

24 Evolving tumors accumulate thousands of mutations. Technological advances have enabled
25 whole genome sequencing of these mutations in large cohorts, such as those from the Pancancer
26 Analysis of Whole Genomes (PCAWG) Consortium. The resulting data explosion has led to
27 many methods for detecting cancer drivers through mutational recurrence and deviation from
28 background mutation rates. However, these methods require a large cohort and underperform
29 when recurrence is low. An alternate approach involves harnessing the variant allele frequency
30 (VAF) of mutations in the population of tumor cells in a single individual. Moreover, ultra-deep
31 sequencing of tumors, which is now possible, allows for particularly accurate VAF
32 measurements, and recent studies have begun to use these to determine evolutionary trajectories
33 and quantify subclonal selection. Here, we developed a method that quantifies tumor growth and
34 driver effects for individual samples based solely on the VAF spectrum. Drivers introduce a
35 perturbation into this spectrum, and our method uses the frequency of "hitchhiking" mutations
36 preceding a driver to measure this perturbation. Specifically, our method applies various growth
37 models to identify periods of positive/negative growth, the genomic regions associated with
38 them, and the presence and effect of putative drivers. To validate our method, we first used
39 simulation models to successfully approximate the timing and size of a driver's effect. Then, we
40 tested our method on 993 linear tumors (i.e. those with linear subclonal expansion, where each
41 parent-subclone has one child) from the PCAWG Consortium and found that the identified
42 periods of positive growth are associated with drivers previously highlighted via recurrence by
43 the PCAWG consortium. Finally, we applied our method to an ultra-deep sequenced AML tumor
44 and identified known cancer genes and additional driver candidates. In summary, our method
45 presents opportunities for personalized diagnosis using deep sequenced whole genome data from
46 an individual.

47

48 **Introduction**

49 Over the past several decades, researchers have proposed different models to explain tumor
50 progression, including stochastic progression, the mutator phenotype, and clonal evolution¹⁻³.

51 Originally suggested about 40 years ago³, Navin and colleagues provided strong evidence that
52 the ‘punctuated clonal evolution’ model constitutes a major force in cancer progression.

53 According to this model, tumor progression is an evolving system subject to selective pressure
54 while accumulating thousands of mutations^{4,5}.

55

56 Advances in technology have allowed scientists to sequence thousands of genomes, revealing
57 millions of variants per individual⁶⁻⁸. In cancer genomics, The Cancer Genome Atlas (TCGA)⁴

58 offers access to thousands of cases encompassing over 30 types of cancer. Similarly, the

59 International Cancer Genome Consortium (ICGC) recently announced ‘data release 26’, which

60 comprises data from more than 17,000 cancer donors and 21 tumor sites. Within ICGC, the

61 Pancancer Analysis of Whole Genomes (PCAWG) study is an international collaboration to

62 identify common patterns of mutations in over 2,800 sequenced whole cancer genomes⁹. As

63 cancer databases continue to expand, the amount of fully sequenced genomes will continue to

64 increase, with future plans setting goals for the storage of more than a million genomes¹⁰.

65 Concurrently, deeper sequencing signifies less noise, more accurate variant allele frequencies

66 (VAFs), and more accurate subclonal and single-nucleotide variant (SNV) identification, while

67 increasing the detection of novel drivers¹¹⁻¹³.

68

69 Recent studies have tackled the effect of selection in tumor progression in the context of clonal
70 evolution, neutral evolution, and selection, providing valuable insights about the clonal
71 progression of the disease^{5,14–16}. By considering tumor progression as an evolutionary process,
72 cancer development follows the trajectory of different evolutionary pathways based on cell and
73 population dynamics, optimization strategies and selective forces. These evolutionary trajectories
74 have been shown to influence primary tumor growth¹⁷ and the timing of landmark events¹⁸.
75 However, the evolutionary and selective mechanisms during tumor progression remain
76 unexplored and strongly debated^{19–22}.

77
78 Accumulated SNVs have been characterized as drivers or passengers, depending on whether or
79 not they provide a selective advantage for the tumor cells. If the selective advantage or their
80 respective effect is weak, the mutations are known as mini-drivers, although the existence and
81 detectability of mini-drivers has been debated^{23,24}. Identifying SNV and gene drivers has been
82 one of the focal points of cancer genomics, where different methods aim to detect driver
83 mutations based on selection, recurrence or changes in mutational density^{23,25}. These methods
84 rely on the deviation from our expectation of the underlying genomic mutation rates, often by
85 considering additional covariates such as replication timing and gene expression^{26–28}. Other
86 methods, characterized as ratiometric, assess the composition of mutations, normalized by the
87 total mutations in a gene²³. This includes the proportion of inactivating mutations, recurrent
88 missense mutations, functional impact bias, mutational composition, or clustering patterns^{29–32}.
89 However, if only a small proportion of mutations within a genomic region (which is potentially
90 under negative selection or functional restrictions) facilitates cancer progression, driver detection
91 requires either a very large sample, a strong effect or otherwise the driver's presence is

92 undetectable²³. Further, mutational heterogeneity in cancer poses an additional problem for large
93 cohorts; as the sample size increases, so does the list of putatively significant genes, producing
94 many false positive driver genes²⁷. More importantly, only a minimal portion of driver mutations
95 are, in fact, true drivers³³. This is particularly important in a clinical context as assessing a cancer
96 gene mutation as a true functional driver is a critical problem for drug selection^{33,34}.

97

98 According to recent studies³⁵ and in agreement with past theories³⁶, a few major genetic hits
99 (strong drivers) can induce tumorigenesis. At the same time, a driver mutation may not actually
100 be the cause of tumorigenesis, but instead only increase growth rate and therefore be under
101 positive selection³⁷. One of the most common and widely used lists of cancer genes is the
102 "Vogelstein list"²⁹, consisting of ~140 oncogenes and tumor-suppressor genes (TSGs). While
103 high-impact mutations in TSGs might favor cancer progression by deactivating tumor
104 suppression, oncogenes need altered expression levels to favor tumor growth. Thus, high-impact
105 mutations such as nonsense mutations in oncogenes might decrease gene expression and burden
106 tumor cells³⁸. Less appreciated is the role of non-coding mutations in tumor progression^{37,39,40}.
107 Interestingly, in the case of TSGs, different studies have reported the role of non-coding intronic
108 mutations that alter correct exon splicing, resulting in faulty tumor suppression⁴¹⁻⁴⁴. Similarly, in
109 the case of oncogenes different studies have reported the potential effect of synonymous
110 mutations^{40,41,45}. For example, Gartner and colleagues showed that the early synonymous
111 mutation F17F in the BLC2-like 12 gene alters the binding affinity of regulatory hsa-miR-671-
112 5p, leading to changes in expression⁴⁵.

113

114 In our study, we developed a framework to model tumor progression and the effect of drivers in
115 individual deep-sequenced tumors. We successfully applied our model using 993 linear tumors
116 (linear subclonal expansion, where each parent-subclone has one child-subclone) from the
117 PCAWG consortium, and found that predicted drivers⁴⁶ are associated with periods of positive
118 growth. Our results suggest that mutations involved in biological processes such as cell
119 development, cell differentiation, and multicellularity appear under strong positive or negative
120 growth enrichment. Missense or nonsense mutations in TSGs were enriched during positive
121 growth. We also identified significant positive enrichment for mutations in the promoter regions
122 of both TSGs and oncogenes. Additionally, in the case of TSGs, we discovered a small but
123 significant signal from intronic mutations. Finally, we applied our framework to a deep-
124 sequenced model AML tumor, where our predicted growth peaks aligned closely with three
125 missense mutations from known cancer genes. Notably, our analysis suggests the potential
126 presence of additional driver candidates.

127

128

129 **Method Overview: Clock-like Hitchhikers, Growth Rates, Local Re-optimization, and**
130 **Driver Effects**

131 When sequencing a cell population or tumor bulk, each mutation is assigned a variant allele
132 frequency (VAF), which corresponds to the mutation's frequency in the resulting pool.

133 According to the infinite sites model⁴⁷, once a mutation occurs it will continue to exist within
134 that cell and its descendants. Therefore, if we assume that there is no selection or chromosomal
135 duplications, the VAF is associated with the time of occurrence and population growth rates.

136 That is, in the presence of a driver (i.e., in cells with higher fitness), non-driver mutations within

137 that cell lineage will also have higher-than-expected VAF and are termed “hitchhikers”²⁹ (Figure
138 1, Supplement). Hitchhikers that initially occurred before the driver mutation but continue to
139 exist within that cell lineage will have a VAF that is higher than or equal to the driver’s
140 frequency. We call these hitchhikers “generational” (g-hitchhikers) because they essentially
141 mark the different generations of an ever-increasing number of tumor cells and thus exhibit a
142 clock-like behavior. Since any non-driver lineage derived from the division of earlier cells will
143 result in a mutation having lower frequency, these pre-driver hitchhiking mutations will indicate
144 generational growth (Figure 1). As the fitness mutation becomes more prevalent over time, so
145 does the prevalence of pre-driver “g-hitchhikers”, but critically at a different pace, which we
146 calculate (see Supplement).

147
148 Our framework’s equations (which we dub “hitchhiker equations”, see Supplement) relate the
149 VAF of generational hitchhiker mutations to the fitness effect of the subclonal driver with which
150 they are hitchhiking, mediated by various growth and population parameters (i.e. the base growth
151 rate r , a scalar multiplier k corresponding to fitness effect of the mutation, the time t_I when the
152 driver mutation is generated, N_{tot} the population size and N_F the driver’s subclone size). The
153 existence and fitness effects of subclonal drivers are not directly observable but are of primary
154 biomedical importance. The VAF of hitchhiker mutations is directly observable, therefore we
155 chose to use these VAFs to infer the presence of subclonal drivers and estimate their fitness
156 effects. Our approach is to fit the known VAFs of the hitchhiker mutations in the hitchhiker
157 equations to estimate the growth pattern and the fitness effect of subclonal drivers. This method
158 requires to simultaneous estimate the various growth and population parameters, which we
159 performed using non-linear least-squares optimization. To address the fact that real tumors differ

160 from idealized behavior, we make use of sliding windows and local timepoint re-optimizations in
 161 the parameter estimation to prevent departures from idealized behavior in one part of the VAF
 162 spectrum from interfering with parameter estimation in other parts of the VAF spectrum. The
 163 details of the growth and population parameters, their estimation, and the use of sliding windows
 164 are described in the Supplement. We derived our estimators for r and k through the
 165 implementation of a deterministic model to a stochastic process with a large final population
 166 N_{tot} .

167

168 *Modeling the frequency of g-hitchhikers using exponential models*

169 We assume a simple and neutral population of cancer cells that grows exponentially with rate r .
 170 For simplicity, we here assign each new daughter cell one new mutation (alternative mutation
 171 rates do not affect the derivation, see Supplement). At time t_1 , a mutation occurs that accelerates
 172 the growth rate of the specific subpopulation by a scalar multiplier k such that the new
 173 population expands with new rate $k \times r$. At the time of biopsy $T = t_1 + t_2$, where the fitness mutation
 174 occurs at t_1 and expands for time t_2 , we expect the frequency of a generational g-hitchhiker
 175 mutation that occurred at time $t_m < t_1$ (see Figure 1, Supplement) to follow a frequency function

176 f_g :

177

178
$$f_g(T, t_m) = \frac{N_R + N_F - N_{RF}}{N_{tot}}$$

179 or

180
$$f_g(T, t_m) = \frac{e^{-rt_m} [N_{tot} - f_{d(T,t_1)} * N_{tot} + \sqrt[k]{f_{d(T,t_1)} * N_{tot}}] + f_{d(T,t_1)} * N_{tot} - \sqrt[k]{f_{d(T,t_1)} * N_{tot}}}{N_{tot}} \quad [1],$$

181

182 where $f_{d(T,t_1)}$ is the frequency of the driver mutation occurring at t_1 and expanding for $t_2=T-t_1$,
183 The terms $\{ e^{-rt_m} * [N_{tot} - f_{d(T,t_1)} * N_{tot}] \}$ and $\{ f_{d(T,t_1)} * N_{tot} \}$ correspond to the growth of
184 regular N_R and fitness N_F populations respectively, while extracting $N_{RF} = \{ \sqrt[k]{f_{d(T,t_1)} * N_{tot}} \}$ for
185 not double-counting the hypothetical regular growth of fitness cells (see Figure 1, Supplement).

186
187 Equation (1) for the m -th hitchhiker implicitly allows one to use the previous $m-1$ potential
188 hitchhikers to refine the estimates of growth rate r and scalar effect k . This estimation is achieved
189 either through a non-linear-least-squares optimization, and/or through the independent
190 calculation of growth r .

191
192 The frequency of g-hitchhiking mutations follows the form of an exponential distribution.
193 Theoretically, this further allows us to estimate growth rate r from consecutive g-hitchhiking
194 mutations m_1, m_2 , and m_3 , which occurred at times t_{m1}, t_{m2} , and t_{m3} (t_{m1}, t_{m2} , and $t_{m3} < t_1$),
195 according to

196
197
$$r = \ln \left(\frac{f_g(T,t_{m1}) - f_g(T,t_{m2})}{f_g(T,t_{m2}) - f_g(T,t_{m3})} \right) \quad [2]$$

198
199 In practice, to obtain more accurate estimates, our default algorithm estimates the growth rate r
200 from three more distant time points $t, t+n$, and $t+m$ ($n < m$ and $t+m < t_1$) with final frequencies
201 $f_g(T, t)$, $f_g(T, t_n)$, and $f_g(T, t_m)$, respectively, as described in the Supplement.

202

203

204 *Optimizing for generational time at any time point during tumor progression*

205

206 In addition to our independent estimate of growth rate r , and in order to avoid previous
207 frequency perturbations in our sample and localize the effect timewise, we also include an extra
208 parameter referred to as ‘*generational time* (\mathbf{t}_g)’, which allows us to calibrate an offset for the
209 number of past generations until that point without considering previous mutations outside our
210 sliding window. Thus, similar to eq. [1], we now have

211

$$212 \quad f_g(T, \mathbf{t}_g, t_i - m) = \frac{e^{-r(\mathbf{t}_g + t_i - m)} * (N_{\text{tot}} - f_{d(T, t_i)} * N_{\text{tot}})^{k_i} \sqrt{f_{d(T, t_i)} * N_{\text{tot}}} + f_{d(T, t_i)} * N_{\text{tot}} - k_i \sqrt{f_{d(T, t_i)} * N_{\text{tot}}}}{N_{\text{tot}}} [3],$$

213 where $f_d(T, t_i)$ is the frequency of the putative driver i occurring at time t_i .

214 This approach allows us to re-optimize \mathbf{t}_g at any time t_i during tumor growth, *independently* of
215 earlier or later calculations.

216

217

218 **Validating Our Model Using Simulations**

219

220 *Birth and death model, Gillespie simulations*

221

222 First, we tested our algorithm on simulated data based on various growth models, including: a)
223 exponential growth, b) exponential growth with delayed cell division, and c) logistic growth
224 (birth and death model). We performed simulation models (a) and (c) using a stochastic Gillespie
225 algorithm, whereas model (b) represents an exponential cell growth model with a lag time for
226 cell division, which prevents a cell from re-dividing immediately. Briefly, for the “*Birth and*
227 *Death*” Gillespie model, which is the workhorse of our simulations, we used a stepwise time-

228 branching process to model the growth of a single transformed cell into a tumor with a dominant
229 subclone. At each time step, an event type is chosen with a probability proportional to the event's
230 prevalence (see supplement) Then, a cell of the eligible type is randomly chosen to undergo that
231 event. In our logistic-growth simulations, the death rate of each cell climbs proportionally as
232 carrying capacity is reached, whereas in our exponential simulations, the death rate of each cell
233 is constant throughout the simulation. The simulation ends randomly, after the driver subclone
234 reaches a critical prevalence (see supplement for more details). The Gillespie algorithm has been
235 frequently used to simulate stochastically dividing cells⁴⁸⁻⁵⁴, although simulations with special
236 attention to cell cycle have also been recommended⁵⁵.

237

238 During simulated growth, we assigned a “driver” mutation with additional propagating effects
239 from nearly neutral to high ($k=1.1, 2, 3,$ and 4), thus leading to faster growth for the respective
240 subpopulation that contains the specific mutation. Using conservative assumptions, these scalar
241 values represent a range of projected selection coefficients s^* from 0.001 to 0.03 in biologically
242 sized populations (see Supplement). For each simulation, we calculated each mutation's
243 frequency in the total population and ordered them based on that frequency. Then, by applying
244 our method we calculated the ranking distance D (as the number of ordered mutations) between
245 the true and our predicted driver (growth peak), as well as the driver's scalar effect k .

246

247 We tested our method's performance in simulated tumors of lower coverage and different
248 effects. Higher sequencing depth and scalar effect k provided more accurate results and improved
249 our method's implementation (Figure 2a,b). Lower coverage was associated with worse k
250 calculations and driver predictions, as well as lower positive predictive values (PPVs). For weak

251 drivers, low sequencing coverage made their identification more difficult. Absolute median
252 ranking distance $|\tilde{D}|$ was 41 for coverage $100x/k=2$, compared to 13 for coverage $1000x/k=2$ and
253 $|\tilde{D}|=11$ for coverage $1000x/k=4$ respectively. In general, driver identification required either a
254 higher than 100x coverage, or a stronger effect (i.e. $k>2$, $s^* >0.01$ for a projected cell population
255 of 1,000,000 cells) (Figure 2i).

256

257 Overall, we were able to well approximate the driver's occurrence and effect (Figure 2). For the
258 birth and death model with simulated coverage 1,000x, the median predicted estimation for
259 simulated effects $k=2$, $k=3$, and $k=4$ was 2.3, 2.9, and 3.8, respectively (Figures 2ii, S6b).

260 Moreover, the median ranking distance \tilde{D} between simulated and predicted drivers with effect
261 $k=1.1$ (nearly neutral), $k=2$, and $k=3$ was 71, 3. 5, and 6, respectively. The corresponding median
262 distances for random mutations were 73, 43, and 41 (Figure S1c). For our nearly neutral
263 simulations ($k = 1.1$, $s^* \sim 0.001$ for a projected cell population of 1,000,000 cells) the median
264 distance \tilde{D} in driver predictions and random predictions was very similar and not significant.

265

266 *Neutral and non-neutral simulations with added stochasticity in mutation rates*

267

268 To further test our model on a separate independent simulation dataset, we applied our method to
269 a) neutral simulations of tumor progression and b) non-neutral simulations for various growth
270 scenarios, as previously developed and described by Williams et al 2016 and Williams et al 2018
271 (see Supplement). These simulations, although also based on the Gillespie growth model,
272 included added stochasticity with varying mutation rates during tumor progression ($\bar{\mu}=10$
273 mutations per cell division). For every simulation, both neutral and non-neutral, we identified our

274 model's highest predicted effect peak, calculated the effect k and absolute median ranking
275 distance $|\widetilde{D}|$ between the simulated and predicted driver in number of ranked mutations. Various
276 scenarios for non-neutral growth included a wide range of simulated selection coefficients s (0 to
277 33, for a population size of 10,000 cells), categorized driver's VAF (small 0.1-0.2 ; medium 0.2-
278 0.3 ; large 0.3-0.4) and larger cell population projections using population genetic models and
279 method adjustments. Corresponding neutral simulations were also generated using the same
280 population parameters. Overall, and in agreement with our previous analyses, our results suggest
281 a small overlap between neutral and non-neutral peaks for weak drivers (figure 2c and S1f) and
282 highly significant driver predictability when the predicted driver effect was larger than our
283 (narrow) neutral-effect distribution (Figure 2c,d and S1g-i). For instance, for simulated
284 populations of 10,000 cell without projection ($0 < \text{simulated } s < 33$) and 1000x coverage our
285 method provided accurate driver detections when the predicted effect was larger than $k=1.29$
286 with $|\widetilde{D}| \sim 50$ mutations compared to 444.5 for random. These results are directly comparable to
287 our previous analyses, considering the new mutation rates. Similarly, for a projected cell
288 population of 1,000,000 cells, our method provided accurate driver detection for projected
289 selection coefficient $s^* > 0.05$ (Figure 2d). Larger population projections typically decreased the
290 predicted effect k^* and selection coefficient s^* , but did not affect our method's ability to detect
291 drivers (Figure S1k) as these projections also decreased the standard deviation of our neutral-
292 effect distribution (predicted k^* for neutral effect peaks). When we combined 140 neutral with
293 360 non-neutral simulations, drivers with medium final VAF showed the highest correlation
294 between simulated selection coefficients and our method's predicted scalar k effects ($r=0.60$,
295 Figure 2c). Drivers with lower final VAFs (small $\sim 0.1-0.2$) provided slightly lower correlation
296 but had the highest driver detectability, with $|\widetilde{D}| = 46$ mutations between the simulated and

297 predicted driver (Figure S11), where random $|\widetilde{D}|$ was 444.5 mutations. A (tenfold) higher $|\widetilde{D}|$
298 here is expected since for these simulations we assumed 10 instead of 1 mutation per cell
299 division.

300

301

302 *Synthetic results using coalescent-based model: estimator \hat{r} for non g-hitchhikers*

303

304 We also tested the behavior of the estimator for r (Eq. 6) on non-g-hitchhiking mutations (i.e.
305 when the assumption that the mutations are generational hitchhikers is not satisfied). For this
306 purpose, we used coalescent theory to estimate the variation in density of mutations across the
307 VAF spectrum for a variety of models (see Supplement). We first analyzed the behavior in a
308 constant-size population, and then in populations with increasing and decreasing exponential
309 growth. Our analysis shows that the growth indicator does not qualitatively change its behavior
310 in this context, so that negative values continue to represent periods of negative growth, and
311 large positive values represent periods of positive growth. However, here we expect a small
312 positive value in the case of zero growth (Figure 2e, S2).

313

314

315 **Growth Patterns and Biological Disruptions in 993 Linear Tumors from the PCAWG**

316

317 Using 993 linear tumors from the PCAWG consortium, we explored the different patterns and
318 dynamics of tumor growth based on our model's assigned growth rates. Tumor "linearity"
319 (where no parental subclone has two or more children subclones) further ensures that tumor

320 subclones do not intermingle and that higher VAF is associated with earlier occurrence. We note
321 that mutational frequency as described in our equations corresponds to $2 \times \text{VAF}$, with correction
322 for purity and copy-number variations. These VAF corrections were obtained from PCAWG and
323 are not implemented in any way by our method, which only considers a final mutational
324 frequency. Using our model, each mutation i from sample in our database is assigned a potential
325 positive or negative growth value r_i and a driver effect k_i . Under ideal conditions, for each
326 sample, a vector of effect-peaks $r_{i-1} \times k_i$ corresponds to potential drivers at position i . However,
327 noise, coverage, and growth stochasticity can cause these peaks to represent the potential
328 presence of a nearby driver, especially in low coverage sequenced tumors (see Figure 3a,b).

329

330 To identify growth patterns across individual tumors, we i) normalized each mutation's growth
331 rate based on the sample's maximum growth value; ii) divided the ordered mutations into 20
332 bins; and iii) applied K-means clustering to the average normalized value per bin. Our results
333 highlighted three main clustering patterns (Figure 3c). As expected, most tumors ($n=525$)
334 showed logistic growth with an increasingly higher growth rate at the beginning and a
335 stabilization at the later stages. For many tumors ($n=366$), an early high growth period was
336 followed by a stagnation and potential reduction in tumor size. This effect could also be
337 artificially enhanced due to sampling errors for mutations with low VAF (during late tumor
338 progression). The last group of tumors ($n=102$) showed relatively steady, continuous growth.
339 However, it is uncertain whether this pattern represents tumors that were sequenced early.
340 Further, some types of cancer seemed to prefer specific growth patterns (Figure 3c).

341

342 By modeling tumor growth, we can find mutations during positive or negative growth periods in
343 single or multiple individual samples. Through positive “*growth enrichment*”, we characterized
344 the degree to which one type of mutation (e.g., TSGs/TP53, nonsynonymous) or region (e.g.,
345 TP53) was significantly enriched and associated with periods of positive growth across multiple
346 samples. We then compared each mutation type to random mutations from their respective
347 samples (see Supplementary Methods for details). To confirm whether we could detect any
348 signal of selection at the gene level, we compared positive *growth enrichment* for mutations
349 between i) the Vogelstein gene list²⁹; ii) a comparable list (in mutational numbers) of randomly
350 selected genes; and iii) a list of assigned drivers from the PCAWG consortium^{33,56}. As expected,
351 PCAWG-assigned driver SNVs clearly showed the highest positive enrichment, followed by
352 SNVs that were not individually called by PCAWG as drivers but that fall within the Vogelstein
353 driver gene list (Figure 3d). We note, however, that our random gene list did show a small
354 positive enrichment, as this list contains several often-mutated genes and potential drivers or
355 mini-drivers. We obtained similar results when we repeated the comparison while considering
356 the difference between additional mutational effect against a random distribution (Figure S3).
357
358 In an effort to better understand the micro-environment of tumor dynamics, the selective forces,
359 and the biological processes that are most keenly affected by tumor progression, we analyzed a
360 list of 1,000 most mutated genes in the PCAWG samples where we identified 293 genes with
361 significant overall association with positive growth (Suppl. Table 1). Then we further tested
362 these genes for Gene Ontology (GO) enrichment. As expected, developmental and differentiation
363 processes were highly enriched during periods of positive growth, showing signals for being

364 under positive selection. Interestingly, we found that genes related to multicellular processes
365 showed the highest enrichment based on raw p-value (Figure 3e, Suppl. Table 2).

366

367 **Tumor-Suppressor Genes vs. Oncogenes**

368

369 Based on each mutation's genomic properties (e.g., genomic position, coding vs. non-coding,
370 TSG vs. oncogene, cancer type, and gene ontology annotation), we can examine whether the
371 specific type of mutation (or "mutation element") is statistically enriched during periods of
372 positive growth when compared to random mutations from their respective samples (see
373 supplementary methods). However, the more specifically that we defined a mutation type, the
374 fewer mutations that corresponded to this category. For example, the Vogelstein TSGs in our
375 dataset contain 321 missense and 103 nonsense mutations, whereas TP53 in our dataset contains
376 71 nonsynonymous mutations and 13 nonsense mutations. Unfortunately, for many tumor genes
377 and cancer types, we currently have a small number of mutations, precluding significance in the
378 results.

379

380 A recent study by Kumar *et al.* suggested that high-impact mutations should have more clear
381 positive effects on tumor growth when they are located in TSGs versus oncogenes³⁸. This is
382 expected, as generally a "defected" oncogene with reduced expression should not favor cancer
383 progression. To better understand the behavior of TSGs and oncogenes, we tested for positive
384 enrichment of synonymous, non-synonymous, premature stop, promoter, and intronic mutations
385 (Figure 4). As expected, our results showed significant enrichment of missense and nonsense
386 mutations in TSG regions. During periods of positive growth, 45 nonsense and 128 missense

387 mutations corresponded to an average of 37.4 and 117.96 random mutations, respectively (100
388 bootstraps replicates, p values=7.823348e-30 and 1.632649e-23). Interestingly, promoter and
389 intronic regions also showed a significant positive effect on tumor growth, suggesting that some
390 non-coding mutations in TSGs might favor positive growth (Figure 4a).

391
392 In the case of oncogenes, we did not find significant enrichment of missense mutations, but we
393 did find significant association between their promoter regions and positive growth (Figures 4b).
394 This might be due to many reasons including the pancancer nature of our analysis, lack of power
395 and small sample size, our modeling assumptions, or the noise due to low sequencing coverage
396 per tumor sample. However, many genes including oncogenes might be under negative selection,
397 with only a small subset of their respective mutations being favorable to cancer growth.
398 Moreover, high-impact mutations in oncogenic regions do not necessarily favor tumor growth.
399 Indeed, our data contain only four nonsense mutations in oncogenic regions. Some oncogenes
400 such as MET and CTNNB1 showed slight overall negative enrichment, but their nonsynonymous
401 mutations, especially in specific cancers, showed enrichment during periods of positive growth
402 (Figure S4).

403
404 To detect mutations during positive growth periods, we applied our model to individual types of
405 mutations (i.e., missense, synonymous, intronic, nonsense, and promoter) for each Vogelstein
406 gene. Overall, our results identified various mutation elements including promoters, nonsense,
407 and missense with significant effects (Figure 4c). Interestingly, synonymous BLC2 mutations
408 that occurred near an early positioned mutational hotspot were significantly associated with
409 positive growth (Figures 4c and S5). Synonymous mutations are not generally considered to be

410 important in cancer; however, previous studies have reported recurrent synonymous F17F
411 mutations in BLC2-like 12, where regulatory hsa-miR-671-5p alters the gene's expression⁴⁵.

412

413 **Predicting Growth Peaks and Driver Effects on a Model Ultra-Deeply-Sequenced AML** 414 **Tumor**

415

416 In addition to the 993 PCAWG low-coverage tumor samples, we implemented our model on an
417 ultra-deeply-sequenced AML (>250x) liquid tumor. A ultra-deeply-sequenced tumor provides
418 more accurate global variant allele frequencies, which should in turn allow for better estimation
419 of model parameters¹²

420

421 In general, the predicted peaks of our model mapped very closely to mutations from known
422 cancer genes (Figure 5). Deep valleys followed by the highest growth peaks corresponded with
423 close approximation to the three missense mutations from known cancer genes (IDH1, IDH2,
424 and FLT3, p-value < 2.2e-16). Thus, in agreement with previous studies^{35,36}, the derived growth
425 patterns suggested three to five major genetic hits from cancer mutations in order to render tumor
426 growth permanent.

427

428 Additionally, we used all the mutations in our previous database to evaluate those in the deeply
429 sequenced AML in order to identify new candidates associated with positive growth. As a result,
430 we further identified five additional candidates from the ultra-deep AML sample that belong to
431 genomic elements associated with positive growth (Figure 5d). These additional candidates
432 consist of four missense mutations (SRCAP, CPS1, GLI1, and COL18A1) and one intronic

433 mutation (MAP3K1), which appeared to align near observed, previously unexplained periods of
434 initial growth. Previous recent studies have also linked CPS1 and GLI1 to various cancers^{57–60}.
435 Finally, based on our PCAWG database, for each driver candidate we detected possible positive
436 enrichment across varying effect ranges [0.9, 1.1, 1.3, 1.5, 1.7, 1.9, and 2.1] (Figure S6).
437 Indicatively, our independent estimation of mutational effect suggested a high correlation when
438 compared to the calculated effect using the deep sequenced model AML tumor (Figure S6).

439

440

441 **Discussion**

442 Most approaches to identify driver candidates are based on recurrent mutations and large
443 cohorts²³. More recently, studies have probed tumor selection either through deviation from
444 background metrics or by using VAF distribution to quantify the subclonal effect^{16,19,22,61,62}.
445 Here, we present a framework that models tumor progression using generational hitchhikers and
446 localized time re-optimizations using mutational frequencies from individual samples to i)
447 determine periods of positive or negative growth, ii) suggest the presence of candidate drivers
448 and estimate their effect on tumor progression, and iii) detect genomic regions or mutation
449 elements that are associated with positive or negative growth periods. Overall, our work
450 highlights the importance of whole genome deep sequencing for modelling tumor progression.

451

452 When we applied our framework to 993 individual tumors from the PCAWG consortium, our
453 growth analysis indicated different growth patterns across cancer types, including steady growth,
454 sigmoidal growth, and modes of stagnation. Determining tumor progression can be useful in
455 understanding each tumor's historic aggressiveness, and the effect of driver mutations on tumor

456 progression (VAFs used by our method typically represent past growth, as latest mutations tend
457 to have undetected frequency in our sample). Additionally, we identified several biological
458 processes significantly affected by tumor progression, including genes involved in
459 multicellularity. These results might indicate an evolutionary transition during tumor progression
460 from multi-cell functionality to single-cell selection.

461
462 As expected, we found significant enrichment of known PCAWG drivers, Vogelstein cancer
463 genes, and nonsense and missense mutation TSGs during periods of positive growth. In
464 accordance with some previous studies⁴¹⁻⁴⁴, our results also suggested that a small proportion of
465 intronic mutations could affect TSGs (but not oncogenes), whereas some synonymous mutations
466 could affect oncogene (but not TSG) expression. Even though defective splicing in TSGs or
467 changes in the negative regulation of oncogenes are not entirely unexpected⁴⁵, non-coding
468 mutations are not generally considered to be major driver events in tumor progression. Thus, it is
469 possible that our results are subject to analytical (e.g., model parametrization, initial parameters,
470 window size selection, low sequencing coverage, sample size) and biological (e.g, hitchhiking)
471 error.

472
473 Using variant allele frequency to quantify driver effects and tumor progression can be
474 challenging. Our analysis might be subject to different types of bias, including sequencing noise,
475 growth stochasticity, model parameterization, low sequencing coverage, tumor ploidy,
476 subclonality, and a low number of tumor samples per cancer or mutational element. Under a
477 neutral model, our method would still detect some growth peaks or suggest the presence of weak
478 drivers. These are false positive predictions, possibly due to noise which results in various signal

479 perturbations in the VAF spectrum, or potential genetic drift. Moreover, our model does not
480 consider the potential effects from deleterious passenger mutations or sequencing errors on the
481 VAF spectrum. However, we consider that -if not depleted- most deleterious mutations should
482 have a small VAF in our sequenced sample. Similarly, we expect that sequencing errors tend to
483 produce spurious mutations of extremely low VAF, which are ignored by our framework.
484 Although some researchers are skeptical of the plausibility of “VAF quantification”^{20,63}, recent
485 analyses have also confirmed that it can be achieved even at low sequencing coverage¹⁶. At the
486 same time, as sequencing cost decreases exponentially, ultra-deep whole genome sequencing for
487 a larger number of samples will become trivially within reach. This is critical for the
488 personalized assessment and parametrization of single samples.

489
490 Similar to previous Darwinian, bacterial, and viral evolution analyses, modeling the variations of
491 cell populations allows us to associate these variations with specific events, even at a single
492 sample level. Our work contributes to our understanding of cancer evolution by directly
493 assessing tumor sample progression at the time of the driver event. This assessment can be very
494 critical for therapeutic strategies and drug selection^{33,34}. Our framework presents opportunities
495 for personalized diagnosis via modeling the tumor’s progression using deep sequenced whole
496 genome data from one single individual.

497

498

499 **Acknowledgements**

500 We thank the PCAWG consortium for the state-of-the-art preliminary analysis and management of the data. We
501 thank the members of the Gerstein lab Cancer Genomics group (SK, SL, PDM) for helpful discussion on the method
502 development and analysis of the data.

503

504 **Author contributions**

505 L.S. conceived of the project, designed and performed the experiments. W.M. designed and developed the
506 simulations. J.W. performed the coalescent analysis and designed the simulations. L.S. drafted the manuscript. L.S.
507 and M.G. wrote the manuscript. All authors read and approved the final manuscript.

508

509 **Competing interests**

510 The authors declare no competing interests.

511

512 **Corresponding authors**

513 Correspondence to Mark Gerstein.

514

515 **Author information**

516 *Affiliations:*

517 1. *Program in Computational Biology and Bioinformatics, Yale University, New Haven, CT, 06520, USA.*

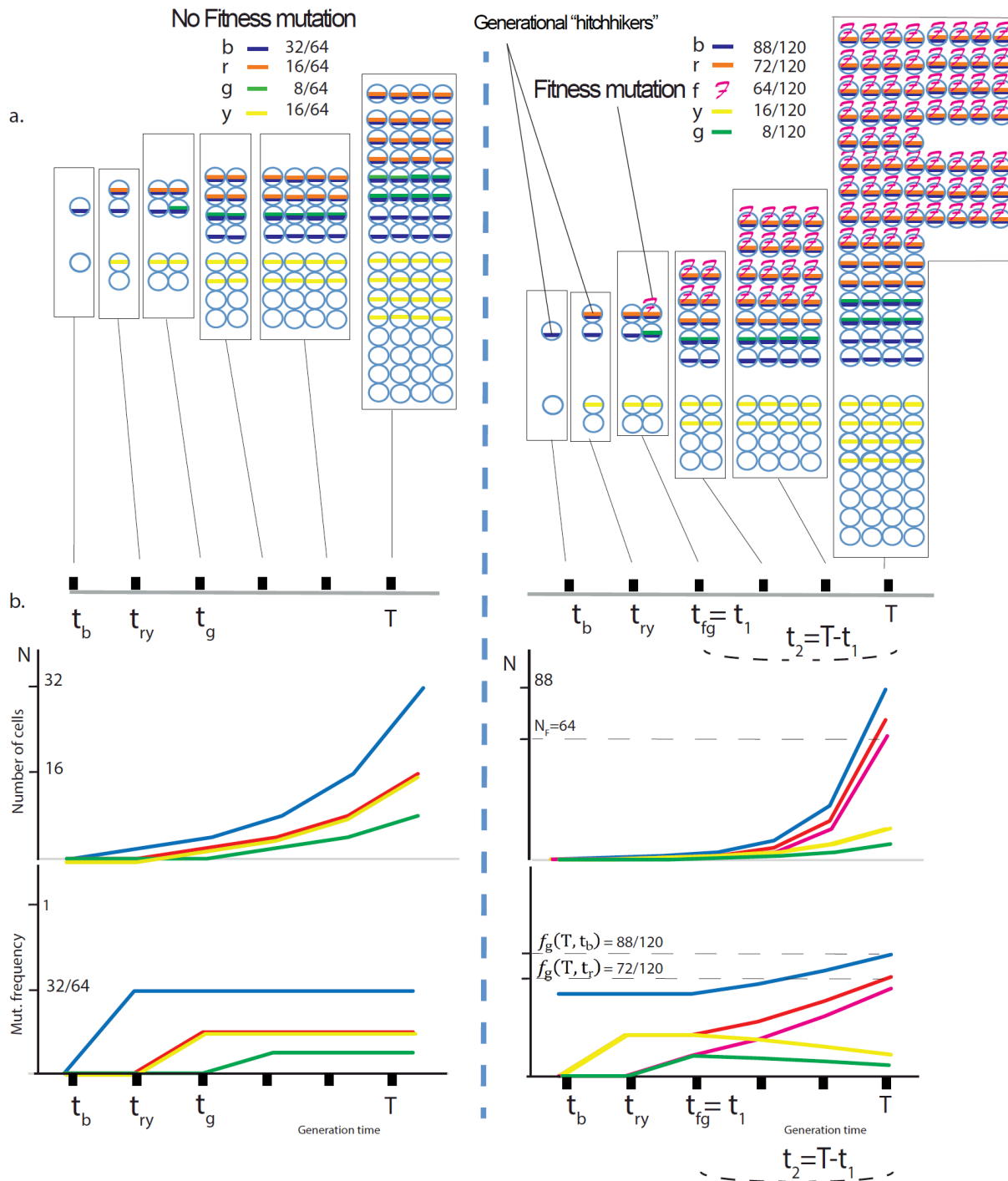
518 2. *Department of Molecular Biophysics and Biochemistry, Yale University, New Haven, CT, 06520, USA.*

519 3. *Department of Computer Science, Yale University, New Haven, CT, 06520, USA.*

520 4. *Center for Biomedical Data Science, Yale University, New Haven, CT, 06520, USA.*

521 **Corresponding author*

Figure 1



522

523

524

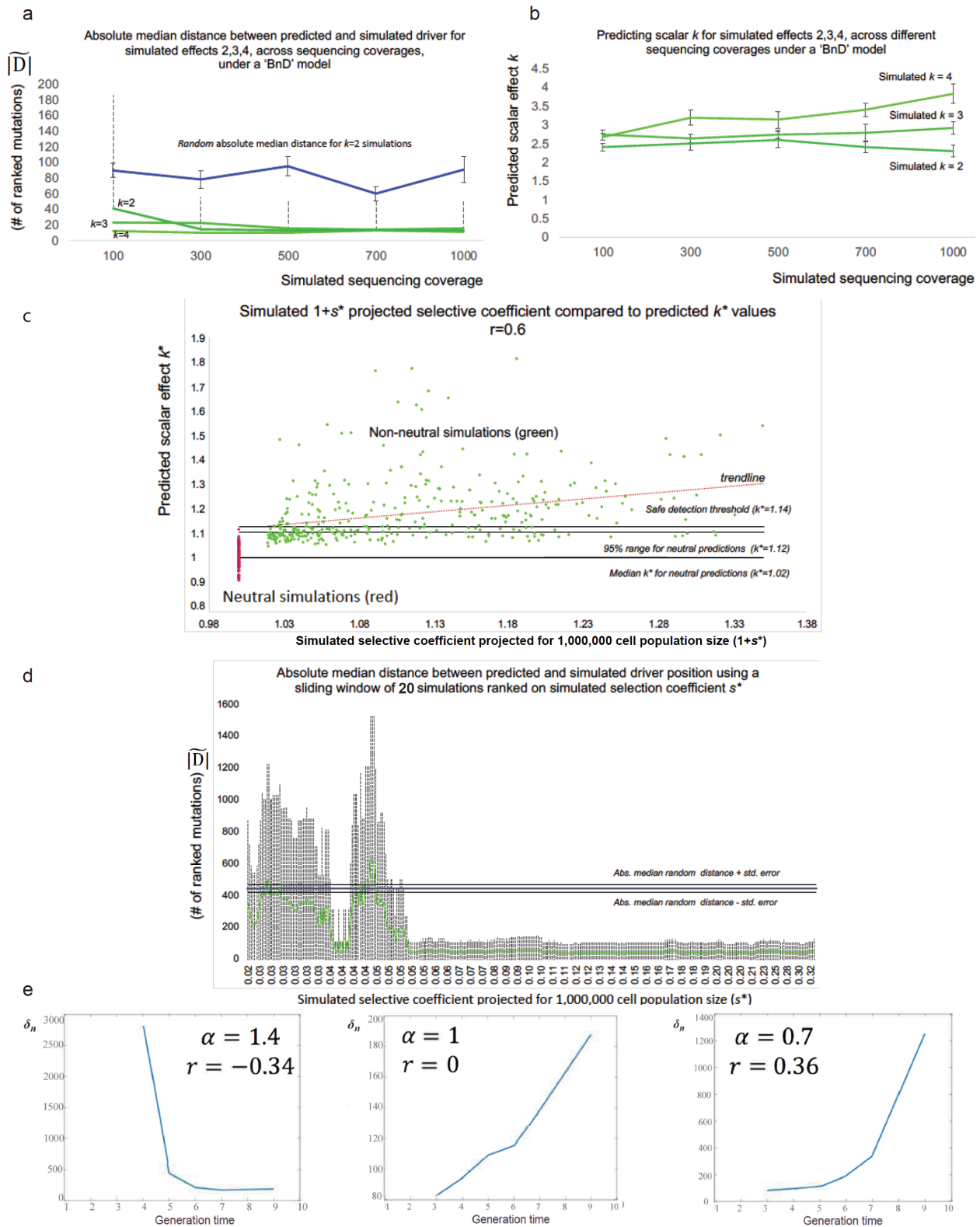
Figure 1. Generational (g-)hitchhikers have increased frequency, which in turn is

525

dependent on the effect of the fitness mutation in the population. We consider a simple

526 population of cancer cells that grows exponentially $N(t) = e^{rt}$; for simplicity, we assign one mutation
527 per cell division. At the time of biopsy T , the frequency of a mutation occurring at time t_n would be equal
528 to $f_n(T, t_n) = \frac{e^{r(T-t_n)}}{e^{rT}} = e^{-rt_n}$. At time t_1 , a mutation occurs that increases the growth rate r of the
529 specific subpopulation by a scalar multiplier k , such that the new population is now expanding as $N_F =$
530 e^{krt_2} . Thus, at the time of biopsy $T=t_1+t_2$, we expect a generational (g-) “hitchhiking” mutation that
531 occurred at time $t_m < t_1$ to have a frequency equal to $f_g(T, t_m) = \frac{e^{r(T-t_m)} + N_F e^{-rt_2}}{N_{tot}}$, where N_{tot} is the total
532 number of cells (or mutations) and N_F is the number of cells that contain the fitness mutation that
533 occurred at t_1 and expanded for t_2 . Therefore $N_F = e^{krt_2}$. In a) we show the mutational frequencies at
534 the time of biopsy T for two growth models; one neutral and one with a fitness mutation
535 occurring at time $t_1=t_{fg}$. Hitchhiking mutations ‘b’ (“blue”), ‘r’ (“red”), as well as passenger
536 mutations ‘g’ (“green”) and ‘y’ (“yellow”), also occur at different time points. **b)** Under an
537 exponential model with a fitness mutation occurring at time $t_1=t_{fg}$, *hitchhikers* ‘b’ and ‘r’ show an
538 increased frequency compared to neutral, subject to time and effect of the fitness mutation.
539 Passenger mutations ‘y’ and ‘g’ that occurred before or with the fitness mutation, but on a
540 different cell lineage, end up with lower frequencies. We characterize mutations ‘b’ and ‘r’
541 as *generational (g-) hitchhikers* since they mark the population’s generational growth.
542

543



544

545

546

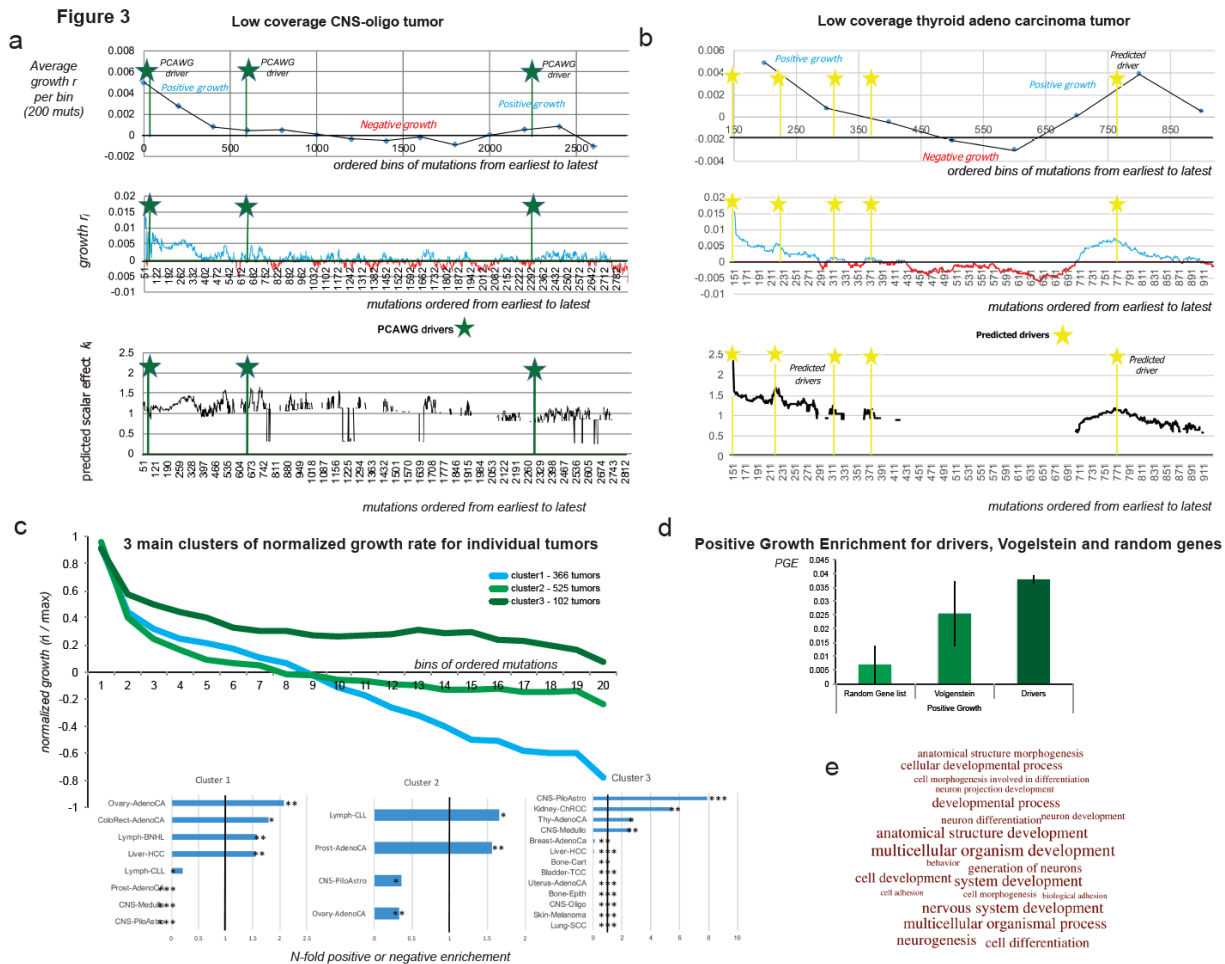
Figure 2.

547 **Figure 2. Higher coverage and stronger drivers improve driver detectability and effect**
548 **prediction.** In **a)** using a total of 541 simulations of tumor growth under a birth and death model
549 (an average of 36 simulations per sequencing depth per simulated effect), we show the absolute
550 median distance $|\widehat{D}|$ as in ‘absolute number of ordered mutations’ between our predicted and the
551 simulated driver for different sequencing depths. With the exception of $k=2$ and sequencing
552 coverage equal to 100x (p value=0.015), we were able to significantly detect the driver’s
553 presence for depth coverages as low as 100x (p value <0.005). Blue line represents the random
554 absolute median distance as derived by selecting a random mutation from each simulation and
555 calculate the absolute distance to the simulated driver. Dotted lines represent the $2 \cdot \sigma$
556 deviation from $|\widehat{D}|$ while capped bars represent the median’s standard error. For convenience and
557 clarity, we only show bars for $k=2$. In **b)** again using a total of 541 simulations of tumor growth
558 under a birth and death model, we show that higher depth coverage provides more accurate k
559 predictions. Low coverage usually results in predicting a lower effect. Capped bars represent the
560 standard error of the median effect prediction. The three lines represent simulations with
561 simulated effect of 2,3 and 4. In **c)** By implementing the Williams et al 2018 algorithm for
562 neutral and non-neutral simulations, we simulated 360 non-neutral and 140 neutral tumor
563 progressions, with a populations size of 10000 cells. Then, we adjusted our effect predictions to
564 account for a larger population with effect size equal to 1,000,000. In addition, we also adjusted
565 the simulated selection coefficient s^* for the same population size. In this figure we show the
566 correlation between the simulated adjusted coefficient ‘ $1+s^*$ ’ against our adjusted predicted k^* .
567 By including both neutral and non-neutral simulations in our sample Pearson correlation was
568 $r=0.6$. In **d)** after ranking simulated driver coefficients s^* for every non-neutral simulation
569 (adapted from Williams et al), we used a sliding window of 20 ranked simulations to estimate the

570 absolute median distance (and 95% deviation) between the simulated and predicted driver within
571 every window of 20 ranked simulations. Dotted lines represent a $2 \times \sigma$ deviation (95%). When
572 our simulated selection coefficient was stronger than 0.05^* our driver detection became highly
573 accurate. Blue line represents absolute median distance for random predictions (444.5), while
574 black lines represent the median standard error for these expectation (24.5). Simulated
575 coefficients s^* have been projected for a population with effect size of 1,000,000. In e) Using
576 Kingman's coalescent theory, for a length of time T_n with n lineages, we show that the growth
577 \hat{r} estimator remains qualitatively unchanged (positive or negative) even for non g-hitchhikers.
578 By approximation, the mutational density δ_n within windows $[1/n \ 1/(n - 1))$, whose lengths
579 are L_n is equal to $\delta_n = \frac{M_n}{L_n} \propto 2\mu n$. As mutational density δ_n increases with n , and hence with
580 time, \hat{r} estimator is predicted to take positive values for both constant and varying size
581 populations. Similarly, for negative growth values, density δ_n decreases with time. A small
582 positive bias is observed in cases of growth $r=0$, as the pattern reverses. Using a population
583 model $N^{t+1} = \alpha N^t$, we let $\alpha > 1$ corresponding to a decreasing (time is indexed in reverse) and α
584 < 1 corresponding to an increasing population.

585
586
587

588



589

590

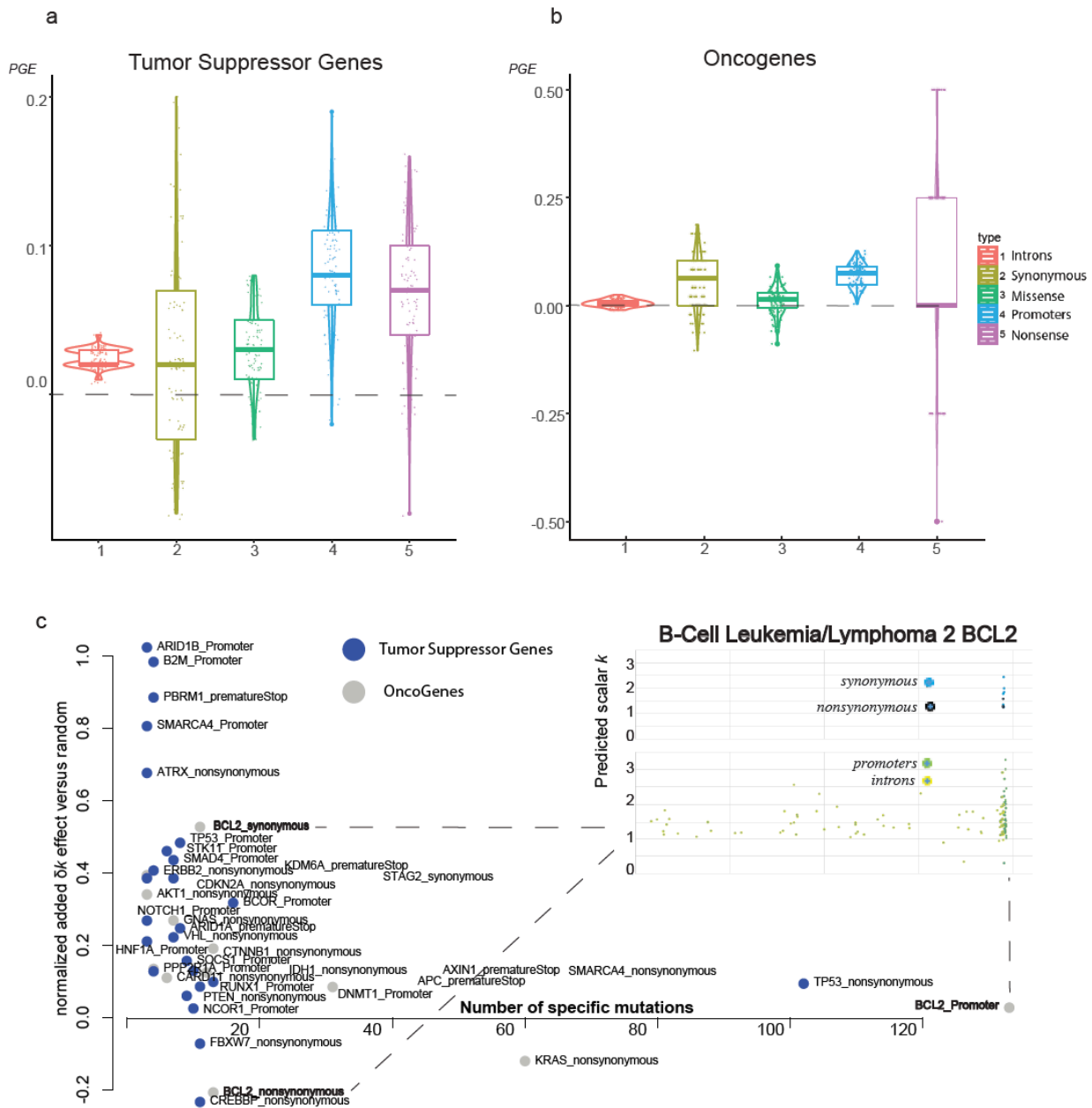
Figure 3. Growth patterns and growth association. Across 993 linear tumors from PCAWG

591 consortium we expect an under-selection mutation to be associated with periods of positive
 592 growth (see supplementary methods). We compared several mutation types (driver mutation,
 593 mutation within geneX, within GO categoryX), to a random distribution from their respective
 594 sample for association with positive growth. **a-b)** we show the i) averaged growth progression,
 595 ii) mutational growth and iii) mutational effect, for a single low coverage CNS-oligo tumor and a
 596 single low coverage thyroid adenocarcinoma tumor without any PCAWG-identified drivers.
 597 Green asterisks denote the ordered position of a PCAWG-predicted driver within the sample.
 598 Yellow asterisks denote a growth peak and putative driver presence. In **c)** we derived three main
 599 growth patterns (steady growth, sigmoid growth, stagnation/shrinkage) for 993 linear tumors, as

600 they were grouped using a k-means clustering algorithm. Various cancer types showed specific
601 enrichment or depletion for the three clusters **d)** PCAWG drivers and Vogelstein genes showed
602 significant positive growth enrichment compared to a list of random highly mutated genes. **e)** We
603 show the GO enrichment for the 20 most affected biological processes, when we use 293 genes,
604 significantly associated with periods of positive growth.

605
606
607
608

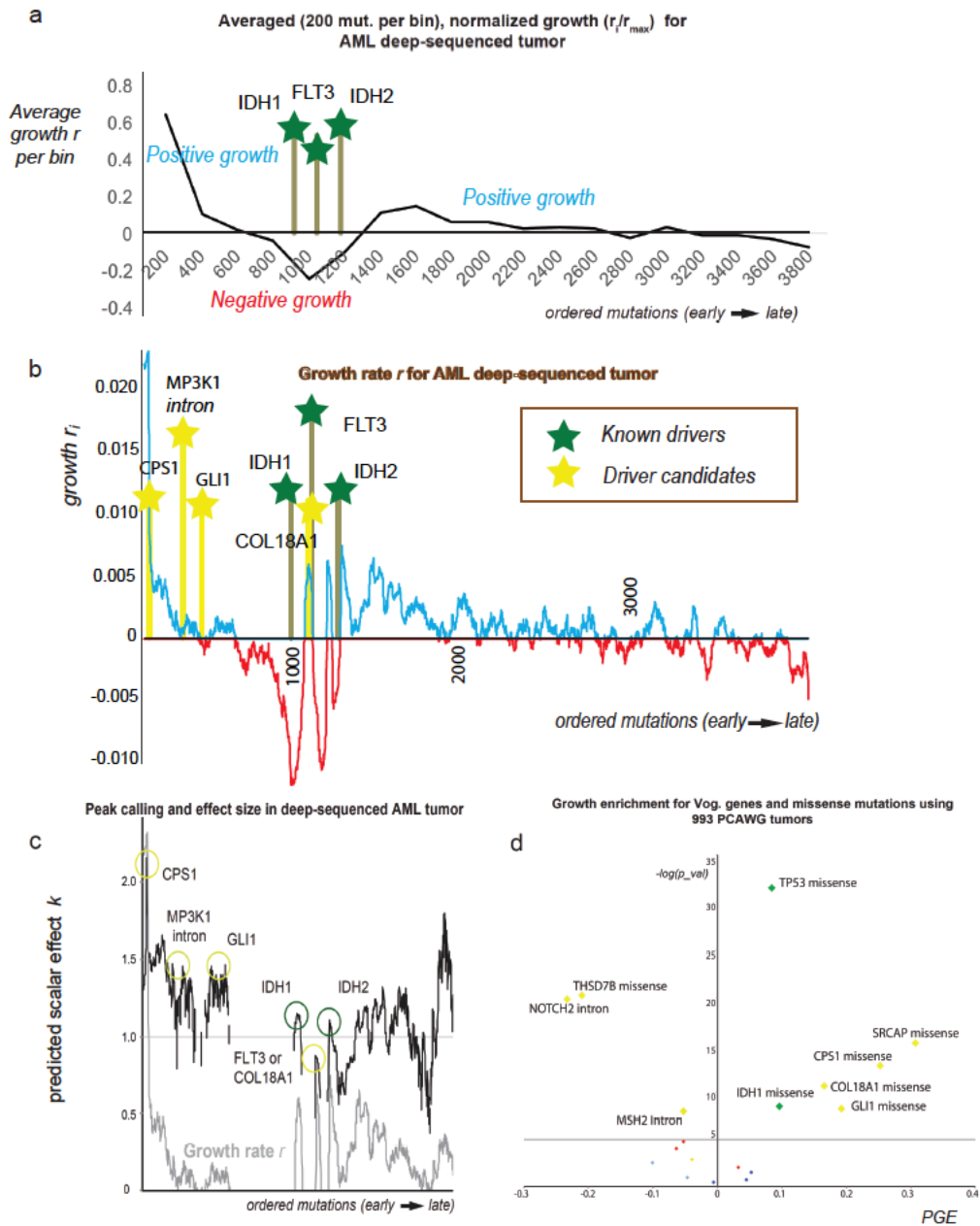
Figure 4



609
 610 **Figure 4. Mutational elements from tumor suppressor genes and oncogenes showing**
 611 **growth enrichment.** We show the positive growth enrichment across different mutation types
 612 (introns, synonymous, missense, nonsense, promoters) for **a)** Vogelstein tumor suppressor genes
 613 and **b)** Vogelstein oncogenes. In **c)** we plot gene elements (e.g. {GeneX_mutation type}) from

614 Vogelstein gene list that showed significant positive or negative enrichment. We further zoom in
 615 to BCL2's genomic region to map missense, nonsynonymous, promoter and intronic mutations.
 616

Figure 5



617
 618
 619

620 **Figure 5. Mapping candidate drivers during tumor progression on an ultra-deeply-**
621 **sequenced AML liquid tumor.**

622 In **a)** we show the averaged growth progression for an AML deep sequenced tumor. We ordered
623 the sample's mutations from highest to lowest frequency and divided them into bins of 200
624 mutations. Three cancer mutations hit the tumor to establish a permanent growth (cancer
625 mutations denoted by green bars). In **b)** we plot the mutational growth r_{i-1} for each mutation
626 across tumor progression. The three cancer genes (IDH1-missense, FLT3-missense, IDH2-
627 missense) aligned well with 3 of our top 5 growth peaks (p-value < 2.2e-16). Candidate driver
628 mutations -denoted by yellow bar- that we identified from our PCAWG database as being
629 associated with positive growth (see also 'd')) aligned well with early –previously unjustified
630 growth peaks. In **c)** we show each mutation's effect in tumor progression. Effect peaks
631 corresponds to putative drivers. **d)** By using our PCAWG database from our previous analysis,
632 we tested which mutations from the deep sequenced sample were associated with positive
633 growth. Overall, we found 6 mutation types that showed positive enrichment across 993
634 PCAWG tumors including TP53-missense (appeared during metastasis), IDH1-missense,
635 COL18A1-missense, CPS1-missense, GLI1-missense and SRCAP-missense. Missense TP53
636 and SRCAP mutations are not included in graph (b) as they were metastatic mutations. For
637 association with positive growth we tested all missense mutations (eg CPS1-missense), and
638 every mutation in the sample from Vogelstein cancer genes (eg. NOTCH2-Intron).

639

640

641

642

643 **References**

644

- 645 1. Heng, H. H. Q. *et al.* Stochastic cancer progression driven by non-clonal chromosome
646 aberrations. *J. Cell. Physiol.* **208**, 461–472 (2006).
- 647 2. Loeb, L. A. Human cancers express mutator phenotypes: origin, consequences and
648 targeting. *Nat. Rev. Cancer* **11**, 450–457 (2011).
- 649 3. Nowell, P. C. The clonal evolution of tumor cell populations. *Science* **194**, 23–28 (1976).
- 650 4. Navin, N. *et al.* Inferring tumor progression from genomic heterogeneity. *Genome Res.*
651 **20**, 68–80 (2010).
- 652 5. Navin, N. *et al.* Tumour evolution inferred by single-cell sequencing. *Nature* **472**, 90–94
653 (2011).
- 654 6. Yngvadottir, B., Macarthur, D. G., Jin, H. & Tyler-Smith, C. The promise and reality of
655 personal genomics. *Genome Biol.* **10**, 237 (2009).
- 656 7. Mu, X. J., Lu, Z. J., Kong, Y., Lam, H. Y. K. & Gerstein, M. B. Analysis of genomic
657 variation in non-coding elements using population-scale sequencing data from the 1000
658 Genomes Project. *Nucleic Acids Res.* **39**, 7058–7076 (2011).
- 659 8. Bernstein, B. E. *et al.* An integrated encyclopedia of DNA elements in the human genome.
660 *Nature* **489**, 57–74 (2012).
- 661 9. Campbell, P. J. *et al.* Pan-cancer analysis of whole genomes. *bioRxiv* (2017).
662 doi:10.1101/162784
- 663 10. Haussler, D. *et al.* *A Million Cancer Genome Warehouse.* (2012).
- 664 11. Campbell, P. J. *et al.* Subclonal phylogenetic structures in cancer revealed by ultra-deep
665 sequencing. *Proc. Natl. Acad. Sci.* (2008). doi:10.1073/pnas.0801523105

- 666 12. Griffith, M. *et al.* Optimizing Cancer Genome Sequencing and Analysis. *Cell Syst.* (2015).
667 doi:10.1016/j.cels.2015.08.015
- 668 13. Krimmel, J. D. *et al.* Ultra-deep sequencing detects ovarian cancer cells in peritoneal fluid
669 and reveals somatic *TP53* mutations in noncancerous tissues. *Proc. Natl. Acad. Sci.*
670 (2016). doi:10.1073/pnas.1601311113
- 671 14. Subramanian, A., Shackney, S. & Schwartz, R. Inference of tumor phylogenies from
672 genomic assays on heterogeneous samples. *J. Biomed. Biotechnol.* **2012**, (2012).
- 673 15. Sottoriva, A. *et al.* A big bang model of human colorectal tumor growth. *Nat. Genet.*
674 (2015). doi:10.1038/ng.3214
- 675 16. Williams, M. J. *et al.* Quantification of subclonal selection in cancer from bulk sequencing
676 data. *Nat. Genet.* (2018). doi:10.1038/s41588-018-0128-6
- 677 17. Turajlic, S. *et al.* Deterministic Evolutionary Trajectories Influence Primary Tumor
678 Growth: TRACERx Renal. *Cell* (2018). doi:10.1016/j.cell.2018.03.043
- 679 18. Mitchell, T. J. *et al.* Timing the Landmark Events in the Evolution of Clear Cell Renal
680 Cell Cancer: TRACERx Renal. *Cell* (2018). doi:10.1016/j.cell.2018.02.020
- 681 19. Tarabichi, M. *et al.* Neutral tumor evolution? *bioRxiv* (2017). doi:10.1101/158006
- 682 20. Williams, M. J., Werner, B., Barnes, C. P., Graham, T. A. & Sottoriva, A. Reply:
683 Uncertainties in tumor allele frequencies limit power to infer evolutionary pressures.
684 *Nature Genetics* (2017). doi:10.1038/ng.3877
- 685 21. Williams, M. J., Werner, B., Barnes, C. P., Graham, T. A. & Sottoriva, A. Reply: Is the
686 evolution of tumors Darwinian or non-Darwinian? *National Science Review* (2018).
687 doi:10.1093/nsr/nwx131
- 688 22. Heide, T. *et al.* Reply: Neutral tumor evolution? *bioRxiv* (2018). doi:10.1101/274142

- 689 23. Tokheim, C. J., Papadopoulos, N., Kinzler, K. W., Vogelstein, B. & Karchin, R.
690 Evaluating the evaluation of cancer driver genes. *Proc. Natl. Acad. Sci.* (2016).
691 doi:10.1073/pnas.1616440113
- 692 24. Castro-Giner, F., Ratcliffe, P. & Tomlinson, I. The mini-driver model of polygenic cancer
693 evolution. *Nature Reviews Cancer* (2015). doi:10.1038/nrc3999
- 694 25. Weghorn, D. & Sunyaev, S. Bayesian inference of negative and positive selection in
695 human cancers. *Nat. Genet.* (2017). doi:10.1038/ng.3987
- 696 26. Parmigiani, G. *et al.* Design and analysis issues in genome-wide somatic mutation studies
697 of cancer. *Genomics* (2009). doi:10.1016/j.ygeno.2008.07.005
- 698 27. Lawrence, M. S. *et al.* Mutational heterogeneity in cancer and the search for new cancer-
699 associated genes. *Nature* (2013). doi:10.1038/nature12213
- 700 28. Dees, N. D. *et al.* MuSiC: Identifying mutational significance in cancer genomes. *Genome*
701 *Res.* (2012). doi:10.1101/gr.134635.111
- 702 29. Vogelstein, B. *et al.* Cancer genome landscapes. *Science* (2013).
703 doi:10.1126/science.1235122
- 704 30. Gonzalez-Perez, A. & Lopez-Bigas, N. Functional impact bias reveals cancer drivers.
705 *Nucleic Acids Res.* (2012). doi:10.1093/nar/gks743
- 706 31. Tamborero, D., Gonzalez-Perez, A. & Lopez-Bigas, N. OncodriveCLUST: Exploiting the
707 positional clustering of somatic mutations to identify cancer genes. *Bioinformatics* (2013).
708 doi:10.1093/bioinformatics/btt395
- 709 32. Davoli, T. *et al.* XCumulative haploinsufficiency and triplosensitivity drive aneuploidy
710 patterns and shape the cancer genome. *Cell* (2013). doi:10.1016/j.cell.2013.10.011
- 711 33. Reiter, J. G. *et al.* Minimal functional driver gene heterogeneity among untreated

- 712 metastases. *Science* (80-.). (2018). doi:10.1126/science.aat7171
- 713 34. Tamborero, D. *et al.* Cancer Genome Interpreter annotates the biological and clinical
714 relevance of tumor alterations. *Genome Med.* (2018). doi:10.1186/s13073-018-0531-8
- 715 35. Tomasetti, C., Marchionni, L., Nowak, M. A., Parmigiani, G. & Vogelstein, B. Only three
716 driver gene mutations are required for the development of lung and colorectal cancers.
717 *Proc. Natl. Acad. Sci. U. S. A.* (2015). doi:10.1073/pnas.1421839112
- 718 36. Knudson, A. G. Two genetic hits (more or less) to cancer. *Nat. Rev. Cancer* (2001).
719 doi:10.1038/35101031
- 720 37. Diederichs, S. *et al.* The dark matter of the cancer genome: aberrations in regulatory
721 elements, untranslated regions, splice sites, non-coding RNA and synonymous mutations.
722 *EMBO Mol. Med.* (2016). doi:10.15252/emmm.201506055
- 723 38. Kumar, S. *et al.* Passenger mutations in 2500 cancer genomes: Overall molecular
724 functional impact and consequences. *bioRxiv* (2018). doi:10.1101/280446
- 725 39. Soussi, T., Taschner, P. E. M. & Samuels, Y. Synonymous Somatic Variants in Human
726 Cancer Are Not Infamous: A Plea for Full Disclosure in Databases and Publications.
727 *Human Mutation* (2017). doi:10.1002/humu.23163
- 728 40. Gotea, V., Gartner, J. J., Qutob, N., Elnitski, L. & Samuels, Y. The functional relevance of
729 somatic synonymous mutations in melanoma and other cancers. *Pigment Cell and*
730 *Melanoma Research* (2015). doi:10.1111/pcmr.12413
- 731 41. Supek, F., Miñana, B., Valcárcel, J., Gabaldón, T. & Lehner, B. Synonymous mutations
732 frequently act as driver mutations in human cancers. *Cell* (2014).
733 doi:10.1016/j.cell.2014.01.051
- 734 42. Chen, Y. T. *et al.* Tumor-associated intronic editing of HNRPLL generates a novel

- 735 splicing variant linked to cell proliferation. *J. Biol. Chem.* (2018).
736 doi:10.1074/jbc.RA117.001197
- 737 43. Jung, H. *et al.* Intron retention is a widespread mechanism of tumor-suppressor
738 inactivation. *Nat. Genet.* (2015). doi:10.1038/ng.3414
- 739 44. Hurst, L. D. & Batada, N. N. Depletion of somatic mutations in splicing-associated
740 sequences in cancer genomes. *Genome Biol.* (2017). doi:10.1186/s13059-017-1337-5
- 741 45. Gartner, J. J. *et al.* Whole-genome sequencing identifies a recurrent functional
742 synonymous mutation in melanoma. *Proc. Natl. Acad. Sci.* (2013).
743 doi:10.1073/pnas.1304227110
- 744 46. Sabarinathan, R. *et al.* The whole-genome panorama of cancer drivers. *bioRxiv* (2017).
745 doi:10.1101/190330
- 746 47. Kimura, M. The number of heterozygous nucleotide sites maintained in a finite population
747 due to steady flux of mutations. *Genetics* (1969). doi:kimura1969
- 748 48. Castellanos-Moreno, A., Castellanos-Jaramillo, A., Corella-Madueño, A., Gutiérrez-
749 López, S. & Rosas-Burgos, R. Stochastic model for computer simulation of the number of
750 cancer cells and lymphocytes in homogeneous sections of cancer tumors. *arXiv e-prints*
751 arXiv:1410.3768 (2014).
- 752 49. Turner, C., Stinchcombe, A. R., Kohandel, M., Singh, S. & Sivaloganathan, S.
753 Characterization of brain cancer stem cells: A mathematical approach. *Cell Prolif.* (2009).
754 doi:10.1111/j.1365-2184.2009.00619.x
- 755 50. Baar, M. *et al.* A stochastic model for immunotherapy of cancer. *Sci. Rep.* (2016).
756 doi:10.1038/srep24169
- 757 51. Figueredo, G. P., Siebers, P. O., Owen, M. R., Reps, J. & Aickelin, U. Comparing

- 758 stochastic differential equations and agent-based modelling and simulation for early-stage
759 cancer. *PLoS One* (2014). doi:10.1371/journal.pone.0095150
- 760 52. Gillespie, D. T. A general method for numerically simulating the stochastic time evolution
761 of coupled chemical reactions. *J. Comput. Phys.* (1976). doi:10.1016/0021-
762 9991(76)90041-3
- 763 53. Székely, T. & Burrage, K. Stochastic simulation in systems biology. *Computational and*
764 *Structural Biotechnology Journal* (2014). doi:10.1016/j.csbj.2014.10.003
- 765 54. Ryser, M. D., Lee, W. T., Ready, N. E., Leder, K. Z. & Foo, J. Quantifying the Dynamics
766 of Field Cancerization in Tobacco-Related Head and Neck Cancer: A Multiscale
767 Modeling Approach. *Cancer Res.* **76**, 7078–7088 (2016).
- 768 55. Yates, C. A., Ford, M. J. & Mort, R. L. A Multi-stage Representation of Cell Proliferation
769 as a Markov Process. *Bull. Math. Biol.* (2017). doi:10.1007/s11538-017-0356-4
- 770 56. Sabarinathan, R. *et al.* The whole-genome panorama of cancer drivers. *bioRxiv* (2017).
771 doi:10.1101/190330
- 772 57. Celiktas, M. *et al.* Role of CPS1 in cell growth, metabolism, and prognosis in LKB1-
773 inactivated lung adenocarcinoma. *J. Natl. Cancer Inst.* (2017). doi:10.1093/jnci/djw231
- 774 58. Kim, J. *et al.* CPS1 maintains pyrimidine pools and DNA synthesis in KRAS/LKB1-
775 mutant lung cancer cells. *Nature* (2017). doi:10.1038/nature22359
- 776 59. Lo, H.-W., Zhu, H., Cao, X., Aldrich, A. & Ali-Osman, F. A novel splice variant of GLI1
777 that promotes glioblastoma cell migration and invasion. *Cancer Res.* (2009).
778 doi:10.1158/0008-5472.CAN-09-0886
- 779 60. Mastrangelo, E. & Milani, M. Role and inhibition of GLI1 protein in cancer. *Lung Cancer*
780 *Targets Ther.* **9**, 35–43 (2018).

- 781 61. Martincorena, I., Raine, K. M., Davies, H., Stratton, M. R. & Campbell, P. J. Universal
782 Patterns of Selection in Cancer and Somatic Tissues. *Cell* (2017).
783 doi:10.1016/j.cell.2017.09.042
- 784 62. Ostrow, S. L., Barshir, R., DeGregori, J., Yeger-Lotem, E. & Hershberg, R. Cancer
785 Evolution Is Associated with Pervasive Positive Selection on Globally Expressed Genes.
786 *PLoS Genet.* (2014). doi:10.1371/journal.pgen.1004239
- 787 63. Noorbakhsh, J. & Chuang, J. H. Uncertainties in tumor allele frequencies limit power to
788 infer evolutionary pressures. *Nat. Genet.* **49**, 1288 (2017).
789
790
791
792