# Queen Mary
## University of London

# Identification of neutral tumor evolution across cancer types.

Williams, MJ; Werner, B; Barnes, CP; Graham, TA; Sottoriva, A

For additional information about this publication click this link.
http://qmro.qmul.ac.uk/xmlui/handle/123456789/13290

# Identification of neutral tumor evolution across cancer types

*Marc J Williams[1,3,4,6], Benjamin Werner[2,6], Chris P Barnes[4,5], Trevor A Graham[1], Andrea Sottoriva[2]*

[1] Evolution and Cancer Laboratory, Barts Cancer Institute, Queen Mary University of London, Charterhouse Square, London, EC1M 6BQ, UK
[2] Centre for Evolution and Cancer, The Institute of Cancer Research, London, SM2 5NG, UK
[3] Department of Cell and Developmental Biology, University College London, London WC1E 6BT, UK
[4] Centre for Mathematics and Physics in Life Sciences and Experimental Biology (CoMPLEX), University College London, London, WC1E 6BT, UK
[5] Department of Genetics, Evolution and Environment, University College London, Gower Street, London, WC1E 6BT, UK

[6] These authors contributed equally to this work

Correspondence should be addressed to T.A.G. (t.graham@qmul.ac.uk) or A.S. (andrea.sottoriva@icr.ac.uk)

## Abstract
Despite extraordinary efforts to profile cancer genomes, interpreting the vast amount of genomic data in the light of cancer evolution remains challenging. Here we demonstrate that neutral tumor evolution results in a power-law distribution of the mutant allele frequencies reported by next-generation sequencing of tumor bulk samples. We find that the neutral power law fits with high precision 323 of 904 cancers from 14 types, selected from different cohorts. In malignancies identified as neutral, all clonal selection occurred prior to the onset of cancer growth and not in later-arising subclones, resulting in numerous passenger mutations that are responsible for intra-tumor heterogeneity. Reanalyzing cancer sequencing data within the neutral framework allowed the measurement, in each patient, of both the *in vivo* mutation rate and the order and timing of mutations. This result provides a new way to interpret existing cancer genomic data and to discriminate between functional and non-functional intra-tumor heterogeneity.

## Introduction
Unraveling the evolutionary history of a tumor is clinically valuable, as prognosis depends on the future course of the evolutionary process[1,2], and therapeutic response is determined by the evolution of resistant subpopulations[3]. In humans, the details of tumor evolution have remained largely uncharacterized as longitudinal measurements are impractical, and studies are complicated by inter-patient variation[4] and intra-tumor heterogeneity (ITH)[5,6]. Several recent studies have begun tackling this complexity[7], revealing patterns of convergent evolution[8], punctuated dynamics[9], and intricate interactions between cancer cell populations[10]. However, the lack of a rigorous theoretical framework able to make predictions on existing data[11] means that results from cancer genomic profiling studies are often difficult to interpret. For example, how much of the

53   detected intra-tumor heterogeneity is actually functional is largely unknown, also
54   because a rigorous 'null model' of genomic heterogeneity is lacking. In particular,
55   interpreting the mutant allele frequency distribution reported by next-generation
56   sequencing (NGS) is problematic because of the absence of a formal model
57   linking tumor evolution to the observed data. Therefore, making sense to the
58   wealth of available sequencing data in cancer remains challenging.
59       Here we show that the subclonal mutant allele frequencies of a significant
60   proportion of cancers of different types and from different cohorts precisely follow
61   a simple power-law distribution predicted by neutral growth. In those neutral
62   cancers, all tumor-driving alterations responsible for cancer expansion were
63   present in the first malignant cell and subsequent tumor evolution was effectively
64   neutral. We demonstrate that under neutral growth, the fundamental parameters
65   describing cancer evolution that have been so far inaccessible in human tumors,
66   such as the mutation rate and the mutational timeline, become measurable.
67   Importantly, this approach allows identifying also non-neutral malignancies, in
68   which ongoing clonal selection and adaption to microenvironmental niches may
69   play a strong role during cancer growth.

## *Results*

70

71

72   **Neutral cancer growth**
73       Recently, we showed that colorectal cancers (CRC) often grow as a single
74   expansion, populated by a large number of intermixed subclones[12].
75   Consequently, we expect that after malignant transformation, individual
76   subclones with distinct mutational patterns grow at similar rates, coexisting within
77   the tumor for long periods of time without overtaking one another. Indeed, only a
78   handful of recurrent driver alterations have been identified in CRC[13], and those
79   are reported to be ubiquitous in multi-region sampling[12] and stable during cancer
80   progression[14], indicating that they all occurred in the "first" cancer cell and that
81   subsequent clonal outgrowths are relatively rare. Consequently, we hypothesized
82   that cancer evolution may often be dominated by neutral evolutionary dynamics.
83       The dynamics of neutral evolutionary processes have been widely studied
84   in the context of molecular evolution and population genetics[15-17] as well as in
85   mouse models of cancer[18]. However, the widely held presumption that subclone
86   dynamics in human cancers are dominated by strong selection has meant these
87   ideas have been neglected in current studies of cancer evolution.
88       Motivated by this, here we present a theoretical model describing the
89   expected pattern of subclonal mutations within a tumor that is evolving according
90   to neutral evolutionary dynamics. The model postulates that, after the
91   accumulation of a "full house" of genomic changes that initiates tumor growth,
92   some tumors expand neutrally, generating a large number of passenger
93   mutations that are responsible for the extensive and common ITH. The
94   parameter-free model is applicable to NGS data from any solid cancer. Here we
95   present the model, and by applying it to large pre-existing cancer genomics
96   datasets, determine which tumors are consistent with neutral growth. When the
97   model applies, we measure new tumor characteristics directly from the patient's
98   data.

99

100  **Model derivation**
101      A tumor is founded by a single cell that has already acquired a significant
102  mutation burden[4]: these "pre-cancer" mutations will be borne by every cell in the
103  growing tumor, and so become "public" or clonal. Mutations that occur within
104  different cell lineages remain "private" or subclonal in an expanding malignancy
105  under the absence of strong selection. We focus on the latter as they contain
106  information on the dynamics of the cancer growth. We denote the number of
107  tumor cells at time $t$ as $N(t)$ which divide at rate $\lambda$ per unit time. During a cell
108  division, somatic mutations may occur with a probability $\mu$. If we consider an
109  average number of $\pi$ chromosome sets in a cancer cell (e.g. the ploidy of the
110  cell), we can calculate the expected number of new mutations per time interval
111  as:

112
$$\frac{dM}{dt} = \mu\pi\lambda N(t) \qquad [1]$$

113
114 Solving this requires integrating over the growth function $N(t)$ in some time
115 interval $[t_0, t]$:

116

117
$$M(t) = \mu\pi\lambda \int_{t_0}^{t} N(t)\,dt \qquad [2]$$

118
119 Since not all cell divisions may be successful in generating two surviving lineages
120 due to cell death or differentiation, we introduce the fraction $\beta$ of "effective" cell
121 divisions in which both resulting lineages survive. In the case of exponential
122 growth, the mean number of tumor cells as a function of time is therefore:
123

124
$$N(t) = e^{\lambda\beta t} \qquad [3]$$

125
126 Substituting into equation [2] gives the explicit solution:
127

128
$$M(t) = \frac{\mu\pi}{\beta}\left(e^{\lambda\beta t} - e^{\lambda\beta t_0}\right) \qquad [4]$$

129
130 This equation describes the total number of subclonal mutations that accumulate
131 within a growing tumor in the time interval $[t_0, t]$. We note that for $t_0=0$ equation [4]
132 corresponds to the Luria-Delbrück model, which describes mutation accumulation
133 in bacteria[19]. In our case, this equation is of limited use as none of the
134 parameters $\mu$, $\lambda$, $\beta$ or the age of the tumor $t$ can be measured directly in humans.
135 However, we do know that for a new mutation occurring at any time $t$, its allelic
136 frequency (the relative fraction) $f$ must be the inverse of the number of alleles in
137 the population:
138

139
$$f = \frac{1}{\pi N(t)} = \frac{1}{\pi e^{\lambda\beta t}} \qquad [5]$$

140
141 For example, if a new mutation arises in a tumor of 100 cells, it will comprise a
142 fraction of 1/100. In the absence of clonal selection (or indeed significant genetic
143 drift), the allelic frequency of a mutation will remain constant during the
144 expansion, as all cells, with and without this mutation, grow at the same rate. In
145 the previous example, after one generation has elapsed we will have 2 cells with
146 that particular mutation, but a total of 200 tumor cells, again a fraction of 1/100.
147 This implies that in the neutral case, tumor age $t$ and mutation frequency $f$ are
148 *interchangeable*. For example, $t_0=0$ in a diploid tumor ($\pi=2$), corresponds to
149 $f_{max}=0.5$ (the expected allelic frequency of clonal variants):
150

151
$$f_{max} = \frac{1}{\pi e^{\lambda\beta t_0}} \qquad [6]$$

152
153 Substituting $t$ for $f$ in equation [4] gives an expression for the cumulative number
154 of mutations in the tumor per frequency $M(f)$:
155

156
$$M(f) = \frac{\mu}{\beta}\left(\frac{1}{f} - \frac{1}{f_{max}}\right) \qquad [7]$$

157
158 thus converging to the solution for expanding populations under neutrality
159 obtained using other approaches[20-23]. Critically, the distribution $M(f)$ is naturally
160 provided by NGS data from bulk sequencing of tumor biopsies and resections,
161 against which the model can be tested. The model predicts that mutations arising
162 during a neutral expansion of a cancer accumulate following a *1/f* power-law
163 distribution. In other words, when neutral evolution occurs in a tumor, the number

3

164 of mutations detected should accumulate linearly with the inverse of their
165 frequency. The *1/f* noise or *pink noise* is common in nature and found in several
166 physical, biological and economic systems[24].
167      Importantly, the coefficient $\mu_e=\mu/\beta$ is the mutation rate per effective cell
168 division, and corresponds to the easily measureable slope of *M(f)*. This model
169 therefore provides a straightforward parameter-free method to measure the *in*
170 *vivo* mutation rate in a patient's tumor using a single NGS sample. We note that
171 the results do not depend on the identity of the alterations considered, since any
172 genomic alteration (mutations, copy number changes or epigenetic modifications)
173 anywhere in the genome that changes the dynamics of tumor growth (e.g. any
174 alteration that is clonally selected) would result in deviation from the neutral *1/f*
175 power law by causing an over- or under-representation of the alleles in that
176 clone. Hence, here we use single nucleotide variants as 'barcodes' to follow
177 clone growth. Stochastic simulations of neutral tumor growth confirm the
178 analytical solution in equation [7] (see Online Methods).
179
**Identification of neutrality in colorectal cancer evolution**
181      A typical allelic frequency distribution of mutations in a tumor measured by
182 NGS whole-exome sequencing is shown in Figure 1A (data from ref [12]).
183 Considering tumor purity and aneuploidy, mutations with high allelic frequency
184 (>0.25) are likely to be public (clonal) while all others are likely subclonal. The
185 same data can be represented as the cumulative distribution *M(f)* of subclonal
186 mutations as in equation [7] (Figure 1B). Remarkably, as reported by the high
187 goodness-of-fit measure $R^2$, these data precisely follow the distribution predicted
188 by the model indicating that this tumor grew with neutral evolutionary dynamics.
189      We next considered our cohort of 7 multi-sampling CRCs[12] and 101 TCGA
190 colon adenocarcinomas[13] selected for high tumor purity (≥70%) that underwent
191 whole-exome sequencing (see Online Methods). The latter were separated
192 between tumors characterized by chromosomal instability (CIN) versus
193 microsatellite instability (MSI). The power-law is remarkably well supported in
194 both these cohorts, with 38/108 (35.1%) of the cases reporting a high $R^2 \geq 0.98$
195 (Figure 1C). These results confirm that in a large proportion of colon cancers,
196 intra-tumor clonal dynamics are not dominated by strong selection but rather
197 follow neutral evolution. In particular, a larger proportion of CIN cancers evolved
198 neutrally (31/82, 37.8%) than MSI cancers (3/19, 15.7%) (Figure 1C), possibly
199 because the latter acquired so many new mutations that some are likely under
200 strong selection. Since *M(f)* is a monotonic growing function, this stringent
201 threshold of $R^2 > 0.98$ was chosen to prevent over-calling neutrality, but we note
202 that we may have therefore misclassified some tumors as non-neutral due to
203 limited sequencing depth or low mutation burden. $R^2$ values were independent
204 from the mean coverage of mutations, the total number of mutations in the
205 sample or the number of mutations within the model range (see Online Methods).
206 See Supplementary Data Set 1 (summary of TCGA data used).
207
**Measurement of the mutation rate in colorectal cancer**
209      Estimating the per-base mutation rate *μ* per division in human
210 malignancies is challenging since direct measurements are not possible.
211 Previous estimates critically depend on assumptions about the cell cycle time
212 and the growth rate *λ*, as well as on the *total* mutational burden of the cancer[25-27].
213 However, accurate measurement of all mutations within a cancer, including
214 heterogeneous subclonal variants, is technically unfeasible since most mutations
215 are present in very small numbers of cells[5]. With our approach it is possible to
216 circumvent this issue by measuring the rate of accumulation of subclonal
217 mutations represented by the slope of *M(f)*. In the case of neutral evolution, this
218 can be done in principle within any (subclonal) frequency range, without the need
219 of detecting extremely rare mutations. We estimated the mutation rate in all
220 samples with $R^2 \geq 0.98$ (Figure 1D) and found that it was more than 15-fold higher
221 in the MSI group (median: $\mu_e=3.65\times10^{-6}$) with respect to the CIN group (median:
222 $\mu_e=2.31\times10^{-7}$; F-test: $p=2.24\times10^{-8}$) and our cohort of CRCs (median: $\mu_e=2.07\times10^{-7}$
223 ), which was comprised of all but one CIN tumors[12]. Different mutational types
224 (e.g. transitions or transversions) are caused by particular mutational
225 processes[28], and so likely occur at different rates and accordingly we found that

226  C>T mutations occurred at median $\mu_{e,C>T}$=2.19×10$^{-7}$, a rate nearly 10-fold higher
227  than any other type of mutation (F-test: $p$=3.13×10$^{-3}$; Supplementary Figure 1A).
228  We stratified according to CIN versus MSI and found that the mutation rate of
229  each mutational type reflected the overall mutation rate for the group
230  (Supplementary Figure 1B). The variation in mutation rates within and between
231  subgroups was remarkably in line with the variation in estimates of mutational
232  burden in colon cancer[4]. We note the mutation rate estimate is scaled by the
233  (unknown) effective division rate $\beta$, which means for example that if only 1 in 100
234  cell divisions leads to two surviving offspring ($\beta$=0.01), then the mutation rate $\mu$ is
235  100 times lower than the effective rate $\mu_e$ reported. Importantly, mutation rates of
236  non-neutral cases ($R^2$<0.98) cannot be estimated, as the model does not fit the
237  dynamics of these tumors.
238          We examined the effect of copy-number changes in the model by
239  performing the analysis using only mutations in diploid regions and found highly
240  similar proportions of neutral tumors and mutation rates (see Online Methods and
241  Supplementary Figure 2). The validity of the variant calls was also corroborated
242  by the consistency of the underlying mutational signature across a range of allelic
243  frequencies; hence the results are unlikely to be influenced by sequencing errors
244  (Supplementary Figure 3).
245          Frequent selection events should induce a higher number of missense and
246  nonsense mutations than expected by chance whereas under neutrality we
247  expect the same rate of silent and non-silent mutations. To test this, we
248  contrasted the estimated rate of synonymous mutations (unlikely to ever be
249  under selection) versus the rate of missense and nonsense mutations (liable to
250  experience selection). Although the latter are more common than the former,
251  after adjustment for the number of potential synonymous and non-synonymous
252  sites in the exome, the two rates were equivalent (Supplementary Figure 4),
253  consistent with neutral evolution.
254
255  **Neutral evolution in coding and non-coding regions**
256          We next tested whether the signature of neutral evolution could be found
257  across the entire genome, not just in coding regions. To do this, we analyzed 78
258  gastric cancers from a recent study[29] subjected to high depth whole-genome
259  sequencing. The large number of mutations detected by WGS accumulated
260  precisely as predicted by the model (example in Figure 2A,B), revealing neutral
261  evolution in 60/78 (76.9%) cases (Figure 2C). A smaller proportion of MSI tumors
262  were neutral (3/10, 30%) than microsatellite stable (MSS) tumors (57/68, 83.8%)
263  consistent with the observation in CRC. A tumor was consistently classified as
264  neutral independently of whether all SNVs or only non-coding SNVs were used to
265  perform the classification (Figure 2C, Venn diagram), whereas due to the limited
266  number of mutations available in the exome alone, fewer tumors were identified
267  as neutral.  Importantly, every case was verified as neutral by at least two
268  different variant sets. These results confirm that neutral evolution can be robustly
269  assessed from mutations anywhere in the genome.
270          Mutation rate analysis of the neutrally evolved gastric cancers revealed
271  that MSI cancers had a more than 4-fold higher mutation rate ($\mu_e$=3.30×10$^{-6}$) with
272  respect to MSS ($\mu_e$=7.82×10$^{-7}$; F-test: $p$=1.35×10$^{-4}$). Results were robust to copy
273  number changes when the analysis was performed only using variants in diploid
274  regions (Supplementary Figure 5). The mutational signature of the variant calls
275  for this cohort was also consistent across the frequency spectrum
276  (Supplementary Figure 6). Synonymous versus nonsynonymous mutation rates
277  were also not consistent with frequent on-going selection (Supplementary Figure
278  7). See Supplementary Data Set 2 (summary of Wang et al. data used).
279
280  **Neutral evolution across cancer types**
281          We then applied our neutral model to a large pan-cancer cohort of 819
282  exome-sequenced cancers from 14 tumor types from the TCGA consortium
283  (which included the 101 colon cancers previously examined). All of these
284  samples had been pre-selected for high tumor purity (≥70%). The fit of the model
285  was remarkably good across types (Figure 3A) with 259/819 (31.6%) cases
286  showing $R^2$≥0.98. We found that neutral evolution was more prominent in some
287  tumor types, such as stomach (validating the WGS analysis), lung, bladder,

288 cervical, and colon. Others showed a consistently poorer fit, indicating that the
289 clonal dynamics in these malignancies were typically not neutral, such as renal,
290 melanoma, pancreatic, thyroid, and glioblastoma. Consistent with these results,
291 "non-neutral" renal carcinoma has been shown to display convergent evolution in
292 spatially disparate tumor regions driven by strong selective forces[8], whereas the
293 same phenomenon was not found in more "neutral" lung cancer[30,31]. Other types
294 displayed mixed dynamics, with some cases that were characterized by neutral
295 evolution and some that were not. We note that a proportion of melanoma
296 samples in this cohort are derived from regional metastases and not primary
297 lesions, and this could potentially explain the lack of neutral dynamics observed.
298      Mutation rate analysis on the neutral cases showed differences of more
299 than an order of magnitude between types (Figure 3B). The highest mutation
300 rates were observed in lung adenocarcinoma (median $\mu_e$=6.79×10$^{-7}$) and in lung
301 squamous cell carcinoma (median $\mu_e$=5.61×10$^{-7}$) and the lowest rates in low
302 grade glioma (median $\mu_e$=9.22×10$^{-8}$) and in prostate (median $\mu_e$=1.04×10$^{-7}$). We
303 stratified the mutation rates into different mutational types (Supplementary Figure
304 8) and found that C>A mutations occurred at a significantly higher rate in lung
305 cancers, consistent with their causation by tobacco smoke[28]. C>T mutation rates
306 were most consistent across cancer types, likely because of their association
307 with normal replicative errors, as opposed to being caused by a particular
308 stochastically-arising defect in DNA replication or repair[28].
309      These results demonstrate that within-tumor clonal dynamics can be
310 neutral, and the classification of tumors based on neutral versus non-neutral
311 growth dynamics leads to new measurements of fundamental tumor biology. See
312 See Supplementary Data Set 1 (summary of TCGA data used).
313
314 ***In silico* validation of the neutral model**
315 To assess the different inherent sources of noise in NGS data (normal
316 contamination, limited sequencing depth, tumor sampling), we designed a
317 stochastic simulation of neutral growth that produced synthetic NGS data from
318 bulk samples (see Online Methods). The simulations produced realistic synthetic
319 NGS data (Supplementary Figure 9) with minimal assumptions and under a
320 range of different scenarios for tumor growth dynamics (variable low mutation
321 rate, variable number of clonal mutations) and sources of assay noise (normal
322 contamination in the sample, sequencing depth, detection limit). For each of
323 these potentially confounding factors, we were able to fit our neutral model to the
324 synthetic NGS data and accurately recover both the underlying neutral dynamics
325 and mutation rate (Supplementary Figure 10). We also validated the prediction
326 that *M(f)* would deviate from the neutral power law in the presence of emerging
327 subclones with a higher fitness advantage (Supplementary Figure 11A,B), as well
328 as in the case of a mixture of subclones (as observed in ref. [32]) emerging either
329 by means of clonal expansions triggered by selection, or by segregating
330 microenvironmental niches (Supplementary Figure 11C-F). Variation of mutation
331 rate between subclones also causes a deviation from neutrality (Supplementary
332 Figure 11G,H). These results confirm the reliability of the conservatively high $R^2$
333 threshold used to call neutrality.
334
335 **Mutational timelines**
336      Under neutral evolution, it is possible to estimate the size of the tumor
337 when a mutation with frequency *f* arose from equation [5]:
338

339 $$N(t) = \frac{1}{\pi f} \qquad\qquad [8]$$

340
341 Figure 4A,B shows the decomposition of the mutational timeline for two
342 illustrative cases: sample TB from[12] and sample TCGA-AA-3712 from[13]. Previous
343 estimates of mutational timelines relied on cross-sectional data[33-36] that are
344 compromised by the extensive heterogeneity, whereas multi-region profiling
345 approaches are instead more accurate but expensive and laborious[8,37,38]. Using
346 our formal model of cancer evolution this timeline information becomes
347 accessible from routinely available genomic data. We found that classical CRC
348 driver alterations, such as in the *APC*, *KRAS* and *TP53* genes, were indeed

349 present in the first malignant cell (likely because they accumulated during
350 previous neoplastic stages). This confirms what we previously reported using
351 single-gland mutational profiling where all these drivers, when present, were
352 found in all glands[12]. However, we also found that when we considered a more
353 extended list of putative drivers, many occurred during the neutral phase of tumor
354 growth, suggesting that the selective advantage conferred by a putative driver
355 alteration may be context-dependent, as demonstrated in a *p53* murine model[39].

## *Discussion*

357      Understanding the evolutionary dynamics of subclones within human
358 cancers is challenging because longitudinal observations are unfeasible and the
359 genetic landscape of cancer is highly dynamic, leading to genomic data that are
360 hard to interpret[40]. In particular, complex non-linear evolutionary trajectories have
361 been observed, such as punctuated evolution and karyotypic chaos[9,40,41]. Here
362 we have presented a formal law that predicts mutational patterns routinely
363 reported in NGS of bulk cancer specimens. Our analysis of large independent
364 cohorts using this framework shows that cancer growth is often dominated by
365 neutral evolutionary dynamics, an observation that is consistent across 14 cancer
366 types. Under neutrality, the clonal structure of a tumor is expected to have a
367 fractal topology characterized by self-similarity (Figure 5). As the tumor grows, a
368 large number of cell lineages are generated and therefore ITH rapidly increases
369 while the allele frequency of the new heterogeneous mutations quickly decreases
370 due to the expansion. This implies that sampling in different parts of the tree
371 leads to the detection of distinct mutations which all show the same *1/f*
372 distribution. Clonal mutations found in a sample (not considered in the model)
373 belong to the most recent common ancestor in the tree.
374      We note that some cancers were dominated by neutral evolution whereas
375 others were not. In non-neutral tumors, strong selection, microenvironmental
376 constrains and non-cell autonomous effects[42] may play a key role. Importantly,
377 our formalization represents the 'null model' of cancer intra-clone heterogeneity
378 that can be used to identify those cases in which complex non-neutral dynamics
379 occur, and to discriminate between functional and non-functional intra-tumor
380 heterogeneity. Furthermore, we speculate that neutral evolutionary dynamics
381 may be favored by the cellular architecture of the tumor (e.g. glandular structures
382 that limit the effects of selection) and/or the anatomical location of the malignancy
383 (e.g. growing in a lumen versus growing in a highly confined space), as well as
384 the presence of potentially selective microenvironmental features of the tumor
385 such as hypoxic regions. Despite the evidence for lack of natural selection during
386 malignant growth, eventual treatment is likely to "change the rules of the game"
387 and strongly select for treatment resistant clones. The same may happen in the
388 context of the purported evolutionary bottleneck preceding metastatic
389 dissemination, wherein treatment-resistance driver alterations that were not
390 under selection during growth may expand due to new selective pressures
391 introduced by therapy. Importantly, this reasoning highlights how 'drivers' can
392 only defined within a context, and so the same 'driver' alteration can be neutral in
393 a certain microenvironmental context (e.g. absence of treatment), and not neutral
394 in another (e.g. during treatment). Moreover, we predict that if a tumor is
395 characterized by different microenvironmental niches but still presents as neutral,
396 it is likely that adaptation will be driven by cancer cell plasticity, rather than clonal
397 selection. Cell plasticity is hard to study in cancer because it implies a change in
398 the cell phenotype that is not caused by any inheritable change (genomic or
399 epigenomic). This means that this phenomenon has been so far largely
400 neglected in cancer. As neutrality can be used as the 'null model' with which to
401 identify clonal selection, this facilitates the study of adaptation through plasticity
402 directly in human malignancies.
403      Furthermore, it is important to note that due to the intrinsic sub-clonal
404 detection limits of sequencing technologies, it is possible to explore only the early
405 expansion of cancer clones (Figure 5) and hence the dynamics of small clones
406 may differ from the tumor bulk as a whole.
407      Importantly, the realization that the within-tumor clonal dynamics are
408 neutral means that the *in vivo* mutation rate per division and the mutational

409  timeline, factors that play a key role in cancer evolution, progression and
410  treatment resistance can be inferred without the need to assume cell division
411  rates. These measurements can be performed in a patient-specific manner and
412  so may be useful for prognostication and the personalization of therapy.
413  Recognizing that the growth of a neoplasm is dominated by neutral clonal
414  dynamics provides an analytically tractable and rigorous method to study cancer
415  evolution and gain clinically relevant insight from commonly available genomic
416  data.

## Accession Codes

## Acknowledgements

## Contributions

439  MJW and BW contributed to the development of the model. MJW designed and
440  performed computational simulations with support from CPB. MJW, AS and TAG
441  analyzed the data. CPB contributed to the analysis. TAG and AS jointly
442  conceived, designed and developed the model, interpreted the results and wrote
443  the manuscript.

## References

446  1.  Basanta, D. & Anderson, A. R. A. Exploiting ecological principles to better
447      understand cancer progression and treatment. *Interface Focus* **3,**
448      20130020 (2013).
449  2.  Almendro, V. *et al.* Inference of tumor evolution during chemotherapy by
450      computational modeling and in situ analysis of genetic and phenotypic
451      cellular diversity. *Cell Rep* **6,** 514–527 (2014).
452  3.  Greaves, M. & Maley, C. C. Clonal evolution in cancer. *Nature* **481,** 306–
453      313 (2012).
454  4.  Vogelstein, B. *et al.* Cancer genome landscapes. *Science* **339,** 1546–1558
455      (2013).
456  5.  Burrell, R. A., McGranahan, N., Bartek, J. & Swanton, C. The causes and
457      consequences of genetic heterogeneity in cancer evolution. *Nature* **501,**
458      338–345 (2013).
459  6.  Marusyk, A., Almendro, V. & Polyak, K. Intra-tumour heterogeneity: a
460      looking glass for cancer? *Nat. Rev. Cancer* **12,** 323–334 (2012).
461  7.  Polyak, K. Tumor Heterogeneity Confounds and Illuminates: A case for
462      Darwinian tumor evolution. *Nat. Med.* **20,** 344–346 (2014).

463 8. Gerlinger, M. *et al.* Intratumor Heterogeneity and Branched Evolution
464 Revealed by Multiregion Sequencing. *N. Engl. J. Med.* **366,** 883–892
465 (2012).
466 9. Baca, S. C. *et al.* Punctuated evolution of prostate cancer genomes. *Cell*
467 **153,** 666–677 (2013).
468 10. Tabassum, D. P. & Polyak, K. Tumorigenesis: it takes a village. *Nat. Rev.*
469 *Cancer* **15,** 473–483 (2015).
470 11. Shou, W., Bergstrom, C. T., Chakraborty, A. K. & Skinner, F. K. Theory,
471 models and biology. *eLife Sciences* **4,** e07158 (2015).
472 12. Sottoriva, A. *et al.* A Big Bang model of human colorectal tumor growth.
473 *Nat. Genet.* **47,** 209–216 (2015).
474 13. The Cancer Genome Atlas. Comprehensive molecular characterization of
475 human colon and rectal cancer. *Nature* **487,** 330–337 (2012).
476 14. Jesinghaus, M. *et al.* Distinctive Spatiotemporal Stability of Somatic
477 Mutations in Metastasized Microsatellite-stable Colorectal Cancer. *The*
478 *American Journal of Surgical Pathology* **8,** 1140–1147 (2015).
479 15. Ohta, T. & Gillespie, J. Development of Neutral and Nearly Neutral
480 Theories. *Theor Popul Biol* **49,** 128–142 (1996).
481 16. P Donnelly, A. & Tavaré, S. Coalescents and Genealogical Structure Under
482 Neutrality. *Annual Review of Genetics* **29,** 401–421 (2003).
483 17. Durrett, R. & Schweinsberg, J. Approximating selective sweeps. *Theor*
484 *Popul Biol* **66,** 129–138 (2004).
485 18. Driessens, G., Beck, B., Caauwe, A., Simons, B. D. & Blanpain, C. Defining
486 the mode of tumour growth by clonal analysis. *Nature* **488,** 527–530
487 (2012).
488 19. Luria, S. E. & Delbrück, M. Mutations of bacteria from virus sensitivity to
489 virus resistance. *Genetics* **28,** 491–511
490 20. Griffiths, R. C. & Tavaré, S. The age of a mutation in a general coalescent.
491 *Communications in Statistics. Part C: Stochastic Models* **14,** 273–295
492 (1998).
493 21. Maruvka, Y. E., Kessler, D. A. & Shnerb, N. M. The Birth-Death-Mutation
494 Process: A New Paradigm for Fat Tailed Distributions. *PLoS ONE* **6,**
495 e26480 (2011).
496 22. Durrett, R. Population genetics of neutral mutations in exponentially
497 growing cancer cell populations. *The Annals of Applied Probability* **23,**
498 230–250 (2013).
499 23. Kessler, D. A. & Levine, H. Large population solution of the stochastic
500 Luria-Delbruck evolution model. *Proc. Natl. Acad. Sci. U.S.A.* **110,** 11682–
501 11687 (2013).
502 24. Bak, P., Tang, C. & Wiesenfeld, K. Self-organized criticality: An explanation
503 of the 1/f noise. *Phys. Rev. Lett.* **59,** 381–384 (1987).
504 25. Jones, S. *et al.* Comparative lesion sequencing provides insights into tumor
505 evolution. *Proc. Natl. Acad. Sci. U.S.A.* **105,** 4283–4288 (2008).
506 26. Bozic, I. *et al.* Accumulation of driver and passenger mutations during
507 tumor progression. *Proc. Natl. Acad. Sci. U.S.A.* **107,** 18545–18550 (2010).
508 27. Sun, S., Klebaner, F. & Tian, T. A new model of time scheme for
509 progression of colorectal cancer. *BMC Syst Biol* **8,** S2 (2014).
510 28. Helleday, T. *et al.* Mechanisms underlying mutational signatures in human
511 cancers. - PubMed - NCBI. *Nat. Rev. Genet.* **15,** 585–598 (2014).
512 29. Wang, K. *et al.* Whole-genome sequencing and comprehensive molecular
513 profiling identify new driver mutations in gastric cancer. *Nat. Genet.* **46,**
514 573–582 (2014).
515 30. de Bruin, E. C. *et al.* Spatial and temporal diversity in genomic instability
516 processes defines lung cancer evolution. *Science* **346,** 251–256 (2014).
517 31. Zhang, J. *et al.* Intratumor heterogeneity in localized lung adenocarcinomas
518 delineated by multiregion sequencing. *Science* **346,** 256–259 (2014).
519 32. Nik-Zainal, S. *et al.* The Life History of 21 Breast Cancers. *Cell* **149,** 994–
520 1007 (2012).
521 33. Attolini, C. S.-O. *et al.* A mathematical framework to determine the
522 temporal sequence of somatic genetic events in cancer. *Proc. Natl. Acad.*
523 *Sci. U.S.A.* **107,** 17604–17609 (2010).
524 34. Gerstung, M. *et al.* The temporal order of genetic and pathway alterations

525       in tumorigenesis. *PLoS ONE* **6,** e27136 (2011).

526  35.  Sprouffske, K., Pepper, J. W. & Maley, C. C. Accurate reconstruction of the

527      temporal order of mutations in neoplastic progression. *Cancer Prev Res*

528      *(Phila)* **4,** 1135–1144 (2011).

529  36.  Guo, J., Guo, H. & Wang, Z. Inferring the temporal order of cancer gene

530      mutations in individual tumor samples. *PLoS ONE* **9,** e89244 (2014).

531  37.  Sottoriva, A. *et al.* Intratumor heterogeneity in human glioblastoma reflects

532      cancer evolutionary dynamics. *Proc. Natl. Acad. Sci. U.S.A.* **110,** 4009–

533      4014 (2013).

534  38.  Wang, Y. *et al.* Clonal evolution in breast cancer revealed by single

535      nucleus genome sequencing. *Nature* **512,** 155–160 (2014).

536  39.  Vermeulen, L. *et al.* Defining stem cell dynamics in models of intestinal

537      tumor initiation. *Science* **342,** 995–998 (2013).

538  40.  Heng, H. H. Q. *et al.* Stochastic cancer progression driven by non-clonal

539      chromosome aberrations. *J. Cell. Physiol.* **208,** 461–472 (2006).

540  41.  Navin, N. *et al.* Tumour evolution inferred by single-cell sequencing. *Nature*

541      **472,** 90–94 (2011).

542  42.  Marusyk, A. *et al.* Non-cell-autonomous driving of tumour growth supports

543      sub-clonal heterogeneity. *Nature* **514,** 54–58 (2014).

544

545 ## *Figure Legends*

546

547 **Figure 1**. **Neutral evolution is common in colon cancer and allows the**
548 **measurement of mutation rates in each tumor. (A)** The output of NGS data,
549 such as whole-exome sequencing, can be summarized as a histogram of mutant
550 allele frequencies, here for sample TB. Considering purity and ploidy, mutations
551 with relatively high frequency (>0.25) are likely to be clonal (public), whereas low
552 frequency mutations capture the tumor subclonal architecture. **(B)** The same data
553 can be represented as the cumulative distribution *M(f)* of subclonal mutations.
554 This was found to be linear with *1/f*, precisely as predicted by our neutral model.
555 **(C)** $R^2$ goodness of fit of our CRC cohort (n=7) and the TCGA colon cancer
556 cohort (n=101) grouped by CIN versus MSI confirmed that neutral evolution is
557 common (38/108, 35.1% with $R^2 \geq 0.98$). **(D)** Measurements of the mutation rate
558 showed that the CIN groups had median mutation rate of $\mu_e = 2.31 \times 10^{-7}$, whereas
559 MSI tumors reported a 15-fold higher rate (median: $\mu_e = 3.65 \times 10^{-6}$, F-test:
560 $p = 2.24 \times 10^{-8}$), as predicted due to their DNA mismatch repair deficiency.

561

562 **Figure 2**. **Neutral evolution across the whole-genome of gastric cancers. (A)**
563 Large number of coding and non-coding mutations can be identified using WGS.
564 **(B)** All detected mutations precisely accumulate as *1/f* following the neutral model
565 in this example. **(C)** Neutral evolution is very common in gastric cancer, with
566 60/78 (76.9%) samples showing goodness of fit of the neutral model $R^2 \geq 0.98$.
567 This was consistent using all, exonic or non-coding subclonal mutations. The
568 same tumors were identified as neutral by all three methods, although limitations
569 in detecting neutrality were present when considering exonic mutations due to
570 the limited number of variants. **(D)** Mutation rates were more than 4 times higher
571 in MSI ($\mu_e = 3.30 \times 10^{-6}$) versus MSS ($\mu_e = 7.82 \times 10^{-7}$; F-test: $p = 1.35 \times 10^{-4}$) cancers,
572 consistently with the underlying biology.

573

574 **Figure 3**. **Neutral evolution and mutation rates across cancer types. (A)** $R^2$
575 values from 819 cancers of 14 different types supported neutral evolution in a
576 large proportion of cases (259/819, 31.6% of $R^2 \geq 0.98$) and across different
577 cancer types, particularly in stomach (validating the WGS analysis), lung,
578 bladder, cervical and colon. On the contrary, renal, melanoma, pancreatic,
579 thyroid, and glioblastoma were characterized by non-neutral evolution. The other
580 types displayed a mixed dynamics. **(B)** The highest mutation rates were found in
581 lung cancer and melanoma. Lower rates were found in thyroid, low grade glioma
582 and prostate.

583

584 **Figure 4**. **Reconstruction of the mutational timeline in each patient.** The
585 frequency of a mutation within the tumor predicts the size of the tumor when the
586 mutation occurred. **(A,B)** The deconvolution of the mutational timeline is
587 illustrated for samples TB and TCGA-AA-3712 respectively. Whereas established
588 CRC drivers (APC, KRAS, TP53) were found to be present from the first
589 malignant cell, several recurrent putative drivers not yet validated were mutated
590 after malignant seeding, despite the underlying neutral dynamics. This suggests
591 that some of these candidate alterations may not be fundamental drivers of
592 growth in all cases. Confidence intervals are calculated using a binomial test on
593 the number of variant reads versus the depth of coverage for each mutation.
594
595 **Figure 5. Neutral evolution and tumor phylogeny.** After the accumulation of
596 genomic alterations, the cancer expansion is likely triggered by a single critical
597 genomic event (the accumulation of a "full house" of genomic changes) followed
598 by neutral evolution that generates a large number of new mutations in ever-
599 smaller subclones. While the tumor heterogeneity rapidly increases, the allele
600 frequency of heterogeneous mutations decreases. In this context, the
601 accumulation of mutations $M(f)$ follows a characteristic $1/f$ distribution. Moreover,
602 the tumor phylogeny displays a characteristic fractal topology that is self-similar.
603 Sampling in different regions of the phylogenetic tree exposes distinct mutations
604 that however show the same $1/f$ distribution. Clonal mutations in a sample (not
605 considered in the model) arose in to the most recent common ancestor of the
606 sampled cells. Due to the large population of cells sampled using bulk
607 sequencing, the overwhelming majority of detected clonal mutations belongs to
608 the trunk of the tree and therefore is found in the first cancer cell. Deviations from
609 the $1/f$ law indicate different dynamics from neutral growth.

## *Online Methods*

610

611

612 **Data analysis**
613     The processing of exome-sequencing data from[1] and TCGA[2] involved
614 variant calling on matched-normal pairs using Mutect[3]. A mutation was
615 considered if the depth of coverage was ≥10 and at least 3 reads supported the
616 variant. Mutations that aligned to a more than one genomic location were
617 discarded. The WGS gastric cancers[4] were processed using VarScan2[5], with
618 minimum depth of coverage for a mutation being 10x and at least 3 reads
619 supporting the variant. Non-CRCs in the TCGA had mutations called using
620 Mutect according to the pipeline described in ref[6]. Microsatellite instability in the
621 TCGA colon cancer samples was called using MSIsensor[7]. Annotation was
622 performed with ANNOVAR[8].
623     To fit the neutral model to allele frequency data we considered only
624 variants with allele frequency in the range [$f_{max}$,$f_{min}$] corresponding to [$t_0$,$t$] in
625 equation [2]. The low boundary $f_{min}$ reflects the limit for the reliable detectability of
626 low-frequency mutations in NGS data, which is in the order of 10%[3]. The high
627 boundary $f_{max}$ is necessary to filter out public mutations that were present in the
628 first transformed cell. In the case of diploid tumors, clonal mutations are expected
629 at $f_{max}$=0.5 (mutations with 50% allelic frequency are heterozygous public or
630 clonal), in the case of triploid tumors, this threshold drops to 0.33 and in the case
631 of tetraploid neoplasms, it drops to 0.25. For all samples we used a boundary of
632 [0.12-0.24] to account only for reliably called subclonal mutations and tumor
633 purity in the samples. All the samples considered in this study were reported to
634 have tumor purity ≥70% and a minimum of 12 reliably called private mutations
635 within the fit boundary. Once these conditions were met in a sample, equation [7]
636 was used to perform the fit as illustrated in Figure 1B and 2B. In particular, for
637 $x=1/f$, equation [7] becomes a linear model with slope $\mu/\beta$ and intercept $-\mu/(\beta$
638 $f_{max})$. We exploited the intercept constraint to perform a more restrictive fit using
639 the model $y=m(x-1/f_{max})+0$.
640     Copy-number changes (allelic deletion or duplication) can alter the
641 frequency of a variant in a manner that is not described by equation [7]. We
642 assessed the impact of copy-number alterations (CNAs) on our estimates of the
643 mutation rate within the TCGA colorectal cancer samples by using the paired

644　publically available segmented SNP-array data to exclude somatic mutations that
645　fell within regions of CNA. CNVs were identified having an absolute log-R-
646　ratio>0.5, and the model fitting was performed only on diploid regions of the
647　genome. In the gastric cancer cohort, regions with copy number changes were
648　identified using Sequenza[9] and removed from the analysis. Mutation rates were
649　adjusted to the size of the resulting diploid genome. Supplementary Figures 2
650　and 5 demonstrate the robustness of our analysis to copy number changes. $R^2$
651　values were independent from the mean coverage of mutations ($p$=0.32), the
652　total number of mutations in the sample ($p$=0.40), the mutation rate ($p$=0.11), or
653　the number of mutations within the model range ($p$=0.65).
654
655　**Stochastic Simulation of Tumor Growth**
656　　　　To further validate our analytical model and to test the robustness to the
657　noise in NGS data, we developed a stochastic simulation of tumor growth and
658　accumulation of mutations that allowed us to generate synthetic datasets. The
659　model was written and analyzed in the Julia programming language.  We then
660　applied the analytical model to the simulated data to confirm that sources of
661　noise in NGS data do not considerably impact our results. In particular, we
662　verified that we could reliably extract input parameters of the simulation (namely
663　the mutation rate) from "noisy" synthetic data. Confounding factors in the data
664　include normal contamination, sampling effects, the detection limit of NGS
665　mutation calling, and variable read depth. We simulate a tumor using a branching
666　process with discrete generations, beginning with a single "transformed" cancer
667　cell that gives rise to the malignancy. Under exponential growth, the population at
668　time $t$ will be given by:
669
670　$$N(t) = R^t = e^{\ln(R)t}$$　　　　　　　　　　　　　　　[9]
671
672　Where R is the average number of offspring per cell and the time $t$ is in units of
673　generations. We will consider primarily the case when $R$=2 (a cell always divides
674　into 2), but we will also consider values <2, noting that $R$ must be greater than 1
675　to have growth. At each division, cells acquire new mutations at a rate $\mu$ and we
676　assume every new mutation is unique (infinite sites approximation). The number
677　of mutations acquired by a newborn cell at division is a random number drawn
678　from a Poisson distribution. Each cell in the population is defined by its mutations
679　and its ancestral history (by recording it's parent cell). Using this information we
680　can then reconstruct the history of the whole tumor and crucially, calculate the
681　variant allele frequency of all mutations in the population. To relate the discrete
682　simulation to the continuous analytical model we will now re-derive equation [7]
683　within the context of our model. As we simulate a growing tumor using discrete
684　generations, both the mutation rate $\mu$ and per capita growth rate $\lambda=ln(R)$ are in
685　units of generations. For an offspring probability distribution $P=(p_0,p_1,p_2)$ where
686　$p_k=P(\# \ of \ OFFSPRING = k)$ where, the average number of offspring  $R$ is simply
687　given by the expected value of $P$:
688
689　$$R = E[P] = p_1 + 2p_2$$　　　　　　　　　　　　　　　[10]
690
691　For example, for $R$=2 we have $P=(p_0=0,p_1=0,p_2=1)$. By choosing different
692　offspring probability distributions we can easily modulate the growth rate. We
693　note that we are now expressing both $\mu$ and $\lambda$ as rates per generation rather than
694　probabilities (all rates are scaled by units of generation). This allows us to write
695　the growth function as $N(t)=exp(\lambda t)$ with $\lambda=ln(R)$. Proceeding as in the main text,
696　our cumulative number of mutations with an allelic frequency $f$ is therefore:
697　$$M(f) = \frac{\mu}{\lambda}\left(\frac{1}{f} - \frac{1}{f_{max}}\right)$$　　　　　　　　　　[11]
698　Therefore, when fitting the model to our stochastic simulation we extract $\mu/\lambda$ from
699　the linear fit, making it straightforward to compare the simulation with the
700　analytical model.
701　　　　NGS data only captures a small fraction of the variability in a tumor, as the
702　resolution is often limited to alleles with frequency >10% due to sequencing

703 depth and limitations in mutation calling. To account for this, we employ a
704 multistage sampling scheme in our simulations. For all simulations reported here
705 we grow the tumor to size 1,024 cells, which gives a minimum allele frequency of
706 ~0.1%, considerably smaller than the 10% attainable in next generation
707 sequencing data. After growing the tumor and calculating the VAF for all alleles,
708 we take a sample of the alleles in the population, noting that we are assuming the
709 population is well mixed and has no spatial structure. We can vary the
710 percentage of alleles we sample, thus allowing us to investigate the effect of the
711 depth of sequencing on our results. As we know the true allelic frequency in the
712 simulated population, we can use the multinomial distribution to produce a
713 sample of the "sequenced" alleles, where the probability of sampling allele $i$ is
714 proportional to its frequency. The probability mass function is given by:

715

716
$$f(x;n,p) = \frac{n!}{x_1! \, \Box \, x_k!} \prod_{i=1}^{k} p_i^{x_i} \quad , \quad x_1 + \Box \, x_k = n \qquad \qquad [12]$$

717

718 where $x_i$ is the sampled frequency of allele $i$, $n$ is the number of trials (the chosen
719 percentage of alleles sampled) and $p_i$ is the probability of sampling allele $i$ (which
720 has frequency $\rho_i$ in the original population):

721
$$p_i = \frac{\rho_i}{\sum_{j=1}^{k} \rho_j} \qquad \qquad [13]$$

722
723 The variant allele frequency VAF is therefore given by:

724
$$VAF = \frac{x_i}{N_i} \qquad \qquad [14]$$

725

726 Where $N_i$ is the total number of sampled cells from which every sampled allele is
727 derived. As we are assuming a constant mutation rate μ, we can assume that the
728 percentage of alleles sampled comes from an equivalent percentage of cells.
729 However, to include an additional element of noise that resembles the variability
730 of read depth, we calculate a new $N_i$ for each allele $i$, which approximates the
731 read depth. For a desired "sequencing" depth $D$ we calculate the corresponding
732 percentage of the population we need to sample that will give us our desired
733 depth. For example, for a desired depth of 100X from a population of 1,000 cells,
734 we would need to sample 10% of the population. To include some variability in
735 depth across all alleles we use Binomial sampling so that $N_i$ is a distribution with
736 mean D.
737 Contamination from non-tumor cells in NGS results in variant allele
738 frequencies being underestimated. To include this effect in our simulation we can
739 modify our $N_i$ by an additional fraction $\varepsilon$, the percentage of normal contamination.
740 Our VAF calculation thus becomes:

741
$$VAF = \frac{x_i}{N_i(1+\varepsilon)}$$

742
743
744 We also include detection limit in our sampling scheme, we only include alleles
745 that have an allelic frequency greater than a specified limit in the original tumor
746 population.

747

748 To include the effects of selection in the simulation we introduce a second
749 population, where on average each cell has a greater number of offspring than
750 the first population. To model this, our second population has a modified offspring
751 probability distribution: the previous offspring probability distribution was
752 $P=(p_0,p_1,p_2)$, and the offspring probability distribution of our second fitter
753 population is defined as $Q=(q_0,q_1,q_2)$, where $q_2>p_2$. The selective advantage of a
754 population – $s$, will be given by the ratio of the expected number of offspring:

755
$$1+s = \frac{E[Q]}{E[P]} = \frac{q_1+2q_2}{p_1+2p_2}$$

756 Therefore given *P*, and a desired selective advantage *s* we can easily calculate
757 the offspring probability distribution of a fitter clone – *Q*.
758
759 Previous studies have detected the presence of mixtures of subclones in breast
760 cancer samples that emerged by means of clonal expansions, thus generating
761 multiple subclonal clusters in the data[10]. We also used our computational model
762 of NGS data to produce similar synthetic data by means of mixing of different
763 clonal clusters and verified that in this scenario (a model of differential selective
764 pressure across subclones), the power law does not hold.
765
766 **Simulation Results**
767    From the simulated data we produced histograms of the allelic frequency
768 and calculated *M(f)* in order to fit the analytical model. We used the same
769 frequency range as applied to empirical data $[f_{max},f_{min}]=[0.12,0.24]$.
770 Supplementary Figure 9A and B shows equivalent plots to Figures 1A and B but
771 with simulated data. These demonstrate that we are able to accurately model the
772 allelic distribution of NGS data with our simple neutral model of tumor growth. We
773 also show the effect of a low mutation rate (Supplementary Figure 9C), a large
774 number of clonal mutations (Supplementary Figure 9D), 30% contamination in
775 the sample (Supplementary Figure 9E) and a low detection limit (Supplementary
776 Figure 9F). Importantly, by fitting the analytical model to the simulated data, we
777 can recover the input mutation rate with high accuracy (Supplementary Figure
778 9G, 10,000 equivalent simulations). The mean percentage error from the fit is
779 1.1%. We also see uniformly high $R^2$ values across all simulations
780 (Supplementary Figure 9H).
781    To test the robustness of the model to the number of clonal mutations, the
782 detection limit and the amount of normal contamination we ran 10,000
783 simulations across the spectrum of these parameters. Supplementary Figures
784 10A-B show that we accurately recover (to within 15%) the mutation rate for 95%
785 of simulations across different numbers of clonal mutations and different
786 detection limits. Differently, we found that levels of normal contamination above
787 30% considerably impact the parameter estimations of the model, hence our
788 decision of only considering samples with ≥70% of tumor content
789 (Supplementary Figure 10C). Indeed, when normal contamination is above 30%,
790 the clonal peak in the allelic frequency distribution interferes significantly with our
791 chosen cumulative sum limit ($f_{max}$ = *0.24*), thus impacting our results.
792 Nevertheless, the estimates are within a factor 2 for normal contamination of up
793 to 50%, which we consider an acceptable level of accuracy. When we consider
794 normal contamination $\varepsilon$ directly within our analytical model, the allelic fraction of a
795 new mutation becomes:
796

$$f = \frac{1}{\pi N(t)} = \frac{1}{\pi e^{\lambda \beta t}(1+\varepsilon)} \qquad [15]$$

798
799 And consequently, *M(f)* is:
800

$$M(f) = \frac{\mu}{\beta(1+\varepsilon)}\left(\frac{1}{f} - \frac{1}{f_{max}}\right) \qquad [16]$$

802
803 Showing that normal contamination alters the measurement of mutation by a
804 factor of *1/(1+ε)*: much lower than one order of magnitude. Furthermore, if normal
805 contamination can be estimated accurately from histopathological scoring or from
806 reliable bioinformatics tools, we would be able to correct the frequency of variants
807 in the data and thus rescue our ability to correctly estimate parameters with up to
808 40-45% normal contamination (Supplementary Figure 10D). We also tested the
809 model with varying read depths and mutation rates. We find that either a low
810 mutation rate or low read depth resulted in a higher proportion of poor model fits
811 ($R^2$<0.98) and inaccurate or higher variance in mutation estimates
812 (Supplementary Figures 10E-H). It is therefore possible that due to our stringent
813 neutrality criteria that the true proportion of tumors that are dominated by neutral

dynamics is higher than reported, and relatedly our gastric cancer cohort covers the whole genome (greater mutation rate per division) and has mean depth of coverage >90X which may explain in part why we see a greater proportion of gastric cancers classified as neutral.

Additionally, we tested the model with simulations using a range of different probability distributions for the number of surviving offspring at each cell division. We simulated a growing tumor 10,000 times with 5 different offspring probability distributions and then reported the distributions of the fitted parameters. Supplementary Figures 10I-J show that as λ decreases the distribution of mutation estimates becomes wider and we see an increase in poorly fitted models (larger number of $R^2$<0.98). Again this suggests that tumor growth may still be neutral even when we classify a tumor as non-neutral due to a poor $R^2$ value. Hence our underestimation of the number of neutral cases may be largely due to a low proportion of cells that successfully produce 2 viable offspring (the β term in equation [7]), rather than the presence of selection.

By introducing a second fitter population early during tumor growth we show that the fitter clone causes an overrepresentation of variants at high frequency compared to what we would expect from our "null" model of neutral tumor growth. This causes the cumulative distribution to bend and deviate from the linear relationship predicted by neutral growth, as shown in Supplementary Figures 11A-B. This is because an overrepresentation of variants at high frequency, as compared to what we would expect from our "null" model, is caused by the clonal selection of the fitter clone, but we note that we do not know what caused this increase (it could be a point mutation, chromosomal aberration or a change in environmental pressures for example). In other words, some passenger mutations are just in the "right clone at the right time" and become overrepresented in the tumour when that "right" clone expands.
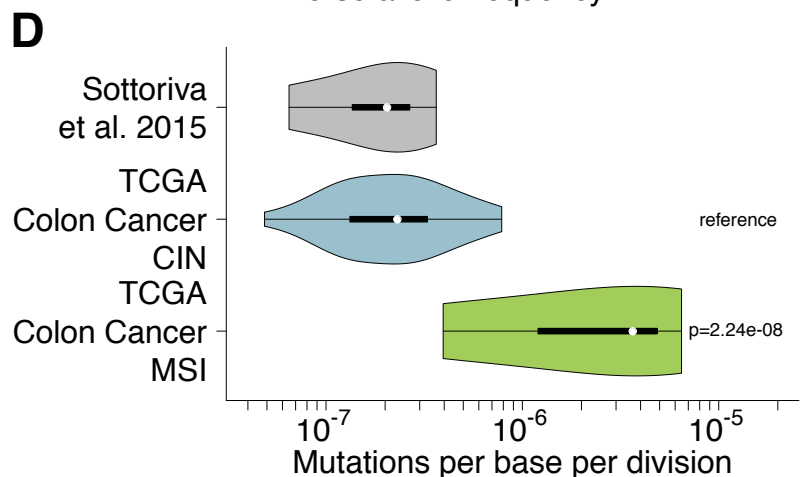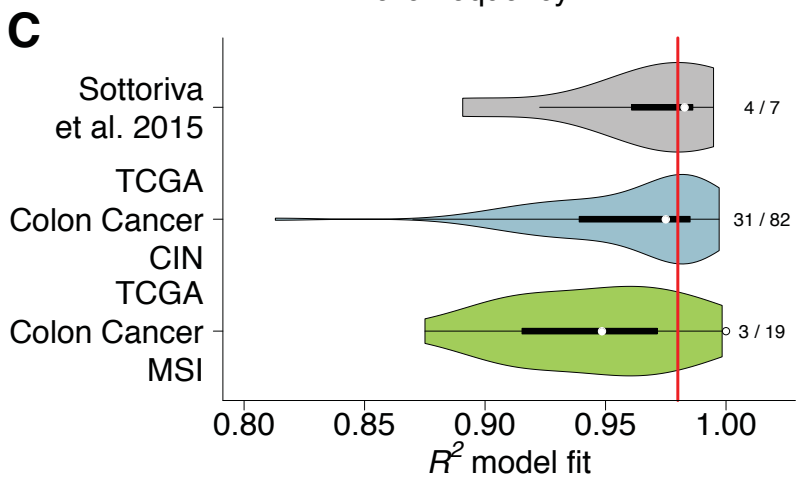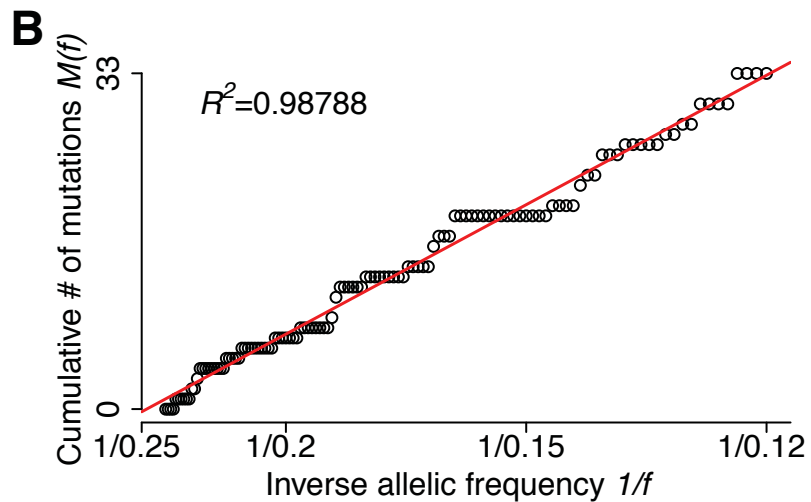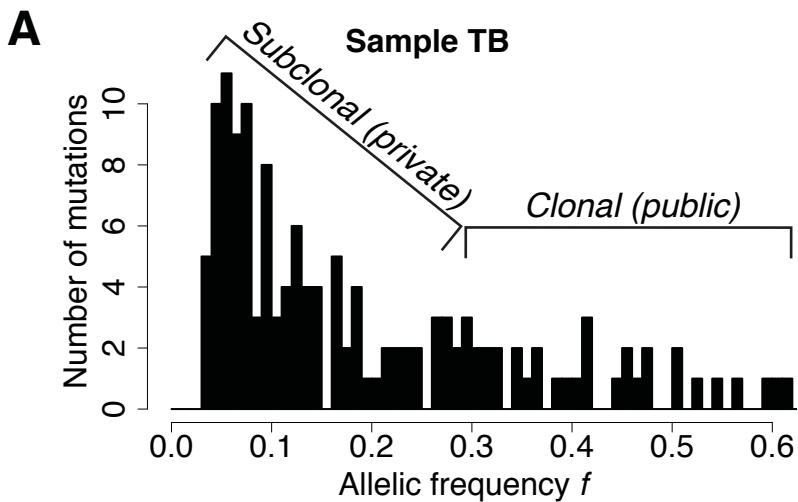
We also show that having multiple subclones that arose by means of clonal expansion, thus producing multiple clonal 'clusters', produces a deviation from the linear relationship we predict (Supplementary Figures 11C-F), as does having a marked increase in the mutation rate early in tumour growth (Supplementary Figures 11G,H).
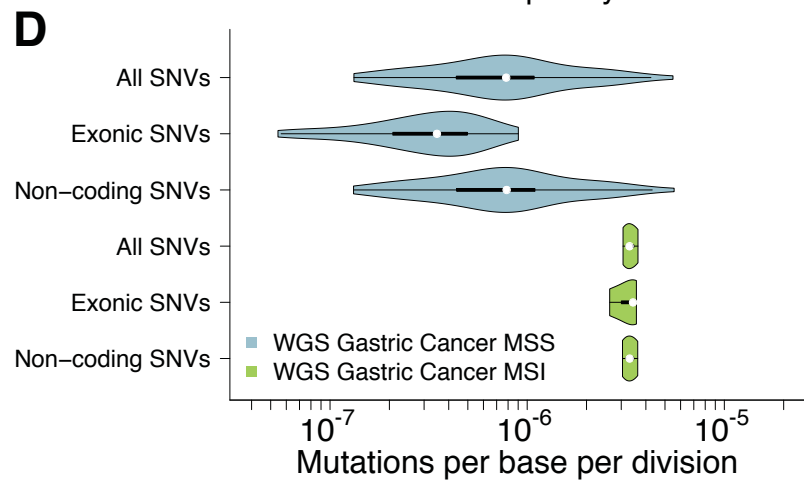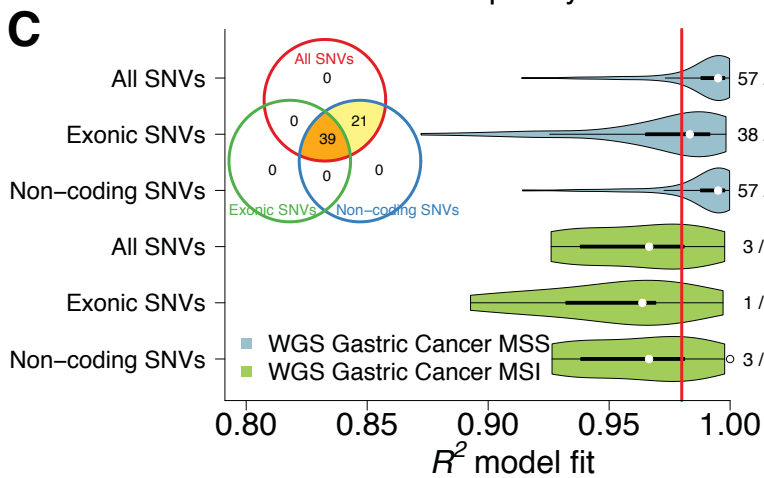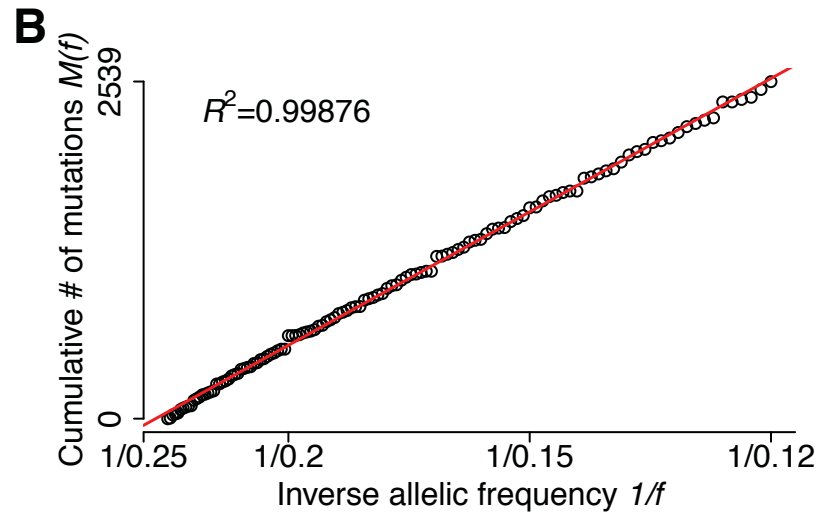
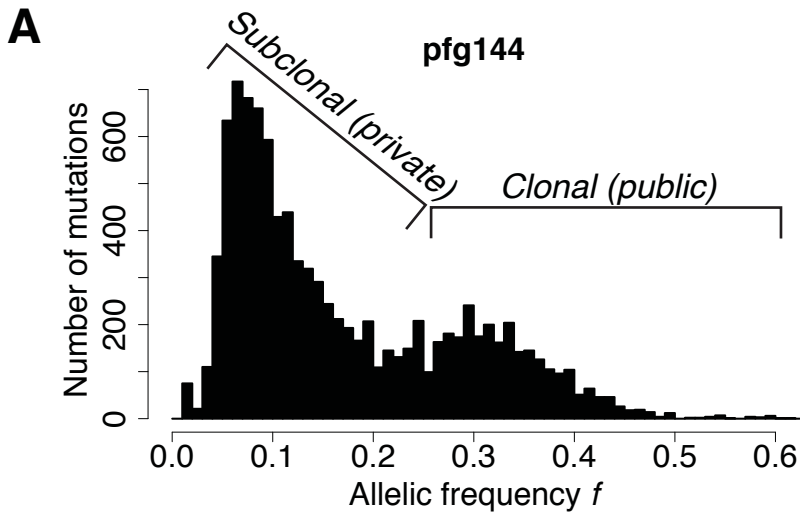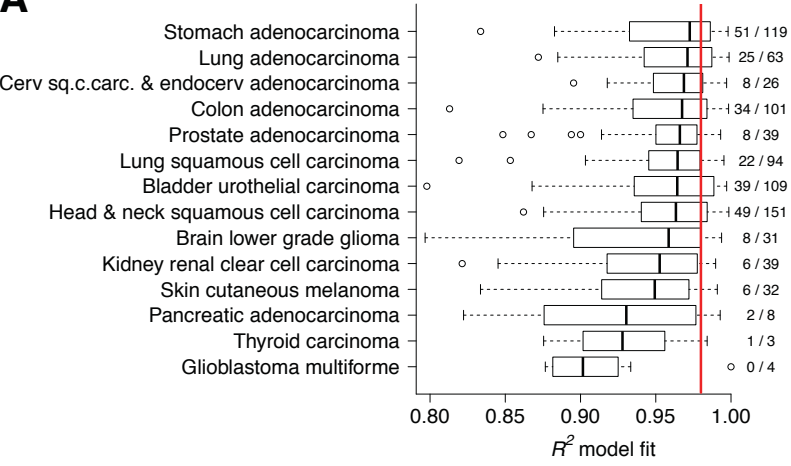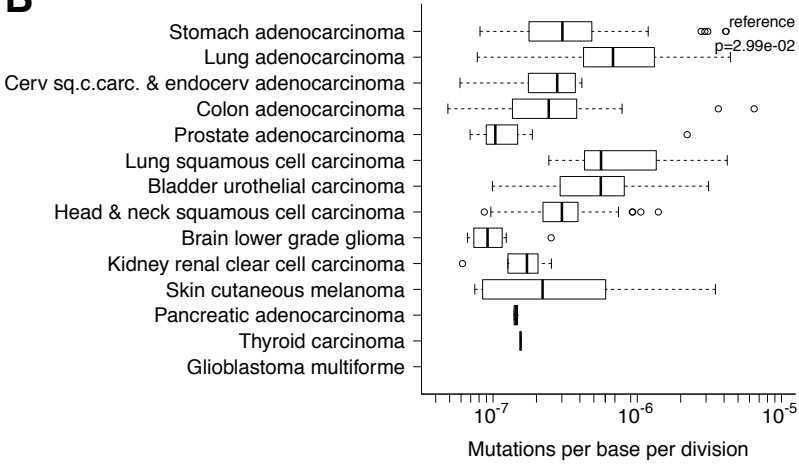## *References*

1.  Sottoriva, A. *et al.* A Big Bang model of human colorectal tumor growth. *Nat. Genet.* **47,** 209–216 (2015).
2.  The Cancer Genome Atlas. Comprehensive molecular characterization of human colon and rectal cancer. *Nature* **487,** 330–337 (2012).
3.  Cibulskis, K. *et al.* Sensitive detection of somatic point mutations in impure and heterogeneous cancer samples. *Nature Biotechnology* **31,** 213–219 (2013).
4.  Wang, K. *et al.* Whole-genome sequencing and comprehensive molecular profiling identify new driver mutations in gastric cancer. *Nat. Genet.* **46,** 573–582 (2014).
5.  Anderson, A. R. A. *et al.* VarScan 2: somatic mutation and copy number alteration discovery in cancer by exome sequencing. *Genome Res.* **22,** 568–576 (2012).
6.  Andor, N. *et al.* Pan-Cancer Analysis of the Causes and Consequences of Intra-Tumor Heterogeneity.
7.  Niu, B. *et al.* MSIsensor: microsatellite instability detection using paired tumor-... - PubMed - NCBI. *Bioinformatics* **30,** 1015–1016 (2014).
8.  Wang, K., Li, M. & Hakonarson, H. ANNOVAR: functional annotation of genetic variants from high-throughput sequencing data. *Nucleic Acids Res.* **38,** e164 (2010).
9.  Favero, F. *et al.* Sequenza: allele-specific copy number and mutation profiles from tumor sequencing data. *Annals of Oncology* **26,** 64–70 (2014).
10. Nik-Zainal, S. *et al.* The Life History of 21 Breast Cancers. *Cell* **149,** 994–1007 (2012).

**Competing financial interests**

874     The authors declare no competing financial interests.

**A** Sample TB

Subclonal (private)

Clonal (public)

Number of mutations

Allelic frequency *f*

**B** $R^2$=0.98788

Cumulative # of mutations $M(f)$

Inverse allelic frequency *1/f*

**C**

Sottoriva et al. 2015 — 4 / 7

TCGA Colon Cancer CIN — 31 / 82

TCGA Colon Cancer MSI — 3 / 19

$R^2$ model fit

**D**

Sottoriva et al. 2015

TCGA Colon Cancer CIN — reference

TCGA Colon Cancer MSI — p=2.24e-08

Mutations per base per division

**A**

| | | |
|---|---|---|
| Stomach adenocarcinoma | | 51 / 119 |
| Lung adenocarcinoma | | 25 / 63 |
| Cerv sq.c.carc. & endocerv adenocarcinoma | | 8 / 26 |
| Colon adenocarcinoma | | 34 / 101 |
| Prostate adenocarcinoma | | 8 / 39 |
| Lung squamous cell carcinoma | | 22 / 94 |
| Bladder urothelial carcinoma | | 39 / 109 |
| Head & neck squamous cell carcinoma | | 49 / 151 |
| Brain lower grade glioma | | 8 / 31 |
| Kidney renal clear cell carcinoma | | 6 / 39 |
| Skin cutaneous melanoma | | 6 / 32 |
| Pancreatic adenocarcinoma | | 2 / 8 |
| Thyroid carcinoma | | 1 / 3 |
| Glioblastoma multiforme | | 0 / 4 |

0.80   0.85   0.90   0.95   1.00

$R^2$ model fit

**B**

| | |
|---|---|
| Stomach adenocarcinoma | reference |
| Lung adenocarcinoma | p=2.99e-02 |
| Cerv sq.c.carc. & endocerv adenocarcinoma | |
| Colon adenocarcinoma | |
| Prostate adenocarcinoma | |
| Lung squamous cell carcinoma | |
| Bladder urothelial carcinoma | |
| Head & neck squamous cell carcinoma | |
| Brain lower grade glioma | |
| Kidney renal clear cell carcinoma | |
| Skin cutaneous melanoma | |
| Pancreatic adenocarcinoma | |
| Thyroid carcinoma | |
| Glioblastoma multiforme | |

$10^{-7}$   $10^{-6}$   $10^{-5}$

Mutations per base per division

**A**  **TB, 152 mutations**

Tumor volume *N*

140,000
100,000
60,000

APC
TP53
LRP1B

TCF7L2

20,000

>0.25   0.2   0.15   0.1   0.05   0

Allelic frequency *f*

**B**  **TCGA−AA−3712, 190 mutations**

Tumor volume *N*

150,000

100,000

TH

50,000

APC
TP53   PIK3CA

>0.25   0.2   0.15   0.1   0.05   0

Allelic frequency *f*

Intra-tumor heterogeneity

Allelic frequency

1/f

M(f)

M(f)

M(f)

M(f)