

1 **Inferring ongoing cancer evolution from single tumour** 2 **biopsies using synthetic supervised learning**

3
4
5

Authors: Tom W. Ouellette^{1,2*} and Philip Awadalla^{1,2*}

6 **Affiliations:**

7 ¹Ontario Institute for Cancer Research, Department of Computational Biology, Toronto, Ontario M5G 0A3,
8 Canada

9 ²Department of Molecular Genetics, Temerty Faculty of Medicine, University of Toronto, Toronto, Ontario
10 M5S 1A8, Canada

11 *Correspondence to: tom.ouellette@oicr.on.ca or philip.awadalla@oicr.on.ca

12 Keywords: cancer evolution, subclonal selection, variant allele frequencies, population genetics,
13 simulation, synthetic data, deep learning, transfer learning

14 **Abstract**

15 Variant allele frequencies (VAF) encode ongoing evolution and subclonal selection in growing tumours.
16 However, existing methods that utilize VAF information for cancer evolutionary inference are compressive,
17 slow, or incorrectly specify the underlying cancer evolutionary dynamics. Here, we provide a proof-of-
18 principle synthetic supervised learning method, TumE, that integrates simulated models of cancer evolution
19 with Bayesian neural networks, to infer ongoing selection in bulk-sequenced single tumour biopsies.
20 Analyses in synthetic and patient tumours shows that TumE significantly improves both accuracy and
21 inference time per sample when detecting positive selection, deconvoluting selected subclonal populations,
22 and estimating subclone frequency. Importantly, we show how transfer learning can leverage stored
23 knowledge within TumE models for related evolutionary inference tasks — substantially reducing data and
24 computational time for further model development and providing a library of recyclable deep learning
25 models for the cancer evolution community. This extensible framework provides a foundation and future
26 directions for harnessing progressive computational methods for the benefit of cancer genomics and, in
27 turn, the cancer patient. TumE is publicly available for use at <https://github.com/tomouellette/TumE>.

28
29
30
31

32 Introduction

33 Cancer is a disease characterized by unrelenting tissue growth and clonal evolution. During evolution,
34 genetic and epigenetic aberrations provide the reservoir for dysfunctional cellular phenotypes that maintain
35 a tumour's replicative advantage, while, over time, fluctuating physiological and ecological properties within
36 the tumour microenvironment drive the need for updated adaptations that sustain immortality¹. Overall, the
37 complex interplay between mutation accumulation and microenvironmental changes leads to a high degree
38 of both cellular and genetic heterogeneity and, by proxy, composite subclonal structure in tumours²⁻⁴.
39 Naturally, the desire to better understand the evolutionary and subclonal dynamics in growing tumour
40 populations has become a major task for cancer genomics - with goals of forecasting tumour progression,
41 developing adaptive evolutionary therapies, and deconvoluting the genetic architecture that drives
42 adaptation^{3,5-8}.

43
44 However, a significant hurdle in understanding cancer evolution *in vivo* are the clinical constraints
45 surrounding serial sequencing, through space or time. For this reason, tumour biopsies are primarily
46 sequenced in bulk from a single site and at a single time point. Although multi-region and single-cell data
47 are becoming increasingly utilized, single time point, bulk sequenced biopsies still represent the major
48 accessible data source for precision genomics guided treatment⁹ and for studying cancer genomics and
49 evolution in patients^{10,11}. Given this limitation, a reasonable strategy for inferring evolution in single tumor
50 biopsies has been to utilize theoretical population genetics to capture signatures of selection from the
51 variant allele frequency (VAF) distribution^{7,12-17}. The premise being that fitness-altering mutations will
52 deterministically change in frequency over time, leading to characteristic and quantifiable deviations in the
53 VAF distribution relative to some neutral evolutionary scenario¹⁸.

54
55 VAF-based methods have been employed to differentiate between positive selection and neutral
56 evolution^{12,13}, to examine growth patterns¹⁹, to quantify subclonal fitness and time subclonal emergence^{7,15},
57 and to build population genetics informed mixture models¹⁶ that account for neutral dynamics, that shape,
58 to some extent, all tumour populations. With that said, existing VAF-based methods used to infer cancer
59 evolution, although mechanistic and useful, have apparent limitations. For example, single statistics^{12,20,21}
60 are maximally compressive and cannot infer complex information, approximate Bayesian computation
61 methods suffer from the curse of dimensionality and can be prohibitively slow due to a rate-limiting
62 simulation step required for each sample^{7,22,23}, and mixture models, used to identify subclonal
63 populations^{16,24,25}, are only implicitly connected to an underlying model of evolution and, until recently¹⁶,
64 have been built under incorrect assumptions that have led to systematic overestimation in the number of
65 subclonal populations in sequenced tumours.

66
67 To address these limitations, we contribute a proof-of-principle synthetic supervised deep learning
68 approach, TumE, for quantifying and classifying the evolutionary and subclonal dynamics in bulk sequenced

69 tumours biopsies using purity-corrected variant allele frequency (VAF) information from diploid genomic
 70 regions. By generating synthetic VAF distributions, as a proxy for evolutionary ground truth, from plausible
 71 simulations of tumour evolution, we were able to build inference models that accurately classify and quantify
 72 evolutionary (e.g. positive selection versus neutral evolution) and subclonal dynamics (e.g. subclone
 73 frequency) in real patient tumours while capturing uncertainty in our estimates, via a form of approximate
 74 Bayesian inference called Monte Carlo dropout^{26,27}. Importantly, our method further highlights the power of
 75 utilizing deep learning for inference - namely exploiting stored knowledge via transfer learning. By recycling
 76 our models for new evolutionary prediction tasks, we reduce the computational burden associated with the
 77 generation of synthetic or simulated data. We validated our synthetic supervised learning approach in
 78 millions of synthetic tumours and applied TumE to 95 copy-number and purity corrected whole-genome
 79 (WGS) and whole-exome (WES) sequenced tumour biopsies.

80

81 Results

82 Inferring cancer evolution using synthetic supervised deep learning

83 Synthetic supervised, or simulation-based, deep learning has been shown to be equivalent to amortized
 84 approximate inference under a generative model²⁸. Therefore, by optimizing a neural network using realistic
 85 synthetic data \mathbf{x} generated from a stochastic generative process $p(\mathbf{x}, \mathbf{z} | \theta)$, where θ indicates the prior or
 86 parameters that define the simulation and \mathbf{z} indicates the latent variables generated during simulation, we
 87 can build inference models that approximate our true posterior of interest $p(\theta, \mathbf{z} | \mathbf{x})$. In our case, by optimizing
 88 a neural network using synthetic VAF distributions sampled from $p(\mathbf{x}, \mathbf{z} | \theta)$, we can build inference models
 89 for evolutionary inference in sequenced tumour biopsies (Figure 1A-C; Methods).

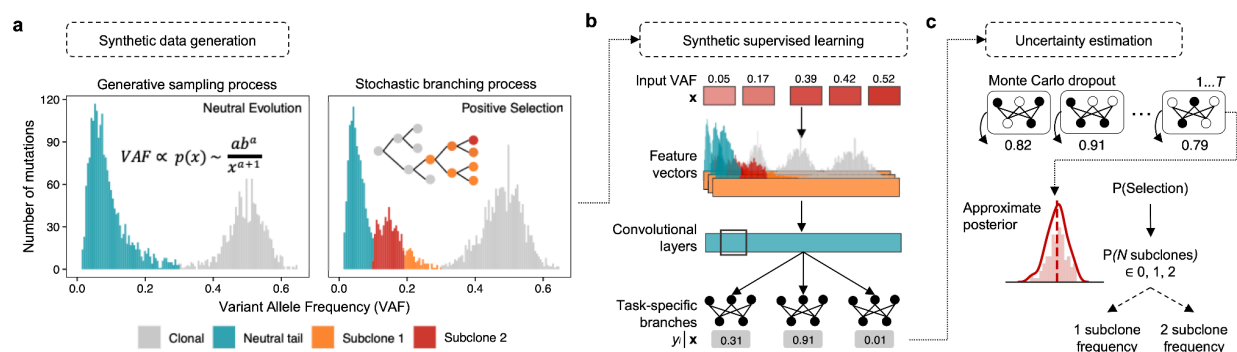


Figure 1. (a) TumE integrates a generative sampling process and stochastic simulation of cancer evolution to build well-specified synthetic variant allele frequency (VAF) distributions with respect to data observed in bulk sequenced tumour biopsies. Assuming copy-neutral diploid regions of tumour genomes, the generative sampling process uses the observation that neutral VAF distributions can be described by a power-law or Pareto neutral ‘tail’^{16,29} in addition to a dispersed clonal peak. By sampling empirically valid Pareto distributions, rapid realizations of the null hypothesis of neutral evolution encoded in the VAF distribution can be created (Methods). The stochastic branching process model of tumour evolution is then used to link parameters and latent states, relevant to positive subclonal selection, back to VAF distributions (Methods). **(b)** Synthetic supervised learning utilizes neural networks capable of handling the complete dimensionality of the simulated VAF distributions, \mathbf{x} , to solve the inverse problem of identifying the evolutionary parameters and latent states, \mathbf{y} , assigned to each synthetic VAF distribution. **(c)** We can then quantify model uncertainty using a computationally efficient form of Bayesian deep learning called Monte Carlo dropout^{26,27}. Approximate posteriors are generated by performing T stochastic passes through the trained neural network.

90 To generate synthetic data that properly captured evolutionary dynamics in patient tumours, we
91 implemented a simulation framework, i.e. a stochastic generative process $p(\mathbf{x}, \mathbf{z} | \Theta)$, combining two
92 complementary approaches to improve the speed and efficiency of synthetic data generation — one for
93 tumours subject to positive selection and one for tumours evolving neutrally (Figure 1A). For growing
94 tumours simulated with positive selection, we utilized a well-established framework of cancer evolution that
95 models exponential tumour growth under a stochastic branching process^{7,12,13,15,19,30} and coupled this with
96 a virtual biopsy procedure to account for sequencing noise observed in real patient tumours (adapted from
97 ref⁷). In our model, we allowed for a completely stochastic arrival of driver mutations that multiplicatively
98 increased the fitness of mutated subclones and tracked the frequency of each subclone until the time of
99 virtual biopsy (Methods). In this study, we define a subclone as a subpopulation of cells with a fitness or
100 growth rate advantage relative to the background population (Methods) and consider subclones detectable
101 if they are between ~10 - 40% VAF (20 - 80% cellular fraction). For tumours that lacked selected subclones
102 (neutrally evolving), we implemented a generative sampling process based on the observation that VAF
103 distributions from tumours without positively selected subclones can be described by a power-law or Pareto
104 distribution^{16,29} in conjunction with a dispersed clonal peak (Supplementary Figure 1). Concisely, this
105 process involved i) sampling allele frequencies from empirically realistic Pareto distributions to generate
106 the neutral power-law 'tail' in the VAF distribution, ii) adding additional diploid clonal heterozygous mutations
107 at 50% VAF, and then iii) injecting additional sequencing noise under a beta-binomial model (Methods). In
108 general, a complete VAF distribution indicative of positive selection, and computed from heterozygous
109 diploid mutations, includes a neutral power-law tail^{12,16}, a heterozygous clonal peak centered at ~50% VAF,
110 and additional subclonal peak(s) in the intermediate frequency ranges (~10 - 40% VAF); whereas a neutrally
111 evolving tumour, or one with undetectable selected subpopulations, lacks the characteristic subclonal
112 peak(s) (Figure 1A). To ensure positively selected and neutrally evolving synthetic tumours were not out of
113 distribution with each other given the alternate data generation approaches, we simulated synthetic tumours
114 in pairs, assigning the neutral VAF distributions with equivalent parameters and mutations with respect to
115 the paired positive selection simulation (Methods; pseudo algorithms and examples provided alongside
116 Supplementary Figure 2 & 3).

117
118 Using this framework, we generated approximately 40 million synthetic tumours across varying mutation
119 rates, selection coefficients, and sequencing noise parameters. We selected broad simulation parameter
120 ranges that were consistent with previous computational studies and empirically estimated values
121 (Methods; Supplementary Table S1). By generating synthetic tumours using well-specified simulations
122 (comparison of real and synthetic data outlined in Methods and Supplementary Figures 3 - 5), we were able
123 to explicitly link each VAF distribution to the parameters and latent states that defined the underlying
124 subclonal and evolutionary dynamics. We then used the millions of annotated synthetic VAF distributions
125 to train hundreds of neural networks using a random hyperparameter search to make inferences on the
126 evolutionary mode (positive selection or neutral evolution), the number of subclones (0, 1, or 2), and the

127 subclone frequency at borderline to optimal sequencing depths (50 - 250X) for evolutionary analysis in
128 cancer genomics (Figure 1B; Methods). To capture model-based uncertainty in our estimates, we
129 implemented a form of Bayesian approximation for deep learning called Monte Carlo (MC) dropout^{26,27}
130 (Figure 1C; Methods). We used MC dropout to mitigate overconfident estimates in cases of high uncertainty
131 or broad approximate posteriors. In general, we structured both neural network training and prediction to
132 favour the more parsimonious explanation of the data (fewer subclones and neutral evolution; Methods).
133 We show how using a classification threshold based on a lower bound of the MC dropout approximate
134 posterior helps mitigate model overconfidence across changing subclone mutation and frequencies in
135 Supplementary Figure 6. Following training, we selected the top scoring models, for predicting the
136 evolutionary mode, number of subclones, and subclone frequency, for further validation (Methods).

137
138 We outline the full synthetic supervised learning pipeline in Methods. In addition, we highlight that even
139 though we model VAF distributions in patient tumours using point mutations from diploid regions, mutations
140 in our framework, as with previous approaches^{7,13,16}, are agnostic to the underlying functional alteration,
141 e.g. missense, silent, driver or copy number driving selection in patient tumours. This is because genome-
142 wide linkage, a by-product of zero recombination, results in hitchhiking of any additional point mutations on
143 the genetic background of any selected clone^{3,18}.

144

145 **Comparison of synthetic supervised learning to existing methods**

146 To evaluate TumE performance on inferred estimates of selection, number of subclones, and subclone
147 frequency, we simulated an additional ~2.8 million synthetic tumours under neutral evolution (0 subclones)
148 and positive selection (1 or 2 detectable selected subclones) assessing the impact of variable sequencing
149 depths (50 - 250x coverage) and read count overdispersion (0 - 0.3 rho) (Methods). We first compared
150 TumE against frequency-based summary statistic approaches for differentiating between neutral evolution
151 and positive selection and found that TumE significantly outperforms recently developed VAF summary
152 statistics¹² (two-sided Wilcoxon test, $p = 2.7 \times 10^{-12}$) as well as common population genetic summary
153 statistics^{20,21} (two-sided Wilcoxon test, $p = 1.9 \times 10^{-8}$), based on AUROC (Figure 2A). Further, TumE
154 outperforms each statistic individually when compared across all sequencing depth and overdispersion
155 combinations analyzed here (ROC analysis; Supplementary Figure 7).

156

157 We next compared TumE against the only mixture model approach, MOBSTER¹⁶, that explicitly and
158 correctly takes into account the neutral dynamics within sequenced tumour VAF distributions to detect
159 subclones. We found that TumE provides comparable or improved performance for predicting the number
160 of subclones (precision-recall, Supplementary Figure 8) and for predicting subclone frequency (Figure 2B;
161 correlation and mean absolute percentage error, Supplementary Figure 9) across all empirically relevant
162 depth (50 - 250x coverage) and read count overdispersions (0 - 0.003 rho) combinations. However, as
163 expected, we found that the performance of TumE and MOBSTER both degrade as sequencing depth

164 decreases ($\leq 75x$ coverage) and overdispersion increases (≥ 0.01) under a beta-binomial sequencing noise
 165 model (Supplementary Figure 8 & 9). Furthermore, additional analysis of subclone frequency estimates in
 166 the 2 subclone setting revealed that as inter-subclone distance increases, i.e. overlap of subclonal peaks
 167 decrease, the mean percentage error for predicting the frequency of both the lowest and highest frequency
 168 subclone decreases towards zero (Supplementary Figures 10 - 12).

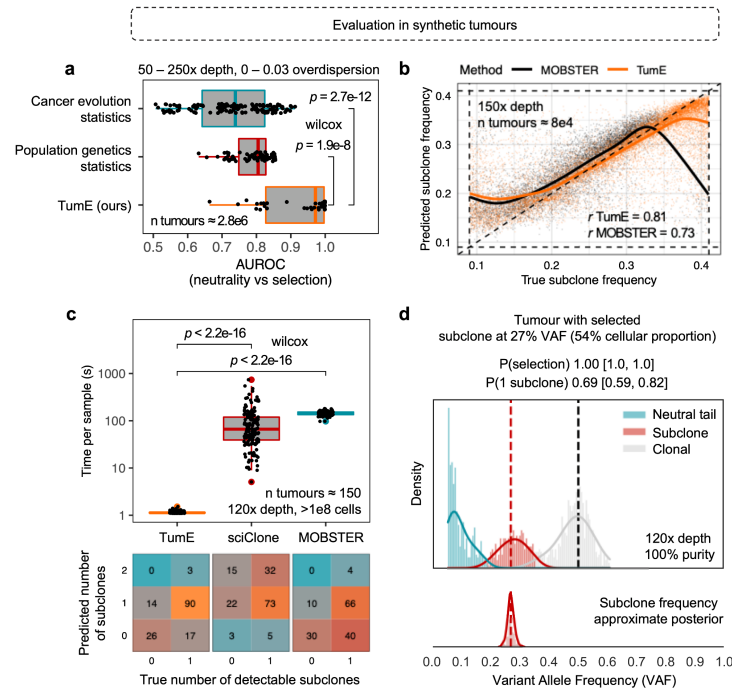


Figure 2 (a) In a cohort of 2.8 million synthetic tumours, TumE outperformed all existing common population genetic^{20,21} and cancer evolution^{7,12} specific summary statistics when differentiating between positive selection and neutral evolution, based on AUROC (two-sided Wilcoxon test). **(b)** Further, for predicting the true frequency of selected subclones, TumE provides comparable or better performance relative to the current state-of-the-art mixture model MOBSTER¹⁶ that properly accounts for neutral dynamics in tumour populations. The panel shows correlation coefficient (r) between the true and predicted subclone frequency in 80,000 synthetic tumours sequenced at 150x mean sequencing depth. **(c)** In an orthogonal dataset of 150 synthetic tumours¹⁶ with either 0 or 1 detectable subclones, TumE was significantly faster at estimating the number of subclones (two-sided Wilcoxon test) than existing mixture model based methods sciClone²⁴ and MOBSTER¹⁶ (measured in inference time per sample). In addition, only TumE and MOBSTER consistently identified the correct number of subclones, as both methods directly account for the neutral dynamics observed in tumour populations. **(d)** TumE estimates in a synthetic tumour sequenced at 120x mean sequencing depth and a subclone at 54% cellular fraction.

169 Given our simulation framework was based on certain approximating assumptions to improve
 170 computational speed and efficiency (namely small population size and no cell death; outlined in Methods),
 171 we sought to perform additional validation of evolutionary estimates in an alternative dataset of synthetic
 172 tumours¹⁶. The orthogonal dataset, described in Caravagna et al. 2020¹⁶, consisted of 150 synthetic
 173 tumours, 40 effectively neutral and 110 with one detectable subclone (between 10 - 45% VAF), sequenced
 174 to 120x depth and grown to a population size of $>10^8$ cells at birth rate of 1 and death rate of 0.2. To frame
 175 our predictions relative to existing methods, we applied TumE, MOBSTER, and a variational Bayesian
 176 mixture model sciClone²⁴ to the synthetic dataset. To make comparisons fair, we limited the maximum

177 number of subclonal cluster assignments to 2 for both MOBSTER and sciClone, as this was the upper
178 bound on TumE estimates (Methods). For sciClone, this meant setting the maximum number of mixture
179 components to 4 (neutral tail, 2 subclones, and a clonal peak) as sciClone doesn't properly account for
180 neutral dynamics (Pareto tail) observed in sequenced tumour populations. Both TumE and MOBSTER
181 consistently identified the correct number of detectable subclones in the majority of cases while sciClone
182 systematically overestimated the number of subclones, even after correcting estimates for the clonal peak
183 and neutral tail (Figure 2C). However, relative to both sciClone and MOBSTER, TumE provided orders-of-
184 magnitude faster estimates (two-sided Wilcoxon test, $p < 2.2 \times 10^{-16}$, Figure 2C), reducing run times per
185 sample from minutes to ~1 second. We provide individual estimates with TumE for each of the 150 synthetic
186 tumours, and an additional 750 synthetic tumours of variable sequencing depth from ref¹⁶, in Supplementary
187 Figures 13 & 14. We provide an example TumE output for a synthetic tumour with a single detectable
188 subclone in Figure 2D.

189
190 In this study, we note that the birth and death rate were set to fixed values (birth rate = $\log(2)$, death rate
191 = 0, in line with ref⁷) to additionally improve the computational efficiency of the stochastic simulations of
192 positively selected tumour populations. Therefore, one additional factor that may impact the accurate
193 detection of selection and subclones with TumE is variable birth and death rates in growing tumours. For
194 example, an elevated cell death can lead to an increase in the number of passenger mutations that are
195 swept to higher frequencies during subclonal selection. In regards to the VAF distribution, this elevated
196 number of mutations 'trailing' the subclonal peak may obscure lower frequency subclones or, alternatively,
197 lead to spurious identification of additional subclones due to an elevated number of neutral mutations
198 entering the subclonal frequency range. To assess the impact of variable growth rates, we generated an
199 additional 6 million synthetic tumours across 26 different birth and death rate combinations (simulation
200 parameters outlined in Supplementary Table S1). Overall, we find that our estimates are robust to changes
201 in tumour growth rates. Any errors that do occur only appear to increase the number of parsimonious
202 explanations of the data (e.g. classifying 2 subclones as 1; Supplementary Figure 15). In addition, the
203 prediction of subclone frequency also remained consistent across all the birth and death rate combinations
204 evaluated here (Supplementary Figure 16).

205

206 **Analysis of whole-genome and exome sequenced tumour biopsies**

207 To make the utility of synthetic supervised learning concrete, we first evaluated TumE in 'gold-standard'
208 tumour biopsies commonly used to evaluate mixture model based approaches, namely a deep sequenced
209 (~320x coverage, 90.7% purity) acute myeloid leukemia (AML) sample from Griffith et al.³¹ and a deep
210 sequenced (~226x coverage, 71.2% purity) breast adenocarcinoma sample retrieved from the pan-cancer
211 analysis of whole genomes (PCAWG)¹¹ but originally from ref³². In both cases, we recovered the correct
212 evolutionary mode, number of subclones, and subclone frequencies (Figure 3A & 3B). In addition, because
213 we provide accurate subclone frequency estimates, we performed heuristic clustering of the clonal,

214 subclonal, and neutral tail mutations by using the expected variance under a binomial sequencing noise
 215 model (Methods). This heuristic approach facilitates subclonal clustering at almost zero additional
 216 computational cost (as observed in the total runtime per sample of ~1s).

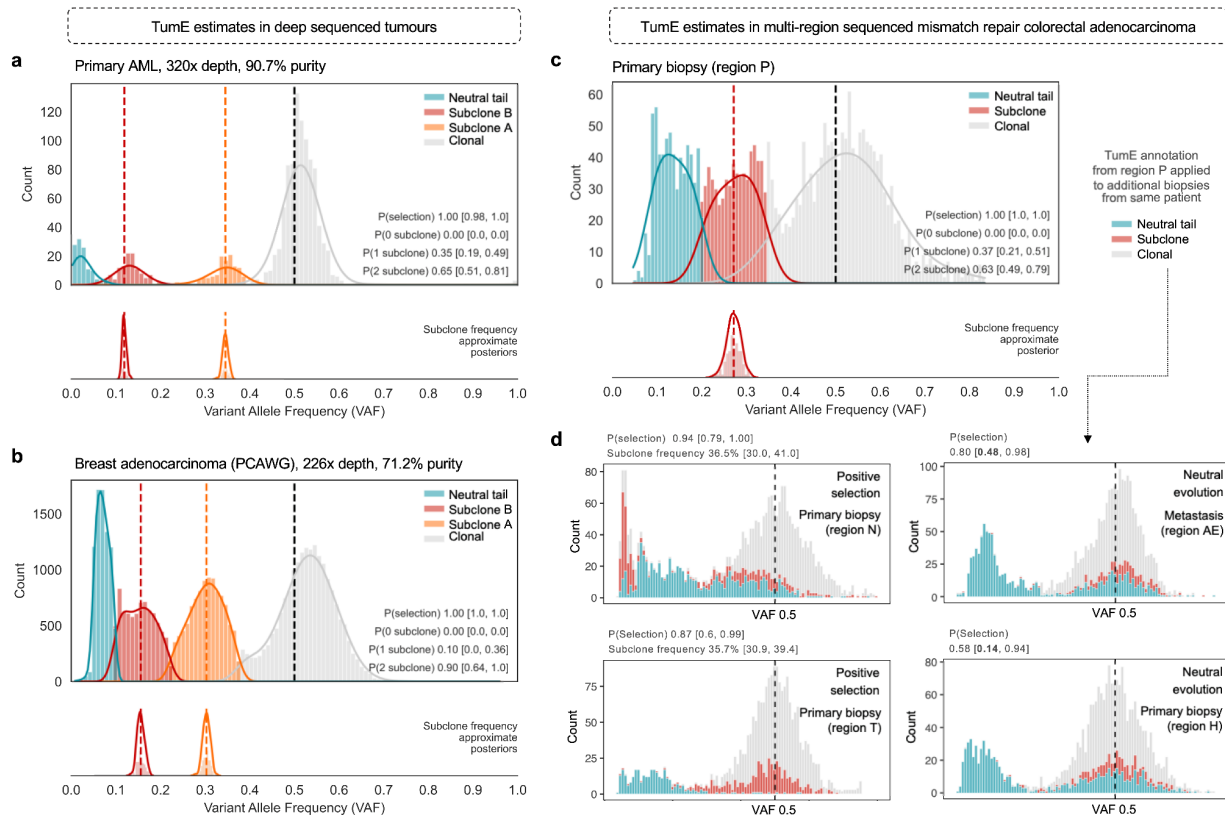


Figure 3. TumE estimates in deep whole-genome or whole exome sequenced tumour biopsies. **(a)** A deep-sequenced primary acute myeloid leukemia (AML) sample from Griffith et al.³¹. TumE estimated two subclones, a neutral tail, and a clonal peak. P(Selection) indicates the probability of selection. P(0, 1, 2 subclone) indicates the probability estimate for the number of subclones. Each probability estimate is provided with the 89% equal-tailed interval generated from 50 Monte Carlo dropout samples. A sample is labeled positive selection if the lower bound of the 89% interval is above $P = 0.5$, and the number of subclones is assigned to a sample if the lower bound of the 89% interval is greater than 0.5 (Methods). Subclone frequency estimates are shown with the complete approximate posterior. **(b)** A deep-sequenced breast adenocarcinoma from the pan-cancer analysis of whole genomes¹¹ (PCAWG). TumE estimated two subclones, a neutral tail, and a clonal peak. **(c)** We applied TumE to a single mismatch repair deficient (MMR) gastro-esophageal tumour sequenced across 5 spatially distinct regions. We first identified an intermediate frequency subclone in region P with TumE. **(d)** Under the hypothesis that TumE could reveal the fixation process of region P subclones in other regions, we annotated each of the remaining regions with the clonal, subclonal, and neutral tail mutations identified in region P. We identified ongoing subclonal selection in 2 out of the 4 remaining regions (N and T) consistent with an increase in frequency of subclonal and neutral tail mutations from region P. In cases where neutral evolution was the most parsimonious explanation, we observed complete fixation of the subclonal region P mutations (region AE and H).

217 We next evaluated TumE in whole-exome sequenced (WES) mismatch repair deficient (MMR) gastro-
 218 esophageal tumours biopsied across multiple spatially distinct regions (collected from von Loga et al.³³). As
 219 evolutionary inference requires high-quality genomes, we only included samples that had a mean effective
 220 coverage (mean sequencing depth * purity) greater than 50x and a minimum purity of 50%. We note that
 221 ~70x mean sequencing depth has been suggested as the minimal threshold for accurate estimates^{7,16}, as

222 we also observed (Supplementary Figures 8 & 9). Following removal of low quality biopsies, we retained
223 biopsies from two tumours with one tumour retaining 5 spatially distinct (WES) biopsies. TumE estimates
224 in the 5 spatially distinct biopsies from a single tumour revealed the fixation process of a positively selected
225 subclone, from intermediate frequency to metastasis fixation (Figure 3C & 3D). In addition, the application
226 of TumE to multi-region samples highlighted the ability of TumE to pick up signatures of selection not directly
227 encoded in distinct subclonal peaks but in the asymmetry of the diploid heterozygous cluster (region N &
228 region T, Figure 3D).

229
230 Finally, we evaluated TumE in 85 whole-genome sequenced (WGS) tumour biopsies, spanning 8 different
231 cancer types, retrieved from the pan-cancer analysis of whole genomes (PCAWG). In total, 38.8% of
232 samples showed evidence for positive, or subclonal, selection whereas the majority, 61.2%, were
233 adequately described by neutral evolutionary dynamics (Supplementary Table S2). Alternative methods
234 applied to large cancer cohorts, including PCAWG, have estimated that as few as 3%¹⁶ to upwards of 96%³⁴
235 of samples show evidence for ongoing subclonal selection. The discrepancy is likely explained in modeling
236 approaches. For example, low estimates are a by-product of utilizing mixture models that rely on distinct
237 and 'clean' subclonal peaks whereas high estimates likely occur from not taking into account the neutral
238 dynamics in tumour evolution. In contrast, TumE generates a non-linear encoding of the VAF distribution,
239 extracting novel representations that increase accuracy while simultaneously accounting for the correct
240 neutrality evolutionary dynamics observed in tumour populations. All samples analyzed in this study,
241 including MMR gastro-esophageal and deep-sequenced AML, are outlined in Methods and Supplementary
242 Figure 16 - 18.

244 **A transfer learning framework to infer additional evolutionary** 245 **parameters**

246 One drawback of simulation-based deep learning approaches is the requirement for the repeated
247 generation of synthetic data for training. Although this allows for fast inference at test time through
248 amortization, altering the models assumptions or changing the parameters being inferred generally requires
249 simulating a completely new set of data and training an entirely new set of models - a computationally
250 expensive process. Practically, overcoming this limitation would provide substantial reductions in the
251 amount of time and data needed to build accurate models and would make simulation-based approaches
252 more accessible to the general user. Therefore, we hypothesized that our trained deep learning models
253 could be used as a source of 'stored' knowledge for related evolutionary inference tasks that also used the
254 VAF distribution as input.

255
256 To explore this possibility, we implemented a transfer learning pipeline, based on domain adaptation^{35,36},
257 to make inferences on additional parameters using a previously developed cancer evolution simulator,
258 TEMULATOR³⁷, that was built under a modified set of assumptions relative to our multiplicative fitness

259 framework (Methods, viable parameter combinations for detectable subclones outlined in Supplementary
 260 Figure 20). In this study, we employ open set domain adaptation³⁶ where the structure of the input space,
 261 i.e. the VAF distribution, is retained whereas the outputs, the evolutionary tasks, are modified. Briefly, this
 262 pipeline involved generating new synthetic tumour sequencing data using TEMULATOR, performing

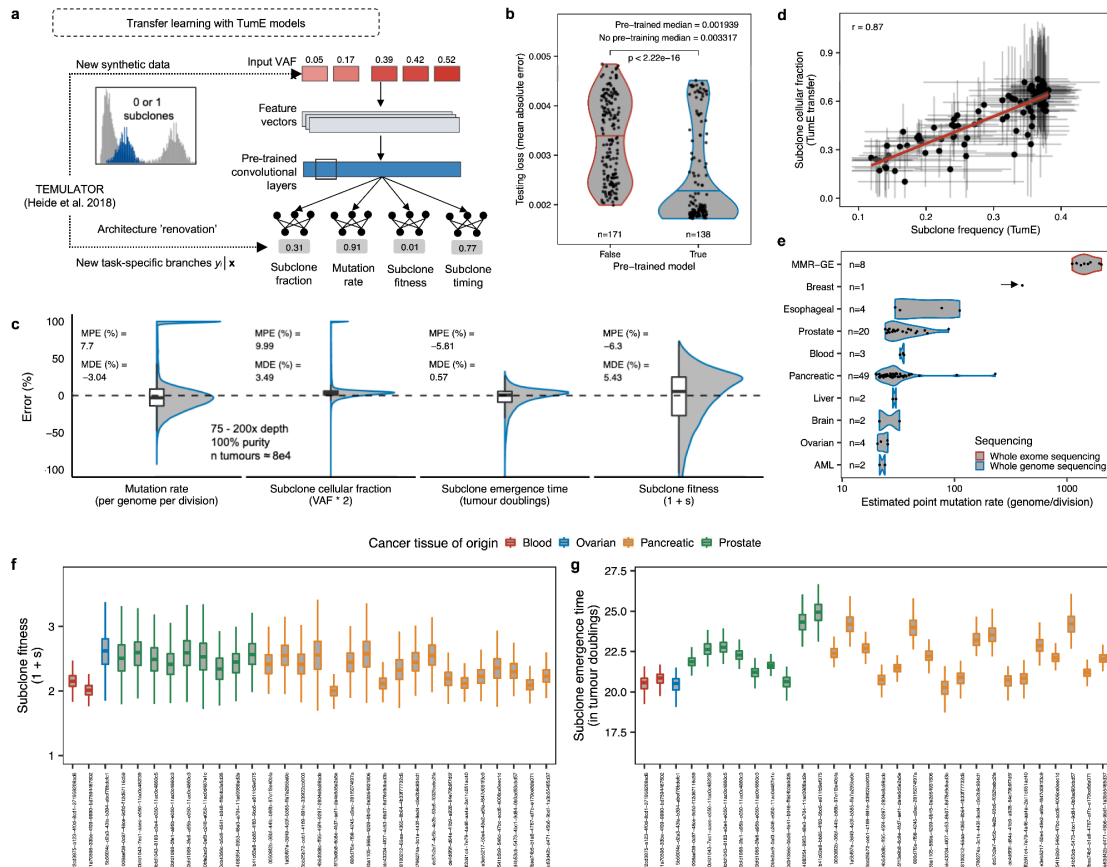


Figure 4. (a) Transfer learning approach utilizing 'renovated' pre-trained neural networks for alternative evolutionary inference tasks in tumour cellular populations. TEMULATOR is an alternative cancer evolution simulator that generates synthetic tumour sequencing data by deterministically initiating subclones at user specified fitnesses and time points³⁷. (b) Pre-trained models provide significant reductions in testing loss, over non-pretrained models, when updating neural network weights on reduced dataset size of 500,000 synthetic VAF distributions (~1.25% of the total dataset size used to originally train TumE). (c) TumE transfer (TumE-T) effectively recovers evolutionary parameters from TEMULATOR simulations (75 - 200x mean sequencing depth, 100% tumor purity) with mean and median percentage errors less than 10% in all cases. A full description of performance across variable sequencing depths, mutation rates, and subclone frequencies is provided in Supplementary Figure 23. (d) We find consistency between the subclone cellular fraction estimated by TumE-T and the subclone frequency (cellular fraction / 2) estimates generated from TumE, indicating nearly identical tasks are easily transferred through pre-training. (e) Per genome per division mutation rate estimates in 95 WES and WGS samples from von Loga et al.³³ (MMR-GE = mismatch deficient repair gastro-esophageal cancer), Griffith et al.³¹ (AML = acute myeloid leukemia), and PCAWG¹¹. (f) Subclone fitness ($1 + s$) estimates (relative growth rate advantage of subclone over background population) and (g) subclone emergence time estimates in 30 tumour biopsies identified with 1 subclone in the PCAWG data. Subclone fitness and emergence time estimates were scaled to a final tumour population size of 10^{10} cells, similar to ref⁷. PCAWG sample identifiers are provided on the x-axis. Boxplots for subclone fitness and emergence time indicate median estimate and 1.5x interquartile range (whiskers) over 500 Monte Carlo dropout samples from TumE-T.

263 architecture 'renovation' on pre-trained TumE neural networks to replace existing task-specific branches
 264 with new ones, and re-tuning the neural network weights and hyperparameters for optimization on the new

265 evolutionary inference tasks (Figure 4A). The evolutionary inference tasks included predicting subclone
266 fitness, subclone emergence time, mutation rate, and subclone cellular fraction (subclone frequency * 2).
267 To highlight the benefit of using pre-trained models on related evolutionary inference tasks, we opted to
268 update network weights with only 500,000 synthetic VAF distributions, representing only a fraction (~1.25%)
269 of the data used in the original training of TumE. Each VAF distribution was generated by simulating
270 synthetic tumours with TEMULATOR at a birth rate of 1, death rate of 0.2, final population size of $\sim 10^4$,
271 and with either 0 or 1 detectable subclone. The remaining parameters, such as mutation rate, were
272 uniformly sampled from empirically plausible ranges (Supplementary Table S3).

273
274 Initially, we used the 500,000 synthetic VAF distributions to compare pre-trained vs non-pretrained models
275 for predicting the evolutionary and subclonal parameters in the presence of 1 subclone. To ensure valid
276 comparisons, we performed a random hyperparameter search, tuning the learning rate and number of fully
277 connected layers in the new task specific branches, across both the pre-trained and non-pretrained model
278 groups. Both groups shared identical neural network architectures. When initially evaluating ~ 300 pre-
279 trained and non-pretrained models on an external test set of 3000 synthetic tumours, we found that
280 pretrained models obtained a significantly lower testing loss (mean absolute error across all tasks, two-
281 sided Wilcoxon test, $p < 2.22 \times 10^{-16}$, Figure 4B). Further, when evaluating the top performing pre-trained
282 and non-pretrained models on an additional 400,000 synthetic tumours, pre-trained models obtained
283 significantly lower mean percentage errors, relative to non-pretrained models, for predicting the mutation
284 rate, subclone emergence time, subclone fitness, and subclone frequency (two-sided Wilcoxon test, $p <$
285 1.7×10^{-8} on all tasks, Supplementary Figure 21).

286
287 Next, we selected the top performing pretrained model, TumE transfer (TumE-T), for further validation. We
288 initially found a modest yet systematic underestimation of the mutation rate (~50% mean percentage error).
289 However, this was easily corrected with a post-hoc adjustment by re-fitting the predicted mutation rate to a
290 set of 1000 synthetic tumours using polynomial regression (degree = 2). Evaluating the updated mutation
291 rate estimates on a holdout set of 100,000 synthetic tumours validated the post-hoc adjustment
292 (Supplementary Figure 22). Overall, we were able to effectively recover all evolutionary parameters in the
293 100,000 synthetic tumours with mean and median percentage errors lower than 10% in all cases (Figure
294 4C). The performance was also consistent across sequencing depths and mutation rates, however, as
295 expected, we could only effectively assign subclonal parameters, such as fitness, at detectable subclone
296 frequencies (~10 - 40% VAF; Supplementary Figure 23).

297
298 Applying TumE-T to the 95 WGS and WES samples described above, we found strong correlation between
299 predicted subclone cellular fraction and subclone frequency estimated by the original TumE models,
300 suggesting nearly identical tasks are easily transferred to new source-target distributions when using pre-
301 trained models (Figure 4D). With respect to mutation rates, estimates were consistent with the general

302 trends observed empirically^{11,33} - with mismatch repair deficient tumours showing extremely high mutation
303 rates (>100 per genome per division) and acute myeloid leukemia showing very few (Figure 3E). For
304 subclone fitness and subclone emergence time estimation, we had to take into account the difference
305 between simulated and true population sizes^{7,19}. In this regard, we rescaled our estimates to account for a
306 true tumour population size of 10^{10} , similar to ref⁷. With rescaling, TumE-T subclone fitness estimates,
307 defined as the relative growth advantage of the selected subpopulation over the background population,
308 ranged from ~1.9 to 2.6 (Figure 4F) while subclone emergence time estimates ranged ~20 to 24 tumour
309 doublings (Figure 4G) in samples with ongoing subclonal selection. We note that emergence times of ~20
310 to 24 tumour doublings represent approximately 0.001 to 0.16% of the final tumour volume, which is
311 consistent with theory and empirical evidence suggesting subclones must arise early during tumour growth
312 to reach detectable frequencies^{7,14,17,38}.

313

314 Discussion

315 In this study, we developed a synthetic supervised learning approach, TumE, for cancer evolutionary
316 inference. Overall, the synthetic supervised learning approach, TumE, provides four major advantages.
317 First, by generating synthetic data using models of cancer evolution, we are able to explicitly account for
318 the neutral and non-neutral evolutionary dynamics observed in tumour VAF distributions^{7,12,16}, thereby
319 avoiding systematic overestimates in the number of subclones due to misclassifying low frequency neutral
320 'tails'. Second, by using neural networks that can naturally handle high-dimensional VAF distributions as
321 input, we avoid information loss that comes with compressing data into a single statistic, or distance metric,
322 prior to inference, improving model accuracy across all evolutionary inference tasks considered here. Third,
323 by separating simulation and model training from prediction, via amortized inference, we significantly
324 decrease inference time per sample, reducing time from minutes to seconds relative to existing methods.
325 Finally, we show how we can use open set domain adaptation^{35,36}, a form of transfer learning, to recycle
326 our models for alternative evolutionary inference tasks that use VAF distributions as input - drastically
327 reducing the number of synthetic samples and computational time required for further model development.
328 Our library of pre-trained models benefits all researchers building inference machines for cancer evolution
329 prediction, even in a non-deep learning setting. For example, providing fast, evolutionary-informed peak
330 initializations for mixture model based methods.

331

332 We mention some current limitations. Firstly, as a neural network requires optimization on a finite, static set
333 of data, estimates, without transfer learning, are constrained to a pre-defined search space. In this study
334 we focused on cancer evolution in the context of 2 detectable selected subpopulations captured from
335 frequency information in diploid genomic regions. Although multiple studies have shown it's rare to detect
336 2, or even 1, subclones^{7,14,16,17} in noisy one-dimensional VAF distributions, it's possible we do not capture
337 extreme cases of selected subclonal heterogeneity. Furthermore, focusing on diploid regions may obscure
338 the detection of ongoing selection if mutations are concentrated in copy number aberrated segments.

339 However, constraining analyses to diploid regions provides a strong baseline for model development, while
340 genome-wide linkage provides biological justification for analyzing diploid regions. Finally, our model of
341 tumour evolution was structured to reflect the biopsy material available here, namely bulk sequenced single
342 site and time point data. We note that tumour growth over space and time can have profound effects on the
343 detectability of selection^{39,40}. In this regard, TumE estimates can still be applied in a localized setting and
344 aggregated globally. Nevertheless, more structured ways of integrating a synthetic supervised learning
345 approach with multi-region data are necessary for maximizing utility.

346
347 Altogether, in this study, we exhibit how coupling well-specified synthetic data with neural networks provides
348 fast and accurate amortized estimates that go beyond the current paradigm of single statistics, mixture
349 models, and approximate Bayesian computation for classifying and quantifying ongoing selection in tumour
350 populations. The integration of generative and simulation-based models of cancer evolution with modern
351 deep learning frameworks facilitates robust and efficient estimates of evolutionary and subclonal dynamics
352 in growing tumour populations. This extensible framework provides future avenues for harnessing
353 progressive computational methods for the benefit of cancer genomics and, as an end goal, the cancer
354 patient.

355

356 **Methods**

357 **Synthetic data generation**

358 We generated synthetic data that encoded the evolutionary dynamics observed in the variant allele
359 frequency (VAF) distribution (namely the neutral tail, subclones, and clonal peaks) using two
360 complementary approaches dependent on the underlying evolutionary mode - one for tumors subject to
361 positive selection and one for tumours evolving neutrally.

362

363 For tumours simulated under positive selection, we utilized a well-established framework of cancer
364 evolution that models exponential tumour growth under a stochastic branching process^{7,12,13,15,19,30} and
365 coupled this with a virtual biopsy procedure to account for sequencing noise observed in whole-
366 genome/whole-exome sequenced tumours from real patient tumours. For implementation, we adapted a
367 previous cancer evolution framework developed by Williams et al.⁷ Briefly, this model simulates
368 exponentially growing tumour populations under a stochastic branching process using a rejection-kinetic
369 Monte Carlo (MC) algorithm, where a given cell accrues mutations at some Poisson-distributed per
370 genome per division rate μ and divides or dies with probabilities proportional to its birth or death rate. This
371 branching process continues by randomly sampling existing cells, weighted by cellular fitness, until a final
372 tumour population size N , sufficient to recapitulate the features of empirical VAF distributions, is reached.
373 Following completion of each simulation, a virtual biopsy procedure to account for sequencing noise
374 observed in real patient VAF distributions is implemented. In this sequencing noise model, the observed

375 frequency for a given mutation (VAF_{obs}) relative to the true underlying frequency (VAF_{true}) in a tumour of
376 population size N is given by

377

378
$$VAF_{obs} = R_{obs} / D_{obs}$$
 where

379

380
$$D_{obs} \sim Bin(n = N, p = \frac{D}{N}), R_{obs} \sim BetaBin(n = D_{obs}, p = VAF_{true}, \rho)$$

381

382 where D total indicates the total observed read depth, R indicates the number of observed reads covering
383 the mutation, VAF_{true} indicates the true population frequency of the mutation, and ρ indicates the
384 overdispersion parameter for the beta-binomial.

385

386 In this study, we modify the Williams et al⁷ framework in two ways. Firstly, we implement a fully stochastic
387 arrival of subclones (driver mutations) rather than deterministically injecting a subclone with a specified
388 fitness at a given time t . Secondly, the fitness of a subclone or cell is dictated by the multiplicative fitness
389 of all driver mutations. Therefore, when a driver mutation does occur, based on some probability p_d , it is
390 assigned a selection coefficient $s > 0$ sampled from an exponential distribution which increases the cell's
391 growth rate ($b - d$) by a factor of $(1 + s)$ i.e. the fitness. In the case of multiple driver mutations, the fitness
392 of a given cell increases multiplicatively i.e. $\prod(1 + s)$. Although this random injection of driver mutations is
393 more computationally intensive, it implicitly captures a wider variety of potential frequency distributions
394 without hard coding additional parameter settings - for example, when additional subclones, beyond 1 or
395 2, are present at undetectable frequencies (e.g. >40% or <10%). In this study, we consider up to 2
396 detectable subclones but allow for up to 3 selected subclones to be present at the time of biopsy (see
397 Simulation Parameter Selection below for more details).

398

399 For tumours simulated under neutral evolution, we use a generative sampling process for producing
400 neutral VAF distributions, rather than using the stochastic simulation framework. We implement this
401 sampling process because we use a small N population size approximation to generate VAF distributions
402 in our stochastic simulations (using a small N allows us to increase simulation speed and efficiency, which
403 makes generating millions of synthetic VAF distributions practically feasible). Although using a small N is
404 reasonable since the VAF distribution contains no information on population size⁷ (a final simulated
405 tumour population size N of $10^3 - 10^4$ has been shown to be sufficient to recapitulate the properties of
406 empirical VAF distributions⁷), neutral stochastic simulations have a higher probability of returning late-
407 occurring spurious subclones due to chance or, in empirical terms, genetic drift. Given the quality of the
408 synthetic data impacts deep learning model performance, we utilize the fully synthetic generative
409 sampling scheme to avoid misspecified data relative to the expected null model of neutral evolution.

410

411 The neutral generative sampling process we implement is based on the observation that neutrally evolving
412 asexual, non-recombining populations, such as cancer populations, have VAF distributions (excluding
413 clonal mutations) that follow a power-law or Pareto distribution^{16,29}. Therefore, a VAF for any mutation i in
414 the neutral tail of a frequency distribution can be realized by sampling

415

$$416 \quad VAF_i \sim f(x | \alpha, m) \text{ with } f(x | \alpha, m) = \alpha m^\alpha x^{-(\alpha+1)}$$

417

418 where α is the shape parameter and m the scale parameter for the Pareto distribution.

419

420 Given that the generative process for neutral tails is known^{16,29}, if we have empirically valid shape and scale
421 parameters that define the Pareto distribution, we can generate realizations of neutral allele frequency
422 distributions that are well-specified. Previously, Caravagna et al¹⁶ fit Pareto distributions (and beta
423 distributions) to thousands of patient tumours extracting both shape and scale parameters. We used these
424 Pareto distribution fits from diploid regions of patient tumours with greater than 50x sequencing coverage
425 to build sampling distributions for the shape and scale parameters. We then used these sampling
426 distributions to generate allele frequencies under a Pareto distribution and, in addition, randomly assigned
427 clonal mutations to each neutral synthetic VAF. In practice, as previously noted¹⁶, the scale parameter can
428 be set to the minimum observed frequency as this is the maximum likelihood estimate for the Pareto
429 distribution.

430

431 We note that we added additional noise to synthetic neutral distributions to better account for variability
432 observed in empirical data in two ways. Firstly, for any synthetically generated neutral distribution, we
433 randomly trimmed the low frequency neutral tail at a frequency f (10 - 30% VAF) with some probability P_{trim}
434 (≤ 0.1). We perform this step as many VAF distributions observed in patient biopsies lack the characteristic
435 neutral tail, even at high sequencing depth¹⁶. By randomly trimming neutral synthetic VAF distributions, we
436 tend to more parsimonious explanations of the data, with respect to positive selection, when assessing
437 incomplete and potentially noisy VAF distributions. Secondly, we randomly shifted the heterozygous, diploid
438 clonal peak (that should be centered at 50% VAF) to between 45 and 50% VAF. We perform this random
439 perturbation of the clonal peak to avoid overestimating positive selection when patient samples have
440 incorrect tumour purity estimates that may have led to spurious elevation in the number of subclonal
441 mutations.

442

443 Finally, to ensure positively selected and neutrally evolving tumours were not out of distribution with each
444 other given the alternate data generation approaches, we built an aggregate simulation framework that
445 generated neutral and positive synthetic tumours in pairs - assigning the neutral VAF distributions with
446 parameter-matched sequencing noise and equivalent clonal and non-clonal mutations with respect to the
447 paired positive selection simulation.

448

449 The synthetic data generation algorithms are outlined in supplementary, and code is available at
450 <https://github.com/tomouellette/CanEvolve.jl>.

451

452 **Simulation parameter selection**

453 Each stochastic simulation described above was parametrized by the mutation rate (per genome per
454 division), the probability a mutation was a driver, the mean for the exponential selection coefficient
455 distribution, the number of clonal mutations in the founder cell, the maximum number of driver mutation
456 events, the final tumour population size, the sequencing depth, and the sequencing overdispersion. We
457 chose simulation parameters that were consistent with previous studies^{7,13,16,30} and that captured the
458 expected qualitative and quantitative attributes of VAF distributions observed empirically (Supplementary
459 Table S1). All non-fixed parameters were uniformly random sampled during the development of the
460 synthetic datasets. To improve computational speed and efficiency in our stochastic simulations, we used
461 similar simulation approximations as ref⁷. Namely, a small N population size approximation (where we
462 simulated data using a final population size of 10^3) and a fixed growth rate (where the birth rate was set to
463 $\log(2)$ and the death rate was set to 0). In addition, as we were focused on differentiating between neutral
464 evolution and selection at effective sequencing depths of $\sim 50 - 250x$, we constrained our search space to
465 1 or 2 detectable subclones present between 10 - 40% VAF. We implemented this constraint as (i) it is
466 extremely rare to detect 3 subclones in a one-dimensional VAF distribution as each subclone has to be $>5-$
467 10% VAF (10-20% cellular fraction) for detection, (ii) most frequency-based methods or studies show limited
468 evidence for detecting >1 subclone at 50 - 250x coverage in a single time point, one-dimensional VAF
469 distribution¹⁶, and (iii) below greater than roughly 10% VAF subclones merge with the neutral tail and above
470 roughly 40% VAF subclones begin to merge with the clonal peak when considering diploid regions.

471

472 **Synthetic supervised learning**

473 As outlined in the results, synthetic or simulation-based deep learning has been shown to be equivalent to
474 amortized approximate inference under a generative model²⁸. Therefore, by optimizing a neural network
475 using synthetic VAF distributions sampled from a stochastic generative process $p(\mathbf{x}, \mathbf{z} | \theta)$ (i.e. the synthetic
476 data generation scheme defined above), we can build fast approximate inference models for evolutionary
477 inference. We describe the synthetic supervised workflow from feature generation to prediction below.

478

479 **Input representation.** For each simulation, we converted mutation frequencies into a VAF distribution
480 (histogram) of length k that had a fixed range between 2% and 50% VAF. To implicitly condition our model
481 on mean sequencing depth (readily available from sequenced tumour biopsies), we only included mutations
482 above a frequency cutoff based on the variance of a binomial sequencing noise model. We note that this
483 conditioning step is arbitrary and simply acts as a way to improve model optimization during training. In
484 general, a simple approach to conditioning a neural network on a measurable variable involvings finding a

485 reasonable encoding within the feature representation. For example, an alternative approach instead of
486 using a frequency cutoff would be to concatenate the sequencing depth to the input feature vector. Overall,
487 each input feature vector was created by counting mutations into k bins where each bin had a width w of
488 $(50 - 2\% \text{ VAF}) / k$ and a lower frequency cutoff defined by $f_{alt} + (2\sqrt{f_{alt}c[1 - f_{alt}]}) / c$ where f_{alt} is the
489 minimum alternative reads to call a mutation divided by mean sequencing depth and c is mean sequencing
490 depth. For all model development and training, we generated and concatenated two feature vectors with k
491 = 64 and $k = 128$ for each simulation to capture varying levels of information depending on the sparsity of
492 mutations in a given synthetic tumour.

493
494 **Model search.** We initially developed neural networks for three single or multi-task inference problems: (i)
495 evolutionary mode (neutral evolution or positive selection) and number of subclones classification (M_{ms}), (ii)
496 frequency prediction for a single subclone (M_{1s}), and (iii) frequency predictions for two subclones (M_{2s}). For
497 each multi-task, we performed a random hyperparameter search using a one dimensional (1D)
498 convolutional neural network (CNN) with task-specific fully connected branches as a base architecture. For
499 the random search, the hyperparameters included the number of convolutional layers (1 - 20) the task-
500 specific branch type (fully connected or global average pooling), the number of feature maps/channels for
501 each convolutional layer (4 - 32), the convolutional kernel width for the left trunk, right trunk, and task-
502 specific branches (1 - 17, odd), the learning rate (10^{-7} - 10^{-3}), and the patience for early stopping (3 - 5). To
503 tend toward higher precision and lower recall for predicting selection, we also tuned a penalty term on the
504 positive class in the binary cross entropy loss. Batch size was fixed to 256. Hardswish activations were
505 used at each hidden layer. Dropout, fixed at a probability of 0.5, was added after each layer to allow for
506 downstream application of uncertainty estimation (see Uncertainty Estimation below). We note that we also
507 explored inferring subclone emergence time under a multiplicative fitness model, but could not effectively
508 recover parameters likely due to a complex non-linear relationship between subclonal fitness and
509 emergence time. However, we provide these estimates as an ‘experimental’ output in the TumE python
510 package (links below).

511
512 **Model training.** We trained over 150 models for each evolutionary inference task(s) using an Adam
513 optimizer, minimizing the cross-entropy loss for classification tasks (M_{ms}) and the L1 loss for regression
514 tasks (M_{1s} and M_{2s}), on approximately 40 million synthetic tumours simulated with parameters outlined in
515 Supplementary Table S1. For training, each batch consisted of 20,000 unique simulations and training was
516 stopped after 4 epochs or when early stopping, updated after each batch, was activated based on specified
517 patience. To avoid overrepresentation of any subclone frequency during training, we re-balanced positive
518 selection simulations before each batch to have an equal number of subclones at each frequency up to two
519 decimal places (e.g. 0.11 or 11% VAF). For two subclone simulations, we re-balanced simulations based
520 on the distance between subclones ($f_{subclone2} - f_{subclone1}$) and only included simulations where the distance

521 between subclones was >5% VAF. Note that only positive selection simulations were used to train M_{1s} and
522 M_{2s} .

523 **Model selection.** Using an independent test set of one hundred thousand simulations, we then selected
524 the top models across each multi-task for further validation. For M_{ms} , we selected models that maximized
525 the mean accuracy across the evolutionary mode, $P(\textit{Selection})$, and number of subclones, $P(N \textit{ subclones})$,
526 classification, and favoured models that assigned a larger penalty term to misprediction of positive selection
527 (i.e. a lower weight to the positive class in the binary cross entropy loss). For the regression models M_{1s}
528 and M_{2s} , we selected models that minimized the mean absolute error between the true and predicted
529 subclone frequency on the test set while also ensuring that predictions properly extrapolated across the
530 entire simulated parameter range (e.g. ~10 - 40% VAF for subclone frequencies).

531 **Uncertainty estimation.** To capture model-based uncertainty in our estimates, we implemented a form of
532 Bayesian approximation for deep learning called Monte Carlo (MC) dropout^{26,27}. Conceptually, MC dropout
533 captures model-based uncertainty by taking advantage of the relationship between model averaging and
534 standard dropout - a network with dropout at every layer encodes 2^n possible network configurations. By
535 keeping dropout on at test time, each prediction is a stochastic pass through a set of randomly activated
536 neurons. More specifically and with a slight abuse of notation w.r.t to ref²⁶ ignoring the variational notation,
537 we make estimates of our target variable y (e.g. subclone frequency) by performing T stochastic forward
538 passes through the network and averaging, $E(y)$, the results:

540

541

$$E(y) \approx \frac{1}{T} \sum_{t=1}^T \hat{y}(x, W_1^t, \dots, W_L^t)$$

542

543 where \hat{y} is the output with respect to the input data x for a neural network with L layers, and W corresponds
544 to a weight matrix for each layer L . For every stochastic pass, each W is assigned a randomly sampled
545 vector of Bernoulli random variables such that each individual neuron is inactivated with a probability equal
546 to the dropout rate. Under this framework, MC sampling over T stochastic passes through the network
547 generates an approximate posterior for our target variables with respect to the input data.

548

549 **Making predictions.** For differentiating between neutral evolution and positive selection, $P(\textit{Selection})$, and
550 predicting the number of subclones, $P(N \textit{ subclones})$, in both synthetic and real patient tumours, we took a
551 conservative, more parsimonious approach to prediction by considering the variance in the approximate
552 posterior. For $P(\textit{Selection})$, we only called positive selection if the lower bound of an 89% equal-tailed
553 interval for the approximate posterior, computed across 50 stochastic passes through M_{ms} , was greater
554 than 0.5. If the lower bound was less than 0.5, we called neutrality and zero subclones, independent of the
555 result of $P(N \textit{ subclones})$. We show the utility of this strategy for mitigating model overconfidence in a
556 synthetic toy example (Supplementary Figure 6). For the regression tasks of predicting subclone frequency

557 and emergence time, we estimated the true value by performing 50 stochastic passes through the networks
558 and averaging the results, while also providing the complete approximate posterior. We describe additional
559 considerations for making estimates in real patient tumour biopsies below.

560

561 All model development and training was done using *pytorch v1.8.1*. We provide a python package, scripts,
562 and all trained neural network models for downloading, use, and modification at
563 <https://github.com/tomouellette/TumE>.

564

565 **Model performance in synthetic tumour sequencing datasets**

566 We simulated or collected 3 different datasets of synthetic tumour sequencing data to study the
567 performance of TumE under changing parameter regimes or changes to model assumptions. The first
568 dataset, generated by our simulation framework described above, consisted of ~2.8 million synthetic
569 tumours simulated across varying sequencing depths and overdispersions (all parameters provided in
570 Supplementary Table S1). Using this dataset, we compared TumE against six frequency-based summary
571 statistics for differentiating between positive selection and neutral evolution. Four of the statistics were
572 cancer evolution statistics developed previously¹² and provided in the R package *neutralitytestr*. For each
573 sample, the parameters of *neutralitytestr* were set as follows: ploidy = 2, cellularity = 1, read_depth =
574 simulated mean sequencing depth, rho = simulated overdispersion (rho). Two of the statistics were common
575 population genetic statistics, Tajima's D^{20} and Fay and Wu's H^21 . Only variant allele frequencies and
576 sequencing depth were required for input to compute these statistics. We provide an implementation of
577 Tajima's D and Fay and Wu's H for tumour sequencing data in the github repository. We additionally
578 evaluated a mixture model based approach MOBSTER¹⁶ for subclone detection and frequency
579 quantification. To enable analysis of ~2.8 million synthetic tumours, we ran MOBSTER with the following
580 parameters: K = 1:3, samples = 2, init = "peaks", tail = c(TRUE, FALSE), epsilon = 1e-6, maxIter = 100,
581 fit.type = "MM", seed = 12345, model.selection = "reICL", pi_cutoff = 0.02, N_cutoff = 10. We defined the
582 number of subclones that MOBSTER detected as follows. If a tail and 3 beta components were fit then we
583 assigned 2 subclones, if a tail and 2 beta components or if no tail and 3 beta components were fit we
584 assigned 1 subclone, and for all remaining fits we assigned 0 subclones or neutrality.

585

586 The second dataset was retrieved from Caravagna et al.¹⁶ and consists of synthetic tumour sequencing
587 data from 150 tumours sequenced to 120x depth using a beta-binomial sequencing model and grown to a
588 size of $>10^8$ at a birth rate of 1 and death rate of 0.2. The complete description is provided in the
589 supplementary of ref¹⁶. We used this dataset to evaluate the small N approximation and to compare TumE
590 to existing mixture model methods. We applied both MOBSTER and a variational Bayesian mixture model
591 sciClone²⁴ to this dataset. MOBSTER was run under default package settings without parallel computation
592 and with K = 1 to 3 beta components. sciClone was run under default package settings with
593 copyNumberCalls fixed to 2 and maximumClusters fixed to 4. To estimate the number of selected subclones
594 with sciClone, which doesn't account for neutral evolutionary dynamics, we took the inferred number of

595 subclones and subtracted 2 (representing the neutral tail and clonal peak). Per-sample runtimes for TumE,
596 MOBSTER, and sciClone were computed on a single machine with 16GB memory and a 2.3GHz quad-
597 core Intel i7 processor.

598

599 The third dataset was used to evaluate variable birth and death rates on TumE estimates for predicting
600 positive selection, determining the number of subclones, and estimating subclone frequency. The dataset
601 consisted of ~6 million synthetic tumours, generated by our simulation framework described above, grown
602 at variable birth and death rate combinations. Mutation rate and mean sequencing depth were both fixed
603 to 100. Other parameters were uniformly sampled and all parameters evaluated are outlined in
604 Supplementary Table S1.

605

606 All statistical analyses comparing methods, including computation of AUROC, correlation coefficients, and
607 performance metrics such as precision and recall, were performed in R *v4.0.3*.

608

609 **Evolutionary parameter estimates in bulk sequenced single tumour** 610 **biopsies**

611 In this study, we used diploid regions of patient tumours for evolutionary inference as we did not have
612 access to accurate phased mutation information for copy number correction of VAFs at non-diploid sites.
613 However, in the absence of complete whole-genome duplication, mutated diploid regions should be
614 sufficient to capture ongoing selection, due to selective sweeps from genome-wide linkage, if a sufficient
615 number of neutral passenger mutations are accumulated during cell division over time^{7,16}.

616

617 In addition to only analyzing diploid regions, we only accepted tumour samples that had at least a 50x mean
618 effective coverage (mean sequencing depth times purity). We set this cutoff as previous studies have shown
619 that tumour genomes sequenced below 50-70x coverage are exceedingly noisy and have insufficient limits
620 of detection relative to low-frequency mutations for proper evolutionary inference^{16,31}.

621

622 Relative to our simulations, VAFs in bulk sequenced single tumour biopsies may also be confounded by
623 impurity, where purity (cellularity) is defined as the percentage of cells in the biopsy that are of malignant
624 or tumour origin. In general, low tumour purity can lead to spurious identification of subclones as it results
625 in lower observed VAFs relative to the true underlying population VAFs. To ensure our inferences weren't
626 biased by impurity, we corrected all VAFs using corresponding tumour purity estimates collected from the
627 study of origin where $VAF_{corrected} = VAF_{observed} / \text{purity}$.

628

629 We also note that some purity estimates may be incorrect - in these cases updating the VAFs with incorrect
630 purity estimates can lead to a heterozygous clonal cluster (that should be centered at approximately 50%
631 VAF) in the subclonal frequency range (~10 - 40% VAF). To ensure clonal clusters were properly centered

632 at 50% VAF following purity correction, we performed additional adjustments to each patients' VAF
633 distribution using the following heuristic. We first computed the density for each VAF distribution and then
634 identified all the locations where the second derivative of the density was zero i.e. peak finding. If the closest
635 peak to 50% VAF (the theoretical diploid clonal cluster) was above 35% VAF, we considered it a
636 misrepresented clonal peak. We made this assumption as analyses in pan-cancer datasets suggest that all
637 tumours are initiated in somatic cells already carrying mutations^{10,11}. We then fit a Gaussian distribution to
638 the identified clonal cluster of each patient VAF distribution and adjusted each VAF by multiplying by 0.5
639 divided by the mean of the fit. Although a Beta distribution is generally used for fitting clonal clusters in
640 cancer genomics, a Gaussian is a reasonable approximation for adjusting VAFs based on incorrect purity
641 estimates as it provides accurate estimates of the cluster mean, and has been used in previous subclonal
642 clustering methods⁴¹. Plots of patient VAF distributions before and after application of this correction are
643 provided in (Supplementary Figures 17 & 18).

644
645 Heuristic clustering using the estimated subclone frequencies was performed either using the expected
646 variance under a binomial sequencing noise model or, alternatively, using the subclone frequency estimates
647 to initialize the means of a gaussian mixture model. Clustering under the binomial framework was performed
648 as follows. Given an estimated subclone frequency q , all mutations within the frequency range of $q \pm$
649 $(\epsilon\sqrt{qc[1-q]}) / c$ were assigned to the subclone, where ϵ scales the cluster width and c is the mean
650 sequencing depth across the tumour genome or exome. We fixed ϵ to 2 in this study.

651 **Transfer learning for inference in alternative synthetic data regimes**

652 Given a pre-trained network with weights optimized to a source domain S , composed of input space X_S ,
653 output space Y_S , transfer learning attempts to use pre-training to improve the performance on another target
654 domain T composed of X_T and Y_T . We employ a variant of transfer learning called open set domain
655 adaptation³⁶ to take advantage of our pre-trained models for additional inference tasks. In this case, the
656 input spaces remain constant ($X_S = X_T$, VAF distribution) but the inferred tasks are allowed to differ. Open
657 set indicates that some tasks may overlap with the output space of both the source and target domains.

658
659
660 To provide a concrete use case for transfer learning in synthetic supervised learning, we aimed to infer
661 additional evolutionary parameters such as subclone fitness, subclone emergence time, mutation rate, and
662 subclone cellular fraction (subclone frequency * 2) using synthetic tumour sequencing data generated by
663 an alternative cancer simulation framework TEMULATOR³⁷. TEMULATOR differs from our synthetic data
664 generation scheme, which was built around a multiplicative fitness driver model, as subclones are
665 deterministically initiated at user specified emergence times and fitnesses. To facilitate transfer between
666 previous and new tasks, we performed architecture renovation on the pre-trained neural networks, retaining
667 all convolutional layers while replacing existing task-specific fully-connected branches with new task-

668 specific fully connected branches (4 in total). To maximize the amount of information transferred to new
669 tasks, we combined the convolutional layers from both the M_{ms} and M_{1s} models described above.

670

671 We then simulated 500,000 synthetic tumours at a birth rate of 1, death rate of 0.2, and final population
672 size of 10^4 (additional parameters such as mutation rate were uniformly sampled and are outlined in
673 Supplementary Table S3). To facilitate efficient simulation, we first fit a noisy Gaussian process (GP)
674 regression to the viable emergence time and fitness parameter combinations (that generated detectable
675 subclones between ~10-40% VAF) and used the GP to sample viable emergence times given a uniformly
676 sampled fitness. We made the GP noisy to facilitate parameter combinations that resulted in subclones
677 across the entire frequency range. The GP was fit using three kernels (RBF with length scale 100, dot-
678 product, and white noise) and an alpha of 10^{-6} in the python package *scikit-learn v1.0*. Next, we used the
679 simulations to re-optimize the pre-trained model weights, using an Adam optimizer to minimize the L1 loss
680 (mean absolute error) for predicting new evolutionary inference tasks. To ensure fair benchmarking
681 between networks with and without pre-trained weights, we performed a random hyperparameter search
682 with ~150 pre-trained and ~150 non-pretrained models, tuning the learning rate and number of fully
683 connected layers in the task-specific branches. Additional synthetic data used for evaluating performance
684 was generated under similar parameter settings. We corrected modest yet systematic overestimates in
685 mutation rate (~50% mean percentage error) in the final transferred model by fitting a polynomial (degree
686 2) ridge regression in *scikit-learn v1.0* to the predicted mutation rates. The mutation rate adjustment was
687 performed using VAF distributions from 1000 synthetic tumours. We validated the correction on an
688 additional 100,000 synthetic tumours. All TEMULATOR synthetic tumours were generated using parameter
689 settings in Supplementary Table S3.

690

691 Predictions in empirical samples were performed by taking 500 Monte Carlo dropout samples and
692 averaging the results. Dropout was only activated at test time on the new task-specific branches. Per-base
693 mutation rate estimates in whole-exome sequenced MMR-GE samples were rescaled based on the 60MB
694 Agilent SureSelectXT Human All Exon v6 kit used in the original study³³. Because subclone fitness and
695 emergence time is impacted by final tumour size, we rescaled our estimates to a realistic tumour size of
696 10^{10} cells, similar to ref⁷. Previous work⁷ has shown that given subclone frequency f_{sub} and an estimated
697 final population size N_{end} , the age of a tumour at time of biopsy can be estimated by $t_{end} = \log_2([1 - f_{sub}] * N_{end})$.
698 Therefore, given that the relationship between emergence time in tumour doublings and log
699 population size is linear, we can generate a rescale fitness estimate w_R as follows.

700

701

$$w_R = I + (w - I) * \frac{t_{end} - t_s}{t_{end_R} - t_{sR}}$$

702

703 where w equals subclone fitness, t_{end} indicates time at tumour biopsy or final population size in tumour
704 doublings, and t_s indicates the time of subclone emergence in tumour doublings. R subscript indicates
705 values rescaled to population size of 10^{10} . The parameters f_{sub} , w , and t_s are all inferred. We approximate
706 the rescaled subclone emergence time t_{sR} as $t_s * \log(N_{endR}) / \log(N_{end})$.

707

708 **Data availability**

709 All TumE predictions in synthetic and empirical datasets, intermediate processing data, data used for
710 generating figures, and fully trained deep learning models can be found at
711 <https://doi.org/10.5281/zenodo.5575877> repository. Whole-genome sequenced AML samples were
712 retrieved from Griffith et al³¹. Multi-region whole-exome sequenced mismatch deficient repair gastro-
713 esophageal samples were retrieved from von Loga et al³³. The remaining whole-genome sequenced
714 samples were retrieved from PCAWG¹¹. We provide hosting of the electronic supplementary at
715 https://tomouellette.gitlab.io/ouellette_awadalla_2021/.

716

717 **Code availability**

718 Scripts for generating figures and analyses can be found at <https://doi.org/10.5281/zenodo.5575877>. Code
719 for generating synthetic tumour sequencing data can be found at
720 <https://github.com/tomouellette/CanEvolve.jl>. Code for performing inference with TumE can be found at
721 <https://github.com/tomouellette/TumE>.

References

1. Hanahan, D. & Weinberg, R. A. Hallmarks of Cancer: The Next Generation. *Cell* **144**, 646–674 (2011).
2. Tabassum, D. P. & Polyak, K. Tumorigenesis: it takes a village. *Nat. Rev. Cancer* **15**, 473–483 (2015).
3. Lipinski, K. A. *et al.* Cancer Evolution and the Limits of Predictability in Precision Cancer Medicine. *Trends Cancer* **2**, 49–63 (2016).
4. Minussi, D. C. *et al.* Breast tumours maintain a reservoir of subclonal diversity during expansion. *Nature* 1–7 (2021) doi:10.1038/s41586-021-03357-x.
5. Fittall, M. W. & Loo, P. V. Translating insights into tumor evolution to clinical practice: promises and challenges. *Genome Med.* **11**, 1–14 (2019).
6. West, J. B. *et al.* Multidrug Cancer Therapy in Metastatic Castrate-Resistant Prostate Cancer: An Evolution-Based Strategy. *Clin. Cancer Res.* **25**, 4413–4421 (2019).

7. Williams, M. J. *et al.* Quantification of subclonal selection in cancer from bulk sequencing data. *Nat. Genet.* **50**, 895–903 (2018).
8. Martincorena, I. *et al.* Universal Patterns of Selection in Cancer and Somatic Tissues. *Cell* **171**, 1029–1041.e21 (2017).
9. Gatenby, R. A. & Brown, J. S. Integrating evolutionary dynamics into cancer therapy. *Nat. Rev. Clin. Oncol.* **17**, 675–686 (2020).
10. Weinstein, J. N. *et al.* The Cancer Genome Atlas Pan-Cancer analysis project. *Nat. Genet.* **45**, 1113–1120 (2013).
11. Campbell, P. J. *et al.* Pan-cancer analysis of whole genomes. *Nature* **578**, 82–93 (2020).
12. Williams, M. J., Werner, B., Barnes, C. P., Graham, T. A. & Sottoriva, A. Identification of neutral tumor evolution across cancer types. *Nat. Genet.* **48**, 238–244 (2016).
13. Sun, R. *et al.* Between-region genetic divergence reflects the mode and tempo of tumor evolution. *Nat. Genet.* **49**, 1015–1024 (2017).
14. Bozic, I., Paterson, C. & Waclaw, B. On measuring selection in cancer from subclonal mutation frequencies. *PLOS Comput. Biol.* **15**, e1007368 (2019).
15. Lee, N. & Bozic, I. Inferring parameters of cancer evolution from sequencing and clinical data. *bioRxiv* 2020.11.18.387837 (2020) doi:10.1101/2020.11.18.387837.
16. Caravagna, G. *et al.* Subclonal reconstruction of tumors by using machine learning and population genetics. *Nat. Genet.* **52**, 898–907 (2020).
17. Tung, H.-R. & Durrett, R. Signatures of neutral evolution in exponentially growing tumors: A theoretical perspective. *PLOS Comput. Biol.* **17**, e1008701 (2021).
18. Crow, J. F. & Kimura, M. *An introduction to population genetics theory.* (New York, Evanston and London: Harper & Row, Publishers, 1970).
19. Salichos, L., Meyerson, W., Warrell, J. & Gerstein, M. Estimating growth patterns and driver effects in tumor evolution from individual samples. *Nat. Commun.* **11**, 732 (2020).
20. Tajima, F. Statistical Method for Testing the Neutral Mutation Hypothesis by DNA Polymorphism. *Genetics* **123**, 585–595 (1989).
21. Fay, J. C. & Wu, C. I. Hitchhiking under positive Darwinian selection. *Genetics* **155**, 1405–1413

- (2000).
22. Sheehan, S. & Song, Y. S. Deep Learning for Population Genetic Inference. *PLOS Comput. Biol.* **12**, e1004845 (2016).
 23. Prangle, D. Summary Statistics in Approximate Bayesian Computation. *ArXiv151205633 Math Stat* (2015).
 24. Miller, C. A. *et al.* SciClone: Inferring Clonal Architecture and Tracking the Spatial and Temporal Patterns of Tumor Evolution. *PLOS Comput. Biol.* **10**, e1003665 (2014).
 25. Roth, A. *et al.* PyClone: statistical inference of clonal population structure in cancer. *Nat. Methods* **11**, 396–398 (2014).
 26. Gal, Y. & Ghahramani, Z. Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. *ArXiv150602142 Cs Stat* (2016).
 27. Gal, Y. & Ghahramani, Z. Bayesian Convolutional Neural Networks with Bernoulli Approximate Variational Inference. *ArXiv150602158 Cs Stat* (2016).
 28. Le, T. A., Baydin, A. G., Zinkov, R. & Wood, F. Using Synthetic Data to Train Neural Networks is Model-Based Reasoning. *ArXiv170300868 Cs Stat* (2017).
 29. Kessler, D. A. & Levine, H. Large population solution of the stochastic Luria–Delbrück evolution model. *Proc. Natl. Acad. Sci.* **110**, 11682–11687 (2013).
 30. McFarland, C. D., Mirny, L. A. & Korolev, K. S. Tug-of-war between driver and passenger mutations in cancer and other adaptive processes. *Proc. Natl. Acad. Sci.* **111**, 15138–15143 (2014).
 31. Griffith, M. *et al.* Optimizing Cancer Genome Sequencing and Analysis. *Cell Syst.* **1**, 210–223 (2015).
 32. Nik-Zainal, S. *et al.* The Life History of 21 Breast Cancers. *Cell* **149**, 994–1007 (2012).
 33. von Loga, K. *et al.* Extreme intratumour heterogeneity and driver evolution in mismatch repair deficient gastro-oesophageal cancer. *Nat. Commun.* **11**, 139 (2020).
 34. Dentre, S. C. *et al.* Characterizing genetic intra-tumor heterogeneity across 2,658 human cancer genomes. *Cell* **184**, 2239–2254.e39 (2021).
 35. Kouw, W. M. & Loog, M. An introduction to domain adaptation and transfer learning. *ArXiv181211806 Cs Stat* (2019).
 36. Busto, P. P. & Gall, J. Open Set Domain Adaptation. in *2017 IEEE International Conference on*

- Computer Vision (ICCV)* 754–763 (IEEE, 2017). doi:10.1109/ICCV.2017.88.
37. Heide, T. *et al.* Reply to ‘Neutral tumor evolution?’ *Nat. Genet.* **50**, 1633–1637 (2018).
 38. Sottoriva, A. *et al.* A Big Bang model of human colorectal tumor growth. *Nat. Genet.* **47**, 209–216 (2015).
 39. Chkhaidze, K. *et al.* Spatially constrained tumour growth affects the patterns of clonal selection and neutral drift in cancer genomic data. *PLOS Comput. Biol.* **15**, e1007243 (2019).
 40. West, J., Schenck, R. O., Gatenbee, C., Robertson-Tessi, M. & Anderson, A. R. A. Normal tissue architecture determines the evolutionary course of cancer. *Nat. Commun.* **12**, 2060 (2021).
 41. Xiao, Y. *et al.* FastClone is a probabilistic tool for deconvoluting tumor heterogeneity in bulk-sequencing samples. *Nat. Commun.* **11**, 4469 (2020).