

# Accounting for errors in data improves divergence time estimates in single-cell cancer evolution

Kylie Chen<sup>1\*</sup>, Jiří C. Moravec<sup>2,3</sup>, Alex Gavryushkin<sup>3</sup>, David Welch<sup>1</sup>, Alexei J. Drummond<sup>1,4</sup>

<sup>1</sup> School of Computer Science, University of Auckland, Auckland, New Zealand

<sup>2</sup> Department of Computer Science, University of Otago, Dunedin, New Zealand

<sup>3</sup> School of Mathematics and Statistics, University of Canterbury, Christchurch, New Zealand

<sup>4</sup> School of Biological Sciences, University of Auckland, Auckland, New Zealand

\* kche309@aucklanduni.ac.nz

## Abstract

Single-cell sequencing provides a new way to explore the evolutionary history of cells. Compared to traditional bulk sequencing, where a population of heterogeneous cells is pooled to form a single observation, single-cell sequencing isolates and amplifies genetic material from individual cells, thereby preserving the information about the origin of the sequences. However, single-cell data is more error-prone than bulk sequencing data due to the limited genomic material available per cell. Here, we present error and mutation models for evolutionary inference of single-cell data within a mature and extensible Bayesian framework, BEAST2. Our framework enables integration with biologically informative models such as relaxed molecular clocks and population dynamic models. Our simulations show that modeling errors increase the accuracy of relative divergence times and substitution parameters. We reconstruct the phylogenetic history of a colorectal cancer patient and a healthy patient from single-cell DNA sequencing data. We find that the estimated times of terminal splitting events are shifted forward in time compared to models which ignore errors. We observed that not accounting for errors can overestimate the phylogenetic diversity in single-cell DNA sequencing data. We estimate that 30-50% of the apparent diversity can be attributed to error. Our work enables a full Bayesian approach capable of accounting for errors in the data within the integrative Bayesian software framework BEAST2.

© The Author(s) 2022. Published by Oxford University Press on behalf of Society for Molecular Biology and Evolution. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<https://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

# Introduction

The growth of cancer cells can be viewed as an evolutionary process where mutations accumulate along cell lineages over time. Within each cell, single nucleotide variants (SNVs) act as markers for the evolutionary process. By sampling and sequencing cells, we can reconstruct the possible evolutionary histories of these cell lineages. This can provide insight into the timing of events and modes of evolution.

Currently, there are two main methods for obtaining genomic sequences, bulk sequencing, and single-cell sequencing. Bulk sequencing data is traditionally used in genomic studies. By pooling the genetic material from many cells to form a single observation, greater coverage and thus genetic signal is retained. However, in the context of cancer phylogenetics, the analysis of bulk data poses challenges. Firstly the intermixing of tumor and normal cells affects the genomic signal. Secondly, the pooled sample may be heterogeneous and thus contain a mixture of different genomic variants (Dagogo-Jack and Shaw, 2018; de Bruin et al., 2014; Liu et al., 2018).

In contrast, single-cell sequencing isolates and amplifies the genetic material within a single cell (Kuipers et al., 2017a). The isolation step alleviates the mixture problem. However, errors are more problematic for single-cell sequencing due to insufficient coverage caused by the limited amount of genetic material. The main sources of errors in single-cell sequencing include: cell doublets, where two

cells are sequenced as one by mistake; allelic dropout (ADO), where one of the alleles fails to be amplified; and sequencing error, where a base is erroneously read as a different base by the sequencing machine (Kuipers et al., 2017a; Woodworth et al., 2017; Lähnemann et al., 2020). Error models proposed to address these issues include models based on false positives and false negatives (Ross and Markowitz, 2016; Jahn et al., 2016; Zafar et al., 2017, 2019), models of allelic dropout and sequencing errors (Kozlov et al., 2022), and models of read count errors (Satas et al., 2020).

To enable easy integration with molecular clock and phylogeography models that are commonly used in other areas of phylogenetics (Meijer et al., 2012; Malmström et al., 2016; Kearns et al., 2018) we implemented two error models within a mature Bayesian evolutionary framework, BEAST2 (Bouckaert et al., 2019). Our motivation is to enable inference and quantify the uncertainty of both the evolutionary history and model parameters for single-cell phylogenetics. Our paper implements: (i) a model for false positive and false negative errors (Ross and Markowitz, 2016; Jahn et al., 2016; Zafar et al., 2017, 2019) and (ii) a model for ADO and sequencing errors (Kozlov et al., 2022).

We show that our implementation is well-calibrated (Dawid, 1982) and demonstrate these models on real and simulated single-cell DNA data. Our simulation studies show that not accounting for errors leads to inac-



independence assumption made by the SiFit model.

SCARLET (Satas et al., 2020) implements a read count model which accounts for false positives and false negatives. This is done by correcting read counts at each site using copy-number variation (CNV) output from another software. Empirical studies have shown ADO is the most significant contributor of errors in single-cell DNA sequencing (Wang and Navin, 2015). CellPhy (Kozlov et al., 2022) explicitly models both ADO and sequencing error on diploid genotypes. Unlike models based on false positives and false negatives where different error types are absorbed into the false positive and false negative parameters, CellPhy is a more realistic model of the errors arising from the sequencing process. Furthermore, CellPhy has been shown to produce the most accurate phylogenetic estimates, followed by SiFit and infSCITE on simulated NGS datasets with ADO, amplification, and doublet errors (Kozlov et al., 2022). Both SCARLET and CellPhy make important advances in using data that is closer to the observed sequencing data than previous methods.

Besides CellPhy, which uses the ML phylogenetic framework RAxML (Stamatakis, 2014), other methods are only available as standalone implementations. The advantage of our work is that it enables easy integration with a wide range of population and clock models. In doing so, making these models available for single-cell phylogenetics. This includes re-

laxed clock models (Drummond et al., 2006), population growth models such as Bayesian skyline plots (Drummond et al., 2005), and phylogeography models such as structured coalescent (Vaughan et al., 2014) and isolation migration models (Nielsen and Wakeley, 2001).

In this paper, we implement error models for binary SNV data and diploid nucleotide SNV data. The binary model accounts for false positive and false negative errors (Ross and Markowitz, 2016; Jahn et al., 2016; Zafar et al., 2017, 2019). The nucleotide model accounts for ADO and sequencing errors (Kozlov et al., 2022). First, we investigate how errors impact the time scale of evolutionary trees inferred from single-cell data. Then, we perform preliminary analyses on real single-cell data to show error models can be used with population growth and molecular clock models. The next section describes the evolutionary models used in this study.

## Materials and Methods

We implement two sets of models: (i) the binary model, which handles mutation presence-absence data and (ii) the GT16 model, which handles diploid nucleotide genotypes.

The mutation process is modeled as a substitution process evolving along the branches of a tree  $\tau$ , with mutation rates defined by the substitution rate matrix  $Q$ . Errors are modeled as a noisy process on tip sequences of the tree, where the true genotype is obfuscated according to error probabilities. To perform

inference on data, we sample the posterior distribution of trees and the model parameters using Markov chain Monte Carlo (MCMC).

## Software and Input format

Our software is available at [www.github.com/bioDS/beast-phylonco](http://www.github.com/bioDS/beast-phylonco). It accepts input files in Nexus, FASTA, or VCF format via a conversion script available at [www.github.com/bioDS/vcf2fasta](http://www.github.com/bioDS/vcf2fasta).

## Binary substitution model

The presence or absence of mutation is represented as a binary state  $\Gamma = \{1, 0\}$ . The rate matrix  $Q$  has a single parameter,  $\lambda$  which is the rate of back-mutation  $1 \rightarrow 0$ , relative to a mutation rate of 1.

The elements of the rate matrix  $Q$  are:

$$Q = \begin{matrix} & \begin{matrix} 0 & 1 \end{matrix} \\ \begin{matrix} 0 \\ 1 \end{matrix} & \begin{pmatrix} -1 & 1 \\ \lambda & -\lambda \end{pmatrix} \end{matrix}.$$

The equilibrium frequencies are:

$$\begin{aligned} \pi_0 &= \lambda / (\lambda + 1), \\ \pi_1 &= 1 / (\lambda + 1). \end{aligned}$$

For data sampled at a single time point, the mutation rate is in units of substitutions per site. Data sampled at multiple time points are required to estimate the mutation rate, which typically has units of substitutions per site per year. Alternatively, if we have prior

information on the mutation rate, such as from empirical experiments, we can also fix the model's mutation rate to the empirical value.

## Binary error model

To account for false positive and false negative errors, we implement the binary error model described in (Jahn et al., 2016; Zafar et al., 2017; Ross and Markowitz, 2016). Let  $\alpha$  be the false positive probability and  $\beta$  be the false negative probability.  $P(x|y)$  is the conditional error probability of observing noisy data  $x$ , given that the true state is  $y$ . For the binary error model, these error probabilities are:

$$\begin{aligned} P(0|0) &= 1 - \alpha, \\ P(1|0) &= \alpha, \\ P(0|1) &= \beta, \\ P(1|1) &= 1 - \beta. \end{aligned} \tag{1}$$

## GT16 substitution model

To model diploid nucleotide sequences, we implement the GT16 substitution and error model described in (Kozlov et al., 2022). The GT16 substitution model is an extension of the four-state general time-reversible nucleotide GTR model (Tavaré et al., 1986) to diploid genotypes:  $\Gamma = \{AA, AC, AG, AT, CA, CC, CG, CT, GA, GC, GG, GT, TA, TC, TG, TT\}$ .

Let  $a, b, c, d$  be alleles chosen from nucleotides  $N = \{A, C, G, T\}$  and  $r_{ab}$  be the rate of going

from allele  $a$  to allele  $b$ .

The elements of the rate matrix  $Q$  are:

$$\begin{aligned} Q_{aa \rightarrow ab} &= r_{ab} \cdot \pi_{ab}, \\ Q_{aa \rightarrow ba} &= r_{ab} \cdot \pi_{ba}, \\ Q_{ab \rightarrow aa} &= r_{ab} \cdot \pi_{aa}, \\ Q_{ab \rightarrow bb} &= r_{ab} \cdot \pi_{bb}. \end{aligned}$$

Other non-diagonal entries not listed above have a rate of zero. The diagonals are the sum of the in-going rates:

$$\begin{aligned} Q_{aa \rightarrow aa} &= - \sum_{b,c \in N \setminus a} Q_{aa \rightarrow bc}, \\ Q_{ab \rightarrow ab} &= - \sum_{\substack{c,d \in N \\ c \neq a \text{ or } d \neq b}} Q_{ab \rightarrow cd}. \end{aligned}$$

The relative rates of the  $Q$  matrix are:

$$\begin{aligned} r_{AC} &= r_{CA} = \alpha, \\ r_{AG} &= r_{GA} = \beta, \\ r_{AT} &= r_{TA} = \gamma, \\ r_{CG} &= r_{GC} = \kappa, \\ r_{CT} &= r_{TC} = \lambda, \\ r_{GT} &= r_{TG} = \mu. \end{aligned}$$

The equilibrium frequencies are:  $\pi = (\pi_{AA}, \pi_{AC}, \pi_{AG}, \pi_{AT}, \pi_{CA}, \pi_{CC}, \pi_{CG}, \pi_{CT}, \pi_{GA}, \pi_{GC}, \pi_{GG}, \pi_{GT}, \pi_{TA}, \pi_{TC}, \pi_{TG}, \pi_{TT})$ .

## GT16 error model

The GT16 error model for diploid nucleotides described in (Kozlov et al., 2022) accounts for amplification errors and biases in single-cell sequencing. This model has two parameters, the amplification, and sequencing error  $\epsilon$ , and

allelic dropout error  $\delta$ . The error probabilities  $P(x|y)$  for genotypes with alleles  $a, b, c$  derived in (Kozlov et al., 2022) are given below:

$$\begin{aligned} P(aa|aa) &= 1 - \epsilon + (1/2) \delta \epsilon, \\ P(ab|aa) &= (1 - \delta) (1/6) \epsilon, \\ P(bb|aa) &= (1/6) \delta \epsilon, \\ P(aa|ab) &= (1/2) \delta + (1/6) \epsilon - (1/3) \delta, \\ P(cc|ab) &= (1/6) \delta \epsilon, \\ P(ac|ab) &= (1 - \delta) (1/6) \epsilon, \\ P(ab|ab) &= (1 - \delta) (1 - \epsilon). \end{aligned} \quad (2)$$

We assume that  $P(ba|aa) = P(ab|aa)$  and  $P(cb|ab) = P(ac|ab)$ . Other combinations not listed above have zero probability. These genotypes can be easily adapted to unphased data by encoding heterozygous states as ambiguities  $P(ab^*) = P(ab) + P(ba)$ , where  $ab^*$  represents  $ab$  without phasing information. Our implementation can handle both phased and unphased data.

## Likelihood calculation

Our input data is a SNV matrix  $D$  with  $n$  sites and  $m$  cells. The cell evolutionary tree  $\tau$  is a rooted binary tree with  $m$  cells at the leaves, and branch lengths  $t_1, t_2, \dots, t_{2(m-1)}$ . These branch lengths are scaled to units of substitutions per site for data sampled at a single time point or years where multiple time points are available. Figure 1 shows an example tree with cells  $a, b, c, d$  sampled at different time points.

The likelihood  $P(D|\tau, M, \theta)$  is the conditional





# Results

## Evaluation on simulated datasets

First, we evaluated our implementations using a well-calibrated study (Dawid, 1982) to test the reliability of the inference when simulating directly from the model. Following the well-calibrated criterion for credible intervals, we expect 95% of the credible interval to cover the true value 95% of the time. Next, we simulated sequences with errors and compared the inference performance with and without modeling the error. Then, we compared the runtime and convergence efficiency of each error model with the baseline non-error substitution model. Lastly, we performed experiments on data simulated with high levels of errors to test the robustness of our methods.

### Simulation 1: Binary data

We performed a well-calibrated study for the binary model using binary sequences with errors. First, we generated trees using a Yule model, then binary sequences were simulated along the branches of the tree, and errors were applied at the tips. Using the sequence data, we jointly estimate the model parameters, tree topology, and branch times using the binary model.

Simulation parameters: We generated 100 trees with 30 leaves from a Yule model, where each tree has a birthrate drawn from  $Normal(\mu = 7.0, \sigma = 1.0)$ . Sequences of length 400 were simulated using the binary model with rate  $\lambda \sim Lognormal(\mu = -1)$ , false positive prob-

ability  $\alpha \sim Beta(1, 50)$  and false negative probability  $\beta \sim Beta(1, 50)$ .

Figure S1 shows the estimates for the model parameters, tree length, and tree height compared to the true simulated values. The estimated 95% highest posterior density (HPD) intervals are shown as bars, where blue indicates the estimate covers the true value, and red indicates otherwise. Our simulations show that the true value of each parameter falls within the estimated 95% HPD interval 91-99% of the time. Figure S2 shows the estimated trees are, on average, 2-5 subtree prune and regraft (SPR) moves away from the true tree.

### Simulation 2: Binary data error vs. no error

To compare the effects of inference with and without error modeling, we used the data from Simulation 1, then performed inference with and without the binary error model. We compared the coverage for each parameter with and without an error model, i.e., how often the estimated 95% HPD covers the true value.

Figure S1 shows the estimated parameters with the binary error model, and Figure S3 shows the estimated parameters without using an error model. Table S1 shows the coverage of each parameter. The coverage of tree length drops from 95% when the error model is used to 39% when no error model is used. Similarly, the coverage of the substitution parameter  $\lambda$  drops 91% to 53%. Furthermore, the tree length tends to be overestimated when no



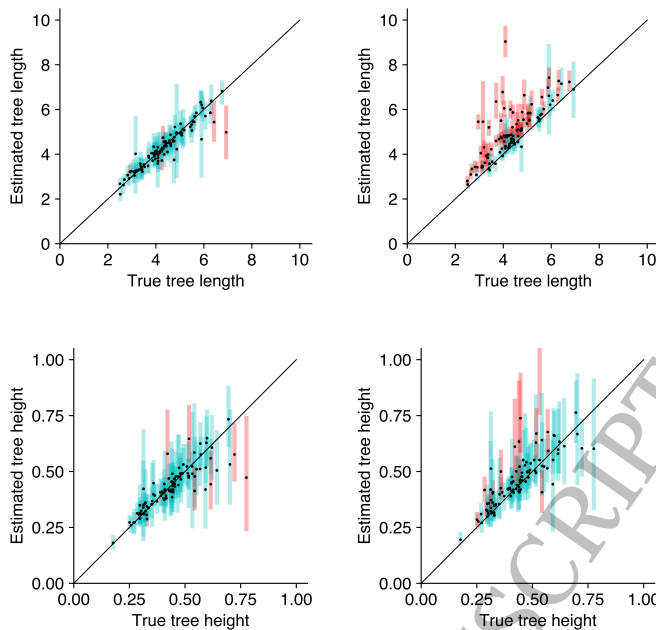


Figure 2: Comparison of branch lengths with and without the binary error model on simulated binary data. Estimated branch lengths using the binary error model (left) and without using the error model (right). Tree length estimates higher than 10, and tree height estimates higher than 1 are truncated on this plot.

error model is used. This suggests that both the tree length and substitution parameters are significantly biased when errors present in the data are not modeled. Figure 2 shows a comparison of the estimated tree length and tree height for these two model configurations. Other parameters such as birthrate and tree height are less biased, with a coverage of 92% and 84% respectively when no error model is used.

### Simulation 3: Diploid nucleotide phased data

We performed a well-calibrated study for the GT16 model using phased sequence data with errors. First, we simulated trees using a coales-

cent model. Sequences were simulated down branches of the tree using the GT16 substitution model, and then errors were applied at the tips. Using these sequences as input, we estimated the tree and model parameters using the GT16 model. The priors on the error probabilities were chosen based on experimental studies for allelic dropout (Huang et al., 2015), amplification and sequencing errors (Gawad et al., 2016; Ross et al., 2013).

Simulation parameters: We generated 100 trees with 16 leaves from a coalescent model, where the population size is drawn from  $\theta \sim \text{LogNormal}(\mu = -2.0, \sigma = 1.0)$ . Sequences of length 200 were simulated using the GT16 model with genotype frequencies  $\pi \sim \text{Dirich-}$

$let(3, 3, \dots, 3)$ , and relative rates  $r \sim Dirichlet(1, 2, 1, 1, 2, 1)$ . Errors were simulated using  $\epsilon \sim Beta(\alpha = 2, \beta = 18)$ , and  $\delta \sim Beta(\alpha = 1.5, \beta = 4.5)$ .

Figures S4-S9 show the 95% HPD estimated for each model parameter and tree branch lengths. The true value of each parameter falls within the 95% HPD interval 91-99% of the time. This shows we are able to accurately estimate the substitution, error, population parameters, and branch lengths for phased data. Next, we computed the accuracy of tree topology by comparing the estimated tree with the true tree. Figure S10 shows the average distance from the estimated trees to the true tree. On average, estimated trees are 2-6 SPR moves away from the true tree.

#### **Simulation 4: Diploid nucleotide unphased data**

We performed a well-calibrated study for the GT16 model using unphased sequencing data. For unphased data, we used the data generated from Simulation 3, with phasing information removed from the sequences. Phasing information was removed by mapping a heterozygous  $ab$  to both states  $ab$  and  $ba$ .

Figures S11-S15 show the estimated model parameters compared to the true simulated values. For each parameter, the estimated 95% HPD interval covers the true value 94-99% of the time, which confirms our implementation is well-calibrated. On average, the estimated trees are 2-6 SPR moves away from the true tree, Figure S16. We note that the

sum of the paired heterozygous frequencies ( $\pi_{ab} + \pi_{ba}$ ) are identifiable, but the individual frequencies ( $\pi_{ab}, \pi_{ba}$ ) are non-identifiable as the data is unphased.

#### **Simulation 5: Diploid nucleotide data error vs. no error**

To compare the effects of inference with and without error modeling for diploid nucleotide data, we used the data from Simulation 3, then performed inference with and without the GT16 error model.

Table S2 shows the coverage of each parameter with and without an error model. Figures S17-S21 show the estimated model parameters when an error model is not used. We observe a similar trend to Simulation 2, where the tree length and substitution parameters are significantly biased without an error model. Although the tree height estimated without an error model are less biased, the tree lengths are overestimated. These differences in the tree heights and tree lengths are highlighted in Figure 3.

#### **Simulation 6: Timing experiments**

We measured the runtime and convergence of the error model compared with the baseline non-error implementation in our framework. Both error models are comparable in computationally runtime efficiency with their baseline non-error substitution models. Runtime comparisons are shown in supplementary figures S24 - S25. The GT16 model takes approximately an hour to reach convergence on

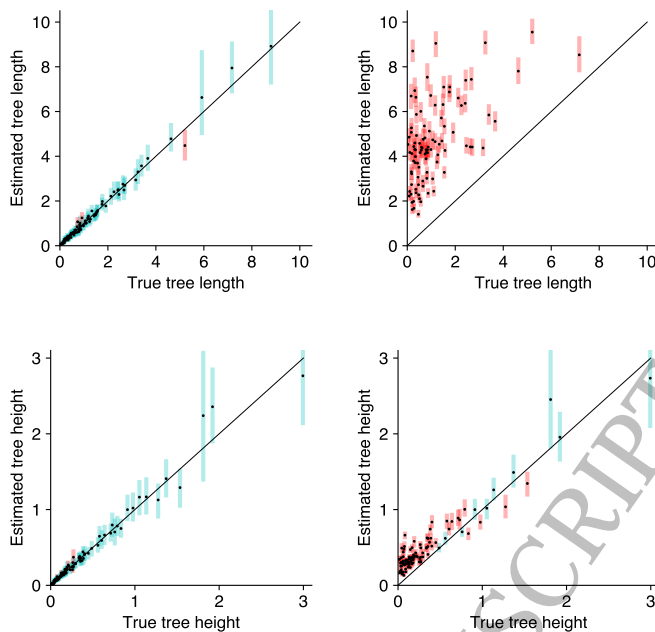


Figure 3: Comparison of branch lengths with and without the GT16 error model on simulated phased nucleotide data. Estimated branch lengths using the GT16 error model (left) and without using the error model (right). Tree length estimates higher than 10, and tree height estimates higher than 3 are truncated on this plot.

simulated datasets with 20 taxa and 500 sites (convergence is measured as the time till the minimum effective sample size is greater than 200). On a similar-sized dataset, the binary model takes less than five minutes to converge. Timing experiments were done on an Intel Xeon E3-12xx v2 virtual machine with 16 processors at 2.7MHz and 32GB RAM hosted by Nectar Research Cloud.

### Simulation 7: Performance on data with high levels of error

Lastly, we performed experiments on extended error ranges based on empirical studies (Huang et al., 2015; Gawad et al., 2016; Ross et al., 2013) to test the robustness of our method.

We used the same simulation parameters as simulation 1 and simulation 3, but with varying levels of error chosen from an extended range:  $\alpha \in [0.001, 0.1]$ ,  $\beta \in [0.1, 0.6]$  for binary data and  $\delta \in [0.1, 0.8]$ ,  $\epsilon \in [0.001, 0.1]$  for diploid genotype data. The priors on the error parameters are  $\alpha \sim \text{Beta}(1, 20)$  and  $\beta \sim \text{Beta}(3, 3)$  for binary data and  $\delta \sim \text{Beta}(1.5, 4.5)$  and  $\epsilon \sim \text{Beta}(2, 18)$  for diploid nucleotide data. Our results in the supplementary materials confirm our methods are robust to high levels of error.

### Evaluation on single-cell datasets

We analyzed two public datasets from previously published studies; L86, a colorectal can-

cer dataset (Leung et al., 2017), and E15, a healthy neurons dataset (Evrony et al., 2015). Preprocessed SNVs from CellPhy (Kozlov et al., 2022) were used for both L86 and E15.

### Colorectal cancer dataset (L86)

L86 contains 86 cells sequenced from a colorectal cancer patient with metastatic spread. The cells were sampled from the primary tumor (colorectal), the secondary metastatic tumor (liver), and matched normal tissue. We used the GT16 model with a relaxed clock to allow for different molecular clock rates in cancer and non-cancer lineages and a coalescent skyline tree prior, which allows changes in population sizes through time.

Model parameters: We used a GT16 substitution model with priors of frequencies  $\pi \sim \text{Dirichlet}(3, 3, \dots, 3)$ , relative rates  $r \sim \text{Dirichlet}(1, 2, 1, 1, 2, 1)$ , and GT16 error model with allelic dropout  $\delta \sim \text{Beta}(1.5, 4.5)$  and sequencing error  $\epsilon \sim \text{Beta}(2, 18)$ . A relaxed clock with a Lognormal prior, and Skyline coalescent tree prior with  $\theta_1 \sim \text{Lognormal}(\mu = -2.3, \sigma = 1.8)$ . We performed two independent repeats of the MCMC chains.

We found that the tree height is similar for both error and non-error models, but the relative ages of terminal branches are shorter for the error model. Figure 4 shows the tree length, treeness (Phillips and Penny, 2003; Lanyon, 1988) and gamma statistics (Pybus and Harvey, 2000) of the tree distributions under different experimental setups: with and without error modeling, and with and with-

out an outgroup constraint. The trees estimated using an error model are more tree-like than ones estimated without an error model. For the default setup without an outgroup, the 95% HPD estimate for tree length is (4.79, 5.85) with the error model and (7.00, 8.19) without the error model. The error parameters estimates are  $\delta \sim (0.62, 0.66)$  and  $\epsilon \sim (7 \cdot 10^{-6}, 1 \cdot 10^{-3})$ . The error estimates are comparable to the estimates reported by CellPhy (Kozlov et al., 2022) which are  $\delta \sim 0.63$  and  $\epsilon \sim 0.00$ .

Figure 5 summarizes the estimated tree with the error model (top) compared to without an error model (bottom). The tips of the tree are colored by cell type. The trees show most cells group together by cell type, which suggests there is signal in the data. However, there is some intermixing of metastatic tumor cells inside the primary tumor clade and mis-sorted normal cells as previously identified by (Leung et al., 2017; Kozlov et al., 2022). We also note that the most recent common ancestor (MRCA) of the normal clade is younger than the MRCA of the two tumor clades for both analyses. This is not what we intuitively expected because we believe the normal ancestral cell should be the ancestor of both tumor and normal cells. Although surprising, this observation is in agreement with the trees estimated by ML algorithms in Kozlov et al. (2022). We believe this issue is closely related to the phylogenetic rooting problem for heterogeneous data (Tian and Kubatko, 2017). Methods by Tian and Kubatko (2017),

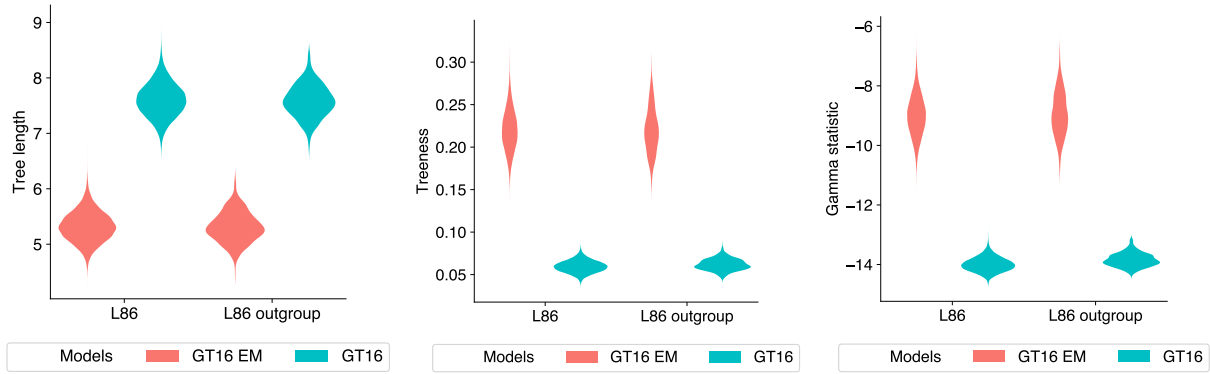


Figure 4: Tree length, treeness and gamma statistics of tree distributions estimated from the L86 dataset. The distributions of each metric is colored by the model used: GT16 error model (red) and GT16 model without error (blue). Two pairs of experiments are shown; L86, which has no tree topology constraints, and L86 outgroup, which has the tree topology constrained to normal cells in the outgroup.

Mai et al. (2017) and Drummond et al. (2006) have provided some partial solutions to the rooting problem, but further research efforts are required to better understand the effects on tree topology.

We also investigated whether constraining the tree topology to have all tumor cells as the ingroup produces different estimates. We repeat the analyses with an outgroup constraint to the tree topology, setting the normal and missorted samples as the outgroup. Figure S22 shows the outgroup constrained summary tree for L86. In this outgroup constrained tree, the age of the MRCA of the normal group is also younger than the MRCA of the primary or metastatic tumors for the error model. The age of the MRCA of the normal group is indistinguishable from the MRCA of the primary or metastatic tumors for the model without error.

### Healthy neurons dataset (E15)

E15 contains 15 neurons and a blood cell taken from the heart region sequenced from a healthy patient.

Model parameters: We used a GT16 substitution with frequencies prior  $\pi \sim \text{Dirichlet}(3, 3, \dots, 3)$  and relative rates  $r \sim \text{Dirichlet}(1, 2, 1, 1, 2, 1)$ , GT16 error model with allelic dropout error  $\delta \sim \text{Beta}(1.5, 4.5)$ , and sequencing error  $\epsilon \sim \text{Beta}(2, 18)$ , with a relaxed clock and Skyline coalescent tree prior with  $\theta_1 \sim \text{Lognormal}(\mu = -2.3, \sigma = 1.8)$ . We performed two independent repeats of the MCMC chains.

We observed that the trees estimated using the error model are more tree-like than ones estimated without an error model, as shown by the tree metrics in figure 6. The 95% HPD estimates of tree length are (1.37, 7.14) with the error model, and (7.60, 13.41) with-





out the error model. We note the tree height for the error model (0.20, 0.80) is lower than that of the non-error model (0.54, 1.01). The estimated interval for the error parameters are  $\delta \sim (0.86, 0.92)$  and  $\epsilon \sim (0.03, 0.17)$ . To test the sensitivity of the error priors, we reran our experiments with adjusted priors  $\epsilon, \delta \sim \text{Beta}(1, 10)$  and  $\epsilon, \delta \sim \text{Beta}(1, 20)$ . We found the error estimates were similar regardless of these adjustments on the error parameter priors.

Figure 7 shows a summary of the estimated trees with the GT16 error model (top) and without an error model (bottom). The tips of the tree are colored by cell types. We expect the blood cell to be placed as an outgroup; however, the estimated trees placed the blood cell inside a clade of neuron cells. To investigate if adding an outgroup constraint to the tree topology can help direct the likelihood in the correct direction, we repeated the analyses with the blood cell as the outgroup. The estimated trees are shown in figure S23. Besides the correct placement of the outgroup enforced by the outgroup constraint, we did not observe any substantial topological discrepancies between the outgroup and non-outgroup analyses.

## Discussion

We demonstrated that incorporating error parameters can affect the relative ages of single-cell datasets. We showed that models incorporating sequencing error could increase the

accuracy of tree branches and model parameters inferred from noisy data. Additionally, we find that using error models is just as fast as the baseline non-error substitution models in our framework. Future work to support multi-threading and add compatibility with the Beagle high-performance library (Ayres et al., 2012) would further increase the computational speed of these models.

From both simulated and real single-cell data, we observed that using an error model tends to shorten the total tree length, as errors explain a portion of the genetic variability within the data. For empirical single-cell data, cells of the same type tend to be placed in the same clade. We believe relaxed clock and local clock models are more suited to heterogeneous data as they allow for changes in mutation rates. The datasets we explored in this paper are sampled at a single time point, so there is no calibration information to allow the mutation rate and time to be disambiguated. Using time sampled data or empirical mutation rate calibrations would improve current analyses and allow node ages to be converted to real time (Drummond et al., 2003, 2002).

Although the effect of filtering strategies in the context of macroevolution shows stringent filtering of sites often leads to worse phylogenetic inference (Tan et al., 2015). The effect of filtering strategies on noisy data such as single-cell phylogenies is yet to be systematically explored. We believe the error parameters in these models can provide increased flexibility, allowing key features of the sequencing and

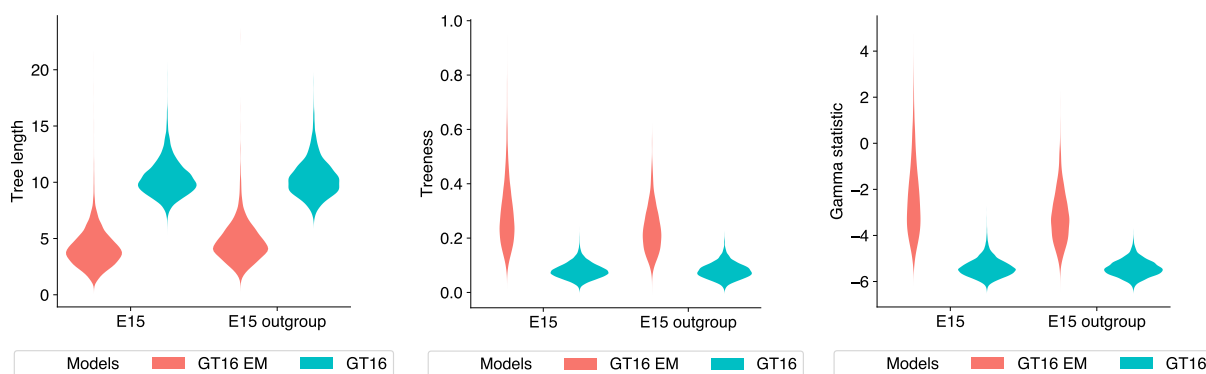


Figure 6: Tree length, treeness and gamma statistics of tree distributions estimated from the E15 dataset. The distributions of each metric is colored by the model used: GT16 error model (red) and GT16 model without error (blue). Two pairs of experiments are shown; E15, which has no tree topology constraints, and E15 outgroup, which has the tree topology constrained with the heart cell as the outgroup.

filtering process to be accounted for during evolutionary inference.

Lastly, incorporating cell biology knowledge during method development would improve the biological significance of model assumptions; and improving the interpretability of tree summarization metrics would enable single-cell phylogenies to be examined in more detail.

## Code and data availability

Our software, Phylonco v0.0.6 is available at [www.github.com/bioDS/beast-phylonco](https://www.github.com/bioDS/beast-phylonco).

Analyses, scripts, and data are available at [www.github.com/bioDS/beast-phylonco-paper](https://www.github.com/bioDS/beast-phylonco-paper).

This paper uses BEAST v2.6.6 (Bouckaert et al., 2019), BeastLabs v1.9.7, LPhy v1.2.0, and LPhyBeast v0.3.0. The following python packages were used: DendroPy (Sukumar and Holder, 2010), lxml (Behnel et al.,

2005), matplotlib (Hunter, 2007), numpy (Harris et al., 2020), and seaborn (Waskom, 2021). The following R packages were used: ggtree (Yu et al., 2017), ggplot2 (Wickham, 2016), tracerR (Rambaut et al., 2018), treeSimGM (Hagen and Stadler, 2018), treeio (Wang et al., 2020), expm, and ape (Paradis et al., 2019).

## Acknowledgements

We thank Prof. David Posada and Dr. Joao Alves for generously providing us with single-cell datasets and advice. We thank Dr. Remco Bouckaert for helpful discussions and Dr. Walter Xie for implementation support. We acknowledge Nectar Research Cloud and Otago Computing Services for providing the computing resources necessary for this study. We thank the reviewers and editors for their valuable comments which greatly improved the quality of our manuscript.

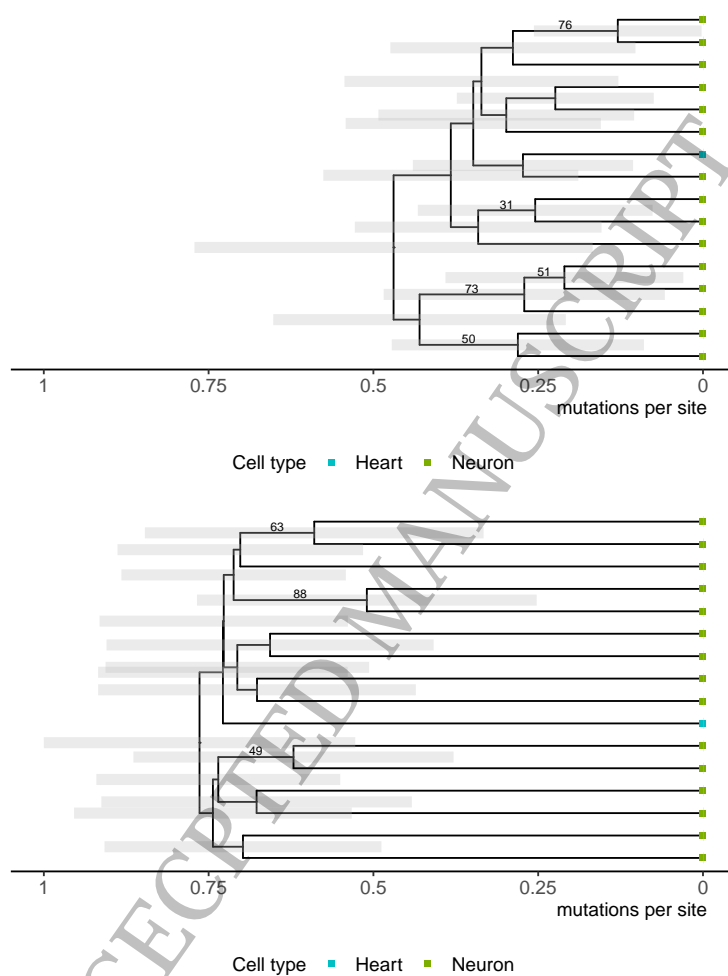


Figure 7: Maximum clade credibility trees for the E15 dataset (healthy patient) using the GT16 model with an error model (top) and without an error model (bottom). Each cell is colored by its cell type: blood cell (blue) and neuron cells (green). The posterior clade support for clades with greater than 30% support are shown on the branches.

AJD and AG acknowledge support from a Data Science Programme grant (UOAX1932); AJD, AG and JCM acknowledge support from an Endeavour Smart Ideas grant (U00X1912); AG was supported by a Rutherford Discovery Fellowship (RDF-U001702); AJD was supported by a James Cook Research Fellowship from the Royal Society of New Zealand; KC was supported by a University of Auckland Doctoral Scholarship.

## References

- Joao M Alves, Sonia Prado-Lopez, Jose Manuel Cameselle-Teijeiro, and David Posada. Rapid evolution and biogeographic spread in a colorectal cancer. *Nature communications*, 10(1):1–7, 2019.
- Daniel L Ayres, Aaron Darling, Derrick J Zwickl, Peter Beerli, Mark T Holder, Paul O Lewis, John P Huelsenbeck, Fredrik Ronquist, David L Swofford, Michael P Cummings, et al. Beagle: an application programming interface and high-performance computing library for statistical phylogenetics. *Systematic biology*, 61(1):170–173, 2012.
- Stefan Behnel, Martijn Faassen, and Ian Bicking. *lxml: Xml and html with python*, 2005.
- Remco Bouckaert, Timothy G. Vaughan, Joëlle Barido-Sottani, Sebastián Duchêne, Mathieu Fourment, Alexandra Gavryushkina, Joseph Heled, Graham Jones, Denise Kühnert, Nicola De Maio, Michael Matschiner, Fábio K. Mendes, Nicola F. Müller, Huw A. Ogilvie, Louis du Plessis, Alex Poppinga, Andrew Rambaut, David Rasmussen, Igor Siveroni, Marc A. Suchard, Chieh-Hsi Wu, Dong Xie, Chi Zhang, Tanja Stadler, and Alexei J. Drummond. Beast 2.5: An advanced software platform for bayesian evolutionary analysis. *PLOS Computational Biology*, 15(4):1–28, 04 2019. doi: 10.1371/journal.pcbi.1006650. URL <https://doi.org/10.1371/journal.pcbi.1006650>.
- Colin S Cooper, Rosalind Eeles, David C Wedge, Peter Van Loo, Gunes Gundem, Ludmil B Alexandrov, Barbara Kremeyer, Adam Butler, Andrew G Lynch, Niedzica Camacho, et al. Analysis of the genetic phylogeny of multifocal prostate cancer identifies multiple independent clonal expansions in neoplastic and morphologically normal prostate tissue. *Nature genetics*, 47(4):367–372, 2015.
- Ibiayi Dagogo-Jack and Alice T Shaw. Tumour heterogeneity and resistance to cancer therapies. *Nature reviews Clinical oncology*, 15(2):81, 2018.
- A Philip Dawid. The well-calibrated bayesian. *Journal of the American Statistical Association*, 77(379):605–610, 1982.
- Elza C de Bruin, Nicholas McGranahan, Richard Mitter, Max Salm, David C Wedge, Lucy Yates, Mariam Jamal-Hanjani, Seema

- Shafi, Nirupa Murugaesu, Andrew J Rowan, et al. Spatial and temporal diversity in genomic instability processes defines lung cancer evolution. *Science*, 346(6206):251–256, 2014.
- Alexei J Drummond, Geoff K Nicholls, Allen G Rodrigo, and Wiremu Solomon. Estimating Mutation Parameters, Population History and Genealogy Simultaneously From Temporally Spaced Sequence Data. *Genetics*, 161(3):1307–1320, 07 2002. ISSN 1943-2631. doi: 10.1093/genetics/161.3.1307. URL <https://doi.org/10.1093/genetics/161.3.1307>.
- Alexei J. Drummond, Oliver G. Pybus, Andrew Rambaut, Roald Forsberg, and Allen G. Rodrigo. Measurably evolving populations. *Trends in Ecology & Evolution*, 18(9):481–488, 2003. ISSN 0169-5347. doi: [https://doi.org/10.1016/S0169-5347\(03\)00216-7](https://doi.org/10.1016/S0169-5347(03)00216-7).
- Alexei J Drummond, Andrew Rambaut, BETH Shapiro, and Oliver G Pybus. Bayesian coalescent inference of past population dynamics from molecular sequences. *Molecular biology and evolution*, 22(5):1185–1192, 2005.
- Alexei J Drummond, Simon Y W Ho, Matthew J Phillips, and Andrew Rambaut. Relaxed phylogenetics and dating with confidence. *PLoS biology*, 4(5):e88, 2006.
- Gilad D Evrony, Eunjung Lee, Bhaven K Mehta, Yuval Benjamini, Robert M Johnson, Xuyu Cai, Lixing Yang, Psalm Haseley, Hillel S Lehmann, Peter J Park, et al. Cell lineage analysis in human brain using endogenous retroelements. *Neuron*, 85(1):49–59, 2015.
- Joseph Felsenstein. Evolutionary trees from dna sequences: a maximum likelihood approach. *Journal of molecular evolution*, 17(6):368–376, 1981.
- Charles Gawad, Winston Koh, and Stephen R Quake. Single-cell genome sequencing: current state of the science. *Nature Reviews Genetics*, 17(3):175, 2016.
- Oskar Hagen and Tanja Stadler. Treesimgm: Simulating phylogenetic trees under general bellman–harris models with lineage-specific shifts of speciation and extinction in r. *Methods in ecology and evolution*, 9(3):754–760, 2018.
- Charles R Harris, K Jarrod Millman, Stéfan J Van Der Walt, Ralf Gommers, Pauli Virtanen, David Cournapeau, Eric Wieser, Julian Taylor, Sebastian Berg, Nathaniel J Smith, et al. Array programming with numpy. *Nature*, 585(7825):357–362, 2020.
- Timon Heide, Angela Maurer, Monika Eipel, Katrin Knoll, Mirja Geelvink, Juergen Veeck, Ruth Knuechel, Julius van Essen, Robert Stoehr, Arndt Hartmann, et al. Multiregion human bladder cancer sequencing reveals tumour evolution, bladder cancer phenotypes and implications for targeted

- therapy. The Journal of pathology, 248(2): 230–242, 2019.
- Lei Huang, Fei Ma, Alec Chapman, Sijia Lu, and Xiaoliang Sunney Xie. Single-cell whole-genome amplification and sequencing: methodology and applications. Annual review of genomics and human genetics, 16: 79–102, 2015.
- J. D. Hunter. Matplotlib: A 2d graphics environment. Computing in Science & Engineering, 9(3):90–95, 2007. doi: 10.1109/MCSE.2007.55.
- Katharina Jahn, Jack Kuipers, and Niko Beerenwinkel. Tree inference for single-cell data. Genome biology, 17(1):86, 2016.
- Yuchao Jiang, Yu Qiu, Andy J Minn, and Nancy R Zhang. Assessing intratumor heterogeneity and tracking longitudinal and spatial clonal evolutionary history by next-generation sequencing. Proceedings of the National Academy of Sciences, 113(37): E5528–E5537, 2016.
- Anna M Kearns, Marco Restani, Ildiko Szabo, Audun Schröder-Nielsen, Jin Ah Kim, Hayley M Richardson, John M Marzluff, Robert C Fleischer, Arild Johnsen, and Kevin E Omland. Genomic evidence of speciation reversal in ravens. Nature Communications, 9(1):1–13, 2018.
- Alexey Kozlov, Joao M Alves, Alexandros Stamatakis, and David Posada. Cellphy: accurate and fast probabilistic inference of single-cell phylogenies from scdna-seq data. Genome Biology, 23(1):1–30, 2022.
- Jack Kuipers, Katharina Jahn, and Niko Beerenwinkel. Advances in understanding tumour evolution through single-cell sequencing. Biochimica et Biophysica Acta (BBA)-Reviews on Cancer, 1867(2):127–138, 2017a.
- Jack Kuipers, Katharina Jahn, Benjamin J Raphael, and Niko Beerenwinkel. Single-cell sequencing data reveal widespread recurrence and loss of mutational hits in the life histories of tumors. Genome research, 27(11):1885–1894, 2017b.
- David Lähnemann, Johannes Köster, Ewa Szczurek, Davis J McCarthy, Stephanie C Hicks, Mark D Robinson, Catalina A Vallejos, Kieran R Campbell, Niko Beerenwinkel, Ahmed Mahfouz, et al. Eleven grand challenges in single-cell data science. Genome biology, 21(1):1–35, 2020.
- Scott M Lanyon. The stochastic mode of molecular evolution: what consequences for systematic investigations? The Auk, 105(3):565–573, 1988.
- Jeongwoo Lee, Daehee Hwang, et al. Single-cell multiomics: technologies and data analysis methods. Experimental & Molecular Medicine, 52(9):1428–1442, 2020.
- Marco L Leung, Alexander Davis, Ruli Gao, Anna Casasent, Yong Wang, Emi Sei, Eduardo Vilar, Dipen Maru, Scott Kopetz,





- Andrew Rambaut, Alexei J Drummond, Dong Xie, Guy Baele, and Marc A Suchard. Posterior summarization in bayesian phylogenetics using tracer 1.7. Systematic biology, 67(5):901, 2018.
- Edith M Ross and Florian Markowetz. Onconem: inferring tumor evolution from single-cell sequencing data. Genome biology, 17(1):1–14, 2016.
- Michael G Ross, Carsten Russ, Maura Costello, Andrew Hollinger, Niall J Lennon, Ryan Hegarty, Chad Nusbaum, and David B Jaffe. Characterizing and measuring bias in sequence data. Genome biology, 14(5):1–20, 2013.
- Gryte Satas, Simone Zaccaria, Geoffrey Mon, and Benjamin J Raphael. Scarlet: Single-cell tumor phylogeny inference with copy-number constrained mutation losses. Cell Systems, 10(4):323–332, 2020.
- Russell Schwartz and Alejandro A Schäffer. The evolution of tumour phylogenetics: principles and practice. Nature Reviews Genetics, 18(4):213–229, 2017.
- Alexandros Stamatakis. Raxml version 8: a tool for phylogenetic analysis and post-analysis of large phylogenies. Bioinformatics, 30(9):1312–1313, 2014.
- Jeet Sukumaran and Mark T Holder. Dendropy: a python library for phylogenetic computing. Bioinformatics, 26(12):1569–1571, 2010.
- Ge Tan, Matthieu Muffato, Christian Ledergerber, Javier Herrero, Nick Goldman, Manuel Gil, and Christophe Dessimoz. Current methods for automated filtering of multiple sequence alignments frequently worsen single-gene phylogenetic inference. Systematic biology, 64(5):778–791, 2015.
- Maxime Tarabichi, Adriana Salcedo, Amit G Deshwar, Máire Ni Leathlobhair, Jeff Wintersinger, David C Wedge, Peter Van Loo, Quaid D Morris, and Paul C Boutros. A practical guide to cancer subclonal reconstruction from dna sequencing. Nature methods, 18(2):144–155, 2021.
- Simon Tavaré et al. Some probabilistic and statistical problems in the analysis of dna sequences. Lectures on mathematics in the life sciences, 17(2):57–86, 1986.
- Yuan Tian and Laura Kubatko. Rooting phylogenetic trees under the coalescent model using site pattern probabilities. BMC evolutionary biology, 17(1):1–11, 2017.
- Timothy G Vaughan, Denise Kühnert, Alex Poppinga, David Welch, and Alexei J Drummond. Efficient bayesian inference under the structured coalescent. Bioinformatics, 30(16):2272–2279, 2014.
- Li-Gen Wang, Tommy Tsan-Yuk Lam, Shuangbin Xu, Zehan Dai, Lang Zhou, Tingze Feng, Pingfan Guo, Casey W Dunn, Bradley R Jones, Tyler Bradley, et al. Treeio: an r package for phylogenetic tree input and output with richly annotated

- and associated data. Molecular biology and evolution, 37(2):599–603, 2020.
- Yong Wang and Nicholas E Navin. Advances and applications of single-cell sequencing technologies. Molecular cell, 58(4):598–609, 2015.
- Michael L Waskom. Seaborn: statistical data visualization. Journal of Open Source Software, 6(60):3021, 2021.
- Hadley Wickham. ggplot2: Elegant Graphics for Data Analysis. Springer-Verlag New York, 2016. ISBN 978-3-319-24277-4. URL <https://ggplot2.tidyverse.org>.
- Mollie B Woodworth, Kelly M Girsakis, and Christopher A Walsh. Building a lineage from single cells: genetic techniques for cell lineage tracking. Nature Reviews Genetics, 18(4):230, 2017.
- Guangchuang Yu, David K Smith, Huachen Zhu, Yi Guan, and Tommy Tsan-Yuk Lam. ggtree: an r package for visualization and annotation of phylogenetic trees with their covariates and other associated data. Methods in Ecology and Evolution, 8(1):28–36, 2017.
- Hamim Zafar, Anthony Tzen, Nicholas Navin, Ken Chen, and Luay Nakhleh. Sifit: inferring tumor trees from single-cell sequencing data under finite-sites models. Genome biology, 18(1):178, 2017.
- Hamim Zafar, Nicholas Navin, Ken Chen, and Luay Nakhleh. Siclonofit: Bayesian inference of population structure, genotype, and phylogeny of tumor clones from single-cell genome sequencing data. Genome research, 29(11):1847–1859, 2019.