

Inferring parameters of cancer evolution from sequencing and clinical data

Nathan Lee¹ and Ivana Bozic^{1,2*}

¹Department of Applied Mathematics, University of Washington, Seattle, WA, USA

² Herbold Computational Biology Program, Fred Hutchinson Cancer Research Center, Seattle, WA, USA

*Correspondence to ibozic@uw.edu

Abstract

As a cancer develops, its cells accrue new mutations, resulting in a heterogeneous, complex genomic profile. We make use of this heterogeneity to derive simple, analytic estimates of parameters driving carcinogenesis and reconstruct the timeline of selective events following initiation of an individual cancer. Using stochastic computer simulations of cancer growth, we show that we can accurately estimate mutation rate, time before and after a driver event occurred, and growth rates of both initiated cancer cells and subsequently appearing subclones. We demonstrate that in order to obtain accurate estimates of mutation rate and timing of events, observed mutation counts should be corrected to account for clonal mutations that occurred after the founding of the tumor, as well as sequencing coverage. We apply our methodology to reconstruct the individual evolutionary histories of chronic lymphocytic leukemia patients, finding that the parental leukemic clone typically appears within the first fifteen years of life.

Introduction

When a cell accrues a sequence of driver mutations – genetic alterations that provide a proliferative advantage relative to surrounding cells – it can begin to divide uncontrollably and eventually develop the complex features of a cancer [1–3]. Thousands of specific driver mutations have been implicated in carcinogenesis, with individual tumors harboring from few to dozens of drivers, depending on the cancer type [4]. Mutations that don’t have a significant effect on cellular fitness also arise, both before and after tumor initiation [5]. These neutral mutations, or “passengers”, can reach detectable frequencies by random genetic drift or the positive selection of a driver mutation in the same cell [6–9]. Mutational burden detectable by bulk sequencing reveals tens to thousands of passengers per tumor [10, 11].

21 Genome sequencing technologies have revealed the heterogeneous, informative genetic profiles produced
22 by the evolutionary process driving carcinogenesis [12, 13]. These genetic profiles have been used to obtain
23 insight into specific features of the carcinogenic process operating in individual patients. For example, the
24 molecular clock feature of passenger mutations has been employed to measure timing of early events in tumor
25 formation, as well as identify stages of tumorigenesis and metastasis [14–22]. Other studies have estimated
26 mutation rates [5, 23, 24], selective growth advantages of cancer subclones [25–28], and the effect of spatial
27 structure on cancer evolution [29–31]. We note that previous approaches typically only estimate one or a
28 few parameters of cancer evolution. In addition, many state of the art methods make use of computationally
29 expensive approaches [24, 30, 32] or simplifying assumptions, such as approximating tumor expansion as
30 deterministic or ignoring cell death [27, 32].

31 Mathematical models of cancer progression, especially when used in conjunction with experimental and
32 clinical data, can provide important insights into the evolutionary history of cancer [9, 19, 33–37]. Branching
33 processes – a type of a stochastic process – can be used to model how different populations of dividing, dying,
34 and mutating cells in a tumor evolve over time [38]. Their theory and applications have been well developed
35 to model the multistage nature of cancer development [25, 29, 35, 38–40]. Here we use a branching process
36 model of carcinogenesis to derive a comprehensive reconstruction of an individual tumor’s evolution.

37 Tumors can grow for many years, even decades, before they reach detectable size [16]. Typically, tumor
38 samples used for sequencing would be obtained at the end of the tumor’s natural, untreated progression.
39 More recently, longitudinal sequencing, where a tumor is sequenced at multiple times during its development,
40 has provided better resolution of tumor growth dynamics and evolution in various cancer types [27, 41–
41 44]. We establish that two longitudinal bulk sequencing and tumor size measurements are sufficient to
42 reconstruct virtually all parameters (mutation rate, growth rates, times of appearance of driver mutations,
43 and time since the driver mutation) of cancer evolution in individual patients. Our analytic approach
44 yields simple formulas for the parameters; thus estimation of the parameters governing cancer growth is not
45 computationally intensive, regardless of tumor size. Our framework makes possible a personalized, high-
46 resolution reconstruction of a tumor’s timeline of selective events and quantitative characterization of the
47 evolutionary dynamics of the subclones making up the tumor.

48 Results

49 Model

50 We consider a multi-type branching process of tumor expansion (Fig. 1a). Tumor growth is started with
51 a single initiated cell at time 0. Initiated tumor cells divide with rate b and die with rate d . These cells
52 already have the driver mutations necessary for expansion, so we assume $b > d$. The population of initiated
53 cells can go extinct due to stochastic fluctuations, or survive stochastic drift and start growing (on average)
54 exponentially with net growth rate $r = b - d$. We will focus only on those populations that survived stochastic
55 drift.

56 At some time $t_1 > 0$ a new driver mutation occurs in a single initiated tumor cell, starting a new
57 independent birth-death process, with birth rate b_1 and death rate d_1 (Fig. 1b). Net growth rate of cells
58 with the new driver is $r_1 = b_1 - d_1$. The new driver increases the rate of growth, i.e., $r_1 > r$. We define
59 the driver's selective growth advantage by $g = (r_1/r - 1)$. In addition, both populations of cells (with and
60 without the driver) accrue passenger mutations with rate u (Fig. 1c).

61 After the driver mutation occurs, an additional time t passes before the tumor is observed. Cells con-
62 taining i new driver mutations, where i is either 0 or 1, will be referred to as type- i cells or simply, clone i .
63 In Materials and Methods we also analyze the more general case of two nested or sibling driver mutations,
64 as well as the fully generalized case of any clonal structure that might arise during tumor expansion.

65 Parameter estimates from two longitudinal measurements

66 We demonstrate that with two longitudinal bulk sequencing measurements, it is possible to accurately
67 estimate net growth rates, time of appearance of a driver mutation, time between a driver mutation and
68 observation, and mutation rate in the tumor. The tumor is first sequenced at time of observation, $t_1 + t$,
69 where both time of driver mutation, t_1 , and time from driver mutation to observation, t , are yet unknown
70 (Fig. 1b). A second bulk sequencing is performed at $t_1 + t + \delta$, a known δ time units after the tumor is first
71 observed (Fig. 1b). From the bulk sequencing data, the fraction of cells carrying the driver mutation, α_1
72 and α_2 , can be measured at the timepoints $t_1 + t$ and $t_1 + t + \delta$, respectively. We denote total number of
73 cells in the tumor at the two bulk sequencing timepoints as M_1 and M_2 . Number of cells in the tumor can
74 be estimated from measurements of tumor volume [45].

75 Equating expected values of the sizes of type-0 and type-1 population at the two bulk sequencing time
76 points with the measured numbers of cells present in clones 0 and 1, we obtain estimates of the net growth

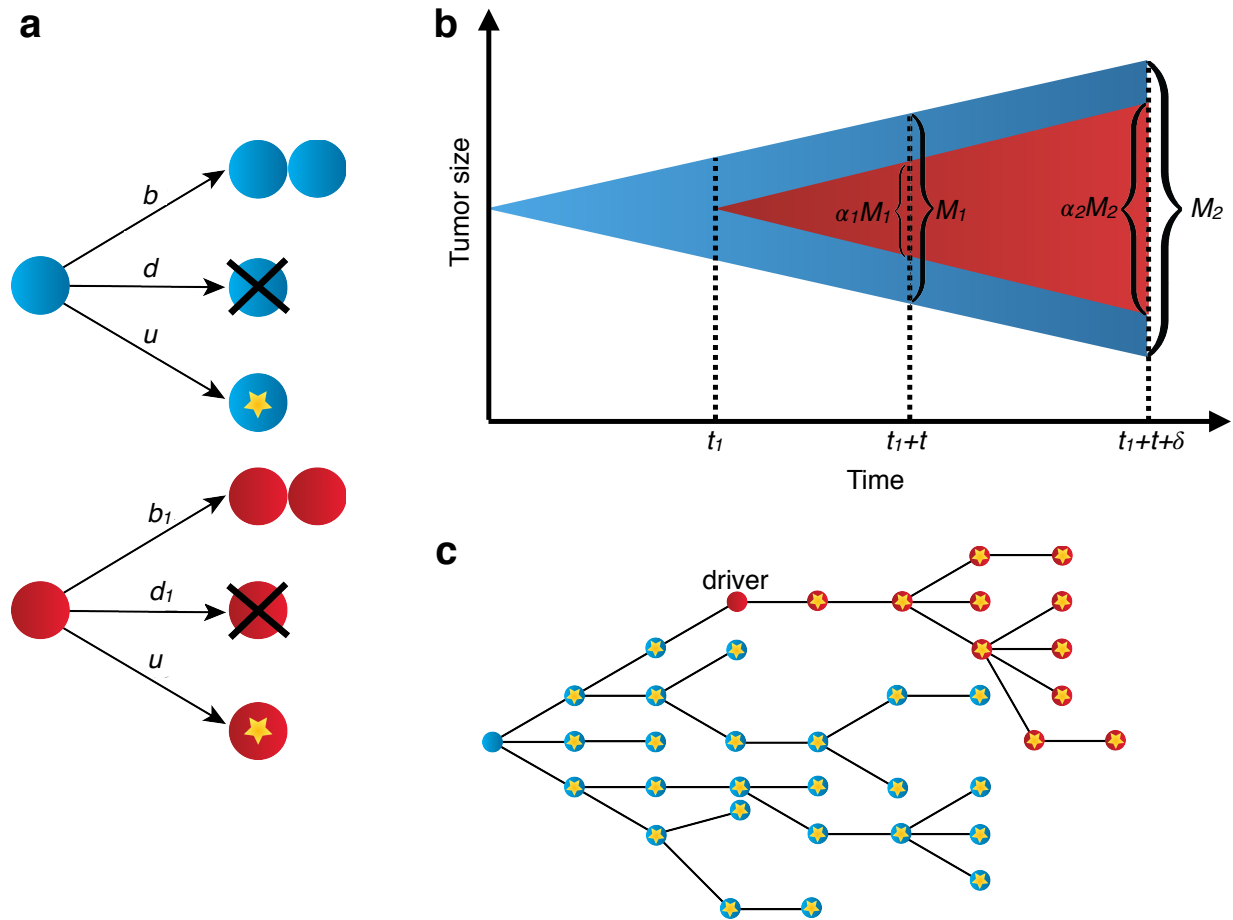


Figure 1: **Stochastic branching process model of tumor evolution.** (a) Stochastic branching process model for tumor expansion. Initiated tumor cells (blue) divide with birth rate b , die with death rate d , and accrue passenger mutations with mutation rate u . Type-1 cells, which carry the driver mutation, divide with birth rate b_1 , die with death rate d_1 , and accrue passenger mutations with mutation rate u . (b) The initiated tumor, or type-0, (blue) population growth is initiated from a single cell. A driver mutation occurs in a single type-0 cell at time t_1 , starting the type-1 population (red). The tumor is bulk sequenced at times $t_1 + t$ and $t_1 + t + \delta$. (c) By the time the tumor is observed, it has a high level of genetic heterogeneity due to the mutations that have accrued in both type-0 (blue) and type-1 populations (red). Each yellow star represents a different passenger mutation.

77 rates of the two subclones:

$$78 \quad r = \frac{1}{\delta} \log \left(\frac{(1 - \alpha_2)M_2}{(1 - \alpha_1)M_1} \right) \quad (1)$$

$$79 \quad r_1 = \frac{1}{\delta} \log \left(\frac{\alpha_2 M_2}{\alpha_1 M_1} \right) \quad (2)$$

80

81 From the growth rate estimates and subclone sizes, we can approximate the expected value of the time a
82 population in a branching process takes to reach an observed size [38]. This yields an estimate of the time t
83 from the appearance of driver mutation until observation:

$$84 \quad t = \frac{1}{r_1} \log(M_1 \alpha_1) \quad (3)$$

85

86 Using the bulk sequencing data from the second timepoint, γ , the number of subclonal passengers between
87 the specified frequencies f_1 and f_2 , can be measured. Using results from previous work [46], we derive the
88 expected value of γ (Materials and Methods), which can be used to estimate the mutation rate u :

$$89 \quad u = \frac{f_1 f_2 r r_1 \gamma}{(f_2 - f_1)(\alpha_2 r + r_1(1 - \alpha_2))} \quad (4)$$

90

91 The m passenger mutations that were present in the original type-1 cell when the driver mutation occurred
92 (Fig. 1c) are present in all type-1 cells. m can be estimated from bulk sequencing data, and used to estimate
93 time of appearance of the driver. We maximize the likelihood function $P(m|t_1)$ with respect to time of
94 appearance of the driver, t_1 , (see Materials and Methods) to obtain the maximum likelihood estimate

$$95 \quad t_1 = \frac{m}{u} \quad (5)$$

96

97 Using formulas (4) and (5), we can now estimate t_1 .

98 **Estimates verified in simulated tumors**

99 To assess the accuracy of the parameter estimates for several modes of tumor evolution, we simulate tumor
100 growth by performing a Monte Carlo simulation, which simulates the birth, death, and accumulation of
101 mutations in the individual cells that make up a tumor. This simulation generates the mutation frequency
102 and tumor size data used by the estimates (see Methods section for details of simulation). We simulate three
103 different types of tumors (slow growing, fast growing, and no cell death), with a high and a low mutation
104 rate for each.

105 In a simulation of a fast growing tumor with a single subclonal driver mutation that confers a strong
 106 selective growth advantage of 100%, we can accurately estimate growth rates, mutation rate, time of driver
 107 event, and time since driver event (Fig. 2). Growth rates of both initiated tumor and driver subclones
 108 can be estimated with a high degree of accuracy, achieving mean percentage error (MPE) of -0.07% and
 109 0.03% for the lower mutation rate ($u = 1$) scenario. The mutation rate u and estimates for time of driver
 110 appearance, t_1 , and time since driver, t , can also be estimated accurately, with MPEs of -0.9% , 3.8% , and
 111 -0.4% , respectively. Estimates for u , t_1 , and t have a somewhat greater degree of variation compared to the
 112 growth rate estimates, due to the inherent randomness of the number of mutations and time to reach the
 113 observed size that occur in each realization of the stochastic process.

114 For the parameter regime with no cell death and the regime for a slow-growing tumor, we again achieve
 115 high accuracies for the net growth rates (Fig. S1, Fig. S2). In the lower mutation rate ($u = 1$) scenario,
 116 parameter estimates for the mutation rate u and time of driver appearance t_1 can be accurately estimated
 117 for both regimes, with MPEs of -1.3% and 4.9% for the no cell death case, and MPEs of -3% and 3.7% for
 118 the slow-growing tumor. The t , time since driver event, estimates have somewhat higher errors, with MPE
 119 of -6.3% for the no cell death case, and MPE of 30.3% for the slow-growing tumor.

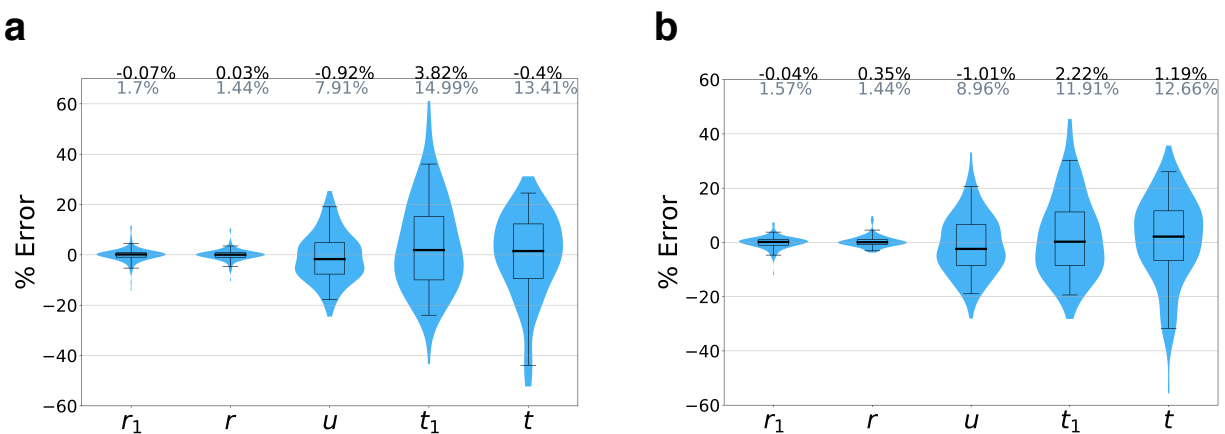


Figure 2: **Accuracy of parameter inferences from simulated data.** We simulated tumor growth by performing a Monte Carlo simulation, which simulates the birth, death, and accumulation of mutations in the individual cells that make up a tumor, and generates the mutation frequency and tumor size data used by the estimates. Mean percent errors (MPEs) of estimates are shown in black above the plots, and mean absolute percent errors (MAPEs) are shown in gray. Boxes contain 25th-75th quartiles, with median indicated by thick horizontal black line. Whiskers of boxplots indicate 2.5 and 97.5 percentiles. Violins are smoothed density estimates of the percent error data points. Ground truth parameter set: $b = b_1 = 0.25$, $d = 0.18$, $d_1 = 0.11$, $t_1 = 70$, $t = 50$, $\delta = 20$, $f_1 = 1\%$, and $f_2 = 20\%$. Mutation rate (a) $u = 1$, (b) $u = 3$. At least 100 Monte Carlo simulation runs with a surviving tumor performed for each parameter combination.

120 Correcting mutation counts observed from genome sequencing data

121 We note that in our estimate for the time of appearance of the driver, t_1 (see formula (5)), used for comparison
122 to simulated data, we employed a correction to m , the number of mutations that were present in the founder
123 type-1 cell at t_1 . From sequencing data, these m mutations are indistinguishable (Fig. 3a) from mutations
124 that occurred after t_1 in type-1 cells, and reached fixation in the type-1 population [46]. Thus, the value of
125 m observed from sequencing data, m_{obs} , will overestimate the true m . In Materials and Methods we show
126 that the expected value of the number of passengers that occurred after t_1 and reached fixation in the type-1
127 population is u/r_1 . We subtract this correction factor from m_{obs} :

$$128 \qquad m = m_{obs} - u/r_1 \qquad (6)$$

129

130 The correction for the m mutations present in the original type-1 cell (6) at time t_1 improves the accuracy
131 of the estimate for time of appearance of driver mutation t_1 . For the fast growing tumor with mutation rate
132 $u = 1$ (Fig. S3a), the correction lowers the mean percent error (MPE) of the t_1 estimate from 14.0% to
133 3.8%. For the slow growing tumor with mutation rate $u = 5$ (Fig. 3b), the correction lowers the MPE of
134 the t_1 estimate from 22.0% to 5.7% (Fig. 3b).

135 Another issue arises from obtaining mutation count γ , number of mutations with frequency between f_1
136 and f_2 , from genome sequencing data. When sequencing data is post-processed by filtering out mutations
137 with L or fewer variant reads, low-frequency mutations will be difficult to detect [35] (Fig. 3c). For a sample
138 with average sequencing coverage of R and tumor purity p , mutations with mutant allele frequency below
139 $L/(pR)$ will typically not be observable. As a result, since mutations with frequencies between f_1 and f_2
140 count towards γ , if $f_1 \leq 2L/(pR)$, the observed number of subclonal mutations between frequencies f_1 and
141 f_2 , γ_{obs} , will underestimate the true value, γ . In the Materials and Methods, we derive a correction for γ ,
142 based on the expected value of the number of subclonal mutations present at cancer cell frequencies (CCFs)
143 between f_1 and $2L/(pR)$:

$$144 \qquad \gamma = \gamma_{obs} \left(\frac{\frac{1}{f_1} - \frac{1}{f_2}}{\frac{pR}{2L} - \frac{1}{f_2}} \right) \qquad (7)$$

145

146 Before applying our methodology to patient sequencing data, we estimated the validity of the above cor-
147 rection applied to observed simulated mutation counts. When we simulate sequencing reads from simulated
148 mutation frequencies (see Materials and Methods) and post-process by removing mutations with $L = 2$ or
149 fewer variant reads, the adjustment we derived for mutation count γ (7) is critical, even for average sequenc-
150 ing coverage of 200x (Fig. 3d). Without any correction, the observed γ has MPE of -53.3% compared to

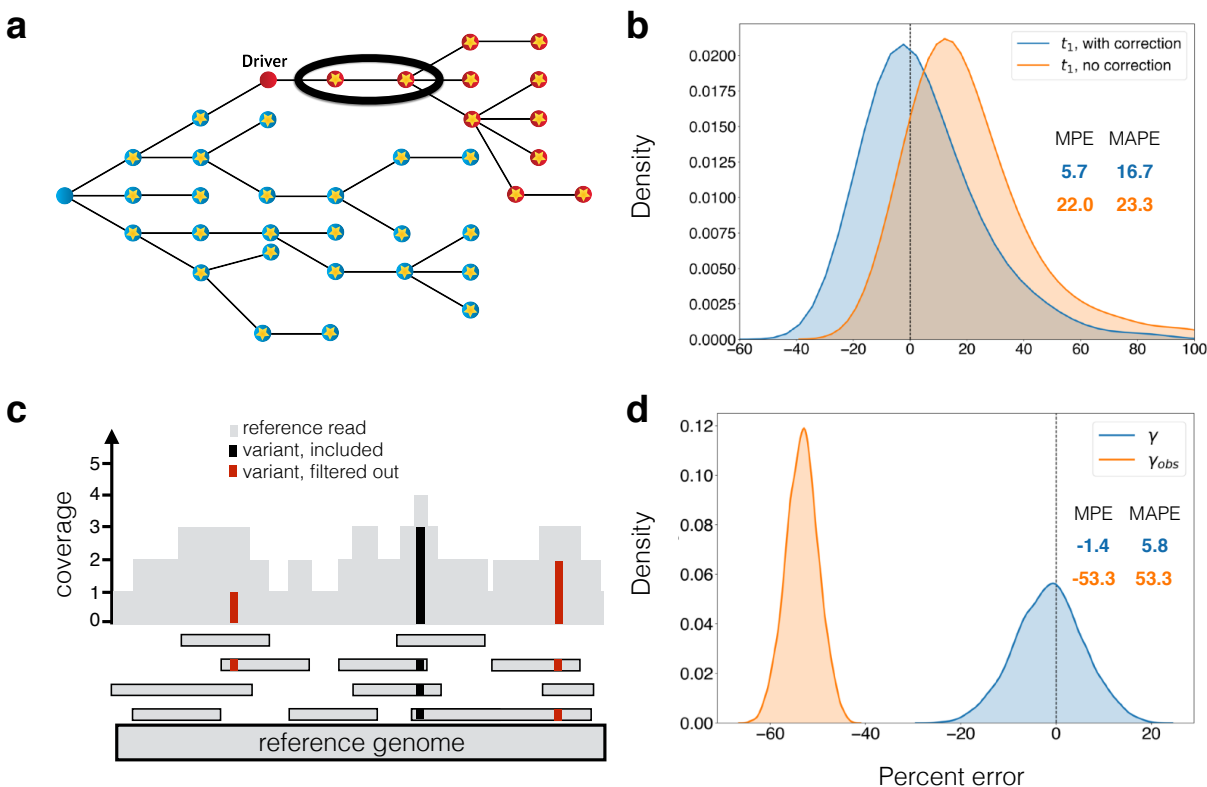


Figure 3: Corrections for observed mutation counts. (a) If passenger mutations (circles with stars) that occur after the driver reach fixation in the driver population (red), then they are indistinguishable from the passengers that were present in the first cell with the driver, which accrued in the type-0 population (blue). The estimate of when the driver occurred needs to account for these mutations (circled). In (b), we compare percent errors of parameter estimates for time from tumor initiation until appearance of a driver subclone, t_1 , with and without this correction (Eq. (6)). Errors for estimate with correction (Eq. (12)) are shown in blue, and for estimate without correction (Eq. (5)) in orange. Errors are plotted as a kernel density estimate for Monte Carlo simulations of slow growing tumor with mutation rate $u = 5$. Mean percent errors (MPEs) and mean absolute percent errors (MAPEs) are listed. (c) Mutations present on two or fewer variant reads (red) are filtered out in post-processing. Mutations with more than two variant reads (black) are included. The number of subclonal mutations between frequencies f_1 and f_2 , γ , which is used in the mutation rate estimate, must be corrected for mutations that are filtered out. In (d), the percent errors for the observed (orange) and corrected (blue) γ (Eq. (7)) are plotted as kernel density estimates. Observed mutations are those that passed post-processing, i.e. those that have more than $L = 2$ mutant reads. True mutation frequencies were generated from 135 surviving runs of a Monte Carlo simulation of a fast growing tumor with mutation rate $u = 1$, from which sequencing reads were simulated with 200x average coverage (see Materials and Methods). Percent errors are calculated relative to the true γ measured from the true mutation frequencies.

151 true γ , but with the correction, the computed γ has MPE of -1.4%. When average coverage is 100x, this
152 correction becomes even more important, as many of the low-frequency mutations are discarded (Fig. S3b).
153 Without any correction, the observed γ has MPE of -79.7%. With the correction the computed γ has MPE
154 of -3.4%. The accuracy of the γ measurement affects our estimate of the mutation rate (4).

155 **Estimating parameters for individual patients with CLL**

156 We use our formulas to infer the patient-specific parameters of cancer evolution for four patients with
157 chronic lymphocytic leukemia (CLL) whose growth patterns and clonal dynamics were analyzed in [27].
158 These CLLs had peripheral white blood cell (WBC) counts measured and whole exome sequencing (WES)
159 performed at least twice before treatment. We consider patients whose WBC counts were classified as
160 having an exponential-like growth pattern, with average $\gamma_{obs} > 2$ and 3 or fewer macroscopic subclones (i.e.
161 subclones with cancer cell fractions of 20% or greater for at least one pre-treatment time point). As in Ref.
162 [27], we perform subclonal reconstruction for each patient using PhylogicNDT [43]. To obtain confidence
163 intervals for our parameter estimates, we utilize a sampling procedure to account for model and measurement
164 uncertainties, including uncertainties in subclone frequencies, fitted growth curves, and the Poisson process
165 for mutation accumulation (see Materials and Methods). For each patient's tumor, we compute estimates of
166 the growth rate of each clone, exome mutation rate, the times that each subclone arose, and how long each
167 subclone expanded before the tumor was detected (Table S1). We reconstruct these histories for tumors
168 with various clonal structures.

169 Patients 3 and 21 are examples of a CLL with a single subclone (Fig. 4). For Patient 3, Clone 0, the
170 most recent common ancestor (MRCA) of this patient's CLL, was initiated when the patient was 14.6 [1.4,
171 26.8] years old (median and [95% confidence interval] of estimate). Clone 0 grew with a net growth rate of
172 0.51 [0.20, 0.85] per year. 18.9 years later, Clone 1 was initiated when the patient was 33.5 [24.1, 39.2] years
173 old. Clone 1 expanded with a growth rate of 0.85 [0.65, 1.04] per year (corresponding to a selective growth
174 advantage of 68.7% over Clone 0), and the patient was diagnosed 29.5 [23.8, 38.9] years later at age 63. We
175 find that the CLL exome mutation rate was 0.48 [0.39, 0.59] mutations per year in this patient.

176 For patient 21, we estimate that the parental clone (MRCA, Clone 0) of this patient's CLL was initiated
177 when the patient was 6.4 [0.3, 16.7] years old, and grew with a net growth rate of 0.79 [0.30, 1.14] per year.
178 Clone 1 appeared when the patient was 19.6 [10.8, 24.0] years old, and grew more quickly than Clone 0, with
179 a growth rate of 1.52 [1.01, 2.04] per year (corresponding to selective growth advantage of 91.4% over Clone
180 0). Clone 1 contained a FGFR1 mutation, which might have been acting as a driver of the increased net
181 proliferation. Clone 1 then grew for 15.4 [11.0, 24.2] years before the patient was diagnosed at age 35. We

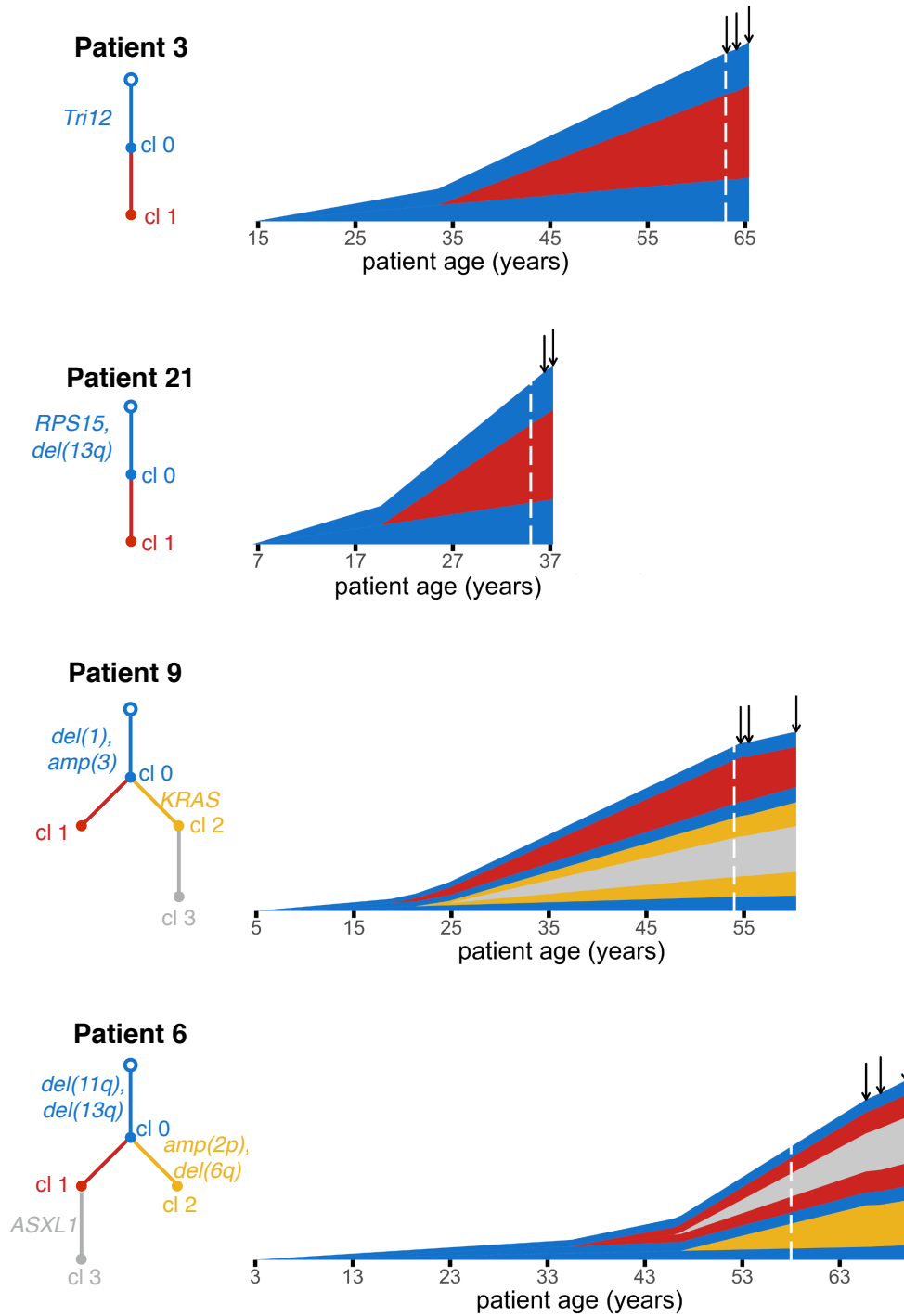


Figure 4: **Reconstructing the timeline of CLL evolution in patients.** We applied our methodology to estimate subclonal growth rates, mutation rates and evolutionary timelines in CLL tumors from Ref. [27]. Vertical height of a clone represents its \log_{10} -scaled size. Phylogenetic trees, colored by clone number, show annotated driver mutations along the trees' edges. For each patient, we show estimates for patient age at CLL initiation and times of appearance of CLL subclones. Dashed white line indicates when the patient was diagnosed. Solid black arrows indicate times of bulk sequencing measurements.

182 estimate that this patient's CLL had an exome mutation rate of 0.20 [0.19, 0.23] mutations per year.

183 Patients 6 and 9 present more complex clonal structures. CLL of Patient 9 contains two sibling subclones,
184 Clones 1 and 2, in addition to the parental population, Clone 0. Clone 2 contains a nested subclone (Clone
185 3). Clone 0 arose when the patient was 4.9 [1.2, 10.8] years old, and had a growth rate of 0.28 [0.17, 0.42]
186 per year. Clone 1 arose when the patient was 18.8 [8.8, 35.1] years old. At the time of sequencing, Clone
187 1 had a negative growth rate of -0.40 [-0.45, -0.19] (/year). Clone 2, containing a KRAS mutation, had the
188 largest net growth rate of the three clones (0.67 [0.49, 0.94] per year), corresponding to a selective growth
189 advantage of 140.9% over the parental clone. Clone 2 arose when the patient was 21.3 [7.7, 31.7] years old.
190 Clone 3 was initiated from within Clone 2 when the patient was 24.8 [10.3, 37.6] years old. We estimate that
191 the CLL exome mutation rate of Patient 9 is 0.36 [0.35, 0.37] mutations/year.

192 CLL of Patient 6 contains two sibling subclones (Clones 1 and 2) descendant from the leukemic MRCA
193 Clone 0. Clone 1 has a nested subclone (Clone 3). We estimate that the CLL was initiated when the patient
194 was 2.8 [0.1, 13.2] years old. Clone 0 then grew at a rate of 0.68 [0.15, 1.30] per year. Approximately 33
195 years after the appearance of Clone 0, when the patient was 35.4 [21.7, 46.1] years old, the first subclone,
196 Clone 1 appeared. Clone 1 had a net growth rate of 0.41 [0.08, 0.73] per year. Clone 3 arose from within
197 Clone 1 when the patient was 45.9 [31.3, 54.6] years old. This clone had net growth rate 1.09 [0.65, 1.78] per
198 year. Clone 3 harbored a driver mutation in ASXL1 and had selective growth advantage of 60.8% over Clone
199 0. Clone 2, nested in parental clone (Clone 0), was initiated when the patient was 46.7 [25.6, 57.5] years
200 old and had growth rate 0.46 [0.08, 0.85] per year. The patient was then diagnosed at age 58, eventually
201 needing treatment 12.0 years after diagnosis. In Patient 6, we estimate a CLL exome mutation rate of 0.15
202 [0.12, 0.19] mutations per year.

203 The average mutation rate in the four CLL patients we analyze is 0.30 mutations/year. This rate is over
204 the exome, which accounts for $\sim 1\%$ of the human genome. Our average estimated mutation rate in CLL
205 exomes is similar to the measured rate of accumulation of mutations in human tissues of 40 mutations per
206 year over the entire genome [47]. Other recent work has estimated a mutation rate of 17 mutations per
207 year in human haematopoietic stem cell/multipotent progenitors [48]. Our estimated mutation rates during
208 CLL progression are on par or higher than the recent estimates in healthy hematopoietic cells [48], in line
209 with the expectation that mutation rates may be increased in cancer. The estimated times of appearance of
210 CLL subclones are very long, on the order of 10 years or more. This finding is agreement with results from
211 Gruber et al. [27], who find few new CLL subclones over years to a decade of evolution. We observe that
212 CLL initiation occurred early in most patients, within the first fifteen years of their lives, consistent with
213 recent work in other cancer types [19, 36].

214 Discussion

215 We use a stochastic branching process model to reconstruct the timing of driver events and quantify the
216 evolutionary dynamics of different subclonal populations of cancer cells. We can accurately estimate the
217 growth rates of tumor subclones, selective growth advantage of individual driver mutations, mutation rate in
218 the tumor, time between tumor initiation and appearance of a subclonal driver mutation, and time between
219 driver mutation and tumor observation. Together, this allows us to estimate the age of the patient at tumor
220 initiation, as well as the age at appearance of a subclonal driver.

221 Previous work has computed relative order of driver events [18, 21, 49], while other studies have given
222 estimates for scaled mutation rates and time of events [24, 32]. However, we present estimates for absolute,
223 unscaled mutation rates and times, which are easily interpretable and don't implicitly depend on unknown
224 parameters. We assume that mutations accrue with time, and not only at cell divisions, which simplifies
225 derivations and is supported by recent experimental data [47].

226 For individual CLLs that underwent bulk sequencing at two time points [27], we infer growth rates of
227 individual subclones, mutation rate in the tumor, the times when cancer subclones began growing, and the
228 time between driver mutations and the patient's diagnosis. Our inferences are limited by the relatively
229 low number of mutations present in CLL, as well as sequencing coverage [27]. The accuracy of estimates
230 presented here is expected to be even higher in cancer types with more mutations, with whole genome
231 sequencing available, or with higher sequencing coverage. Our methodology is in principle applicable to any
232 cancer type, not only CLL or liquid cancers. We note, however, that in the case of solid tumors, multiple
233 biopsies would potentially be needed to fully account for the existing heterogeneity.

234 Our model and derivations assume a fixed mutation rate u after transformation and fixed growth rates of
235 cancer subclones, similar to previous approaches [24, 30, 35]. Using an exponential model of cancer growth
236 with constant mutation and growth rate to estimate parameters of cancer evolution has its weaknesses: some
237 cancer subclones (such as Clone 1 from Pt. 9) not only do not grow exponentially, they actually decline in
238 absolute cell numbers. Sudden genomic instability events, or a change in cancer mutation and/or growth
239 rate over time could also introduce errors into our parameter inferences. Recent sequencing data points to
240 mutational processes that change over time during cancer evolution [20, 50]; incorporating possible changes
241 in the mutation and/or growth rate into the model would require much higher density of sequencing and
242 clinical data [37].

243 **Materials and Methods**

244 **Branching process model of tumor evolution**

245 We employ a continuous, multi-type branching process model of cancer evolution. Tumor expansion is
246 initiated by a single type-0, or initiated tumor cell. Type-0 cells divide with rate b and die with rate d ,
247 yielding a net growth rate of $r = b - d$. At time t_1 , a single driver mutation is introduced into a randomly
248 selected cell in the type-0 population, founding a new type-1 population of cells. This type-1 population
249 undergoes its own independent branching process. They divide with rate b_1 , die with rate d_1 , and have
250 net growth rate $r_1 = b_1 - d_1$. If the driver mutation gives type-1 cells a selective growth advantage over
251 the type-0 population, then $r_1 > r$. With the ratios of the growth rates denoted as $s = r_1/r$, the growth
252 advantage can be quantified as $g = (s - 1) \cdot 100\%$. In the case of neutral evolution, $g = 0$. If there is a selective
253 advantage, $g > 0$. Neutral mutations, or passengers, have no effect on the cell's fitness, and accrue according
254 to a Poisson process with rate u . We assume an infinite alleles model such that there is no back mutation
255 and an infinite sites model such that every new passenger mutation is unique. Only surviving populations
256 are considered. All derivations below will condition on survival. The type-0 and type-1 populations at time
257 t will be denoted as $X_0(t)$ and $X_1(t)$, respectively.

258 **Measurements sufficient to determine evolutionary history**

259 We derive estimates for parameters describing the carcinogenic process using measurements taken from two
260 timepoints late in the tumor's development. We require sequencing of the tumor at the two timepoints,
261 when the tumor is first observed at the unknown time $t_1 + t$ and a specified δ later, at $t_1 + t + \delta$. From these
262 two bulk sequencing measurements, we obtain measurements of α_1 and α_2 , the fraction of cells carrying the
263 driver mutation at $t_1 + t$ and $t_1 + t + \delta$, respectively. In addition, from the bulk sequencing at $t_1 + t + \delta$,
264 we obtain measurements of m , the number of mutations present in the founder type-1 cell, as well as γ , the
265 number of mutations with frequency between the specified f_1 and f_2 . The total population size at these
266 times, M_1 and M_2 , is also measured.

267 **Expected value of γ , number subclonal mutations**

268 For a population consisting of a single clone with birth and death rates b and d , the expected number of
269 subclonal mutations present at a frequency larger than f is shown to be [46]

$$270 \frac{\bar{u}(1-f)}{(1-\delta)f} \tag{8}$$

271

272 where $\delta = d/b$ and \bar{u} is the probability that a daughter cell gains a new passenger mutation at cell division.
 273 In this paper, we allow mutations to occur at any point in time and consider the absolute mutation rate per
 274 cell, u , which is equal to $\bar{u}b$. Then the expected number of subclonal mutations between f_1 and f_2 , $\mathbb{E}\gamma$, is

$$275 \quad \mathbb{E}\gamma = \frac{u(1-f_1)}{b(1-\delta)f_1} - \frac{u(1-f_2)}{b(1-\delta)f_2} \quad (9)$$

$$276 \quad = \frac{u}{r}(1/f_1 - 1/f_2) \quad (10)$$

278 where $r = b - d > 0$.

279 Now we derive $\mathbb{E}\gamma$ in the case of clones 0 through k , each clone with growth rate $r_i > 0$ and fraction α_i^c .
 280 Each clone i has $\alpha_i^c \frac{u}{r_i}(1/f_1 - 1/f_2)$ expected subclonal passengers between frequencies f_1 and f_2 . Thus, the
 281 total expected number of passengers with frequencies between f_1 and f_2 is

$$282 \quad \mathbb{E}\gamma = (1/f_1 - 1/f_2) \sum_{i=0}^k \frac{u\alpha_i^c}{r_i} \quad (11)$$

284 For the simplest case we consider, a tumor with a single driver mutation occurring in the initiated tumor
 285 population, there is a type-0 population with growth rate r and a type-1 population with growth rate r_1 .
 286 Equation (11) reduces to

$$287 \quad \mathbb{E}\gamma = \left(\frac{u\alpha}{r_1} + \frac{u(1-\alpha)}{r} \right) \left(\frac{1}{f_1} - \frac{1}{f_2} \right) \quad (12)$$

289 where α is the fraction of cells having the driver mutation.

290 Derivation of estimates of evolutionary parameters

291 With the two bulk sequencings at $t_1 + t$ and $t_1 + t + \delta$, we are able to derive estimates for t_1 , t , r , r_1 , and u .
 292 First we solve for r and r_1 , based on the estimated cell counts at $t_1 + t$ and $t_1 + t + \delta$. The observed type- i
 293 cell count is equated to the expected value of the type- i population size, conditioned on survival. For the
 294 type-0 population,

$$295 \quad \mathbb{E}[X_0(t_1 + t) | X_0(t_1 + t) > 0] = \frac{b}{r} e^{r(t_1+t)} = (1 - \alpha_1)M_1 \quad (13)$$

$$296 \quad \mathbb{E}[X_0(t_1 + t + \delta) | X_0(t_1 + t + \delta) > 0] = \frac{b}{r} e^{r(t_1+t+\delta)} = (1 - \alpha_2)M_2 \quad (14)$$

298 Proceeding similarly for the type-1 population, we obtain

$$299 \quad r_1 = \frac{1}{\delta} \log \left(\frac{\alpha_2 M_2}{\alpha_1 M_1} \right) \quad (15)$$

$$300 \quad r = \frac{1}{\delta} \log \left(\frac{(1 - \alpha_2) M_2}{(1 - \alpha_1) M_1} \right) \quad (16)$$

301

302 The expected value of the first time a population of type-1 cells in a branching process reaches the observed
303 size $\alpha_1 M_1$ is [38]

$$304 \quad \mathbb{E}t = \frac{1}{r_1} \log \left(\frac{\alpha_1 M_1 r_1}{b_1} \right) - \frac{1}{r_1} \int_0^\infty e^{-z} \log z dz \quad (17)$$

$$305 \quad = \frac{1}{r_1} \log \left(\frac{\alpha_1 M_1 r_1}{b_1} \right) + \frac{0.5772}{r_1} \quad (18)$$

306

307 which we approximate as

$$308 \quad \mathbb{E}t \approx \frac{1}{r_1} \log(\alpha_1 M_1) \quad (19)$$

309

310 We make use of two approximations to arrive at (19). First, we neglect the second term in (18), which serves
311 as a small correction to the first term. This term will be dominated by the first term as it increases with
312 logarithm of the cancer size. For $r_1 = 0.5$, $\alpha_1 M_1 \sim 10^{11}$, and $r_1 \approx b_1$, the second term (1.2) will be only
313 2.3% of the first term (50.7). For any growth rate, the second term will be 2.3% of the first term. Second,
314 we assume r_1 is similar in magnitude to b_1 .

315 With the measurement of γ , the number of subclonal passengers with frequency between f_1 and f_2 , we
316 can estimate the mutation rate u . In the previous section we derive the expected value of γ as

$$317 \quad \mathbb{E}\gamma = \left(\frac{u\alpha}{r_1} + \frac{u(1-\alpha)}{r} \right) \left(\frac{1}{f_1} - \frac{1}{f_2} \right) \quad (20)$$

318

319 Using the estimates of r and r_1 from (15) and (16), and the measured value of γ from the second bulk
320 sequencing, equation (20) can be solved for the mutation rate u ,

$$321 \quad u = \frac{f_1 f_2 r r_1 \gamma}{(f_2 - f_1)(\alpha_2 r + r_1(1 - \alpha_2))} \quad (21)$$

322

323 When estimating mutation rate for the CLL patients from Ref. [27], for which there is bulk sequencing at
324 two or more timepoints, we average the mutation rate calculated at each of these timepoints. (21) is applied
325 for each timepoint with the respective CCFs and observed γ values for each timepoint.

326 To derive the maximum likelihood estimates of t_1 , we consider the likelihood function $P(m|t_1)$. The
327 number of passenger mutations present in the founder type-1 cell that appeared at time t_1 is a Poisson
328 process with rate u . Thus,

$$329 \quad P(m|t_1) \propto \frac{(ut_1)^m e^{-ut_1}}{m!} \quad (22)$$

330

331 Maximizing the logarithm of the likelihood function with respect to t_1 yields a MLE for t_1 in terms of
332 estimated or measured quantities:

$$333 \quad t_1 = m/u \quad (23)$$

334

335 **Estimating number of unobserved subclonal mutations from sequencing data**

336 When sequencing data is post-processed by filtering out any mutations with L or fewer variant reads, the
337 number of mutations between f_1 and f_2 will likely be underestimated if $2L/(Rp) > f_1$, where R is average
338 sequencing coverage and p is tumor purity. Define γ_{obs} as the observed number of mutations between
339 frequencies f_1 and f_2 , after post-processing has been performed that filtered out any mutations with L or
340 fewer variant reads. The expected number of subclonal mutations between frequencies f_1 and x is given by

$$341 \quad \gamma(x) = c(1/f_1 - 1/x) \quad (24)$$

342

343 where c is a constant that will vary depending on the patient and sample. It can be fit on the sequencing
344 data by noting

$$345 \quad \gamma_{obs} = \gamma(f_2) - \gamma(2L/(Rp)) \quad (25)$$

$$346 \quad = c(Rp/(2L) - 1/f_2) \quad (26)$$

347

348 Therefore, c can be estimated from the sequencing data as

$$349 \quad c = \frac{\gamma_{obs}}{Rp/(2L) - 1/f_2} \quad (27)$$

350

351 Then, we can estimate γ as

$$352 \quad \gamma = \gamma_{obs} \left(\frac{\frac{1}{f_1} - \frac{1}{f_2}}{\frac{Rp}{2L} - \frac{1}{f_2}} \right) \quad (28)$$

353

354 **Number of passengers reaching fixation after t_1**

355 We estimate the number of passengers that occurred after t_1 and reached fixation in the type-1 population
356 in order to adjust the m_{obs} mutation count. From [46], when mutations occur at cell division, the expected
357 number of clonal passengers is $\delta\bar{u}/(1-\delta)$. \bar{u} is the probability that a daughter cell gains a new passenger
358 mutation at cell division, so the mutation rate is $u = \bar{u}b_1$. For the type-1 population, $\delta = d_1/b_1 < 1$. When
359 mutations accrue over time, and not only at divisions, the expected number of clonal passengers is thus

$$360 \quad \bar{u}/(1-\delta) = u/r_1 \quad (29)$$

361

362 Similarly, for a clone i , the expected number of passengers that occur after time t_i and reach fixation is

$$363 \quad u/r_i \quad (30)$$

364

365 where $r_i = b_i - d_i > 0$.

366 **Simulation of tumor evolution and sequencing data**

367 To assess the accuracy of the analytic results, we perform a continuous time Monte Carlo simulation to
368 model tumor evolution and collection of sequencing data with an implementation of the Gillespie algorithm
369 [51]. Simulations are written in C/C++.

370 The type- j population has division rate b_j , death rate d_j , and mutation rate u . Mutations can occur at
371 any point of the cell cycle, not just during division. z_n is the number of type- j cells with passenger n as
372 their most recent passenger mutation. The type-0 population is initiated with a single cell at time 0, and
373 the type-1 population is initiated with a single cell at time t_1 . Let a be the vector recording the ancestor of
374 new mutations. Element a_i is the subclonal ancestor of the i th passenger mutation. Repeat 1-4 while time
375 is less than $t_1 + t + \delta$.

- 376 1) Set $\Gamma = N_j(b_j + d_j + u)$. Time increment to next event time is randomly sampled from $\text{Exp}[\Gamma]$.
 - 377 • For type-0, if time is greater than or equal to t_1 for first time, randomly select type-0 subclone i
378 to have driver mutation, remove one cell from type-0 population count, and set $N_1 = 1$. Record
379 the true value of m , the number of passenger mutations present in the founder type-1 cell.
- 380 2) Randomly select cell, with most recent passenger mutation i , to have the event.
- 381 3) Determine which type of event and update population and mutation frequencies. Sample Y from
382 $\text{Uniform}[0, \Gamma]$ to determine event type:

- 383 i) $y \in (0, b_j) \rightarrow$ birth. $N_j += 1, z_i += 1$.
- 384 ii) $y \in (b_j, b_j + d_j) \rightarrow$ death. $N_j -= 1, z_i -= 1$.
- 385 iii) $y \in (b_j + d_j, b_j + d_j + u) \rightarrow$ passenger mutation. Suppose it's the k th passenger, $z_i -= 1, z_k = 1$.
- 386 Update ancestor: $a_k = i$.
- 387 4) For type-0, if time is less than t_1 and population goes extinct, restart simulation. For type-1, if time
- 388 is greater than t_1 and population goes extinct, restart type-1 simulation at t_1 with a single cell.
- 389 5) Reindex to remove extinct passenger mutations, and traverse back through ancestor vector \mathbf{a} to sum
- 390 total number of cells with each passenger.

391 Measurements are taken at bulk sequencing times $t_1 + t$ and $t_1 + t + \delta$. If time is greater than or equal to

392 $t_1 + t$, we measure $M_1 = N_0 + N_1$ and $\alpha_1 = N_1 / (N_0 + N_1)$. Then an additional bulk sequencing measurement

393 is taken at the final time $t_1 + t + \delta$, where we measure $M_2 = N_0 + N_1$ and $\alpha_2 = N_1 / (N_0 + N_1)$. At $t_1 + t + \delta$,

394 we measure γ , the number of mutations with frequency between f_1 and f_2 .

395 To measure m_{obs} , the observed number of passengers in the founder type-1 cell, we count the number of

396 passengers present in all type-1 cells. We also save the true value of m .

397 For when we calculate a percent error of corrected and observed γ values in Figure 3d and Supplementary

398 Figure 3b, we simulate sequencing data by sampling from the mutation frequencies obtained in the Monte

399 Carlo simulation, outlined above, using the approach of [35]. Define average sequencing coverage as R ,

400 number of cells at time of sequencing as M , Z_i as the number of cells with mutation i , R_i as read coverage,

401 and χ_i as the true mutation frequency from Monte Carlo simulation. For each saved Monte Carlo simulation

402 run, repeat the following 100 times:

- 403 1) Generate read coverage: $R_i \sim \text{Binomial}[M, R/M]$
- 404 2) Generate number of cells carrying mutation i : $Z_i \sim \text{Binomial}[R_i, \chi_i/2]$
- 405 3) Post-processing. If there are $L = 2$ or fewer variant reads, discard mutation.
- 406 4) Measure γ_{obs} , the observed number of subclonal mutations between frequencies f_1 and f_2 : $\gamma_{obs} =$
- 407 $\sum_i I(f_1 \leq 2Z_i/R \leq f_2, Z_i > L)$
- 408 5) Calculate the truth, γ_{true} , from the true mutation frequencies: $\gamma_{true} = \sum_i I(f_1 \leq \chi_i \leq f_2)$

409 Parameter values for simulations

410 For the simulation we consider three parameter sets corresponding to three modes of tumor evolution: a fast

411 growing tumor, slow growing tumor, and tumor with no cell death. For each parameter regime we have a low

412 and high mutation rate. Mutation rate parameter values lie within observed genome wide point mutation
413 rates per day [52]. For the fast growing tumor $b = b_1 = 0.25$, $d = 0.18$, $d_1 = 0.11$, $t_1 = 70$, $t = 50$, $\delta = 20$,
414 and $u = 1, 3$. For the slow growing tumor $b = 0.25$, $b_1 = 0.25$, $d = 0.225$, $d_1 = 0.2125$, $t_1 = 180$, $t = 135$,
415 $\delta = 45$, and $u = 1, 5$. For the parameter regime with no cell death $b = 0.25$, $b_1 = 0.375$, $d = d_1 = 0.0$,
416 $t_1 = 23$, $t = 17$, $\delta = 6$, and $u = 1, 10$. The fast growing tumor dynamics are from [34]. The slower growing
417 tumor parameter regime has a reduced net growth of $r = 0.025$, compared to the fast growing tumor's net
418 growth rate of $r = 0.07$.

419 Subclonal reconstruction of CLL sequencing data

420 The sequencing data from all CLLs analyzed is from Ref. [27], Supplementary Tables 2-4. As in that
421 publication, we use PhylogicNDT [43] to perform subclonal reconstruction. We run the Cluster and BuildTree
422 modules of PhylogicNDT on the longitudinal mutation data from Supplementary Table 3 of [27], using
423 mutation alternate/reference counts, copy number, and tumor purity at all pre-treatment time points. Then
424 for each patient, PhylogicNDT outputs a clonal reconstruction, which includes a phylogenetic tree of the
425 subclones and posterior distributions of subclone CCFs. Additionally, it clusters mutations and assigns them
426 to clones. We directly use subclone assignments and posteriors generated from PhylogicNDT. In our analysis
427 we focus on estimating timing and growth rates of macroscopic subclones whose CCFs are greater than 20%
428 for at least one pre-treatment timepoint.

429 Accounting for uncertainties in subclone frequencies and growth rates

430 Our estimates for parameters of cancer evolution require as input the information on the number of subclonal
431 populations in the tumor, their CCFs and their phylogenetic relationships. In order to obtain this informa-
432 tion, we use PhylogicNDT [43], which performs subclonal reconstruction of longitudinal cancer sequencing
433 data. The uncertainty in subclone CCFs reported by PhylogicNDT affects our estimates for subclone growth
434 rates, which in turn affect the estimates of mutation rate and and time t between driver(s) and diagnosis.
435 We account for this uncertainty by drawing from the CCF posterior distributions that are output by Phy-
436 logicNDT. Using these sampled CCF values, we then calculate growth rates, mutation rate u , and time t
437 between driver(s) and diagnosis, thereby generating confidence intervals for these parameters due to CCF
438 uncertainty.

439 To estimate subclonal growth rates, we fit an exponential growth curve to subclonal sizes measured at
440 two or more time points. This regression yields fitted values for each clone's growth rate and age. To account
441 for uncertainty in the curve fit (in the case of more than two longitudinal samples), we sample the growth

442 rates and age of clone from a bivariate normal distribution with mean equal to the fitted parameters and
443 variance equal to the covariance matrix of the fitted parameters. In line with recent findings [53], we found
444 that sometimes the estimated growth rate is smaller than minimal possible growth rate necessary to reach
445 the observed clone size. In that case, for calculating mutation rate, time of the driver(s), and time between
446 driver(s) and diagnosis we use the minimal growth rate.

447 **Accounting for model uncertainty**

448 The largest source of model uncertainty is the Poisson process for how mutations accumulate, which is used
449 to estimate the time t_1 of the driver mutation. In the simulation experiments, the time t_1 had the largest
450 error and variation (Fig. 2). The estimate for t_1 depends on the m mutations present in all cells in the
451 driver subclone. The observed m is a single random sample from a Poisson distribution. To account for
452 the uncertainty in t_1 arising from m in the CLLs analyzed, we sample t_i from the posterior distributions
453 $P(t_1|m)$. This source of model uncertainty due to the Poisson process will be most significant for cancers
454 like CLL with a smaller number of mutations.

455 The time t between driver mutation and diagnosis (t) is a random variable due to the stochasticity of
456 cancer cell growth, and will naturally have a certain amount of variation. Time between driver event and
457 diagnosis in a branching process follows a Gumbel distribution [38], and will have a constant variance. The
458 mean, however, will increase with the logarithm of the cancer cell counts, which for the CLLs analyzed are
459 $\sim 10^{11}$. The simulations of cancer evolution grow to smaller tumor sizes ($\sim 10^5$) and, as a result, the estimate
460 for t has a significant amount of uncertainty (Fig. 2). However, for time scales necessary to generate a tumor,
461 the estimate for t will be quite accurate. For commonly observed tumor sizes, the stochastic fluctuations in
462 the time for the cancer to reach that size will be smaller relative to the magnitude of the time. For a cancer
463 with cell count $\sim 10^{11}$, the standard deviation of the time t will be less than 5% of its expected value.

464 **Tumor with two nested driver subclones**

465 Here we consider the case where there are two nested driver subclones (Fig. S4a). “Nested” means that all
466 cells carrying the second driver mutation also carry the first. Type-0, or initiated tumor, cells have birth
467 rate b_0 , death rate d_0 , and net growth rate $r_0 = b_0 - d_0$. Type 1 cells, which only have the first driver, have
468 birth rate b_1 , death rate d_1 , and net growth rate $r_1 = b_1 - d_1$. Type-2 cells, which carry both drivers, have
469 birth rate b_2 , death rate d_2 , and net growth rate $r_2 = b_2 - d_2$. The first driver occurred in a type-0 cell at
470 time t_1 . The second driver occurred in a type-1 cell at $t_2 = t_1 + t'_2$. The mutation rate u is the same for all
471 subclones.

472 At times $t_1 + t'_2 + t$ and $t_1 + t'_2 + t + \delta$, the tumor is bulk sequenced. The bulk sequencing allows the
 473 measurement of the fraction of cells with driver 1 at time $t_1 + t'_2 + t$, α_1 ; the fraction of cells with driver 2
 474 at $t_1 + t'_2 + t$, α_2 ; fraction of cells with driver 1 at time $t_1 + t'_2 + t + \delta$, β_1 ; the fraction of cells with driver
 475 2 at $t_1 + t'_2 + t + \delta$, β_2 ; and the observed number of subclonal passenger mutations between frequencies f_1
 476 and f_2 , γ_{obs} . Note that the fraction of the population that is a type-1 cell at the two times is $\alpha_1 - \alpha_2$ and
 477 $\beta_1 - \beta_2$. The fraction of type-0 cells at the two bulk sequencing timepoints are $1 - \alpha_1$ and $1 - \beta_1$. The
 478 number of total cells at bulk sequencing timepoints are M_1 and M_2 . Equating the estimated cell counts to
 479 the expected value of the type- i population size X_i , conditioned on survival,

$$480 \quad \mathbb{E}\left[X_i\left(t_1 + t'_2 + t\right) \mid X_i\left(t_1 + t'_2 + t\right) > 0\right] = \begin{cases} \frac{b_0}{r_0} e^{r_0(t_1+t'_2+t)} & i = 0 \\ \frac{b_1}{r_1} e^{r_1(t'_2+t)} & i = 1 \\ \frac{b_2}{r_2} e^{r_2 t} & i = 2 \end{cases} \quad (31)$$

$$481 \quad = \begin{cases} (1 - \alpha_1)M_1 & i = 0 \\ (\alpha_1 - \alpha_2)M_1 & i = 1 \\ \alpha_2 M_1 & i = 2 \end{cases} \quad (32)$$

$$484 \quad \mathbb{E}\left[X_i\left(t_1 + t'_2 + t + \delta\right) \mid X_i\left(t_1 + t'_2 + t + \delta\right) > 0\right] = \begin{cases} \frac{b_0}{r_0} e^{r_0(t_1+t'_2+t+\delta)} & i = 0 \\ \frac{b_1}{r_1} e^{r_1(t'_2+t+\delta)} & i = 1 \\ \frac{b_2}{r_2} e^{r_2(t+\delta)} & i = 2 \end{cases} \quad (33)$$

$$485 \quad = \begin{cases} (1 - \beta_1)M_2 & i = 0 \\ (\beta_1 - \beta_2)M_2 & i = 1 \\ \beta_2 M_2 & i = 2 \end{cases} \quad (34)$$

487 Solving the above equations for r_i , we obtain the growth rate estimates:

$$488 \quad r_0 = \frac{1}{\delta} \log\left(\frac{(1 - \beta_1)M_2}{(1 - \alpha_1)M_1}\right) \quad (35)$$

$$489 \quad r_1 = \frac{1}{\delta} \log\left(\frac{(\beta_1 - \beta_2)M_2}{(\alpha_1 - \alpha_2)M_1}\right) \quad (36)$$

$$490 \quad r_2 = \frac{1}{\delta} \log\left(\frac{\beta_2 M_2}{\alpha_2 M_1}\right) \quad (37)$$

491

492 The expected value of the first time a population of type-2 cells in a branching process reaches the observed
 493 size $\alpha_2 M_1$ [38],

$$494 \quad \mathbb{E}t = \frac{1}{r_2} \log\left(\frac{\alpha_2 M_1 r_2}{b_2}\right) - \frac{1}{r_2} \int_0^\infty e^{-z} \log z dz \quad (38)$$

$$495 \quad = \frac{1}{r_2} \log\left(\frac{\alpha_2 M_1 r_2}{b_2}\right) + \frac{0.5772}{r_2} \quad (39)$$

497 can be approximated as

$$498 \quad \mathbb{E}t \approx \frac{1}{r_2} \log(\alpha_2 M_1) \quad (40)$$

500 We make use of two approximations to arrive at (40). First, we neglect the second term in (39), which serves
 501 as a small correction to the first term. Second, we assume r_2 is similar in magnitude to b_2 .

502 By (11),

$$503 \quad \mathbb{E}\gamma = u \left(\frac{1 - \beta_1}{r_0} + \frac{\beta_1 - \beta_2}{r_1} + \frac{\beta_2}{r_2} \right) \left(\frac{1}{f_1} - \frac{1}{f_2} \right) \quad (41)$$

505 Using the estimates for r_0 , r_1 , and r_2 from (35)-(37), and setting (41) equal to the value of γ obtained from
 506 (28) and the second bulk sequencing, u can be estimated:

$$507 \quad u = \frac{f_1 f_2 \gamma}{(f_2 - f_1) \left(\frac{1 - \beta_1}{r_0} + \frac{\beta_1 - \beta_2}{r_1} + \frac{\beta_2}{r_2} \right)} \quad (42)$$

509 When estimating mutation rate for the CLL patients from Ref. [27], for which there is bulk sequencing at
 510 two or more timepoints, we average the mutation rate calculated at each of these timepoints. (42) is applied
 511 for each timepoint with the respective CCFs and observed γ values for each timepoint.

512 Every type-1 cell carries the m_1 passenger mutations that were present in the original type-1 cell when
 513 the first driver mutation mutation occurred at t_1 . Similarly, every type-2 cell carries the m_2 passengers that
 514 were present in the founder type-2 cell when the second driver mutation occurred at t_2 . Note, none of the
 515 m_1 mutations are counted towards m_2 . Now we consider the likelihood function

$$516 \quad P(m_1, m_2 | t_1, t_2) \quad (43)$$

517

518
$$P(m_1, m_2 | t_1, t'_2) \propto P(m_1 | t_1) P(m_2 | t'_2) \quad (44)$$

519
$$\propto \frac{(ut_1)^{m_1} e^{-ut_1}}{m_1!} \frac{(ut'_2)^{m_2} e^{-ut'_2}}{m_2!} \quad (45)$$

520

521 Now, maximizing the logarithm of (45) with respect to t_1 and t'_2 ,

522
$$t_1 = \frac{m_1}{u} \quad (46)$$

523
$$t'_2 = \frac{m_2}{u} \quad (47)$$

524

525 The number of passengers present in the founder type- i cell cannot be directly observed, but we can
 526 measure $m_{i\text{ obs}}$, the number of passengers present in all type- i cells. An expected u/r_1 passengers occurring
 527 after t_1 in type-1 cells and reaching fixation in the type-1 subclone will be incorrectly included in $m_{1\text{ obs}}$,
 528 rather than in $m_{2\text{ obs}}$ (see Methods). Similarly, an expected u/r_2 passengers occurring after t_2 in type-2 cells
 529 and reaching fixation in the type-2 subclone will be incorrectly included in $m_{2\text{ obs}}$. Thus,

530
$$m_1 = m_{1\text{ obs}} - u/r_1 \quad (48)$$

531
$$m_2 = m_{2\text{ obs}} - u/r_2 + u/r_1 \quad (49)$$

532

533 Tumor with two sibling driver subclones

534 Here we consider a tumor with two “sibling” driver mutations (Fig. S4b). Sibling driver mutations are
 535 drivers that occur in separate subclones. In this case, cells are either initiated tumor cell (type-0), carry
 536 driver 1 (type-1), or carry driver 2 (type-2). No cells contain both drivers. Driver 1 occurred in a type-0
 537 cell at time t_1 . Driver 2 occurred in a type-0 cell at t_2 . Type-0 cells have birth rate b_0 , death rate d_0 , and
 538 net growth rate $r_0 = b_0 - d_0$. Type-1 cells, which carry driver 1, have birth rate b_1 , death rate d_1 , and
 539 net growth rate $r_1 = b_1 - d_1$. Type-2 cells, which carry driver 2, have birth rate b_2 , death rate d_2 , and net
 540 growth rate $r_2 = b_2 - d_2$. The mutation rate u is the same for all subclones.

541 Suppose time τ_i elapses between driver mutation i and tumor observation. Bulk sequencing of the tumor
 542 is performed at $t_1 + \tau_1$ (or equivalently $t_2 + \tau_2$), and a known δ later. Sequencing the tumor allows the
 543 measurement of the fraction of cells with driver 1 at the first sequencing, α_1 ; the fraction of cells with driver
 544 2 at the first sequencing, α_2 ; fraction of cells with driver 1 at the second sequencing, β_1 ; the fraction of
 545 cells with driver 2 at the second sequencing, β_2 ; and the number of subclonal passenger mutations between

546 frequencies f_1 and f_2 , γ . The fraction of type-0 cells at the two bulk sequencing timepoints are $1 - \alpha_1 - \alpha_2$
 547 and $1 - \beta_1 - \beta_2$. The number of total cells at the two sequencing timepoints are M_1 and M_2 .

548 Equating the estimated cell counts to the expected value of the type- i population size X_i , conditioned
 549 on survival,

$$550 \quad \mathbb{E}\left[X_i(t_i + \tau_i) \mid X_i(t_i + \tau_i) > 0\right] = \begin{cases} \frac{b_0}{r_0} e^{r_0(t_1 + \tau_1)} & i = 0 \\ \frac{b_i}{r_i} e^{r_i(\tau_i)} & i = 1, 2 \end{cases} \quad (50)$$

$$551 \quad = \begin{cases} (1 - \alpha_1 - \alpha_2)M_1 & i = 0 \\ \alpha_i M_1 & i = 1, 2 \end{cases} \quad (51)$$

553

$$554 \quad \mathbb{E}\left[X_i(t_i + \tau_i + \delta) \mid X_i(t_i + \tau_i + \delta) > 0\right] = \begin{cases} \frac{b_i}{r_i} e^{r_i(t_1 + \tau_1 + \delta)} & i = 0 \\ \frac{b_i}{r_i} e^{r_i(\tau_i + \delta)} & i = 1, 2 \end{cases} \quad (52)$$

$$555 \quad = \begin{cases} (1 - \beta_1 - \beta_2)M_2 & i = 0 \\ \beta_i M_2 & i = 1, 2 \end{cases} \quad (53)$$

556

557 Solving the above equations for r_i , we obtain

$$558 \quad r_0 = \frac{1}{\delta} \log\left(\frac{(1 - \beta_1 - \beta_2)M_2}{(1 - \alpha_1 - \alpha_2)M_1}\right) \quad (54)$$

$$559 \quad r_i = \frac{1}{\delta} \log\left(\frac{\beta_i M_2}{\alpha_i M_1}\right) \quad i = 1, 2 \quad (55)$$

560

561 The expected value of the first time a population of type- i cells in a branching process reaches the
 562 observed size $\alpha_i M_1$ is [38]

$$563 \quad \mathbb{E}\tau_i = \frac{1}{r_i} \log\left(\frac{\alpha_i M_1 r_i}{b_i}\right) - \frac{1}{r_i} \int_0^\infty e^{-z} \log z dz \quad (56)$$

$$564 \quad = \frac{1}{r_i} \log\left(\frac{\alpha_i M_1 r_i}{b_i}\right) + \frac{0.5772}{r_i} \quad (57)$$

565

566 which we approximate as

$$567 \quad \mathbb{E}\tau_i \approx \frac{1}{r_i} \log(\alpha_i M_1) \quad i = 1, 2 \quad (58)$$

568

569 We use two approximations to arrive at (58). We neglect the second term in (57), which serves as a small

570 correction to the first term. Second, we assume r_i is similar in magnitude to b_i .

571 By (11),

$$572 \mathbb{E}\gamma = u \left(\frac{1 - \beta_1 - \beta_2}{r_0} + \frac{\beta_1}{r_1} + \frac{\beta_2}{r_2} \right) \left(\frac{1}{f_1} - \frac{1}{f_2} \right) \quad (59)$$

574 Using the estimates for r_0 , r_1 , and r_2 from (54) and (55), and setting (59) equal to the value of γ obtained
575 from (28) and the second bulk sequencing, u can be estimated:

$$576 u = \frac{f_1 f_2 \gamma}{(f_2 - f_1) \left(\frac{1 - \beta_1 - \beta_2}{r_0} + \frac{\beta_1}{r_1} + \frac{\beta_2}{r_2} \right)} \quad (60)$$

578 When estimating mutation rate for the CLL patients from Ref. [27], for which there is bulk sequencing at
579 two or more timepoints, we average the mutation rate calculated at each of these timepoints. (60) is applied
580 for each timepoint with the respective CCFs and observed γ values for each timepoint.

581 Every type-1 cell carries the m_1 passenger mutations that were present in the original type-1 cell when
582 the first driver mutation mutation occurred at t_1 . Similarly, every type-2 cell carries the m_2 passengers that
583 were present in the founder type-2 cell when the second driver mutation occurred at t_2 . We assume that
584 m_1 and m_2 don't contain any shared mutations. In the CLL dataset we use, this is true. We consider the
585 likelihood function $P(m_1, m_2 | t_1, t_2)$

$$586 P(m_1, m_2 | t_1, t_2) \propto P(m_1 | t_1) P(m_2 | t_2) \quad (61)$$

$$587 \propto \frac{(ut_1)^{m_1} e^{-ut_1}}{m_1!} \frac{(ut_2)^{m_2} e^{-ut_2}}{m_2!} \quad (62)$$

589 Maximizing the logarithm of (62) with respect to t_1 and t_2 yields the maximum likelihood estimates:

$$590 t_1 = \frac{m_1}{u} \quad (63)$$

$$591 t_2 = \frac{m_2}{u} \quad (64)$$

593 Using the same approach as in the case of a single driver, we obtain the corrections for the observed number
594 of mutations present in all cells of each subclone:

$$595 m_1 = m_{1\text{ obs}} - u/r_1 \quad (65)$$

$$596 m_2 = m_{2\text{ obs}} - u/r_2 \quad (66)$$

597

598 Fully generalized estimates for any phylogeny of k drivers

599 Here we derive estimates for a completely general tumor phylogeny. Suppose a tumor has k driver mutations.
 600 In this general case, define a type- i cell as a cell where its most recent driver mutation was driver i . Note
 601 that a type- i cell can have between 0 and $k - 1$ other driver mutations. A phylogenetic reconstruction of the
 602 k driver mutations is necessary for the completely general case. From this phylogenetic tree, the ancestor
 603 of each subclone can be obtained. Define the function $a(i)$ as the ancestor of the type- i population. That
 604 is, if all driver mutations contained in the type- i population are ordered, $a(i)$ gives the driver mutation
 605 that occurred prior to i . Define t_i as the time between when driver i occurred and when the type- i cells'
 606 previous driver mutation occurred. At time of observation, assume the type- i population has κ_i total driver
 607 mutations, where $1 \leq \kappa_i \leq k$ for all $1 \leq i \leq k$. Denote the time between the type- i 's κ_i , or last, driver
 608 mutation and when the tumor is observed as τ_i . This is the time between the founder type- i cell's birth
 609 and tumor observation. Then the tumor is first observed and bulk sequenced at $T_1 \equiv (\sum_{j=0}^{\kappa_i-1} t_{a^j(i)}) + \tau_i$
 610 (equivalently τ_0 for $i = 0$), where we denote a^j as the j th iterate of the function a :

$$611 \quad a^0(i) \equiv i \quad (67)$$

$$612 \quad a^j(i) \equiv a(a^{j-1}(i)) \quad \forall j \geq 1 \quad (68)$$

614 The tumor is also bulk sequenced at $T_2 \equiv (\sum_{j=0}^{\kappa_i-1} t_{a^j(i)}) + \tau_i + \delta$ (equivalently $\tau_0 + \delta$ for $i = 0$). These
 615 assumptions allow for any subclone phylogeny, including combinations of the previously discussed sibling
 616 and nested subclone types.

617 The bulk sequencing allows the measurement of the fraction of cells with driver i at T_1 , α_i ; the fraction
 618 of cells with driver i at time T_2 , β_i ; and the number of subclonal passenger mutations between frequencies
 619 f_1 and f_2 , γ . Again, the number of total cells at measurement times T_1 and T_2 are M_1 and M_2 . To write the
 620 type- i frequencies, α_i^c and β_i^c , in terms of the driver frequencies, we subtract the fraction of cells descending
 621 from type- i cells but gaining additional driver mutation(s) after i , from the fraction of cells containing driver
 622 i :

$$623 \quad \alpha_i^c = \begin{cases} \alpha_i - \sum_{j=1}^k \delta_{i,a(j)} \alpha_j & 1 \leq i \leq k \\ 1 - \sum_{j=1}^k \alpha_j & i = 0 \end{cases} \quad (69)$$

$$624 \quad \beta_i^c = \begin{cases} \beta_i - \sum_{j=1}^k \delta_{i,a(j)} \beta_j & 1 \leq i \leq k \\ 1 - \sum_{j=1}^k \beta_j & i = 0 \end{cases} \quad (70)$$

625

626 where $\delta_{i,a(j)}$ is the Kronecker delta, defined as

$$627 \quad \delta_{i,a(j)} = \begin{cases} 0 & \text{if } i \neq a(j) \\ 1 & \text{if } i = a(j) \end{cases}$$

628

629 Equating the estimated cell counts at the first bulk sequencing timepoint to the expected value of the type- i
630 population size X_i , conditioned on survival,

$$631 \quad \mathbb{E}[X_i(T_1)|X_i(T_1) > 0] = \frac{b_i}{r_i} e^{r_i \tau_i}$$

$$632 \quad \quad \quad = \alpha_i^c M_1 \quad (71)$$

633

634 And similarly, at the second bulk sequencing timepoint,

$$635 \quad \mathbb{E}[X_i(T_2)|X_i(T_2) > 0] = \frac{b_i}{r_i} e^{r_i(\tau_i + \delta)}$$

$$636 \quad \quad \quad = \beta_i^c M_2 \quad (72)$$

637

638 Solving the above equations for r_i , we obtain

$$639 \quad r_i = \frac{1}{\delta} \log \left(\frac{\beta_i^c M_2}{\alpha_i^c M_1} \right) \quad \forall i = 0, 1, \dots, k \quad (74)$$

640

641 By (11)

$$642 \quad \mathbb{E}\gamma = \left(u \sum_{i=0}^k \frac{\beta_i^c}{r_i} \right) \left(\frac{1}{f_1} - \frac{1}{f_2} \right) \quad (75)$$

643

644 Now, using the growth rate estimates r_i and the subclone sizes, we can estimate each τ_i . The expected
645 value of the first time a population of type- i cells in a branching process reaches the observed size $\alpha_i^c M_1$ is
646 [38]

$$647 \quad \mathbb{E}\tau_i = \frac{1}{r_i} \log \left(\frac{\alpha_i^c M_1 r_i}{b_i} \right) - \frac{1}{r_i} \int_0^\infty e^{-z} \log z dz \quad (76)$$

$$648 \quad \quad \quad = \frac{1}{r_i} \log \left(\frac{\alpha_i^c M_1 r_i}{b_i} \right) + \frac{0.5772}{r_i} \quad (77)$$

649

650 which we approximate as

$$651 \quad \mathbb{E}\tau_i \approx \frac{1}{r_i} \log(\alpha_i^c M_1) \quad (78)$$

652

653 We make use of two approximations to arrive at (78). First, we neglect the second term in (77), which serves
654 as a small correction to the first term. Second, we assume r_i is similar in magnitude to b_i .

655 Using the $(k+1)$ r_i estimates from (74), and setting (75) equal to the value of γ obtained at the second
656 bulk sequencing from (28), u can be estimated:

$$657 \quad u = \frac{f_1 f_2 \gamma}{(f_2 - f_1) \left(\sum_{i=0}^k \frac{\beta_i^c}{r_i} \right)} \quad (79)$$

658

659 When estimating mutation rate for the CLL patients from Ref. [27], for which there is bulk sequencing at
660 two or more timepoints, we average the mutation rate calculated at each of these timepoints. (79) is applied
661 for each timepoint with the respective CCFs and observed γ values for each timepoint.

662 The number of passengers present in the original type i founder cell cannot be directly observed, but we
663 can measure m_i , the number of clonal passengers present in the type i population, only including passengers
664 not present in other clones. We will assume that the m_i don't contain any shared mutations, which is true
665 for the CLL dataset we consider. The likelihood function $P(m_1, \dots, m_k | t_1, \dots, t_k)$ is proportional to

$$666 \quad \prod_{i=1}^k P(m_i | t_i) \propto \prod_{i=1}^k \frac{(ut_i)^{m_i} e^{-ut_i}}{m_i!} \quad (80)$$

667

668 Then, maximizing the logarithm of (80) with respect to t_1, t_2, \dots, t_k ,

$$669 \quad t_i = \frac{m_i}{u} \quad \forall i = 1, \dots, k \quad (81)$$

670

671 The observed clonal passengers in the founder type- i cell will incorrectly include passengers that reached
672 fixation in the type- i population after driver mutation i occurred, instead of correctly being counted toward
673 the descendant of clone i . As a result, we again correct for the expected number of these passengers, u/r_i .

674 That is,

$$675 \quad m_i = m_{i,obs} - u/r_i + u/r_{a(i)} \quad \forall i = 1, \dots, k \quad (82)$$

676

⁶⁷⁷ **Availability of data and materials**

⁶⁷⁸ All simulated data generated during this study are included in this published article and its supplementary
⁶⁷⁹ information files. CLL data analyzed is publicly available in Supplementary Tables from Ref. [27]. Code
⁶⁸⁰ can be found at https://github.com/nathanlee543/Cancer_Inf_Sims

References

- 681
- 682 [1] Nowell, P. C. The Clonal Evolution of Tumor Cell Populations. *Science* **194**, 23–28 (1976). URL
683 <http://www.jstor.org/stable/1742535>.
- 684 [2] Stratton, M. R., Campbell, P. J. & Futreal, P. A. The cancer genome. *Nature* **458**, 719–724 (2009).
685 URL <https://www.nature.com/articles/nature07943>.
- 686 [3] Hanahan, D. & Weinberg, R. A. Hallmarks of Cancer: The Next Generation. *Cell* **144**, 646–674 (2011).
687 URL [https://www.cell.com/fulltext/S0092-8674\(11\)00127-9](https://www.cell.com/fulltext/S0092-8674(11)00127-9).
- 688 [4] Bailey, M. H. et al. Comprehensive Characterization of Cancer Driver Genes and Mutations. *Cell* **173**,
689 371–385.e18 (2018). URL [https://www.cell.com/cell/abstract/S0092-8674\(18\)30237-X](https://www.cell.com/cell/abstract/S0092-8674(18)30237-X).
- 690 [5] Williams, M. J., Werner, B., Barnes, C. P., Graham, T. A. & Sottoriva, A. Identification of neutral
691 tumor evolution across cancer types. *Nature Genetics* **48**, 238–244 (2016). URL <https://www.nature.com/articles/ng.3489>.
- 692
- 693 [6] Kimura, M. Evolutionary Rate at the Molecular Level. *Nature* **217**, 624–626 (1968). URL <https://www.nature.com/articles/217624a0>.
- 694
- 695 [7] Kimura, M. Genetic variability maintained in a finite population due to mutational
696 production of neutral and nearly neutral isoalleles*. *Genetics Research* **11**, 247–270
697 (1968). URL [http://www.cambridge.org/core/journals/genetics-research/article/
698 genetic-variability-maintained-in-a-finite-population-due-to-mutational-production-of-neutral-and-
699 A74BD3A5D72ED2C52444FD99DFE483EF](http://www.cambridge.org/core/journals/genetics-research/article/genetic-variability-maintained-in-a-finite-population-due-to-mutational-production-of-neutral-and-A74BD3A5D72ED2C52444FD99DFE483EF).
- 700 [8] Smith, J. M. & Haigh, J. The hitch-hiking effect of a favourable gene. *Genetics Research*
701 **23**, 23–35 (1974). URL [http://www.cambridge.org/core/journals/genetics-research/article/
702 hitchhiking-effect-of-a-favourable-gene/918291A3B62BD50E1AE5C1F22165EF1B](http://www.cambridge.org/core/journals/genetics-research/article/hitchhiking-effect-of-a-favourable-gene/918291A3B62BD50E1AE5C1F22165EF1B).
- 703 [9] Turajlic, S., Sottoriva, A., Graham, T. & Swanton, C. Resolving genetic heterogeneity in can-
704 cer. *Nature Reviews Genetics* **20**, 404–416 (2019). URL [https://www.nature.com/articles/
705 s41576-019-0114-6](https://www.nature.com/articles/s41576-019-0114-6).
- 706 [10] Vogelstein, B. et al. Cancer Genome Landscapes. *Science* **339**, 1546–1558 (2013). URL [https://
707 science.sciencemag.org/content/339/6127/1546](https://science.sciencemag.org/content/339/6127/1546).
- 708 [11] Lawrence, M. S. et al. Mutational heterogeneity in cancer and the search for new cancer-associated
709 genes. *Nature* **499**, 214–218 (2013). URL <https://www.nature.com/articles/nature12213>.

- 710 [12] Merlo, L. M. F., Pepper, J. W., Reid, B. J. & Maley, C. C. Cancer as an evolutionary and ecological pro-
711 cess. Nature Reviews Cancer **6**, 924–935 (2006). URL <https://www.nature.com/articles/nrc2013>.
- 712 [13] Pepper, J. W., Findlay, C. S., Kassen, R., Spencer, S. L. & Maley, C. C. SYNTHESIS: Cancer research
713 meets evolutionary biology. Evolutionary Applications **2**, 62–70 (2009). URL <https://onlinelibrary.wiley.com/doi/abs/10.1111/j.1752-4571.2008.00063.x>.
- 714
- 715 [14] Tsao, J.-L. et al. Genetic reconstruction of individual colorectal tumor histories.
716 Proceedings of the National Academy of Sciences **97**, 1236–1241 (2000). URL <https://www.pnas.org/content/97/3/1236>.
- 717
- 718 [15] Jones, S. et al. Comparative lesion sequencing provides insights into tumor evolu-
719 tion. Proceedings of the National Academy of Sciences of the United States of America **105**, 4283–
720 4288 (2008). URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2393770/>.
- 721 [16] Yachida, S. et al. Distant metastasis occurs late during the genetic evolution of pancreatic cancer.
722 Nature **467**, 1114–1117 (2010). URL <https://www.nature.com/articles/nature09515>.
- 723 [17] Naxerova, K. et al. Hypermutable DNA chronicles the evolution of human colon cancer.
724 Proceedings of the National Academy of Sciences **111**, E1889–E1898 (2014). URL <http://www.pnas.org/cgi/doi/10.1073/pnas.1400179111>.
- 725
- 726 [18] McGranahan, N. et al. Clonal status of actionable driver events and the timing of mutational processes
727 in cancer evolution. Science Translational Medicine **7**, 283ra54–283ra54 (2015). URL <https://stm.sciencemag.org/content/7/283/283ra54>.
- 728
- 729 [19] Mitchell, T. J. et al. Timing the Landmark Events in the Evolution of Clear Cell Renal Cell Cancer:
730 TRACERx Renal. Cell **173**, 611–623.e17 (2018). URL <http://www.sciencedirect.com/science/article/pii/S0092867418301648>.
- 731
- 732 [20] PCAWG Evolution & Heterogeneity Working Group et al. The evolutionary history of 2,658 cancers.
733 Nature **578**, 122–128 (2020). URL <http://www.nature.com/articles/s41586-019-1907-7>.
- 734 [21] Sundermann, L. K., Wintersinger, J., Rättsch, G., Stoye, J. & Morris, Q. Reconstructing tumor evolu-
735 tionary histories and clone trees in polynomial-time with SubMARine. PLOS Computational Biology
736 **17**, e1008400 (2021). URL <https://dx.plos.org/10.1371/journal.pcbi.1008400>.
- 737 [22] PCAWG Evolution and Heterogeneity Working Group et al. Reconstructing evolutionary trajectories of
738 mutation signature activities in cancer using TrackSig. Nature Communications **11**, 731 (2020). URL
739 <http://www.nature.com/articles/s41467-020-14352-7>.

- 740 [23] Tomasetti, C. & Bozic, I. The (not so) immortal strand hypothesis. Stem Cell Research **14**, 238–241
741 (2015).
- 742 [24] Werner, B. et al. Measuring single cell divisions in human tissues from multi-region sequenc-
743 ing data. Nature Communications **11**, 1–9 (2020). URL [https://www.nature.com/articles/](https://www.nature.com/articles/s41467-020-14844-6)
744 [s41467-020-14844-6](https://www.nature.com/articles/s41467-020-14844-6). Number: 1 Publisher: Nature Publishing Group.
- 745 [25] Bozic, I. et al. Accumulation of driver and passenger mutations during tumor progression.
746 Proceedings of the National Academy of Sciences of the United States of America **107**, 18545–18550
747 (2010). URL <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2972991/>.
- 748 [26] Sun, R. et al. Between-region genetic divergence reflects the mode and tempo of tumor evolution.
749 Nature Genetics **49**, 1015–1024 (2017). URL <https://www.nature.com/articles/ng.3891>.
- 750 [27] Gruber, M. et al. Growth dynamics in naturally progressing chronic lymphocytic leukaemia. Nature
751 **570**, 474–479 (2019). URL <https://www.nature.com/articles/s41586-019-1252-x>.
- 752 [28] Salichos, L., Meyerson, W., Warrell, J. & Gerstein, M. Estimating growth patterns and driver effects
753 in tumor evolution from individual samples. Nature Communications **11**, 1–14 (2020). URL <https://www.nature.com/articles/s41467-020-14407-9>.
754 [//www.nature.com/articles/s41467-020-14407-9](https://www.nature.com/articles/s41467-020-14407-9).
- 755 [29] Noble, R. et al. Spatial structure governs the mode of tumour evolution. Nature Ecology & Evolution
756 (2021). URL <https://www.nature.com/articles/s41559-021-01615-9>.
- 757 [30] Chkhaidze, K. et al. Spatially constrained tumour growth affects the patterns of clonal selection and
758 neutral drift in cancer genomic data. PLOS Computational Biology **15**, e1007243 (2019). URL <https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1007243>.
759 [//journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1007243](https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1007243).
- 760 [31] Fu, X. et al. Spatial patterns of tumour growth impact clonal diversification in a computational model
761 and the TRACERx Renal study. Nature Ecology & Evolution (2021). URL [https://www.nature.com/](https://www.nature.com/articles/s41559-021-01586-x)
762 [articles/s41559-021-01586-x](https://www.nature.com/articles/s41559-021-01586-x).
- 763 [32] Williams, M. J. et al. Quantification of subclonal selection in cancer from bulk sequencing data.
764 Nature Genetics **50**, 895 (2018). URL <https://www.nature.com/articles/s41588-018-0128-6>.
- 765 [33] Avanzini, S. et al. A mathematical model of ctDNA shedding predicts tumor detection size.
766 Science Advances **6**, eabc4308 (2020). URL [https://www.science.org/doi/10.1126/sciadv.](https://www.science.org/doi/10.1126/sciadv.abc4308)
767 [abc4308](https://www.science.org/doi/10.1126/sciadv.abc4308).

- 768 [34] Bozic, I. *et al.* Evolutionary dynamics of cancer in response to targeted combination therapy. *eLife* **2**,
769 e00747 (2013). URL <https://elifesciences.org/articles/00747>.
- 770 [35] Dinh, K. N., Jaksik, R., Kimmel, M., Lambert, A. & Tavaré, S. Statistical Inference for the Evolutionary
771 History of Cancer Genomes. *Statistical Science* **35**, 129–144 (2020). URL <https://projecteuclid.org/euclid.ss/1583226033>.
772
- 773 [36] Lahouel, K. *et al.* Revisiting the tumorigenesis timeline with a data-driven generative model.
774 *Proceedings of the National Academy of Sciences* **117**, 857–864 (2020). URL <http://www.pnas.org/lookup/doi/10.1073/pnas.1914589117>.
775
- 776 [37] Bozic, I. & Wu, C. J. Delineating the evolutionary dynamics of cancer from theory to reality.
777 *Nature Cancer* **1**, 580–588 (2020). URL <http://www.nature.com/articles/s43018-020-0079-6>.
- 778 [38] Durrett, R. Branching Process Models of Cancer. In Durrett, R. (ed.)
779 *Branching Process Models of Cancer*, Mathematical Biosciences Institute Lecture Series, 1–63 (Springer
780 International Publishing, Cham, 2015). URL https://doi.org/10.1007/978-3-319-16065-8_1.
- 781 [39] Tavaré, S. The linear birth–death process: an inferential retrospective. *Advances in Applied Probability*
782 **50**, 253–269 (2018). URL https://www.cambridge.org/core/product/identifier/S0001867818000848/type/journal_article.
783
- 784 [40] Heyde, A., Reiter, J. G., Naxerova, K. & Nowak, M. A. Consecutive seeding and transfer of genetic
785 diversity in metastasis. *Proceedings of the National Academy of Sciences* **116**, 14129–14137 (2019).
786 URL <https://www.pnas.org/content/116/28/14129>.
- 787 [41] Griffith, M. *et al.* Optimizing Cancer Genome Sequencing and Analysis. *Cell Systems* **1**, 210–223 (2015).
788 URL <http://www.sciencedirect.com/science/article/pii/S2405471215001131>.
- 789 [42] Haber, D. A. & Velculescu, V. E. Blood-Based Analyses of Cancer: Circulating Tumor Cells and
790 Circulating Tumor DNA. *Cancer Discovery* **4**, 650–661 (2014). URL <https://cancerdiscovery.aacrjournals.org/content/4/6/650>.
791
- 792 [43] Leshchiner, I. *et al.* Comprehensive analysis of tumour initiation, spatial and temporal progression under
793 multiple lines of treatment. preprint, Bioinformatics (2018). URL <http://biorxiv.org/lookup/doi/10.1101/508127>.
794
- 795 [44] Myers, M. A., Satas, G. & Raphael, B. J. CALDER: Inferring Phylogenetic Trees from Longitudinal
796 Tumor Samples. *Cell Systems* **8**, 514–522.e5 (2019). URL <http://www.sciencedirect.com/science/article/pii/S2405471219301917>.
797

- 798 [45] Carlsson, G., Gullberg, B. & Hafström, L. Estimation of liver tumor volume using different formulas? An
799 experimental study in rats. Journal of Cancer Research and Clinical Oncology **105**, 20–23 (1983). URL
800 <http://link.springer.com/10.1007/BF00391826>.
- 801 [46] Bozic, I., Gerold, J. M. & Nowak, M. A. Quantifying Clonal and Subclonal Passenger Mutations in
802 Cancer Evolution. PLOS Computational Biology **12**, e1004731 (2016). URL <http://dx.plos.org/10.1371/journal.pcbi.1004731>.
803
- 804 [47] Blokzijl, F. et al. Tissue-specific mutation accumulation in human adult stem cells during life. Nature
805 **538**, 260–264 (2016). URL <http://www.nature.com/articles/nature19768>.
- 806 [48] Mitchell, E. et al. Clonal dynamics of haematopoiesis across the human lifespan. preprint, Genomics
807 (2021). URL <http://biorxiv.org/lookup/doi/10.1101/2021.08.16.456475>.
- 808 [49] Auslander, N., Wolf, Y. I. & Koonin, E. V. In silico learning of tumor evolution through mutational
809 time series. Proceedings of the National Academy of Sciences **116**, 9501–9510 (2019). URL <https://www.pnas.org/content/116/19/9501>.
810
- 811 [50] Petljak, M. et al. Characterizing Mutational Signatures in Human Cancer Cell Lines Reveals Episodic
812 APOBEC Mutagenesis. Cell **176**, 1282–1294.e20 (2019). URL <https://linkinghub.elsevier.com/retrieve/pii/S0092867419301618>.
813
- 814 [51] Gillespie, D. T. Exact stochastic simulation of coupled chemical reactions.
815 The Journal of Physical Chemistry **81**, 2340–2361 (1977). URL <https://doi.org/10.1021/j100540a008>.
816
- 817 [52] Bozic, I., Paterson, C. & Waclaw, B. On measuring selection in cancer from subclonal mutation fre-
818 quencies. PLOS Computational Biology **15**, e1007368 (2019). URL <https://journals.plos.org/ploscompbiol/article?id=10.1371/journal.pcbi.1007368>.
819
- 820 [53] Fabre, M. A. et al. The longitudinal dynamics and natural history of clonal haematopoiesis. preprint,
821 Cancer Biology (2021). URL <http://biorxiv.org/lookup/doi/10.1101/2021.08.12.455048>.