensitivity' for mutation detection. In these cases, where all mutations were
350 resolved, we measured dN/dS≈1 (and hence infer a selection coefficient of 0), despite some
351 lineages being positively selected (Figure S9). If only higher frequency variants were
352 analysed, then the measured dN/dS > 1 and the correct selective coefficient is inferred
353 (Figure S9). We note that at very low frequencies the detected mutations are newly arisen
354 in the population, and so are as yet 'unfiltered' by selection. Consequently the ratio of non-
355 synonymous to synonymous mutations is expected to be proportional to the respective
356 mutation rates of the two mutation types. The abundance of low frequency mutations also
357 increases exponentially with decreasing clone frequency, and so including very low-
358 frequency variants 'drowns out' the effects of selection (Figure S9C). We note that the
359 limited sequencing data of the majority of currently available cancer genomic data means
360 that typically only high frequency variants are detected.
361
362 **Currently available cancer sequencing data is insufficient to infer selective advantages**
363 Limitations in the quality of currently available sequencing data meant that the theoretically
364 predicted frequency dependence of dN/dS values could not be assessed in cancer genomics
365 data (Figure S8). Limited sequencing depth introduces uncertainty into the determination
366 of variant allele frequencies ("sequencing noise") which can result in incorrect classification
367 of mutation clonality. Visual inspection of the mutation copy number histogram for TCGA
368 data (Figure 4a) showed a very broad dispersion of MCNs, and the resolution at lower
369 (subclonal) frequencies was particularly poor. Issues arising due to sequencing noise are
370 exacerbated in the setting of dN/dS analysis where pooling the data from multiple patients
371 with different sequencing depth and purities is required. Consequently, the range of
372 subclonal frequencies where interval dN/dS could be calculated was severely restricted.
373
374 We tested whether or not looking at individual genes (rather than individual mutations)
375 allowed for measurement of the DFE. However, the lack of recurrent subclonal mutations
376 on a gene-by-gene basis precluded this approach. Power calculations predicted that a
377 minimum of 30 subclonal mutations in a given gene were required to accurately fit the
378 interval dN/dS model (Figure 5b). This level of subclonal recurrence of individual mutations
379 was not seen in the data: for example, the average number of subclonal mutations in *TP53*
380 per cancer type, as well as the number of subclonal VHL mutations (which has been
381 reported to occur subclonally at an appreciable frequency [41]) were both well below this
382 cutoff (Figure 5B). Consequently, large cohorts of tumours sequenced to higher depth are
383 required to apply this approach.
384
385 Aside, we note that the traditional dN/dS approach, and also our modelling framework,
386 assumes that mutations are independent, and consequently the possibility of hitchhiking of
387 mutations (e.g. nested driver mutations within clones) is neglected. In simulated data, we
388 observed high mutation rates for both driver and passenger mutations led to hitchhiking
389 being common, and subsequent obscuring of the signal of selection (Figure 5c). In extreme
390 cases this led to dN/dS = 1 (apparent neutral evolution) even in the presence of multiple
391 selected lineages. For most cancers, the number of driver mutations per cancer is thought
392 to be low (<10)[3], but nevertheless in hypermutator cancers the hitchhiking effect is likely to
393 be common. Thus, despite hypermutator tumours tending to have fewer copy-number

394  alterations and hence less problematic estimation of MCNs, the prevalence of hitchhiking
395  precludes analysis of these tumours.
396
397  **Discussion**
398
399  Here we have shown that the combination of dN/dS values with mutation frequency-based
400  information provides additional quantitative insight into dynamics of somatic evolution than
401  either method alone. Specifically, the combined approach enables direct inference of the
402  selection coefficients of mutations in somatic tissues.
403
404  Using this methodology we have begun the construction of the distribution of fitness effects
405  (DFE) in somatic evolution (Figure 3c,d & Figure S6f). In histologically normal epithelium,
406  mutations of most genes considered showed minimal effects on fitness (near-neutral
407  evolution), though selection coefficients for some loci, foremost *NOTCH1* and *TP53* were
408  considerable (>1% and >5% respectively), and consequently the DFE has most mass close to
409  s=0 with a long right-tail of highly-selected variants.  We observed that values of selective
410  coefficients of individual genes varies between patients, likely because of inter-patient
411  difference in the precise location of point mutations, but potentially also because of inter-
412  patient variation in selective pressure from the microenvironment. Nevertheless, the
413  comparative rank of per-gene fitness coefficients was broadly consistent across patients.
414  This consistency in selective coefficients is in agreement with the observation highly
415  recurrent gene mutations in cancer[42] and evidence of repeatability in cancer evolution[43].
416
417  We have previously measured fitness effects in individual cancers (but were unable to
418  ascribe fitness changes to individual genes) finding increases in growth rate in a selected
419  clone approaching 100% in some cases[27].  Care must be taken when comparing selective
420  coefficients between normal and cancer populations, because in the former we quantify
421  selection as tilt away from homeostasis and towards net growth of a lineage, whereas in
422  cancer we infer the relative growth rate of a clone within the tumour as a whole. With this
423  important caveat in mind, nevertheless the fitness increases observed in cancer appear to
424  be much larger than for normal tissues. We hypothesise that this is because the effect of
425  selection is weaker in expanding populations like cancer, wherein the generation of a
426  subclonal expansion requires very large increases in fitness[44].
427
428  On a cautionary note, our theoretical work shows that the clonality of mutations strongly
429  determine the observed value of dN/dS, and so a misleading picture of the selective forces
430  operating in a tumour (or healthy tissue) will be produced if dN/dS frequency-dependent
431  effects are not corrected for. The accuracy of any estimate of evolutionary dynamics from
432  dN/dS values is of course dependent of the underlying accuracy of the dN/dS measure itself,
433  which is compromised by uncharacterised variability in the mutation rate across the
434  genome[45] and in the uncertain pathogenicity of individual single nucleotide variants
435  (extensions to estimate site level selection coefficients may circumvent some of these
436  issues[46,47]). Finally, we note that dN/dS measures cannot elucidate evolutionary pressures in
437  individual samples as insufficient (subclonal) mutations will be found at any individual locus.
438  dN/dS cohort measurements are sensitive to outliers, where a few patients with high
439  selection can drive the results [48]. Other approaches, such as using the site frequency
440  spectrum, are likely more powerful for these types of questions.

441

442 Combining population genetics methods with comparative genomics is a powerful way to
443 infer selection pressures in human somatic evolution, giving new insight into the
444 fundamental parameters that determine evolutionary dynamics in health and disease.
445

446 **Acknowledgements**

453

**Methods**

**TCGA Data Processing**

MAF (Mutation Annotation Format) files from the Mutect2 mutation calling algorithm and copy number segmentation data for 9950 cancers from 26 cancer types were downloaded from the genomic data commons portal using the TCGAbiolinks R package [49]. Cellularity and ploidy estimates derived from ASCAT were obtained from COSMIC (https://cancer.sanger.ac.uk/cosmic/download). We then filtered for >2 reads reporting the variant and >9 reads coverage at each locus in both the tumour and normal sample. We removed samples where the effective depth (defined as cellularity times depth) was < 50X and those that had likely undergone genome doubling (ploidy > 2.5). This left 2619 samples from 17 cancer types which we deemed suitable for analysis.

Copy number (CN) segmentations together with cellularity estimates were used to correct the variant allele frequency and produce mutation copy number estimates. We assume that the observed CN state $(\overline{CN})$ was a combination of signals from the tumour sample and contamination from normal cells (with two copies) assuming tumour purity c.

$$\overline{CN} = c \times CN + 2(1 - c)$$

With this, log(R) ratios were transformed into copy number states using the following formula:

$$CN = \frac{2\left(2^{\log(R)} - 1 + c\right)}{c}$$

Using these corrected copy number states, mutation copy number (MCN) values were calculated. Given mutation $i$ with variant allele frequency $VAF_i$, copy number $CN_i$ at the locus and cellularity estimate of the tumour c, the MCN was calculated as follows:

$$MCN_i = \frac{CN_i \times VAF_i}{c}$$

Visual inspection of the MCN histograms (Figure 4a) show a dominant peak at MCN = 1 representing clonal mutations present in a single copy, confirming that the corrections we applied work as intended.

**Oesophagus and skin data**

For the oesophagus and skin data we used mutation calls provided by the original studies. In the oesophagus data when a mutation was present in multiple adjacent biopsies we used the sum of the mutation frequency times the area of the biopsies ($2mm^2$) as our readout of clone size and performed the dN/dS analysis on a patient by patient basis.

**dN/dS calculations**

For calculating dN/dS ratios the dndscv R package was used which calculates both global dN/dS ratios across the whole exome or a panel of genes as well as per gene dN/dS ratios using a covariate based model to infer dN/dS values with a limited number of mutations [3]. In an attempt to enrich for positive selection in some of our analysis we calculated dN/dS for a subset of 198 high confidence driver genes [50].

499

500 Over or under filtering of possible germline SNPs is known to influence dN/dS values in

501 somatic genomes[3]. We previously found that mutation calls provided by TCGA are likely

502 over stringent on filtering germline SNPs resulting in inflated dN/dS values [48]. To circumvent

503 this issue, we calculated a baseline dN/dS value by randomly selecting 1,000 genes

504 (excluding drivers) and then running dndscv across the whole TCGA cohort, reasoning that

505 this should on average return dN/dS = 1, and any deviation from this would be due to

506 under/over filtering of SNPs . Repeating this procedure 50 times and then taking the mean

507 value gave us our baseline value which we could then subtract from further dN/dS values

508 we calculate in our analysis.  To confirm this procedure produces the expected result of

509 dN/dS = 1 in the absence of selection, we repeated the procedure and again, randomly

510 selected 1,000 genes 100 times and then applied the correction (subtracting the calculated

511 deviation from 1). As would be expected the mean of this distribution was dN/dS = 1,

512 validating our approach, Figure S10.

513

514 To calculate the interval dN/dS measure we took our corrected mutation frequency data

515 and determined a low cutoff $f_{min}$ based on the minimum mutation frequency. We then

516 created a vector of frequencies $f_{max}$ that covered the total range of mutation frequencies

517 and calculated dN/dS between $f_{min}$ and all values of $f_{max}$. This allowed us to plot dN/dS as

518 a function of $f_{max}$ and fit our interval dN/dS models.

519

520 **Model fitting**

521 We used a maximum likelihood approach to fit our models to the data. Defining the

522 observed interval dN/dS as $y$ and the model dN/dS as $\hat{y}(\theta) = \frac{\mu_p}{\mu_d} \frac{\int_{f_{min}}^{f_{max}} g(\theta,\mu_d,s,f)df}{\int_{f_{min}}^{f_{max}} g(\theta,\mu_p,s=0,f)df}$ . First

523 of all we define the residuals between the data and the model as $R = y - \hat{y}$.  Assuming that

524 the residuals are normally distributed with mean 0 we can write down the negative log

525 likelihood (NLL) as

526 $$NLL(\theta) = -\sum_{y-\hat{y}(\theta)} \log\left(N(y - \hat{y}(\theta), \mu = 0, \sigma)\right)$$

527 where $N$ denotes the normal probability density function. We can then find the parameters

528 $\theta$ that minimize the NLL and calculate confidence intervals on these estimates using the

529 Fisher information matrix.

530

531 **Interval dN/dS models**

532 For the stem cell model, using equations [2]-[6] in the main text, interval dN/dS is given by:

533 $$i\text{-}\frac{dN}{dS} = \frac{1}{1+\Delta} \frac{\left[E_i\left(-\frac{\rho A_{max}}{N_\Delta(t)}\right)-E_i\left(-\frac{\rho A_{min}}{N_\Delta(t)}\right)+\frac{1}{2}\left(\frac{e^{-\frac{\rho A_{max}}{N_\Delta(t)}}}{\rho A_{max}}+\frac{e^{-\frac{\rho A_{min}}{N_\Delta(t)}}}{\rho A_{min}}\right)\right]}{\left[E_i\left(-\frac{\rho A_{max}}{N(t)}\right)-E_i\left(-\frac{\rho A_{max}}{N(t)}\right)+\frac{1}{2}\left(\frac{e^{-\frac{\rho A_{max}}{N(t)}}}{\rho A_{max}}+\frac{e^{-\frac{\rho A_{min}}{N(t)}}}{\rho A_{min}}\right)\right]}$$

534

535 Where $E_i$ is the exponential integral $E_i(x) = -\int_x^\infty \frac{e^{-n}}{n} dn$. Given that the data is in terms of

536 area, A we made the transformation $f = \rho A$, where $\rho$ is density of stem cells per mm$^2$,

537 which we set to 5,000 cells /mm$^2$ for fitting.

538

539    For the cancer model, interval dN/dS is given by:

540

$$i\text{-}\frac{dN}{dS} = \frac{\mu_p \int_{f_{min}}^{f_{max}} C_{selection} df}{\mu_d \int_{f_{min}}^{f_{max}} C_{neutral} df} = N^{\frac{s}{1+s}}(1+s)\frac{\beta_p}{\beta_d^{\frac{1}{1+s}}}\frac{b_p}{b_d}\Gamma\left(\frac{2+s}{1+s}\right)\frac{f_{min}^{\frac{-1}{1+s}} - f_{max}^{\frac{-1}{1+s}}}{\frac{1}{f_{min}} - \frac{1}{f_{max}}}$$

541

542    We note that in the cancer setting because the final population size N is generally unknown

543    we fit the model $\hat{y}(\theta = \{A, s\}) = A \times \frac{f_{min}^{\frac{-1}{1+s}} - f_{max}^{\frac{-1}{1+s}}}{\frac{1}{f_{min}} - \frac{1}{f_{max}}}$.

544

545    For a detailed description of the mathematical background of the clone size distribution in

546    these models and comparison with simulation see the supplementary Jupyter notebooks.
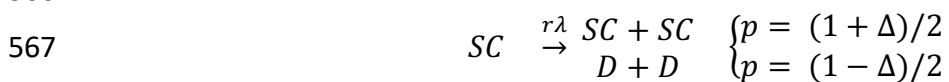
547

548    **Simulations**

549

550    To confirm our analytical models and investigate the influence of uncertainty in mutation

551    frequencies due to sequencing noise and to challenge some of the underlying assumptions

552    of our theoretical approach, we developed 2 simulation based models. The first one models

553    cancer evolution and the second models stem cell evolution under homeostasis. For the

554    cancer evolution model, we adapted our previously described model[27] so that mutations

555    can be one of two types, neutral passengers or mutations that have an effect on fitness of

556    cells (either positive or negative). We model cancer growth as a continuous time branching

557    process. At each division, daughter cells acquire mutations with a fitness effect s at rate $\mu_d$

558    and passenger mutations (which are neutral) at rate $\mu_p$. This is implemented by drawing a

559    Poisson random variable with mean given by $\mu_d$ or $\mu_p$. Fitness of passenger mutations is 0,

560    while driver mutations have fitness advantage s, where s is defined by equation [8]. We also

561    implemented a model where fitness was a random exponentially distributed variable with

562    mean s.

563

564    For the stem cell model we seed a population of $N_s$ stem cells that then undergo

565    loss/replacement as described by the following rate equations

566

567

$$SC \xrightarrow{r\lambda} \begin{matrix} SC + SC \\ D + D \end{matrix} \quad \begin{cases} p = (1+\Delta)/2 \\ p = (1-\Delta)/2 \end{cases}$$

568

569    As only the stem cells are long lived the differentiated cells are not explicitly modelled such

570    that when a stem cell "differentiates" it is effectively lost from the population. As in the

571    cancer model, during division, daughter cells acquire mutations with a fitness effect at rate

572    $\mu_d$ and passenger mutations at rate $\mu_p$. Fitness increases the bias toward self-proliferation $\Delta$

573    of a stem cell lineage. Additional driver mutations do not further increase the fitness of

574    stem cells.

575

576    To calculate dN/dS across a cohort of simulated tumours or tissue biopsies we count the

577    number of driver mutations $N_d$ and the number of passenger mutations, $N_p$ and then

578    normalize by their respective mutation rates. In our model drivers = non-synonymous and

579    thus every driver has an effect on fitness. Then the ratio of these two numbers gives us the

580    excess or deficit of mutations due to selection – ie the dN/dS ratio.

16

581

$$\frac{dN}{dS} = \frac{N_d/\mu_d}{N_p/\mu_p}$$

582

583

584 For the interval dN/dS we simply calculate the $N_x$ between $f_{min}$ and $f_{max}$.

585

586 To introduce uncertainty into mutation frequencies we perform a process of empirically
587 motivated sampling to the true underlying frequency $f$. Firstly, we specify the average
588 depth of sequencing D, then the depth of sequencing for mutation i is given by

589 $$D_i = Po(D)$$

590 The sampled number of read counts is then

591 $$n_s = Bo(n = D_i, p = f)$$

592 And the sampled variant frequency is then $f_s = n_s/D_i$

593

594 **Code and data availability**
595 Code used for the analysis are included as a snakemake pipeline which will reproduce all the
596 analysis and generate all the figures. Julia [51] was used for the majority the simulations and R
597 [52] was used to analyse the data and generate the figures. Some of the analysis rely in
598 bespoke packages written for this which are freely available under and open source licence.
599 Code is available at github.com/marcjwilliams1/dnds-clonesize.

600

601
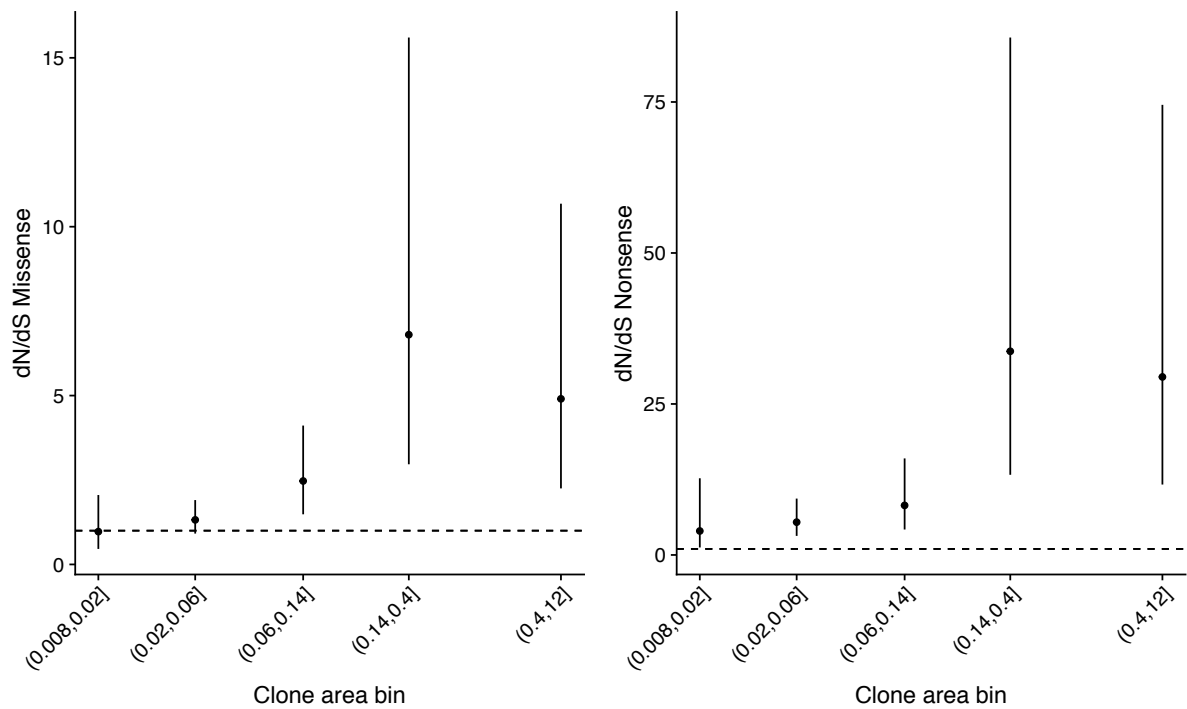
602

**References**

1. Williams, M. J., Sottoriva, A. & Graham, T. Measuring Clonal Evolution in Cancer with Genomics. *Annu. Rev. Genom. Hum. Genet. In Press*

2. Eyre-Walker, A. & Keightley, P. D. The distribution of fitness effects of new mutations. *Nat Rev Genet* **8,** 610–618 (2007).

3. Martincorena, I. *et al.* Universal Patterns of Selection in Cancer and Somatic Tissues. *Cell* 1–35 (2017). doi:10.1016/j.cell.2017.09.042

4. Vermeulen, L. *et al.* Defining stem cell dynamics in models of intestinal tumor initiation. *Science* **342,** 995–998 (2013).

5. Rogers, Z. N. *et al.* Mapping the in vivo fitness landscape of lung adenocarcinoma tumor suppression in mice. *Nature Genetics* **50,** 483–486 (2018).

6. Watson, C. J. *et al.* The Evolutionary Dynamics and Fitness Landscape of Clonal Haematopoiesis. BioRxiv 1–34 (2019). doi:10.1101/569566

7. Körber, V. *et al.* Evolutionary Trajectories of IDHWT Glioblastomas Reveal a Common Path of Early Tumorigenesis Instigated Years ahead of Initial Diagnosis. *Cancer Cell* 1–37 (2019). doi:10.1016/j.ccell.2019.02.007

8. Bailey, M. H. *et al.* Comprehensive Characterization of Cancer Driver Genes and Mutations. *Cell* **173,** 371–385.e18 (2018).

9. Weghorn, D. & Sunyaev, S. Bayesian inference of negative and positive selection in human cancers. *Nature Genetics* **49,** 1–8 (2017).

10. Zapata, L. *et al.* Negative selection in tumor genome evolution acts on essential cellular functions and the immunopeptidome. *Genome Biology* **19,** 1–17 (2018).

11. Wu, C.-I., Wang, H.-Y., Ling, S. & Lu, X. The Ecology and Evolution of Cancer: The Ultra-Microevolutionary Process. *Annu. Rev. Genet.* **50,** 347–369 (2016).

12. Greenman, C., Wooster, R., Futreal, P. A., Stratton, M. R. & Easton, D. F. Statistical analysis of pathogenicity of somatic mutations in cancer. *Genetics* **173,** 2187–2198 (2006).

13. Yang, Z., Ro, S. & Rannala, B. Likelihood models of somatic mutation and codon substitution in cancer genes. *Genetics* **165,** 695–705 (2003).

14. Martincorena, I. *et al.* Somatic mutant clones colonize the human esophagus with age. *Science* **57,** eaau3879–14 (2018).

15. Lee-Six, H. *et al.* Population dynamics of normal human blood inferred from somatic mutations. *Nature* **14,** 213–478 (2018).

16. Nielsen, R. & Yang, Z. Estimating the distribution of selection coefficients from phylogenetic data with applications to mitochondrial and viral DNA. *Mol Biol Evol* **20,** 1231–1239 (2003).

17. McGranahan, N. & Swanton, C. Clonal Heterogeneity and Tumor Evolution: Past, Present, and the Future. *Cell* **168,** 613–628 (2017).

18. Sottoriva, A. *et al.* A Big Bang model of human colorectal tumor growth. *Nature Genetics* **47,** 209–216 (2015).

19. Kryazhimskiy, S. & Plotkin, J. B. The Population Genetics of dN/dS. *PLOS Genet* **4,** e1000304–10 (2008).

20. Mugal, C. F., Wolf, J. B. W. & Kaj, I. Why Time Matters: Codon Evolution and the Temporal Dynamics of dN/dS. *Mol Biol Evol* **31,** 212–231 (2013).

21. Simons, B. D. Deep sequencing as a probe of normal stem cell fate and preneoplasia in human epidermis. *PNAS* **113,** 128–133 (2016).
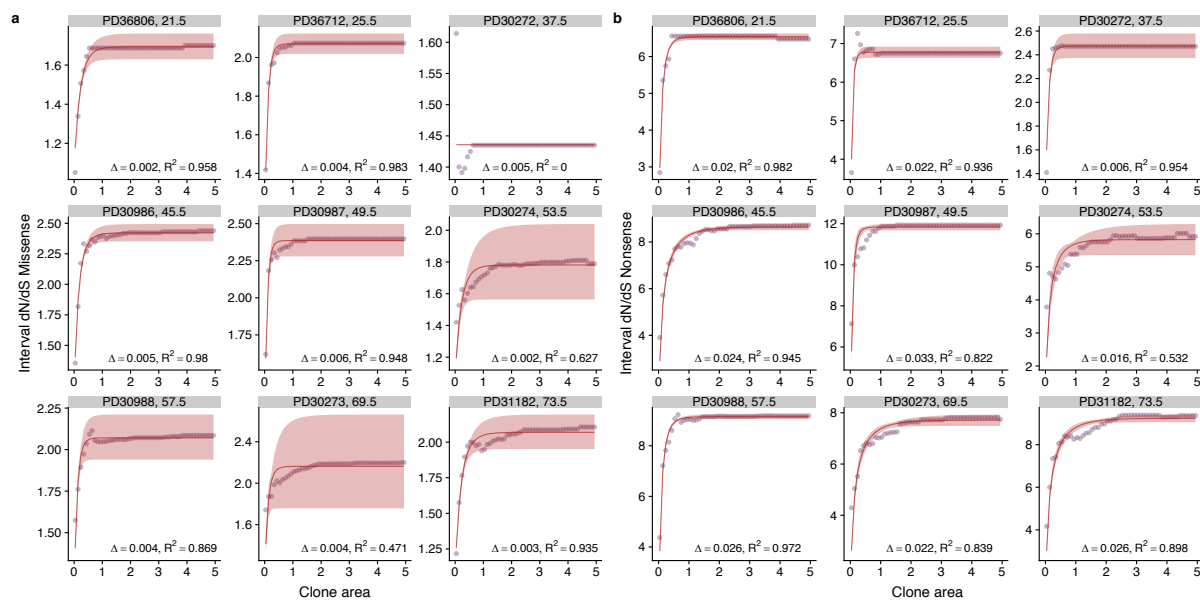
650  22.  Durrett, R. Population genetics of neutral mutations in exponentially growing cancer
651       cell populations. *The Annals of Applied Probability* **23,** 230–250 (2013).
652  23.  Klein, A. M., Brash, D. E., Jones, P. H. & Simons, B. D. Stochastic fate of p53-mutant
653       epidermal progenitor cells is tilted toward proliferation by UV B during preneoplasia.
654       *Proc. Natl. Acad. Sci. U.S.A.* **107,** 270–275 (2010).
655  24.  Lopez-Garcia, C., Klein, A. M., Simons, B. D. & Winton, D. J. Intestinal stem cell
656       replacement follows a pattern of neutral drift. *Science* **330,** 822–825 (2010).
657  25.  Vermeulen, L. *et al.* Defining stem cell dynamics in models of intestinal tumor
658       initiation. *Science* **342,** 995–998 (2013).
659  26.  Williams, M. J., Werner, B., Barnes, C. P., Graham, T. A. & Sottoriva, A. Identification
660       of neutral tumor evolution across cancer types. *Nature Genetics* **48,** 238–244 (2016).
661  27.  Williams, M. J. *et al.* Quantification of subclonal selection in cancer from bulk
662       sequencing data. *Nature Genetics* **50,** 895–903 (2018).
663  28.  Bozic, I., Gerold, J. M. & Nowak, M. A. Quantifying Clonal and Subclonal Passenger
664       Mutations in Cancer Evolution. *PLoS Comput Biol* **12,** e1004731 (2016).
665  29.  Ling, S. *et al.* Extremely high genetic diversity in a single tumor points to prevalence of
666       non-Darwinian cell evolution. *Proc. Natl. Acad. Sci. U.S.A.* **112,** E6496–505 (2015).
667  30.  Simons, B. D. Reply to Martincorena et al.: Evidence for constrained positive selection
668       of cancer mutations in normal skin is lacking. *Proc. Natl. Acad. Sci. U.S.A.* **113,** E1130–
669       E1131 (2016).
670  31.  Martincorena, I., Jones, P. H. & Campbell, P. J. Constrained positive selection on
671       cancer mutations in normal skin. *Proc. Natl. Acad. Sci. U.S.A.* **113,** E1128–E1129
672       (2016).
673  32.  Klein, A. M. & Simons, B. D. Universal patterns of stem cell fate in cycling adult
674       tissues. *Development* **138,** 3103–3111 (2011).
675  33.  Doupé, D. P. *et al.* A single progenitor population switches behavior to maintain and
676       repair esophageal epithelium. *Science* **337,** 1091–1093 (2012).
677  34.  Alcolea, M. P. *et al.* Differentiation imbalance in single oesophageal progenitor cells
678       causes clonal immortalization and field change. *Nature Cell Biology* **16,** 612–619
679       (2014).
680  35.  Nicholson, M. D. & Antal, T. Universal Asymptotic Clone Size Distribution for General
681       Population Growth. *Bull. Math. Biol.* **78,** 2243–2276 (2016).
682  36.  Ewens, W. J. Mathematical Population Genetics. 1–435 (2012).
683  37.  Martincorena, I. *et al.* Tumor evolution. High burden and pervasive positive selection
684       of somatic mutations in normal human skin. *Science* **348,** 880–886 (2015).
685  38.  Bielski, C. M. *et al.* Widespread Selection for Oncogenic Mutant Allele Imbalance in
686       Cancer. *Cancer Cell* 1–24 (2018). doi:10.1016/j.ccell.2018.10.003
687  39.  Zheng, Q. Progress of a half century in the study of the Luria–Delbrück distribution.
688       *Math Biosci* (1999). doi:10.1016/S0025-5564(99)00045-0
689  40.  Kessler, D. A. & Levine, H. Scaling Solution in the Large Population Limit of the
690       General Asymmetric Stochastic Luria–Delbrück Evolution Process. *J Stat Phys* **158,**
691       783–805 (2014).
692  41.  Gerlinger, M. *et al.* Genomic architecture and evolution of clear cell renal cell
693       carcinomas defined by multiregion sequencing. *Nature Genetics* **46,** 225–233 (2014).
694  42.  Lawrence, M. S. *et al.* Mutational heterogeneity in cancer and the search for new
695       cancer-associated genes. *Nature* **499,** 214–218 (2013).

696    43.    Caravagna, G. *et al.* Detecting repeated cancer evolution from multi- region tumor
697           sequencing data. *Nat Methods* **15,** 1–13 (2018).
698    44.    Korolev, K. S. *et al.* Selective sweeps in growing microbial colonies. *Phys Biol* **9,**
699           026008 (2012).
700    45.    Van den Eynden, J. & Larsson, E. Mutational Signatures Are Critical for Proper
701           Estimation of Purifying Selection Pressures in Cancer Somatic Mutation Data When
702           Using the dN/dS Metric. *Front. Genet.* **8,** 415–9 (2017).
703    46.    Cannataro, V. L., Gaffney, S. G. & Townsend, J. P. Effect Sizes of Somatic Mutations in
704           Cancer. *JNCI Journal of the National Cancer Institute* **110,** 1171–1177 (2018).
705    47.    Temko, D., Tomlinson, I. P. M., Severini, S., Schuster-Böckler, B. & Graham, T. A. The
706           effects of mutational processes and selection on driver mutations across cancer
707           types. *Nat Commun* **9,** 1857 (2018).
708    48.    Heide, T. *et al.* Reply to 'Neutral tumor evolution?'. *Nature Genetics* **48,** 1–9 (2018).
709    49.    Colaprico, A. *et al.* TCGAbiolinks: an R/Bioconductor package for integrative analysis
710           of TCGA data. *Nucleic Acids Research* **44,** e71–e71 (2016).
711    50.    Martincorena, I. & Campbell, P. J. Somatic mutation in cancer and normal cells.
712           *Science* **349,** 1483–1489 (2015).
713    51.    Bezanson, J., Edelman, A., Karpinski, S. & Shah, V. B. Julia: A Fresh Approach to
714           Numerical Computing. *SIAM Review (2017)*
715    52.    R Core Team, *R: A Language and Environment for Statistical Computing, R Foundation*
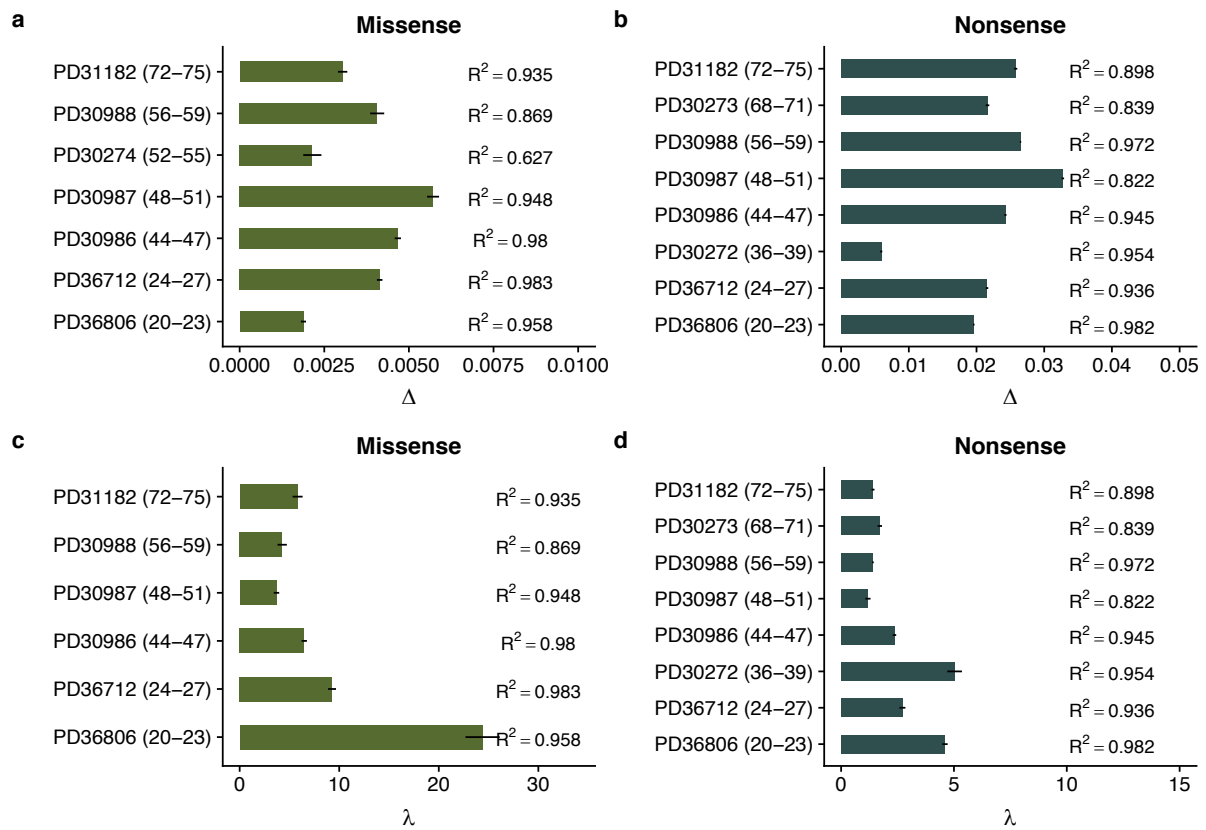716           *for Statistical Computing, Vienna.* 2018
717

**Figure S1**

Global dN/dS values in different frequency bins for patient PD31182 showing that the values depend on the frequency of mutations.
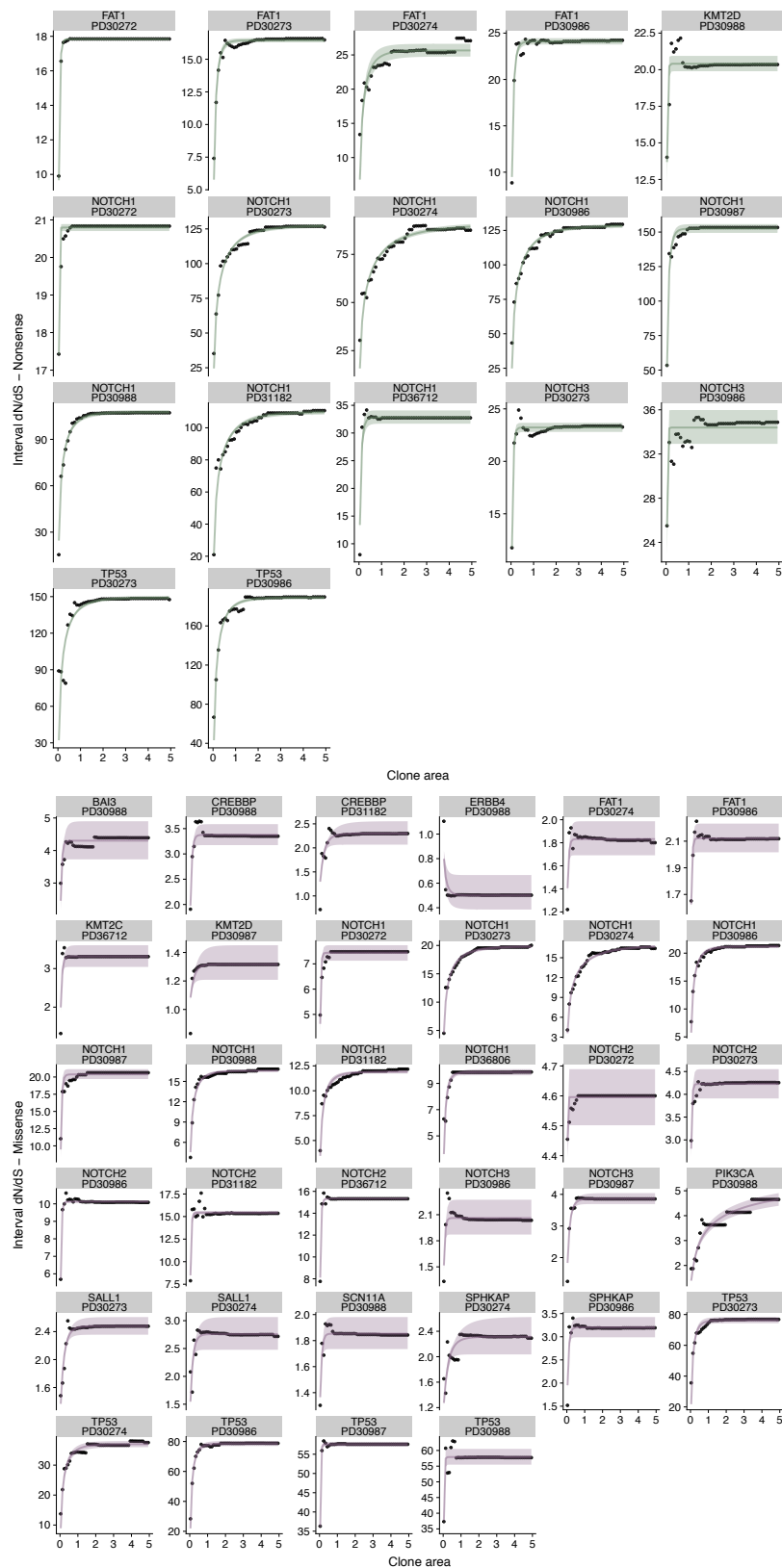


**Figure S2**

Model fits for all patients in the oesophagus data set. Purple points are data and red lines model fits. Fits were performed separately for missense, **a** and nonsense mutations, **b.** Each plot is annotated with the inferred bias $\Delta$ and the $R^2$ value.
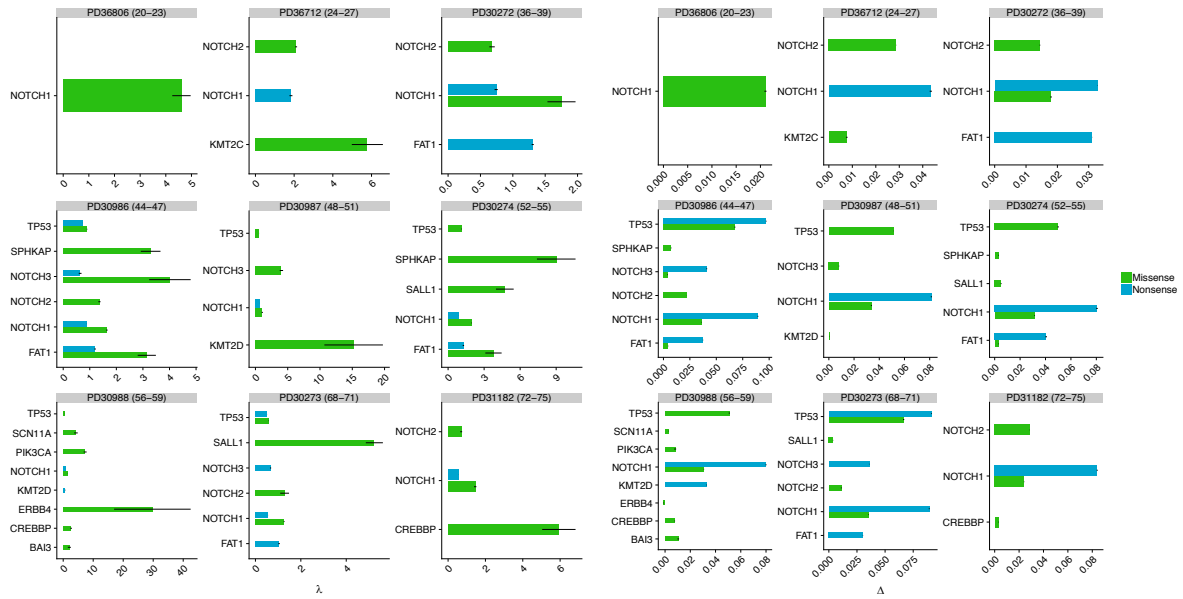
**Figure S3**

Inferred biases for for each patient in the oesophagus dataset based on missense , **a** and nonsense mutations, **b**. Inferred loss replacement rates, $\lambda$ for each patient based on missense, **a** and nonsense mutations, **b**.

**Figure S4**

Individual fits for each gene in each patient in the oesophagus dataset. Points are data and lines are model fits. Analysis performed separately for nonsense, **a** and missense, **b**.
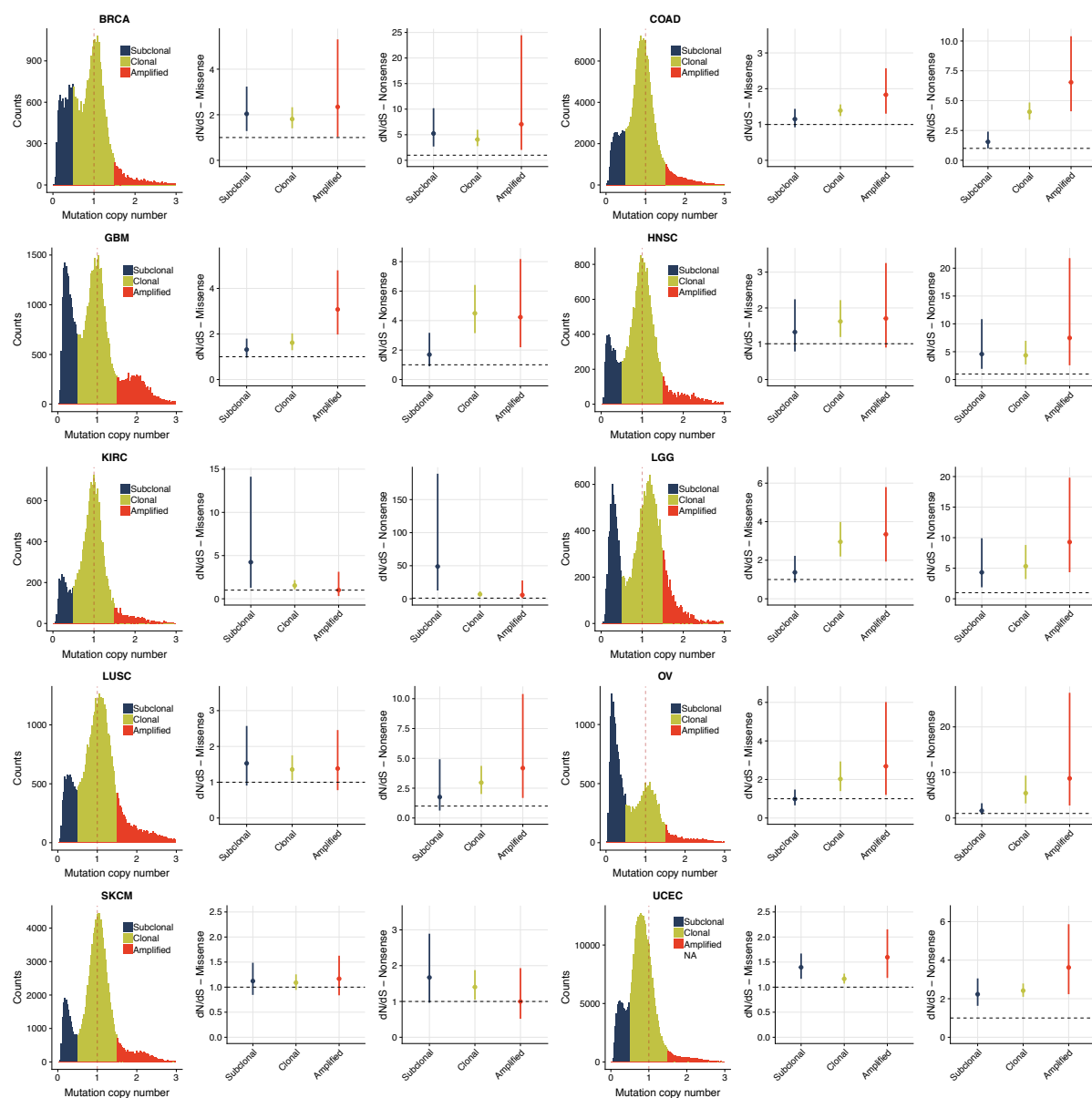
**Figure S5** Inferred parameters for each gene in each patient in the oesophagus dataset where there were sufficient mutations to perform the analysis. Left hand plot shows inferred loss replacement rates $\lambda$ and right hand plot inferred biases $\Delta$.
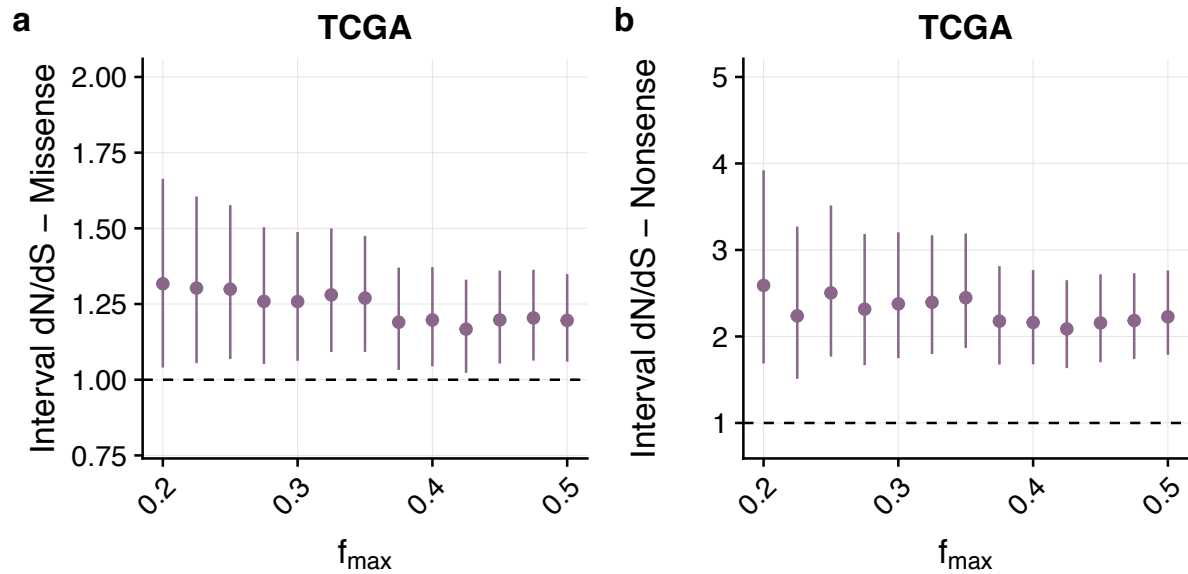


**Figure S6**
Model fits per patient and per gene per patient when there were sufficient mutations in the skin dataset. Points are data and lines are model fits, **a-e. f** shows the distributions of fitness effects for missense mutations across the cohort. There were insufficient nonsense mutations in the majority of genes to draw the equivalent plot for nonsense mutations.
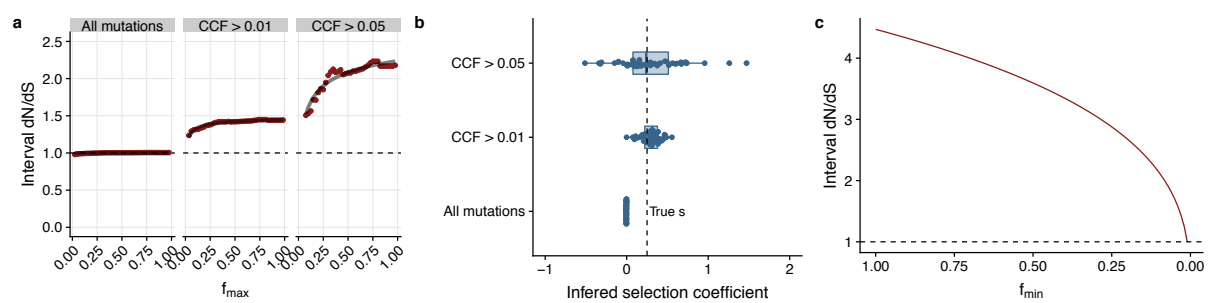
**Figure S7**

Mutation copy number histograms and dN/dS values for different cancer types with >100 samples (post filtering) in TCGA. Histograms and dN/dS plots coloured by mutation clonality.
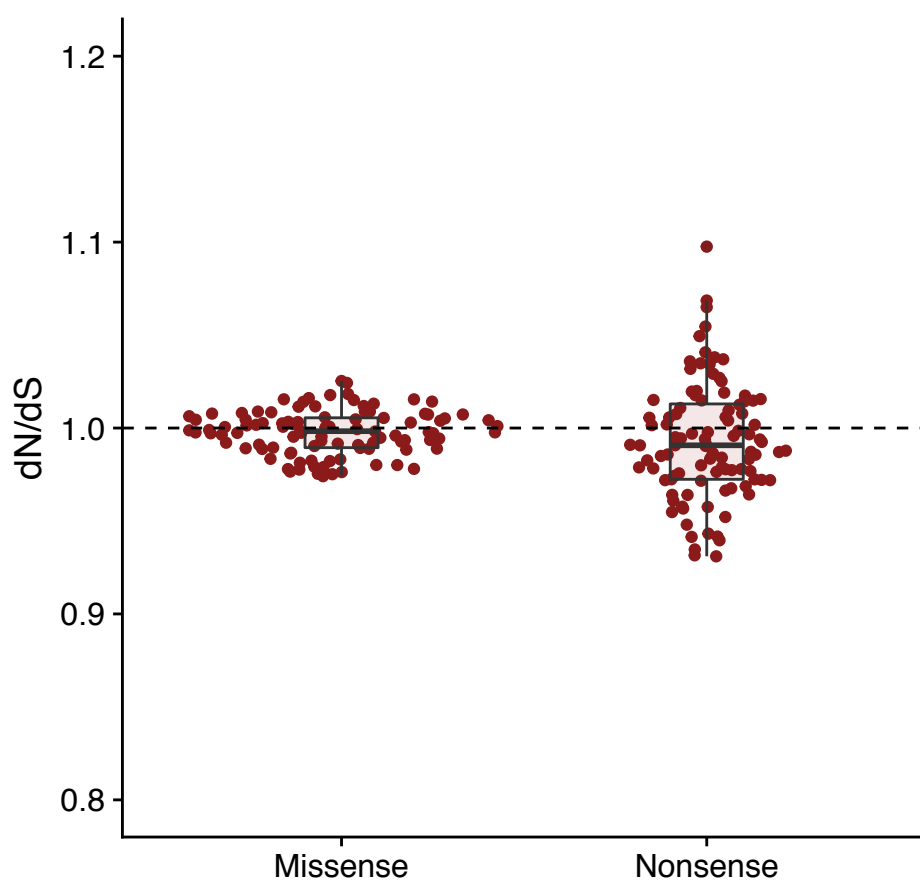
**Figure S8**

Interval dN/dS for 192 high confidence driver mutations. We observe no patterns that are predicted by our theoretical model.



**Figure S9**

Generating a synthetic cohort with selection and using all mutations to infer dN/dS values shows that in this case dN/dS~1, while if we restrict our attention to high frequency variants dN/dS>1, **a**. Inferred selection coefficients are accurate only when using high frequency variants, **b**. Using our theoretical interval model equation we see that fixing $f_{min} = 1$ and taking the limit $f_{min} \rightarrow 0$ results in dN/dS = 1.

**Figure S10**

Corrected dN/dS values from 100 sets of 1000 randomly samples genes. Average dN/dS ~ 1 as would be expected.