

(2021), 0, 0, pp. 1–37
doi:10.1093//output

A Phylogenetic Approach to Inferring the Order in Which Mutations Arise during Cancer Progression

YUAN GAO*

Division of Biostatistics, The Ohio State University, 1958 Neil Ave, Columbus, OH 43210, US

gao.957@osu.edu

JEFF GAITHER

Institute for Genomic Medicine, Nationwide Children's Hospital, 700 Childrens Dr., Columbus, OH 43205,

US

JULIA CHIFMAN

Dept of Mathematics and Statistics, American University, 3501 Nebraska Ave NW, Washington 20016, US

LAURA KUBATKO

Mathematical Biosciences Institute, The Ohio State University, 1735 Neil Ave, Columbus, OH 43210, US

Depts of Statistics and Evolution, Ecology, and Organismal Biology, The Ohio State University, 1958 Neil Ave, Columbus, OH 43210, US

1 SUMMARY

2 Although the role of evolutionary processes in cancer progression is widely accepted, increasing attention
3 is being given to evolutionary mechanisms that can lead to differences in clinical outcome. Recent studies
4 suggest that the temporal order in which somatic mutations accumulate during cancer progression is im-
5 portant. Single-cell sequencing provides a unique opportunity to examine the mutation order during cancer
6 progression. However, the errors associated with single-cell sequencing complicate this task. We propose a
7 new method for inferring the order in which somatic mutations arise within a tumor using noisy single-cell
8 sequencing data that incorporates the errors that arise from the data collection process. Using simulation, we

*To whom correspondence should be addressed.

9 show that our method outperforms existing methods for identifying mutation order in most cases, especially
10 when the number of cells is large. Our method also provides a means to quantify the uncertainty in the
11 inferred mutation order along a fixed phylogeny. We apply our method to empirical data for colorectal and
12 prostate cancer.

13 *Key words:* Bayesian inference; Cancer evolution; Error effects quantification; Mutation order; Single-cell sequencing.

14

1. INTRODUCTION

15 Cancer progression is a dynamic evolutionary process that occurs among the individual cells within each
16 patient's tumor. Cancer develops from a single cell in normal tissue whose genetic alterations endow a growth
17 advantage over the surrounding cells, allowing that cell to replicate and to expand, resulting in the formation
18 of a clonal population of identical cells. Cells within this clonal population may then undergo their own
19 somatic mutations, followed by replication and formation of subclones. During this complex process, many
20 competitive and genetically diverse subpopulations may be formed, resulting in intratumoral heterogeneity
21 (ITH) depicted in Fig. 1(a) (O'Sullivan *and others*, 2003; Ishwaran *and others*, 2009; Jamal-Hanjani *and*
22 *others*, 2017; Ascolani and Liò, 2019). Ortmann *and others* (2015) demonstrate that the type of malignancy
23 and the response to treatment of myeloproliferative neoplasm patients are affected by the order in which
24 somatic mutations arose within the patients' tumors. Though this study is specific to one type of cancer,
25 the timing and organization of somatic mutations are crucial to clinical outcomes for other cancers as well.
26 Determining the temporal order of mutations required for tumor progression is thus critical, especially in
27 the field of targeted therapeutics. However, this information cannot be observed directly, since genomic data
28 is most often collected at one snapshot in time. Consequently, use of computational methods that infer the
29 order of mutations from DNA sequence data is the approach of choice.

30 Most studies on cancer phylogenetics utilize bulk high-throughput sequencing data, but signals from
31 bulk sequencing only reflect the overall characteristics of a population of sequenced cells, rather than the
32 characteristics of individual cells. Variation in the mutational signatures among different cells in a tumor is
33 thus difficult to evaluate from bulk sequencing data. Single-cell sequencing (SCS) data is promising because
34 it enables sequencing of individual cells, thus providing the highest possible resolution available on the
35 mutational history of cancer. However, the high error probabilities associated with SCS data complicate the

36 development of methods for inference of the mutational history. The whole-genome amplification (WGA)
37 process used to produce SCS data results in a variety of errors, including allelic dropout (ADO) errors, false
38 positives (FP), non-uniform coverage distribution, and low coverage regions. ADO contributes a considerable
39 number of false negatives (FN) in point mutations (Navin, 2014).

40 Recently, several studies have proposed various mathematical methods to infer mutation order (Fig. 1(c)
41 - Fig. 1(e)) from data arising from single-cell somatic mutations. Of particular interest are the methods of
42 Jahn *and others* (2016) and Zafar *and others* (2017), called SCITE and SiFit, respectively. SiFit uses an
43 MCMC approach as a heuristic to find the maximum likelihood tree from imperfect SCS data. Based on the
44 inferred tumor phylogenetic tree, SiFit estimates the mutation order by estimating the most likely mutation
45 status of the tips and the internal nodes using a dynamic programming algorithm. Although both SCITE and
46 SiFit by default output only the order of the mutations, both can be used to account for uncertainty in the
47 inferred order. For example, because SCITE uses an MCMC algorithm for inference, the posterior probability
48 associated with various mutation orders can be obtained by examining the frequency with which these orders
49 are sampled by the MCMC algorithm. Similarly, the authors of SiFit recently developed a method called
50 SiCloneFit (Zafar *and others*, 2019) that utilizes MCMC to sample trees, and thus the algorithm from SiFit
51 for inferring mutation order on a fixed tree could be applied to a posterior sample of trees to measure the
52 uncertainty in the mutation order that results from uncertainty in the true tumor phylogeny.

53 In this paper, we propose a novel method for inferring the order in which mutations arise within an
54 individual tumor given SCS data from the tumor at a single time point. Our approach utilizes models for
55 both the mutational process within the tumor and the errors that arise during SCS data collection in a
56 Bayesian framework, thus allowing us to quantify the uncertainty in the inferred mutation orders along
57 a fixed tumor phylogeny. Our approach thus represents a conceptually distinct and practically important
58 extension of earlier methods.

59 2. METHODS

60 We assume that we are given a phylogenetic tree with branch lengths that displays the evolutionary rela-
61 tionships among a sample of J cells within a tumor. To infer the locations (branches) on which a set of
62 somatic mutations are acquired in the tree, we need to model the evolutionary process of the somatic mu-

63 tations and quantify the technical errors that arise from the SCS data collection process. We assume that
64 during the evolutionary process, somatic mutations evolve independently across sites, and each mutation
65 evolves independently on different branches. We also assume that each somatic mutation occurs once along
66 the phylogeny and that no back mutation occurs, so that all descendant cells linked by the mutation branch
67 will harbor the corresponding mutation. When quantifying the effect of errors, we assume that SCS technical
68 errors for mutations are independent of one another.

69 2.1 Notation and terminology

70 Consider somatic mutations of interest at I loci across the genome for a sample of J single cells. The J
71 single cells are sampled from different spatial locations (clones) within the tumor. The mutation data can
72 be either binary or ternary. For binary data, 0 denotes the absence of mutation and 1 means that mutation
73 is present, while for ternary data, 0, 1 and 2 represent the homozygous reference (normal), heterozygous
74 (mutation present) and homozygous non-reference (mutation present) genotypes, respectively.

75 The I somatic mutations evolve along the tumor evolutionary tree \mathcal{T} . Each tip in \mathcal{T} represents one single
76 cell C_j , where $j = 1, \dots, J$. Let $C = \{C_1, \dots, C_J\}$ be the set of the J single cells under comparison. $\mathcal{T} = (T, \mathbf{t})$
77 includes two parts: the tree topology T and a vector of branch lengths \mathbf{t} . The tree topology $T = (V, E)$ is
78 a connected graph without cycles and is composed of nodes and branches, where V is the set of nodes and
79 E is the set of branches. Trees are rooted, and the root r represents the common ancestor (a normal cell
80 without somatic mutations) for all the single cells under comparison. In the context of this paper, all the
81 definitions in the following sections will apply to rooted bifurcating trees. There are $2J - 2$ branches in a
82 rooted bifurcating tree with J tips, i.e., $E = \{e_1, e_2, \dots, e_{2J-2}\}$. Let v and w be two nodes in the node set V
83 that are connected by the branch x in the branch set E (i.e., $x = \{v, w\}$: v is the immediate ancestor node of
84 w , and x connects v and w). Then the set $U^x(w)$, which includes the node w and all nodes descended from
85 w in \mathcal{T} , is called the *clade induced by w* . The branch x connects the ancestor node v and the clade induced
86 by w , and we define branch x as the *ancestor branch of clade $U^x(w)$* . $E^x(w)$ is a subset of E that includes
87 branches connecting nodes in $U^x(w)$, and $C^x(w)$ are the tips in $U^x(w)$.

88 Let G_{ij} denote the true genotype for the i^{th} genomic site of cell C_j . The i^{th} genomic site will then have
89 a vector $\mathbf{G}_i \in \{0, 1\}^J$ (for binary data) or $\{0, 1, 2\}^J$ (for ternary data) representing its true genotype for all

100 the J cells represented by the tips in the tree, where $i = 1, \dots, I$. Let S_{ij} denote the observed data for the
101 i^{th} genomic site of cell C_j . Due to the technical errors associated with SCS sequencing, the observed data
102 S_{ij} does not always equal the true genotype G_{ij} . For both binary and ternary data, the observed state S_{ij}
103 might be flipped with respect to the true mutation G_{ij} due to FP or FN. Missing states (“-”) or low-quality
104 states (“?”) may be present for some genomic sites, as well. Fig. 2 shows an example of true and observed
105 binary genotype data for the mutations in Fig. 1. In Fig. 2, the observed state is highlighted in red color
106 if it is not consistent with the true genotype. The red numbers are those mutations with flipped observed
107 mutation states relative to the true mutation states. The red dash (“-”) indicates a missing value and the red
108 question mark (“?”) represents a low-quality value.

109 Mathematically, we represent the observed mutation states of the J single cells at I different genomic
110 sites by an $I \times J$ mutation matrix \mathbf{S} for convenience,

$$\mathbf{S} = \begin{pmatrix} \mathbf{S}_1 \\ \mathbf{S}_2 \\ \vdots \\ \mathbf{S}_I \end{pmatrix} = \begin{pmatrix} S_{11} & \dots & S_{1J} \\ S_{21} & \dots & S_{2J} \\ \vdots & \ddots & \vdots \\ S_{I1} & \dots & S_{IJ} \end{pmatrix}. \quad (2.1)$$

111 Each entry (i, j) denotes the state observed for mutation i in cell C_j , so \mathbf{S}_i gives the observed data for genomic
112 site i as a vector with J values corresponding to the J single cells. Column j represents the mutations of
113 interest for cell C_j . In \mathcal{T} , let \mathcal{B} be the vector of locations (branches) on which the I mutations occur, i.e.,
114 $\mathcal{B} = \{B_1, \dots, B_I\}$, where B_i is the branch on which mutation i is acquired. Note that B_i takes values in
115 $\{e_1, e_2, \dots, e_{2J-2}\}$.

116 2.2 Somatic mutation process

117 To model the somatic mutation process, we consider continuous-time Markov processes, which we specify by
118 assigning a rate to each possible transition between states. We consider point mutations. Once a mutation i
119 is acquired on a branch $x \in E$, all the branches in the set $E^x(w)$ will harbor mutation i but those branches
120 in the set $E \setminus (x \cup E^x(w))$ will not carry this mutation. Specification of the rates of mutation among states
121 allows for flexibility in the modeling procedure.

122 2.2.1 **Binary genotype data** For binary genotype data, the mutation process can be modeled by the
123 2×2 instantaneous rate matrix

$$\mathcal{Q}_\lambda = \begin{matrix} & 0 & 1 \\ \begin{matrix} 0 \\ 1 \end{matrix} & \begin{pmatrix} -\lambda & \lambda \\ 0 & 0 \end{pmatrix} \end{matrix}, \quad (2.2)$$

114 where λ denotes the instantaneous transition rate per genomic site. The transition probability matrix $P(t)$
 115 along a branch of length t is then computed by matrix exponentiation of the product of \mathcal{Q}_λ and the branch
 116 length t , which gives

$$P(t) = \begin{matrix} & 0 & 1 \\ \begin{matrix} 0 \\ 1 \end{matrix} & \begin{pmatrix} P_{00}(t) & P_{01}(t) \\ P_{10}(t) & P_{11}(t) \end{pmatrix} \end{matrix} = \begin{matrix} & 0 & 1 \\ \begin{matrix} 0 \\ 1 \end{matrix} & \begin{pmatrix} \exp(-\lambda t) & 1 - \exp(-\lambda t) \\ 0 & 1 \end{pmatrix} \end{matrix}. \quad (2.3)$$

117 Note that $P_{01}(t)$ is the probability that mutation i is acquired along a branch of length t . Under this
 118 model and recalling that each mutation evolves independently along different branches in \mathcal{T} , the marginal
 119 probability that mutation i is acquired on branch $x \in E$, denoted by $P(B_i = x | \mathcal{T}, \mathcal{Q}_\lambda)$, is thus given by

$$P(B_i = x | \mathcal{T}, \mathcal{Q}_\lambda) = \frac{\left[\prod_{B \in [E \setminus (x \cup E^x(w))]} P_{00}(t_B) \right] P_{01}(t_x) \left[\prod_{B \in E^x(w)} P_{11}(t_B) \right]}{\sum_{z \in E} \left(\left[\prod_{B \in [E \setminus (z \cup E^z(h))]} P_{00}(t_B) \right] P_{01}(t_z) \left[\prod_{B \in E^z(h)} P_{11}(t_B) \right] \right)}, \quad (2.4)$$

120 where t_B is length of branch B . In the numerator, the first term is a product of probabilities over all branches
 121 without the mutation, the second term is the probability that the mutation is acquired on branch x , and the
 122 third term is a product of probabilities over all branches with the mutation, i.e., all branches in $E^x(w)$. The
 123 denominator is needed to create a valid probability distribution over all possible branches, and is obtained
 124 by summing the numerator over all valid branches $z \in E$. The $P(B_i = x | \mathcal{T}, \mathcal{Q}_\lambda)$ term is normalized by the
 125 denominator because we exclude two possible outcomes: a mutation is not acquired on any branch in \mathcal{T} , or
 126 a mutation is acquired more than once on different branches in \mathcal{T} .

127 As an example, Fig. 3 depicts the observed and true binary genotype for mutation $i = 1$ shown in
 128 Fig. 2. The set of branches is $E = \{e_1, \dots, e_8\}$ and the corresponding set of branch lengths would be
 129 $\mathbf{t} = \{t_1, \dots, t_8\}$. If mutation i is acquired on branch e_1 , the cell descending along branch e_8 will not carry
 130 the mutation, while those descending from the blue branches would carry this mutation. The marginal
 131 probability that mutation $i = 1$ is acquired on branch e_1 would be proportional to its numerator, i.e.,
 132 $P(B_1 = e_1 | \mathcal{T}, \mathcal{Q}_\lambda) \propto P_{00}(t_8)P_{01}(t_1)[P_{11}(t_2)P_{11}(t_3)P_{11}(t_4)P_{11}(t_5)P_{11}(t_6)P_{11}(t_7)]$.

133 2.2.2 **Ternary genotype data** The mutation model for ternary data is complex and includes three pos-
 134 sible ways that mutation i originates on a branch x in \mathcal{T} :

135 1. The status of mutation i transitions from $0 \rightarrow 1$ on a branch x and there is no further mutation at this
 136 genomic site in \mathcal{T} .

137 2. The status of mutation i transitions directly from $0 \rightarrow 2$ on a branch x in \mathcal{T} .

138 3. The status of mutation i transitions from $0 \rightarrow 1$ on a branch x and then transitions from $1 \rightarrow 2$ on a
 139 branch $y \in E^x(w)$ in \mathcal{T} .

140 We let B_i be the location at which mutation i originates, $B_i^{0 \rightarrow 1}$ would be the branch on which mutation
 141 status transitions from 0 to 1, $B_i^{0 \rightarrow 2}$ is the branch on which mutation status transitions from 0 to 2, and
 142 $B_i^{1 \rightarrow 2}$ is the branch on which mutation status transitions from 1 to 2. If the mutation i occurs on branch x ,
 143 all cells belonging to $C^x(w)$ will carry 1 or 2 mutations. In other words, $G_{ij} = 1$ or 2 for all $C_j \in C^x(w)$ and
 144 $G_{ij} = 0$ for all $C_j \in C \setminus C^x(w)$. We define the instantaneous rate matrix \mathcal{Q}_λ as

$$\mathcal{Q}_\lambda = \begin{matrix} & \begin{matrix} 0 & 1 & 2 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \end{matrix} & \begin{pmatrix} -(\lambda_1 + \lambda_1 \lambda_2) & \lambda_1 & \lambda_1 \lambda_2 \\ 0 & -\lambda_2 & \lambda_2 \\ 0 & 0 & 0 \end{pmatrix} \end{matrix}, \quad (2.5)$$

145 where λ_1 and λ_2 denote the instantaneous transition rate per genomic site of the transitions $0 \rightarrow 1$ and
 146 $1 \rightarrow 2$, respectively. Studies have provided evidence that direct mutation of $0 \rightarrow 2$ at rate $\lambda_1 \lambda_2$ is possible
 147 in principle, although it is extremely rare (Iwasa *and others*, 2004). If λ_2 is 0 in Expression (2.5), the model
 148 will be reduced to the infinite sites diploid model. The transition probability matrix $P(t) = \exp(\mathcal{Q}_\lambda t)$ is then
 149 given by

$$P(t) = \begin{matrix} & \begin{matrix} 0 & 1 & 2 \end{matrix} \\ \begin{matrix} 0 \\ 1 \\ 2 \end{matrix} & \begin{pmatrix} \exp(-(\lambda_1 + \lambda_1 \lambda_2)t) & \frac{\lambda_1(\exp(-(\lambda_1 + \lambda_1 \lambda_2)t) - \exp(-\lambda_2 t))}{\lambda_2 - (\lambda_1 + \lambda_1 \lambda_2)} & \frac{(\lambda_1 \lambda_2 - \lambda_2) \exp(-(\lambda_1 + \lambda_1 \lambda_2)t) + \lambda_1 \exp(-\lambda_2 t)}{\lambda_2 - (\lambda_1 + \lambda_1 \lambda_2)} + 1 \\ 0 & \exp(-\lambda_2 t) & 1 - \exp(-\lambda_2 t) \\ 0 & 0 & 1 \end{pmatrix} \end{matrix}. \quad (2.6)$$

150 The marginal probability that mutation i originates on branch $x \in E$ for these three possible conditions is
 151 thus given by

$$P(B_i^{0 \rightarrow 1} = x | \mathcal{T}, \mathcal{Q}_\lambda) = \frac{Q(B_i^{0 \rightarrow 1} = x)}{\sum_{z_1 \in E} [Q(B_i^{0 \rightarrow 1} = z_1) + Q(B_i^{0 \rightarrow 2} = z_1) + \sum_{z_2} Q(B_i^{0 \rightarrow 1} = z_1, B_i^{1 \rightarrow 2} = z_2)]}, \quad (2.7)$$

$$P(B_i^{0 \rightarrow 2} = x | \mathcal{T}, \mathcal{Q}_\lambda) = \frac{Q(B_i^{0 \rightarrow 2} = x)}{\sum_{z_1 \in E} [Q(B_i^{0 \rightarrow 1} = z_1) + Q(B_i^{0 \rightarrow 2} = z_1) + \sum_{z_2} Q(B_i^{0 \rightarrow 1} = z_1, B_i^{1 \rightarrow 2} = z_2)]}, \quad (2.8)$$

$$P(B_i^{0 \rightarrow 1} = x, B_i^{1 \rightarrow 2} = y | \mathcal{T}, \mathcal{Q}_\lambda) = \frac{Q(B_i^{0 \rightarrow 1} = x, B_i^{1 \rightarrow 2} = y)}{\sum_{z_1 \in E} [Q(B_i^{0 \rightarrow 1} = z_1) + Q(B_i^{0 \rightarrow 2} = z_2) + \sum_{z_2} Q(B_i^{0 \rightarrow 1} = z_1, B_i^{1 \rightarrow 2} = z_2)]}, \quad (2.9)$$

152 where

$$Q(B_i^{0 \rightarrow 1} = x) = \left[\prod_{B \in [E \setminus (x \cup E^x(w))]} P_{00}(t_B) \right] P_{01}(t_x) \left[\prod_{B \in E^x(w)} P_{11}(t_B) \right], \quad (2.10)$$

153

$$Q(B_i^{0 \rightarrow 2} = x) = \left[\prod_{B \in [E \setminus (x \cup E^x(w))]} P_{00}(t_B) \right] P_{02}(t_x) \left[\prod_{B \in E^x(w)} P_{22}(t_B) \right], \quad (2.11)$$

154

$$Q(B_i^{0 \rightarrow 1} = x, B_i^{1 \rightarrow 2} = y) = \left[\prod_{B \in [E \setminus (x \cup E^x(w))]} P_{00}(t_B) \right] P_{01}(t_x) \left[\prod_{B \in [E^x(w) \setminus (y \cup E^y(b))]} P_{11}(t_B) \right] P_{12}(t_y) \left[\prod_{B \in E^y(b)} P_{22}(t_B) \right]. \quad (2.12)$$

155 As for binary data, we normalize the marginal probabilities to exclude scenarios in which mutations are
 156 acquired more than once or in which mutations are not acquired in \mathcal{T} . As an example, Fig. S1 in the
 157 Supplementary Material depicts the same mutation as in Fig. 3, but considers ternary data, leading to the
 158 following:

159 1. The marginal probability that mutation i transitions from $0 \rightarrow 1$ on branch e_1 is $P(B_i^{0 \rightarrow 1} = e_1 | \mathcal{T}, \mathcal{Q}_\lambda) \propto$

160 $P_{00}(t_8) P_{01}(t_1) [P_{11}(t_2) P_{11}(t_3) P_{11}(t_4) P_{11}(t_5) P_{11}(t_6) P_{11}(t_7)].$

161 2. The marginal probability that mutation i transitions from $0 \rightarrow 2$ on branch e_1 is $P(B_i^{0 \rightarrow 2} = e_1 | \mathcal{T}, \mathcal{Q}_\lambda) \propto$

162 $P_{00}(t_8) P_{02}(t_1) [P_{22}(t_2) P_{22}(t_3) P_{22}(t_4) P_{22}(t_5) P_{22}(t_6) P_{22}(t_7)].$

163 3. The marginal probability that mutation i transitions from $0 \rightarrow 1$ on e_1 , and from $1 \rightarrow 2$ on e_3 is

164 $P(B_i^{0 \rightarrow 1} = e_1, B_i^{1 \rightarrow 2} = e_3 | \mathcal{T}, \mathcal{Q}_\lambda) \propto P_{00}(t_8) P_{01}(t_1) P_{11}(t_2) P_{12}(t_3) P_{22}(t_4) P_{22}(t_5) P_{22}(t_6) P_{22}(t_7).$

165 The probability $P(B_i^{0 \rightarrow 1} = e_1, B_i^{1 \rightarrow 2} = e_3 | \mathcal{T}, \mathcal{Q}_\lambda)$ is the marginal probability that two mutations at the same
 166 locus along the genome mutate on two branches e_1 and e_3 , respectively. After the first mutation occurs on
 167 branch e_1 , the second mutation can occur on any branch except e_1 and e_8 .

168

2.3 Quantification of SCS errors

169 To account for FPs and FNs in the observed SCS data, our method applies the error model for binary and
 170 ternary data from Kim and Simon (2014), Jahn *and others* (2016), and Zafar *and others* (2017). Let α_{ij} be
 171 the probability of a false positive error and β_{ij} be the probability of a false negative error for genomic site i
 172 of cell C_j .

173 For binary data, if the true genotype is 0, we may observe a 1, which is a false positive error. If the
 174 true genotype is 1, we may observe a 0, which is a false negative error. The conditional probabilities of the
 175 observed data given the true genotype at genomic site i of cell C_j are

$$\mathbf{N}^{ij} = \begin{matrix} & S_{ij} = 0 & S_{ij} = 1 \\ \begin{matrix} G_{ij} = 0 \\ G_{ij} = 1 \end{matrix} & \begin{pmatrix} 1 - \alpha_{ij} & \alpha_{ij} \\ \beta_{ij} & 1 - \beta_{ij} \end{pmatrix} \end{matrix}, \quad (2.13)$$

176 where $\mathbf{N}_{01}^{ij} = P(S_{ij} = 1 | G_{ij} = 0) = \alpha_{ij}$, and other entries are defined similarly. Under the assumption that
 177 sequencing errors are independent, if mutation i is acquired on branch x , we can precisely quantify the effect
 178 of SCS technical errors for mutation i as

$$P(\mathbf{S}_i | B_i = x, \mathcal{T}, \mathbf{N}^i) = \prod_{j=1}^J P(S_{ij} | G_{ij}), \quad (2.14)$$

179 where $\mathbf{N}^i = \{\mathbf{N}^{i1}, \dots, \mathbf{N}^{iJ}\}$. Using the example in Fig. 3, the error probability of the observed genotype condi-
 180 tioning on the mutation $i = 1$ occurring on branch e_1 would be $P(\mathbf{S}_1 | B_1 = e_1, \mathcal{T}, \mathbf{N}^1) = \mathbf{N}_{11}^{11} \mathbf{N}_{10}^{12} \mathbf{N}_{11}^{13} \mathbf{N}_{10}^{14} \mathbf{N}_{00}^{15}$,
 181 where $\mathbf{N}^1 = \{\mathbf{N}^{11}, \dots, \mathbf{N}^{15}\}$ for this binary data example.

182 For ternary data, the conditional probabilities of the observed data given the true genotype are given by

$$\mathbf{N}^{ij} = \begin{matrix} & S_{ij} = 0 & S_{ij} = 1 & S_{ij} = 2 \\ \begin{matrix} G_{ij} = 0 \\ G_{ij} = 1 \\ G_{ij} = 2 \end{matrix} & \begin{pmatrix} 1 - \alpha_{ij} - \alpha_{ij}\beta_{ij}/2 & \alpha_{ij} & \alpha_{ij}\beta_{ij}/2 \\ \beta_{ij}/2 & 1 - \beta_{ij} & \beta_{ij}/2 \\ 0 & 0 & 1 \end{pmatrix} \end{matrix}, \quad (2.15)$$

183 where $\mathbf{N}_{01}^{ij} = P(S_{ij} = 1 | G_{ij} = 0) = \alpha_{ij}$, and the other entries are defined similarly. Under the same assump-

184 tions as for binary genotype data, we can precisely quantify the effect of SCS technical errors as in Equation
 185 (2.14) if mutation i is acquired on branch x . Using the example in Fig. S1 in the Supplementary Material,
 186 the error probabilities for the three possible ways that mutation $i = 1$ may arise on branch e_1 are

187 1. The error probability under the condition that the true mutation transitions from $0 \rightarrow 1$ on branch e_1
 188 is $P(\mathbf{S}_1|B_i^{0 \rightarrow 1} = e_1, \mathcal{T}, \mathbf{N}^1) = \mathbf{N}_{12}^{11}\mathbf{N}_{10}^{12}\mathbf{N}_{11}^{13}\mathbf{N}_{10}^{14}\mathbf{N}_{00}^{15}$.

189 2. The error probability under the condition that the true mutation transitions from $0 \rightarrow 2$ on branch e_1
 190 is $P(\mathbf{S}_1|B_i^{0 \rightarrow 2} = e_1, \mathcal{T}, \mathbf{N}^1) = \mathbf{N}_{22}^{11}\mathbf{N}_{20}^{12}\mathbf{N}_{21}^{13}\mathbf{N}_{20}^{14}\mathbf{N}_{00}^{15}$.

191 3. The error probability under the condition that the true mutation transitions from $0 \rightarrow 1$ on branch e_1 ,
 192 and transitions from $1 \rightarrow 2$ on branch e_3 is $P(\mathbf{S}_1|B_i^{0 \rightarrow 1} = e_1, B_i^{1 \rightarrow 2} = e_3, \mathcal{T}, \mathbf{N}^1) = \mathbf{N}_{12}^{11}\mathbf{N}_{20}^{12}\mathbf{N}_{21}^{13}\mathbf{N}_{20}^{14}\mathbf{N}_{00}^{15}$.

193 And $\mathbf{N}^1 = \{\mathbf{N}^{11}, \dots, \mathbf{N}^{15}\}$ for this ternary data example. The term $P(\mathbf{S}_1|B_i^{0 \rightarrow 1} = e_1, B_i^{1 \rightarrow 2} = e_3, \mathcal{T}, \mathbf{N}^1)$
 194 gives the error probability for the case in which the two mutations at the same locus occur on branches e_1
 195 and e_3 .

196 2.4 *Missing and low-quality data*

197 In real data, missing and low-quality states are observed and must be taken into account. For each mutation
 198 i , we exclude cells with missing states, and a subtree \mathcal{T}_i from \mathcal{T} is extracted. The number of tips J_i in subtree
 199 \mathcal{T}_i is less than or equal to J . Let E_i be the set of branches in subtree \mathcal{T}_i . The probability that mutation
 200 i occurs on branch x is then given by $P(B_i = x|\mathcal{T}, \mathcal{Q}_\lambda) = P(B_i = x|\mathcal{T}_i, \mathcal{Q}_\lambda)$, where $P(B_i = x|\mathcal{T}_i, \mathcal{Q}_\lambda)$ is
 201 computed based on branches in the subtree \mathcal{T}_i , and $P(B_i = x|\mathcal{T}_i, \mathcal{Q}_\lambda)$ is 0 for those branches $x \in E \setminus E_i$. We
 202 quantify the effect of the SCS technical errors as

$$P(\mathbf{S}_i|B_i = x, \mathcal{T}, \mathbf{N}^i) = \prod_{j=1}^{J_i} \left(\sum_{S_{ijk}} w_{ijk} P(S_{ijk}|G_{ijk}) \right), \quad (2.16)$$

203 where w_{ijk} is the weight for each possible genotype state at a mutation site. For a site with an observed
 204 state that is not missing or ambiguous, w_{ijk} is 1 for the observed state and 0 for all other states. For
 205 an ambiguous site, we can assign equal weight for each possible state, or we can assign weight based on
 206 sequencing information or other biological characteristics.

207 2.5 *Inferring the location of a mutation in \mathcal{T}*

208 Once the observed status matrix $\mathbf{S} = [\mathbf{S}_1 \dots \mathbf{S}_I]^T$ of the I mutations has been collected, the next step is to
 209 infer the branch on which mutation i takes place, conditioning on \mathbf{S} . Given the observed data matrix \mathbf{S} , the

210 tumor phylogenetic tree \mathcal{T} , the error probability matrix $\mathbf{N} = \{\mathbf{N}^{ij} | 1 \leq i \leq I, 1 \leq j \leq J\}$, and the mutation
 211 process \mathcal{Q}_λ , we can assign a posterior probability distribution $P(B_i | \mathbf{S}, \mathcal{T}, \mathbf{N}, \mathcal{Q}_\lambda)$ to the location of mutation
 212 i using Bayes' Theorem,

$$P(B_i = x | \mathbf{S}, \mathcal{T}, \mathbf{N}, \mathcal{Q}_\lambda) = \frac{P(\mathbf{S}_i | B_i = x, \mathcal{T}_i, \mathbf{N}^i) P(B_i = x | \mathcal{T}_i, \mathcal{Q}_\lambda)}{P(\mathbf{S}_i | \mathcal{T}_i, \mathbf{N}^i, \mathcal{Q}_\lambda)}. \quad (2.17)$$

213 For mutation i , $P(B_i = x | \mathbf{S}_i, \mathcal{T}_i, \mathbf{N}^i, \mathcal{Q}_\lambda)$ is computed for all x in set E_i . For example, there are 8 branches in
 214 the tree in Fig. 3, so the branch on which mutation $i = 1$ occurs, B_1 , can be any of the 8 branches. For the
 215 binary example, the posterior probability that mutation $i = 1$ occurs on e_1 is $P(B_1 = e_1 | \mathbf{S}_1, \mathcal{T}_1, \mathbf{N}^1, \mathcal{Q}_\lambda) \propto$
 216 $P_{00}(t_8) P_{01}(t_1) [P_{11}(t_2) P_{11}(t_3) P_{11}(t_4) P_{11}(t_5) P_{11}(t_6) P_{11}(t_7)] \cdot \mathbf{N}_{11}^{11} \mathbf{N}_{10}^{12} \mathbf{N}_{11}^{13} \mathbf{N}_{10}^{14} \mathbf{N}_{00}^{15}$. In this way, the posterior
 217 probability that the mutation occurs on each of the 8 branches can be computed, giving the probability
 218 distribution for the location of mutation $i = 1$, i.e. $P(B_1 = x | \mathbf{S}_1, \mathcal{T}_1, \mathbf{N}^1, \mathcal{Q}_\lambda)$ for $x \in \{e_1, \dots, e_8\}$.

219 To summarize this probability distribution, we construct a $(1 - \theta) \times 100\%$ credible set for the location of
 220 mutation i as follows. First, the branches are ranked by their posterior probabilities, and then branches are
 221 added to the credible set in the order of decreasing posterior probability until the sum of their probabilities
 222 reaches $(1 - \theta)$. The number of branches in the credible set is informative about the level of certainty
 223 associated with the inferred location for the mutation. To obtain a point estimate, we pick the branch that
 224 maximizes the posterior probability, i.e., the maximum a posteriori (MAP) estimate. The MAP estimator
 225 for the location of mutation i is given by

$$\hat{B}_{i_{MAP}} = \operatorname{argmax}_{B_i \in \{e_1, \dots, e_{2J-2}\}} P(B_i | \mathbf{S}_i, \mathcal{T}, \mathbf{N}^i, \mathcal{Q}_\lambda). \quad (2.18)$$

226 For the example in Fig. 3, the branch with the largest posterior probability is $\hat{B}_{1_{MAP}}$ for mutation $i = 1$.

227 2.6 Inferring the mutation order in \mathcal{T}

228 We now consider the joint posterior probability distribution of the locations for the I mutations in the
 229 sample of J single cells, which is a distribution on a set of cardinality $(2J - 2)^I$. Based on the assumption
 230 of independence among the I mutations being considered, the posterior distribution for \mathcal{B} is given by

$$P(\mathcal{B} | \mathbf{S}, \mathcal{T}, \mathbf{N}, \mathcal{Q}_\lambda) = \prod_{i=1}^I P(B_i | \mathbf{S}_i, \mathcal{T}, \mathbf{N}^i, \mathcal{Q}_\lambda), \quad (2.19)$$

231 where $\mathbf{N}^i = \{\mathbf{N}^{i1}, \dots, \mathbf{N}^{iJ}\}$. From this distribution, we can extract information on the ordering of mutations
 232 of interest. For example, if we are interested in the order of mutation $i = 1$ and mutation $i = 2$ in Fig. 2,

233 the joint posterior probability distribution that mutation $i = 1$ occurs on branch $x \in E$ and mutation $i = 2$
 234 occurs on branch $y \in E$ can be used to find the probability that mutation $i = 1$ occurs earlier in the tree than
 235 mutation y . Note that $P^{B_1=x, B_2=y} = P(B_1 = x, B_2 = y | \mathbf{S}, \mathcal{T}, \mathbf{N}, \mathcal{Q}_\lambda) = P(B_1 = x | \mathbf{S}, \mathcal{T}, \mathbf{N}, \mathcal{Q}_\lambda) \cdot P(B_2 =$
 236 $y | \mathbf{S}, \mathcal{T}, \mathbf{N}, \mathcal{Q}_\lambda)$. This joint distribution can be represented in a matrix given by

$$\begin{matrix}
 & B_1 = e_1 & B_1 = e_2 & \dots & B_1 = e_8 \\
 B_2 = e_1 & \left(P^{B_1=e_1, B_2=e_1} & P^{B_1=e_2, B_2=e_1} & \dots & P^{B_1=e_8, B_2=e_1} \right) \\
 B_2 = e_2 & \left(P^{B_1=e_1, B_2=e_2} & P^{B_1=e_2, B_2=e_2} & \dots & P^{B_1=e_8, B_2=e_2} \right) \\
 \vdots & \left(\vdots & \vdots & \ddots & \vdots \right) \\
 B_2 = e_8 & \left(P^{B_1=e_1, B_2=e_8} & P^{B_1=e_2, B_2=e_8} & \dots & P^{B_1=e_8, B_2=e_8} \right)
 \end{matrix}$$

237 Adding entries of the matrix for which branch e_1 is earlier in the tree than branch e_2 thus gives the probability
 238 that mutation 1 occurs before mutation 2. To measure the uncertainty of the ordering of the mutations, we
 239 rank all possible mutation orders by their posterior probabilities, and construct a $(1 - \theta) \times 100\%$ credible
 240 set by adding orders with decreasing probability until their sum exceeds $1 - \theta$. The MAP estimator for the
 241 order of I mutations is thus given by

$$\hat{B}_{MAP} = \underset{\mathcal{B} \in \{e_1, \dots, e_{2J-2}\}^I}{\operatorname{argmax}} P(\mathcal{B} | \mathbf{S}, \mathcal{T}, \mathbf{N}, \mathcal{Q}_\lambda). \tag{2.20}$$

242

3. SIMULATION STUDY

243 To evaluate the ability of our method, which we call MO (Mutation Order), to correctly identify the locations
 244 and the order of a set of mutations under different conditions, we conduct a series of simulation studies with
 245 data simulated under different assumptions. The goal is to assess the effect of data quality (complete or
 246 incomplete, high or low error probabilities), number of cells, branch lengths, number of mutations and type
 247 of genotype data on the performance of our method. We consider a total of 12 scenarios, with 100 replicates
 248 for each setting within each scenario. Scenarios 1 - 4 involve data generated under our model for either 10
 249 cells (scenarios 1 and 2) or 50 cells (scenarios 3 and 4) for either long branch lengths (scenarios 1 and 3) or
 250 short branch lengths (scenarios 2 and 4) data. Scenarios 5 - 8 consider data simulated under various models
 251 implemented in the CellCoal software (Posada, 2020). Scenarios 9 and 10 involve data generated under
 252 our model, but with mutations placed on branches with varying (rather than equal) probabilities. Finally,
 253 scenarios 11 and 12 consider data simulated under the finite sites assumption (all other simulation settings

254 used the infinite sites assumption). The methods used to simulate data under these different scenarios are
255 described in detail in Sections A.1 to A.4 of the Supplementary Materials, and Section D of the Supplementary
256 Materials provides information about computational requirements.

257 *3.1 Accuracy of MAP estimates*

258 We assess the accuracy of the MAP estimates in MO across the 100 trees within each simulation setting in
259 several ways, including whether the mutation is inferred to occur on the correct branch (“location accuracy”),
260 whether any pair of mutations are inferred to occur in the correct order (“order accuracy”), and whether a
261 pair of mutations that occur on adjacent branches are inferred to occur in the correct order (“adjacent order
262 accuracy”). In evaluating both the order accuracy and adjacent order accuracy, if two sequential mutations
263 are inferred to occur on the same branch, then it is counted as ordering the mutations incorrectly. In addition,
264 pairs of mutations that occur on the same branch are also included in the computation of order accuracy
265 and adjacent order accuracy. The details of how the MAP estimates are assessed are given in Section B of
266 the Supplementary Material. Tables 1 to 4 in the Supplementary Material show the location accuracy for
267 scenarios 1 to 4 with each cell entry corresponding to a unique setting of α , β , type of genotype and missing
268 data percentage. In most cases, the location accuracy of MO is high except when the error probabilities are
269 high. Comparing the effect of the size of the tree (i.e., the number of cells), the accuracy of settings with
270 50 cells (Tables 3 and 4) is as good as for settings with 10 cells (Tables 1 and 2) in most cases. When error
271 probabilities are large, the accuracy of settings with 50 cells in Tables 3 and 4 are slightly lower than those
272 with 10 cells in Tables 1 and 2. Within one table, the accuracy of MAP estimation with ternary genotype
273 data tends to be higher than that of binary genotype data when fixing other parameters. With the same type
274 of genotype and same error probability setting, the accuracy decreases as the percentage of missing values
275 increases. When α (or β) is fixed, accuracy tends to decrease as β (or α) increases. Comparing Tables 1 and
276 2, the accuracy of MAP estimation in Table 1 (each tree has 10 cells and longer branch lengths) tends to
277 be slightly lower than that in Table 2 (each tree has 10 cells and shorter branch lengths) within the same
278 setting with a few exceptions for large α and β , although the difference is very small. Comparing Tables 3
279 and 4, the effects of branch lengths are flipped and the accuracy of MAP estimation in Table 3 tends to be
280 higher than that in Table 4.

281 The results for order accuracy (Tables 5 to 8 in the Supplementary Material) and adjacent order accuracy
282 (Tables 9 to 12 in the Supplementary Material) are similar. In addition to the same overall trends due to
283 number of cells, data type, percentage of missing data and error probabilities, the magnitudes of the order
284 accuracies are higher than the corresponding adjacent order accuracies.

285 The results for location accuracy, order accuracy and adjacent order accuracy of MO in scenarios 5 to
286 10 have similar patterns to those observed for scenarios 1 to 4. The accuracy in scenarios 5 to 10 is not
287 affected by the number of mutations. In addition to the same overall trends due to the number of cells, type
288 of genotype, missing data percentage and error probabilities, the magnitudes of the accuracies in scenarios
289 5 to 10 are higher than the corresponding accuracies in scenarios 1 to 4. Especially when error probabilities
290 are low, the accuracies can be as high as 99%.

291 3.2 *Credible set accuracy*

292 The credible set accuracy of the inferred mutation branch is assessed as well. If the true mutation branch
293 is within the credible set, we count this as correct; otherwise, it is incorrect. We use 95% credible set for
294 computation (Tables 13 to 16 in the Supplementary Material). The credible set accuracies have the same
295 overall trends as the accuracies of MAP estimates due to the number of cells, type of genotype, missing
296 data percentage and error probabilities, though the accuracies are much higher than for the corresponding
297 MAP estimates, especially for settings with large error probabilities and higher missing data percentages.
298 As for the MAP estimates, the overall trends for scenarios 5 to 10 are similar to scenarios 1 to 4, but the
299 corresponding magnitudes of the accuracies in scenarios 5 to 10 are higher than those in scenarios 1 to 4.

300 3.3 *Comparison with competing approaches*

301 To further assess the performance of MO, we compare its performance with the methods SCITE (Jahn *and*
302 *others*, 2016) and SiFit (Zafar *and others*, 2017) for the simulation data in scenarios 1 to 12. SCITE can
303 estimate the order of mutations for either binary or ternary genotype data. We use the maximum likelihood
304 mutation order inferred by SCITE with 1,000,000 iterations given the true error probabilities. SiFit can
305 use either binary or ternary genotype data when inferring the phylogenetic tree, but it can only use binary
306 genotype data when inferring mutation order. We estimate the most likely mutational profiles for the tips,
307 the internal nodes, and the mutation locations by SiFit given the true phylogenetic tree, error probabilities
308 and mutation rates. We then extract the mutation order information from the output. The three methods are

309 compared with respect to the order accuracy and adjacent order accuracy for the above simulation settings.
310 3.3.1 **Scenarios 1 to 4** Fig. 4 to Fig. 7 plot the order accuracy and the adjacent order accuracy for the
311 three methods for scenarios 1 to 4, respectively. In each figure, the top row shows the results for binary data
312 and the bottom row shows the results for ternary data. In each panel, different methods are highlighted with
313 different colors.

314 In scenarios 1 to 4, order accuracy and adjacent order accuracy show general decreasing trends as data
315 quality becomes worse for all three methods. For results estimated from the trees with 10 cells (scenarios 1
316 and 2), MO is comparable to SCITE in terms of order accuracy estimated from binary and ternary data.
317 Only when both α and β are large does SCITE have higher order accuracy rates than MO. Comparing
318 adjacent order accuracy when there are 10 cells in each tree, MO has comparable adjacent order accuracy
319 when estimated from ternary data. MO has lower adjacent order accuracy than SCITE when estimated from
320 binary data but the discrepancies of adjacent order accuracies between MO and SCITE are only 7% on
321 average. When there are 50 cells in each tree (scenarios 3 and 4), MO is superior to SCITE in all settings
322 in terms of order accuracy and adjacent order accuracy estimated from both binary and ternary genotype
323 data. Specifically, the order accuracy for MO is 25% higher than SCITE on average, and the adjacent order
324 accuracy for MO is 20% higher than SCITE on average. In all settings, SiFit has the worst performance since
325 only a subset of the input mutations are inferred to occur on the tree. Although the output partial mutation
326 orders from SiFit are mostly correct, the accuracy is low due to the small number of inferred mutation orders.
327 MO thus dominates SiFit when assessing the performance using order accuracy and adjacent order accuracy.
328 Comparing between settings with 10 cells and those with 50 cells, the performance of MO is consistently
329 good, and the accuracy is slightly higher as the number of cells increases. SiFit performs better as the
330 number of cells increases as well. However, the performance of SCITE becomes worse when the number of
331 cells increases. Although the number of correct pairs inferred by SCITE increases, the accuracy decreases
332 because the total number of true pairs increases.

333 3.3.2 **Scenarios 5 to 8** Fig. S2 and Fig. S3 in the Supplementary Material plot the order accuracy and
334 adjacent order accuracy for scenarios 5 and 6, respectively. In scenarios 5 and 6 where mutations evolve by
335 the infinite sites diploid model, order accuracy and adjacent order accuracy show general decreasing trends
336 as data quality becomes worse for all three methods, as is observed for scenarios 1 to 4. MO is superior to
337 SCITE in all settings in terms of adjacent order accuracy and order accuracy for both the complete and

338 missing data settings. In all the settings, SiFit has the worst performance with respect to order accuracy
339 when there are 10 cells in each tree. However, SiFit has comparable adjacent order accuracy to SCITE when
340 error probabilities are small. Similar to scenarios 1 to 4, only a proportion of mutations are inferred to occur
341 on the tree. MO thus dominates SiFit in scenarios 5 and 6. In all settings for scenarios 5 and 6, the number
342 of mutations and number of tips in the tree do not affect the order accuracy or adjacent order accuracy of
343 MO and SiFit very much. However, the performance of SCITE is affected by the number of mutations. As
344 the number of mutations increases, the accuracy of SCITE becomes worse. In addition, the adjacent order
345 accuracy of SCITE increases as the number of cells increases.

346 Fig. S4 and Fig. S5 in the Supplementary Material plot the order accuracy and the adjacent order accuracy
347 for scenarios 7 and 8, respectively. In scenarios 7 and 8, mutations arise by the infinite sites diploid model,
348 as was the case for scenarios 5 and 6, but now a small proportion of the mutations are lost. Compared to the
349 complete settings in scenarios 5 and 6, the performance of all the three methods becomes worse. However,
350 the performance of the three methods is comparable to settings with missing values in scenarios 5 and 6.

351 In addition to the above comparisons, we also apply MO to data from scenarios 5 and 6 when transition
352 rates are misspecified. Fig. S9 and Fig. S10 show the order accuracy and adjacent order accuracy when MO is
353 applied with misspecified transition rates $\lambda_1 = 1$ and $\lambda_2 = 10^5$. In each panel, red, blue, and green correspond
354 to MO, SCITE, and SiFit, respectively, when the true transition rates ($\lambda_1 = 1$ and $\lambda_2 = 0$) are used, as in
355 the initial analysis in scenarios 5 and 6. Purple color corresponds to MO when the misspecified transition
356 rates are used. The performance of SCITE is not affected by misspecified transition rates. Comparing the
357 plots, we see that when binary data are used, the effect of misspecified transition rates are ignorable, and
358 the accuracy with either the correct or the incorrect transition rates are nearly identical. However, when
359 using ternary data, the differences are noticeable. In scenario 5, the order accuracy for MO with misspecified
360 transition rates is comparable to SCITE when error probabilities are small and higher than SCITE when
361 error probabilities are large. In scenario 6, the order accuracy inferred from ternary genotype data for MO
362 with misspecified transition rates is lower than SCITE. Comparing the adjacent order accuracy with ternary
363 data, the performance of MO with the misspecified transition rates is worse than when the transition rates
364 are correctly specified in MO, but MO still performs better than SCITE.

365 3.3.3 *Scenarios 9 to 10* In scenarios 9 and 10, mutations are simulated under the mutation process
366 defined in Section 2.2. Although the transition rates are the same as in scenarios 1 to 4, each mutation is
367 not equally likely to occur on all of the branches. In Fig. S6 and Fig. S7, we observe that MO has higher
368 accuracy than SCITE and SiFit in all settings in terms of both order accuracy and adjacent order accuracy.

369 3.3.4 *Scenarios 11 to 12* In scenarios 11 and 12, mutations are simulated under the finite sites assump-
370 tion. Because it is unclear how mutation order should be defined when mutations can arise multiple times
371 along a phylogeny, we instead plot the location accuracy of MO and SiFit in Fig. S8. When there are only
372 10 tips in the tree, most simulated mutations occur only once along the tree and MO has higher accuracy
373 than SiFit. However, when there are 50 tips, most are back mutations and/or parallel mutations. SiFit per-
374 forms better than MO when the data are complete and the missing percentage is low. When the missing
375 percentage is high (e.g., 20%), neither MO nor SiFit identify the correct mutation location. MO is limited
376 by its assumption that all mutations occur only once on the tree. Although SiFit can infer parallel/back
377 mutations, it is not able to identify all the locations on which the mutations occur for the simulated data.

378

4. EMPIRICAL EXAMPLES

379 We apply MO to two experimental single-cell DNA sequencing datasets, one for prostate cancer (Su *and*
380 *others*, 2018) and one for metastatic colon cancer patients (Leung *and others*, 2017). For the prostate cancer
381 dataset, we retrieve publicly available data from the single-cell study of Su *and others* (2018), which includes
382 10 single-cell genomes for each patient. For the colon cancer dataset, we use the somatic single nucleotide
383 variants (SNVs) after variant calling provided in the original study (16 SNVs for patient CRC1 and 36 SNVs
384 for patient CRC2) of Leung *and others* (2017).

385

4.1 *Prostate cancer data*

386 4.1.1 *Data analysis* To infer tumor evolutionary trees for patients 1 and 2 (labeled P1 and P2), we
387 use the SVDQuartets method of Chifman and Kubatko (2014) as implemented in PAUP* (Swofford, 1999)
388 using the aligned DNA sequences for all somatic mutations as input with the expected rank of the flattening
389 matrix set to 4. We specify the normal cell sample as the outgroup. We use the maximum likelihood method
390 to estimate the branch lengths.

391 We select common tumor suppressor genes and oncogenes for both P1 and P2 identified by Su *and others*
392 (2018). In addition to these common cancer-associated genes across different cancers, we map mutations in

393 prostate cancer-specific genes (genes that are more commonly mutated in prostate cancer patients) suggested
394 by Barbieri *and others* (2013) and Tate *and others* (2018). For both binary and ternary genotype data for
395 these genes, we use MO to compute the posterior probability of mutation on each branch of the tumor
396 phylogeny for each of the two patients. Su *and others* (2018) estimated the error probabilities to be $(\alpha, \beta) =$
397 $(0.29, 0.02)$ for P1, and $(\alpha, \beta) = (0.31, 0.02)$ for P2. Although our method in Section 2 allows the assignment
398 of varying error probabilities across genomic sites and cells, here we use same probabilities for all sites.

399 To examine the effect of informativeness of the prior on the resulting inference, we consider two priors for
400 each parameter with mean equal to the estimated error probability from the empirical data and with either
401 a large or a small variance as described in Section C in the Supplementary Materials. For P1, we consider
402 $\alpha|\mathbf{S}_i \sim \text{Beta}(0.29, 0.71)$ (larger variance) and $\alpha|\mathbf{S}_i \sim \text{Beta}(2.9, 7.1)$ (smaller variance). For P2, we consider
403 $\alpha|\mathbf{S}_i \sim \text{Beta}(0.31, 0.69)$ (larger variance) and $\alpha|\mathbf{S}_i \sim \text{Beta}(3.1, 6.9)$ (smaller variance). For β for both P1
404 and P2, we consider $\beta|\mathbf{S}_i \sim \text{Beta}(0.02, 0.98)$ (larger variance) and $\beta|\mathbf{S}_i \sim \text{Beta}(0.2, 9.8)$ (smaller variance).

405 According to Iwasa *and others* (2004), the mutation rates for the first and second mutation are estimated
406 to be $\lambda_1 = 10^{-7}$ and $\lambda_2 = 10^{-2}$, respectively. We use these values to specify the priors for the transition
407 rates. Similar to the sequencing error parameter priors, we set two priors for each transition rate with equal
408 means but different variances. The distribution of the transition rate λ_1 ($0 \rightarrow 1$ for ternary genotype) is set as
409 $\lambda_1|\mathbf{S}_i \sim \text{Gamma}(2, 5.0 \times 10^{-8})$ (larger variance) and $\lambda_1|\mathbf{S}_i \sim \text{Gamma}(5, 2.0 \times 10^{-8})$ (smaller variance). The
410 distribution of the transition rate λ_2 ($1 \rightarrow 2$ for ternary genotype) is set as $\lambda_2|\mathbf{S}_i \sim \text{Gamma}(2, 5.0 \times 10^{-3})$
411 (larger variance) and $\lambda_2|\mathbf{S}_i \sim \text{Gamma}(5, 2.0 \times 10^{-3})$ (smaller variance). The estimated probabilities of
412 mutation do not vary substantially when the priors with larger or smaller variance are used for any of these
413 parameters. The heatmaps of estimated probabilities with different priors (larger or smaller variance) are in
414 the Supplementary Material.

415 4.1.2 **Results** Fig. 8 and Fig. 9 show the tumor evolutionary tree estimated for P1 and P2, respectively.
416 In both tumor trees, the trunk connects the tumor clone to the normal clone. We annotate the genes on their
417 inferred mutation branches. The uncertainty in the inferred mutation locations is highlighted with colors.
418 Mutations with strong signal (defined to be a probability larger than 0.7 that the mutation occurred on a
419 single branch) are colored red, while mutations with moderate signal (defined to be a total probability larger
420 than 0.7 on two or three branches) are colored blue.

421 We also compare the estimated posterior probability distributions for each mutation of common cancer-
422 associated genes for patients P1 and P2, which are used to construct credible sets and to measure the
423 uncertainty of the inferred mutation order. Fig. S11 to Fig. S14 in the Supplementary Material are heatmaps
424 for the posterior probability distribution of each mutation for patients P1 and P2 with different priors (larger
425 or smaller variance).

426 Fig. 12 shows heatmaps of the estimated posterior probabilities for prostate cancer-specific genes when the
427 variance of the prior distributions are large for P1 and P2. The corresponding heatmaps for the cases in which
428 the prior distribution has small variance are shown in Fig. S15 and Fig. S16 in the Supplementary Material.
429 In agreement with the results of Su et al. (2018), we find that TP53, a mutation commonly associated with
430 tumor initiation in many cancers (see, e.g., Yu *and others* (2014)), is inferred to occur on the trunk of the
431 tree with high probability in patient P1, but not in patient P2. Gene *ZFHX3* has a high probability of having
432 mutated on the trunk of the tree in both patients. In addition, the data for patient P1 shows strong signal
433 that *FOXP1* mutates on the trunk of the tumor tree, while *BRCA2* has a high probability of having mutated
434 on the trunk of the tree for patient P2. Comparing the heatmaps of common cancer-associated genes with
435 the prostate cancer-specific genes, mutations inferred to have occurred on the trunk of the tree tend to be
436 those that are common across cancer types, while mutations known to have high frequency within prostate
437 cancer are generally found closer to the tips of the tree in both patients.

438 4.2 *Metastatic colorectal cancer data*

439 4.2.1 *Data analysis* The SVDQuartets method of Chifman and Kubatko (2014) is also applied to these
440 data to estimate each colorectal patient's tumor phylogeny. The normal cells in each patient are merged into
441 one normal sample and used as the outgroup. We also merge collections of cells with high similarity (similar
442 mutations). We use the maximum parsimony method to compute the number of changes on each branch and
443 scale the number of changes on each branch by the total number of changes to estimate the branch lengths.

444 The original study of Leung *and others* (2017) reported 16 and 36 SNVs for patients CRC1 and CRC2
445 after variant calling. Leung *and others* (2017) reported error probabilities of $(\alpha, \beta) = (0.0152, 0.0789)$ and
446 $(\alpha, \beta) = (0.0174, 0.1256)$ for CRC1 and CRC2, respectively. For each patient, we use these values to specify
447 the same priors across all sites. For CRC1, we consider $\alpha|\mathbf{S}_i \sim \text{Beta}(0.015, 0.985)$ (larger variance) and
448 $\alpha|\mathbf{S}_i \sim \text{Beta}(0.15, 9.85)$ (smaller variance); and $\beta|\mathbf{s}_i \sim \text{Beta}(0.078, 0.922)$ (larger variance) and $\beta|\mathbf{S}_i \sim$

449 $Beta(0.78, 9.22)$ (smaller variance). For CRC2, we consider $\alpha|\mathbf{S}_i \sim Beta(0.0174, 0.9826)$ (larger variance)
450 and $\alpha|\mathbf{S}_i \sim Beta(0.174, 9.826)$ (smaller variance); and $\beta|\mathbf{S}_i \sim Beta(0.1256, 0.8744)$ (larger variance) and
451 $\beta|\mathbf{S}_i \sim Beta(1.256, 8.744)$ (smaller variance). The priors for the transition rates for CRC1 and CRC2 are
452 same as for P1 and P2. As was found for the prostate cancer patients, the estimated probabilities do not
453 vary substantially when we use priors with small or large variance.

454 4.2.2 **Results** The inferred tumor tree and mutation order are depicted in Fig. 10 and Fig. 11. The poste-
455 rior probabilities of the inferred mutation locations are indicated with colors as for the prostate cancer data,
456 and agree overall with the findings of Leung *and others* (2017). Fig. S17 and Fig. S18 in the Supplementary
457 Material are heatmaps for the posterior probability distribution of each mutation for patients CRC1 and
458 CRC2 with different priors. For patient CRC1, mutations in *APC*, *KRAS* and *TP53* are inferred to have been
459 acquired on the trunk of the tumor phylogeny with high posterior probability, in agreement with Leung *and*
460 *others* (2017) and in agreement with past studies. The studies of Fearon and Vogelstein (1990) and Powell
461 *and others* (1992) have shown that the mutation order of these genes appears to be fixed in initializing col-
462 orectal cancer, providing further support for our findings. In addition, we identify the five mutations specific
463 to metastatic cells that are found by Leung *and others* (2017), with three (*ZNF521*, *TRRAP*, *EYS*) inferred
464 to occur on branch 21 in Fig. 10 and the remaining two (*RBFOX1*, *GATA1*) inferred to occur on branch 29.

465 For CRC2, we identify strong signals on branch 2 in Fig. 11 for 7 genes reported by Leung *and others*
466 (2017) that are shared by primary and metastatic cells, including driver mutations in *APC*, *NRAS* and *TP53*.
467 We also identify an independent lineage of primary diploid cells (colored in pink in Fig. 11) that evolved in
468 parallel with the rest of the tumor with moderate to strong signals for mutations in *ALK*, *ATR*, *EPHB6*,
469 *SPEN* and *NR3C2* and that do not share the mutations listed in the previous sentence. Our analysis further
470 agrees with that of Leung *and others* (2017) in that we also identify the subsequent formation of independent
471 metastatic lineages. For example, on branch 124 we find strong support for mutations in *FUS*; on branch
472 125 we find strong support for mutations in *ATP7B* and *NR4A3*; and on branch 133 mutations in *HELZ*
473 and *PRKCB* are strongly supported. Many of the genes showing weaker or moderate support for mutation
474 in these metastatic lineages agree with those identified by Leung *and others* (2017). The primary difference
475 between our result and that of Leung *and others* (2017) is that we identify mutation in *ATP7B* along a
476 second major metastatic lineage, rather than in the primary tumor.

477

5. DISCUSSION

478 Development of computational tools based on a phylogenetic framework for use in studying cancer evolution
479 has the potential to provide tremendous insight into the mechanisms that lead to ITH, especially the role
480 of the temporal order of mutations in cancer progression. For example, Ortmann et al. (2015) have shown
481 differences in clinical features and the response to treatment for patients with different mutation orders,
482 indicating that inference of the order in which mutations arise within an individual's cancer may have direct
483 implications in clinical oncology, both for diagnostic applications in measuring the extent of ITH and for
484 improving targeted therapy. SCS data provide an unprecedented opportunity to estimate mutation order at
485 the highest resolution. However, such data are subject to extensive technical errors that arise during the
486 process of whole-genome amplification.

487 To analyze such data, we introduce MO, a new Bayesian approach for reconstructing the ordering of
488 mutational events from the imperfect mutation profiles of single cells. MO is designed to infer the temporal
489 order of a collection of mutations of interest based on a phylogeny of cell lineages that allows modeling of
490 the errors at each tip. MO can infer the mutation order that best fits single-cell data sets that are subject
491 to the technical noise that is common for SCS data, including ADO, false positive errors, low-quality data,
492 and missing data. The assumption of independence of mutations made by MO is the same as that made in
493 other methods developed for inferring mutation order (e.g., Zafar *and others* (2017), Zafar *and others* (2019),
494 and Jahn *and others* (2016)). Thus, MO does not presently account for possible interactions between the
495 occurrences of mutations, though it could be extended to accommodate this if biological information about
496 these interactions is available. However, recent work (Canisius *and others*, 2016) indicates that observed
497 dependence typically takes the form of mutual exclusivity (i.e., only one gene in the group will be mutated
498 in any given patient) rather than positive association, making the independence assumption of less concern
499 here, as the set of mutations we study are assumed to be present within an individual patient. MO could
500 also be extended to work on clonal trees and models that include errors in observed data for multiple cells
501 in a tip instead of a single cell. In addition, MO could be modified to account for the accelerated mutation
502 rates common in late-stage cancers, or to allow for back or parallel mutation.

503 An important difference between MO and existing methods, such as SCITE (Jahn *and others*, 2016) and
504 SiFit (Zafar *and others*, 2017), is the mechanism for quantifying uncertainty in the inferred order. Options

505 available within SCITE (Jahn *and others*, 2016) allow for estimation of the posterior probability distribution
506 across orders. SiFit (Zafar *and others*, 2017), on the other hand, could be modified to account for uncertainty
507 in the orders because the true tumor phylogeny is unknown and must first be estimated. In contrast, because
508 MO uses a probabilistic model for inferring mutation locations along a fixed tree, it is able to provide an
509 estimate of uncertainty in the inferred locations conditioning on the correct tumor phylogeny, thus capturing
510 a source of uncertainty that differs from what SCITE and SiFit provide. MO performs accurately, as is
511 evident from a comprehensive set of simulation studies that take into account different aspects of modern
512 SCS data sets by examining a wide range of error probabilities, fractions of missing data, branch lengths,
513 and numbers of cells in each tree. The simulation studies also demonstrate that MO outperforms the state-
514 of-the-art methods when the number of cells is large and performs comparably to other methods when the
515 number of cells is small. MO is robust to the technical errors that arise during whole-genome amplification.
516 When applied to data from two prostate cancer patients and from two colorectal cancer patients, MO is able
517 to not only provide insight into the locations of cancer-associated mutations, but also the level of certainty
518 in the locations. However, MO does not provide estimates of transition rates and error probabilities as do
519 SiFit and SCITE, but rather integrates over uncertainty in these parameters.

520 The methodology underlying MO could be enhanced by incorporating models for copy number alterations,
521 as well as by considering mutations that affect the same allele more than once. As SCS data collection
522 becomes more advanced, enabling hundreds of cells to be analyzed in parallel at reduced cost and increased
523 throughput, MO is poised to analyze the resulting large-scale data sets to make meaningful inference of
524 the mutation order during tumor progression for individual patients. MO thus represents an important step
525 forward in understanding the role of mutation order in cancer evolution and as such may have important
526 translational applications for improving cancer diagnosis, treatment, and personalized therapy. If inferred
527 mutation orders can be associated with clinical outcomes, future research can explore the cause of clinical
528 outcomes given specific mutation orders with the goal of developing novel, targeted treatments. This will
529 allow clinical providers to make decisions concerning treatment based on the mutation landscapes of patients.
530 Although the current study focuses on cancer, MO can potentially also be applied to single-cell mutation
531 profiles from a wide variety of fields. These applications are expected to provide new insights into our
532 understanding of cancer and other human diseases.

REFERENCES

23

533

6. SOFTWARE

534 MO has been implemented in R and is available at <https://github.com/lkubatko/MO>.

535

7. SUPPLEMENTARY MATERIAL

536 Supplementary material is available.

537

8. ACKNOWLEDGMENTS

538 The simulation experiments and data analyses were carried out using the ASC Unity Cluster at The Ohio
539 State University, USA. The authors thank two anonymous reviewers for helpful comments on an earlier draft
540 of this manuscript.

541 *Conflict of Interest:* None

542

REFERENCES

543 ASCOLANI, GIANLUCA AND LIÒ, PIETRO. (2019). Modeling breast cancer progression to bone: how driver
544 mutation order and metabolism matter. *BMC Medical Genomics* **12**(6), 106.

545 BARBIERI, CHRISTOPHER E, BANGMA, CHRIS H, BJARTELL, ANDERS, CATTO, JAMES WF, CULIG, ZO-
546 RAN, GRÖNBERG, HENRIK, LUO, JUN, VISAKORPI, TAPIO AND RUBIN, MARK A. (2013). The mutational
547 landscape of prostate cancer. *European Urology* **64**(4), 567–576.

548 CANISIUS, SANDER, MARTENS, JOHN W. M. AND WESSELS, LODEWYK F. A. (2016). A novel independence
549 test for somatic alterations in cancer shows that biology drives mutual exclusivity but chance explains most
550 co-occurrence. *Genome Biology* **17**, 261.

551 CHIFMAN, JULIA AND KUBATKO, LAURA. (2014). Quartet inference from SNP data under the coalescent
552 model. *Bioinformatics* **30**(23), 3317–3324.

553 FEARON, ERIC R AND VOGELSTEIN, BERT. (1990). A genetic model for colorectal tumorigenesis. *cell* **61**(5),
554 759–767.

555 ISHWARAN, HEMANT, BLACKSTONE, EUGENE H, APPERSON-HANSEN, CAROLYN AND RICE, THOMAS W.
556 (2009). A novel approach to cancer staging: application to esophageal cancer. *Biostatistics* **10**(4), 603–620.

- 557 IWASA, YOH, MICHOR, FRANZISKA AND NOWAK, MARTIN A. (2004). Stochastic tunnels in evolutionary
558 dynamics. *Genetics* **166**(3), 1571–1579.
- 559 JAHN, KATHARINA, KUIPERS, JACK AND BEERENWINKEL, NIKO. (2016). Tree inference for single-cell data.
560 *Genome Biology* **17**(1), 86.
- 561 JAMAL-HANJANI, MARIAM, WILSON, GARETH A, MCGRANAHAN, NICHOLAS, BIRKBAK, NICOLAI J,
562 WATKINS, THOMAS BK, VEERIAH, SELVARAJU, SHAFI, SEEMA, JOHNSON, DIANA H, MITTER,
563 RICHARD, ROSENTHAL, RACHEL *and others*. (2017). Tracking the evolution of non-small-cell lung cancer.
564 *New England Journal of Medicine* **376**(22), 2109–2121.
- 565 KIM, KYUNG IN AND SIMON, RICHARD. (2014). Using single cell sequencing data to model the evolutionary
566 history of a tumor. *BMC Bioinformatics* **15**(1), 27.
- 567 LEUNG, MARCO L., DAVIS, ALEXANDER, GAO, RULI, CASASENT, ANNA, WANG, YONG, SEI, EMI, VILAR,
568 EDUARDO, MARU, DIPEN, KOPETZ, SCOTT AND NAVIN, NICHOLAS E. (2017). Single-cell DNA sequenc-
569 ing reveals a late-dissemination model in metastatic colorectal cancer. *Genome Research* **27**, 1287–1299.
- 570 NAVIN, NICHOLAS E. (2014). Cancer genomics: one cell at a time. *Genome Biology* **15**(8), 452.
- 571 ORTMANN, CHRISTINA A, KENT, DAVID G, NANGALIA, JYOTI, SILBER, YVONNE, WEDGE, DAVID C,
572 GRINFELD, JACOB, BAXTER, E JOANNA, MASSIE, CHARLES E, PAPAEMMANUIL, ELLI, MENON, SURAJ
573 *and others*. (2015). Effect of mutation order on myeloproliferative neoplasms. *New England Journal of*
574 *Medicine* **372**(7), 601–612.
- 575 O’SULLIVAN, FINBARR, ROY, SUPRATIK AND EARY, JANET. (2003). A statistical measure of tissue hetero-
576 geneity with application to 3D PET sarcoma data. *Biostatistics* **4**(3), 433–448.
- 577 POSADA, DAVID. (2020). CellCoal: coalescent simulation of single-cell sequencing samples. *Molecular Biology*
578 *and Evolution* **37**(5), 1535–1542.
- 579 POWELL, STEVEN M, ZILZ, NATHAN, BEAZER-BARCLAY, YASMIN, BRYAN, TRACY M, HAMILTON, STAN-
580 LEY R, THIBODEAU, STEPHEN N, VOGELSTEIN, BERT AND KINZLER, KENNETH W. (1992). Apc muta-
581 tions occur early during colorectal tumorigenesis. *Nature* **359**(6392), 235–237.

REFERENCES

25

- 582 SU, FEI, ZHANG, WEI, ZHANG, DALEI, ZHANG, YAQUN, PANG, CHENG, HUANG, YINGYING, WANG, MIAO,
583 CUI, LUWEI, HE, LEI, ZHANG, JINSONG *and others.* (2018). Spatial intratumor genomic heterogeneity
584 within localized prostate cancer revealed by single-nucleus sequencing. *European Urology* **74**(5), 551–559.
- 585 SWOFFORD, DL. (1999). Phylogenetic analysis using parsimony, PAUP* 4.0, beta version 4.0 b2. *Sinauer*
586 *Associates, Boston, Mass.*
- 587 TATE, JOHN G, BAMFORD, SALLY, JUBB, HARRY C, SONDKA, ZBYSLAW, BEARE, DAVID M, BINDAL,
588 NIDHI, BOUTSELAKIS, HARRY, COLE, CHARLOTTE G, CREATORE, CELESTINO, DAWSON, ELISABETH
589 *and others.* (2018). COSMIC: the catalogue of somatic mutations in cancer. *Nucleic Acids Re-*
590 *search* **47**(D1), D941–D947.
- 591 YU, CHANG, YU, JUN, YAO, XIAOTIAN, WU, WILLIAM K. K., LU, YOUYONG, TANG, SENWEI, LI,
592 XIANGCHUN, BAO, LI, LI, XIAOXING, HOU, YONG, WU, RENHUA, JIAN, MIN, CHEN, RUOYAN, ZHANG,
593 FAN, XU, LIXIA, FAN, FAN, HE, JUN, LIANG, QIAOYI, WANG, HONGYI, HU, XUEDA, HE, MINGHUI,
594 ZHANG, XIANG, ZHENG, HANCHENG, LI, QIBIN, WU, HANJIE, CHEN, YAN, YANG, XU, ZHU, SHIDA,
595 XU, XUN, YANG, HUANMING, WANG, JIAN, ZHANG, XIUQING, SUNG, JOSEPH J. Y., LI, YINGRUI *and*
596 *others.* (2014). Discovery of biclonal origin and a novel oncogene SLC12A5 in colon cancer by single-cell
597 sequencing. *Cell Research* **24**, 701–712.
- 598 ZAFAR, HAMIM, NAVIN, NICHOLAS, CHEN, KEN AND NAKHLEH, LUAY. (2019). SiCloneFit: Bayesian infer-
599 ence of population structure, genotype, and phylogeny of tumor clones from single-cell genome sequencing
600 data. *Genome Research* **29**, 1–13.
- 601 ZAFAR, HAMIM, TZEN, ANTHONY, NAVIN, NICHOLAS, CHEN, KEN AND NAKHLEH, LUAY. (2017). SiFit:
602 inferring tumor trees from single-cell sequencing data under finite-sites models. *Genome Biology* **18**(1),
603 178.

604

||

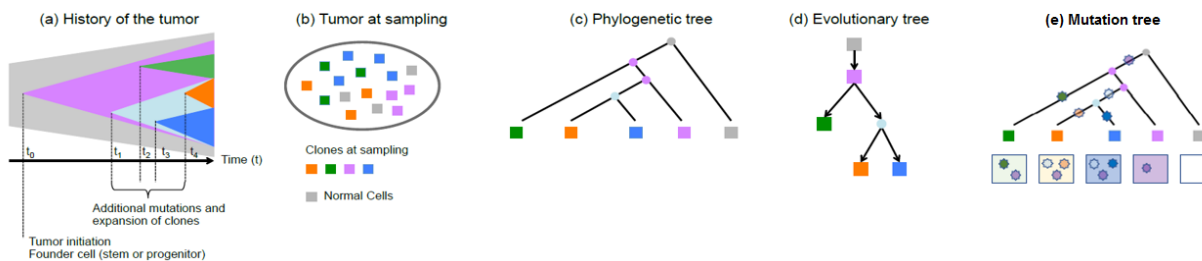


Fig. 1: Pictorial representation of tumor evolution. (a) - (b) A pictorial representation of the evolution of a tumor from the first initiating mutation to the heterogeneous tissue at the time of sampling, which consists of four different clones and normal tissue. (c) A phylogenetic tree with single cells as the tips. (d) A clonal lineage tree inferred from sampled cells where each node represents a subclone (cluster of cells). (e) A mutation tree inferred from sampled cells where each star represents the occurrence of one mutation. Boxes underneath each tip show which mutations are present in the cell represented by the tip.

(a) True binary genotype						(b) Observed binary genotype					
	C_1	C_2	C_3	C_4	C_5		C_1	C_2	C_3	C_4	C_5
G_1						S_1					
G_2						S_2					
G_3						S_3					
G_4						S_4					
G_5						S_5					

Fig. 2: True and observed binary data. (a) True binary mutation matrix representing the mutation status of the sequenced tumor cells in the mutation tree in Fig. 1(e). Each row represents true genotypes for one genomic site in all cells and each column represents the true genotypes of multiple genomic sites for one single cell. (b) Observed mutation matrix with missing and ambiguous values (red), as well as genotypes that are misrecorded with respect to the true mutation matrix (red numbers; these are either false positives or false negatives). The red dash indicates a missing value since the sequencing process does not return signal at this site of this cell, and the red question mark represents an ambiguous value. Each row represents observed genotypes for one genomic site in all cells and each column represents the observed genotypes of multiple genomic sites for one single cell.

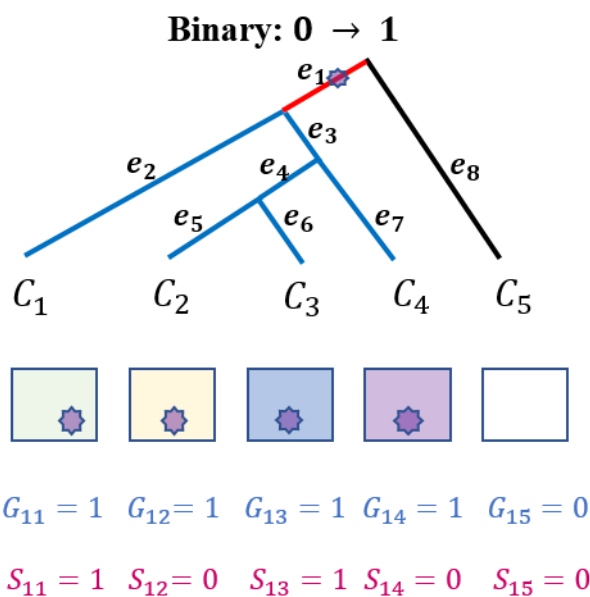


Fig. 3: Binary mutation process example. Example in which a mutation is acquired on branch e_1 (marked with red color). The cell descending from branch e_8 (marked with black color) does not carry the mutation, while the cells descending from the blue branches carry the mutation.

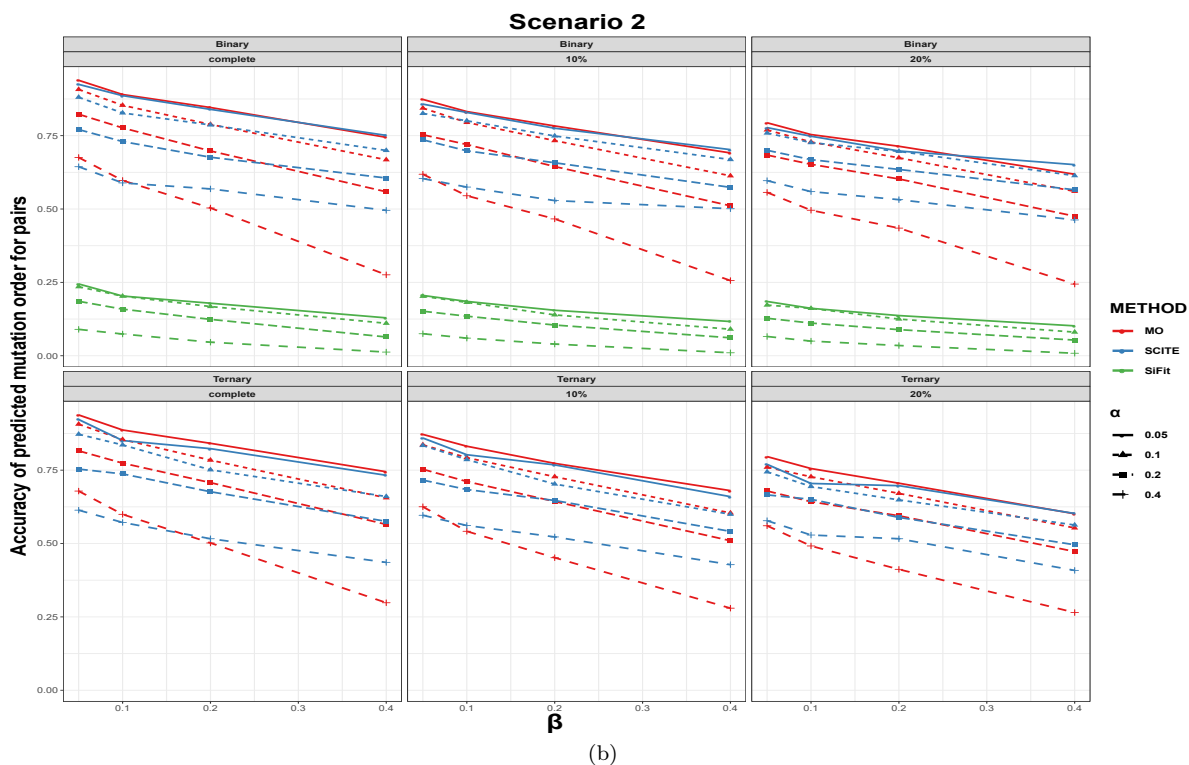
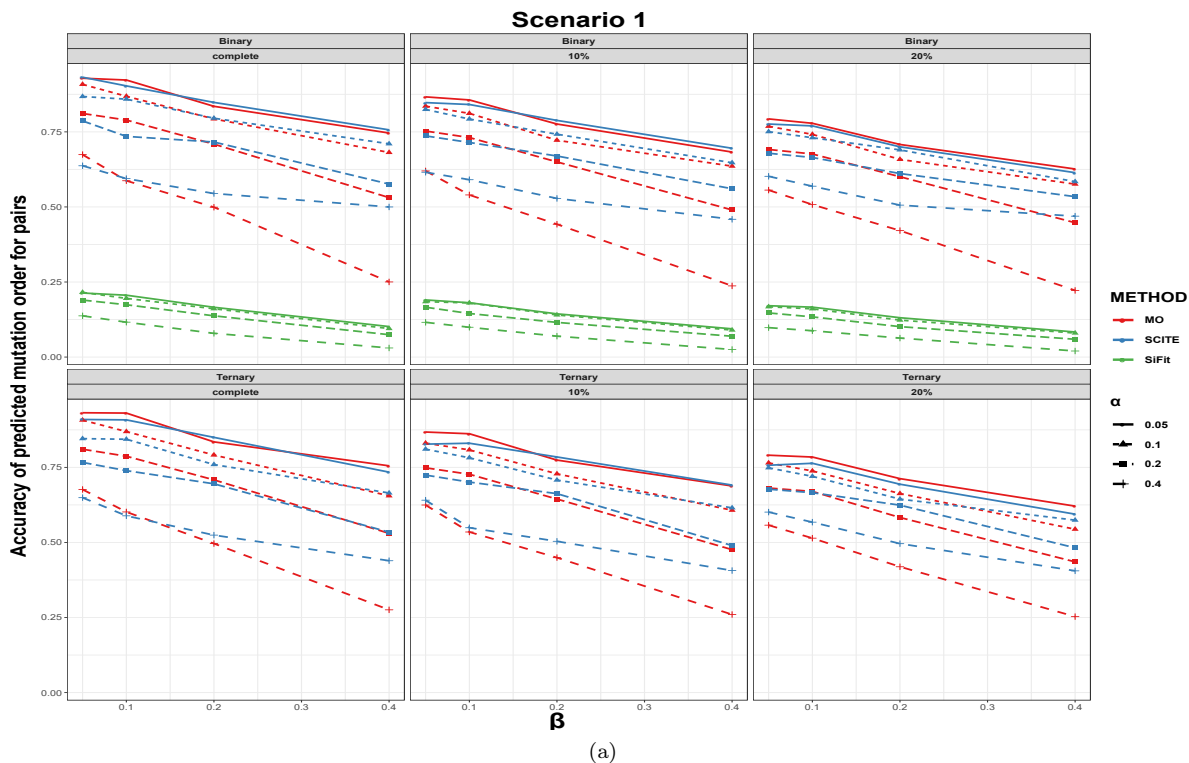


Fig. 4: Order accuracy in scenarios 1 and 2 for MO, SCITE and SiFit. Each panel includes the results from the specific type of genotype and missing data percentage. In each panel, red, blue and green colors correspond to MO, SCITE and SiFit, respectively. Each plotting symbol on the line represents a different probability of a false positive error, α . The x-axis is the probability of a false negative error, β , and the y-axis is order accuracy.

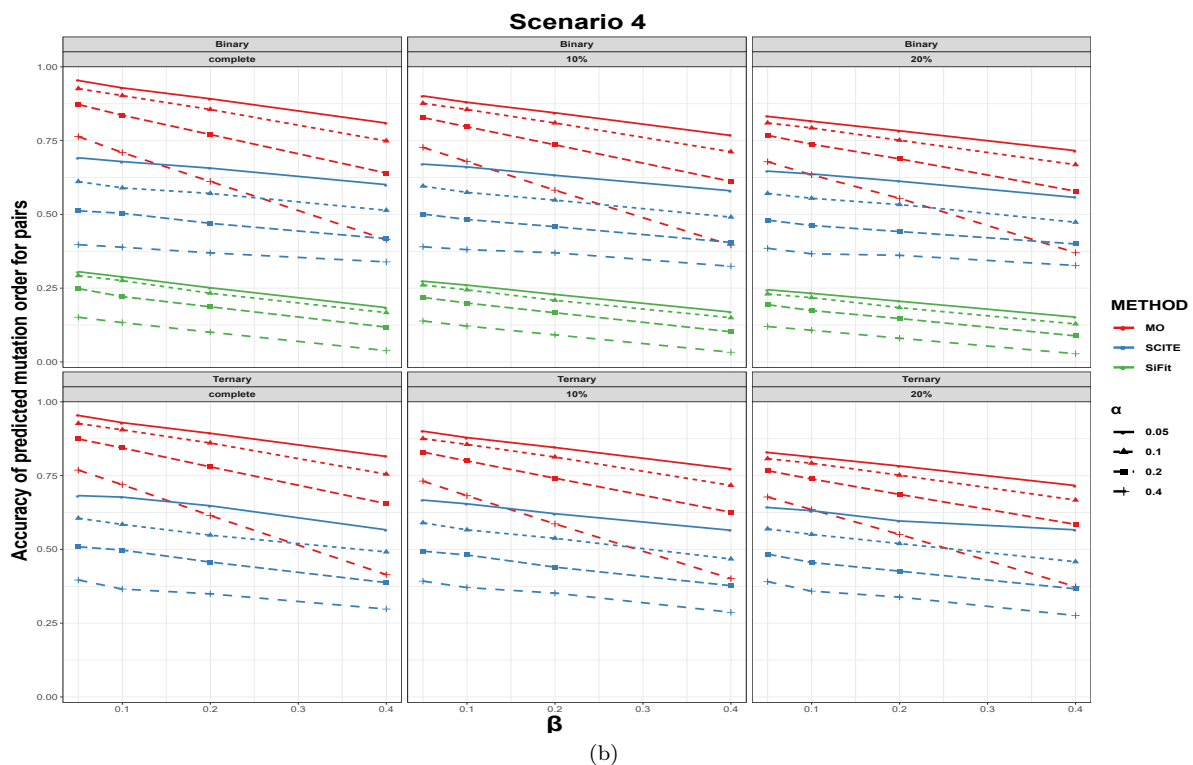
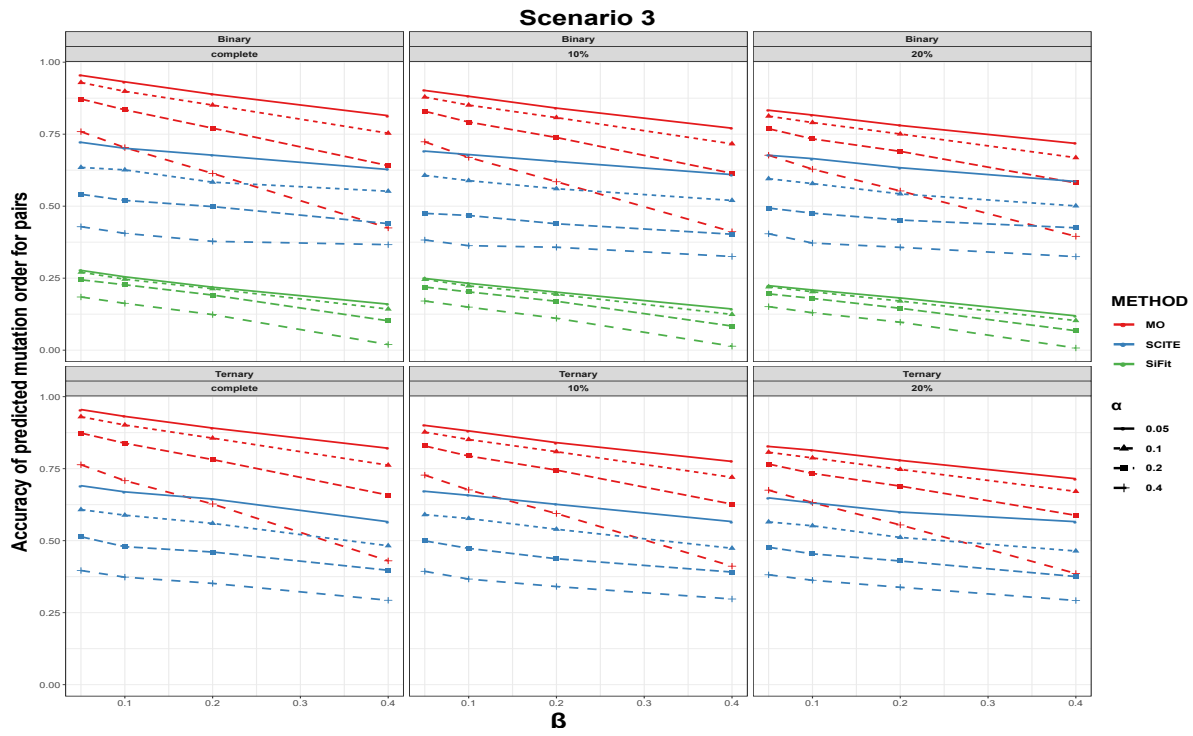


Fig. 5: Order accuracy in scenarios 3 and 4 for MO, SCITE and SiFit. Each panel includes the results from the specific type of genotype and missing data percentage. In each panel, red, blue and green colors correspond to MO, SCITE and SiFit, respectively. Each plotting symbol on the line represents a different probability of a false positive error, α . The x-axis is the probability of a false negative error, β , and the y-axis is order accuracy.

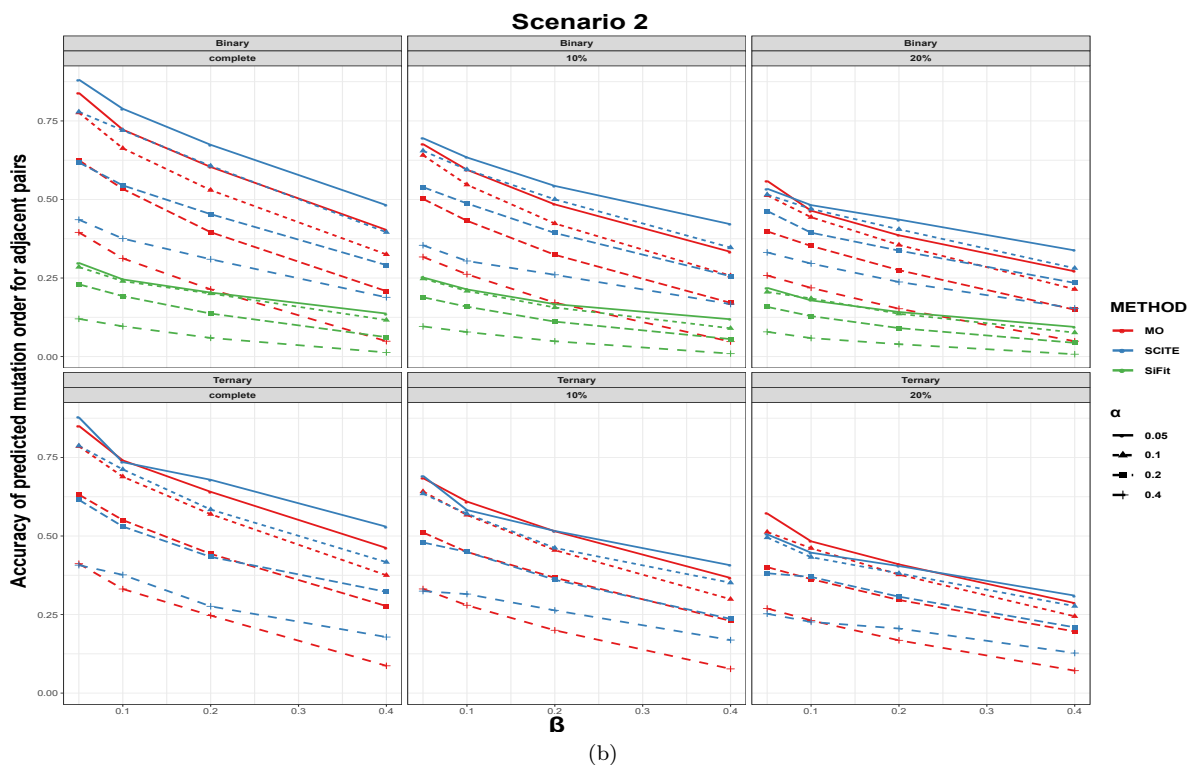
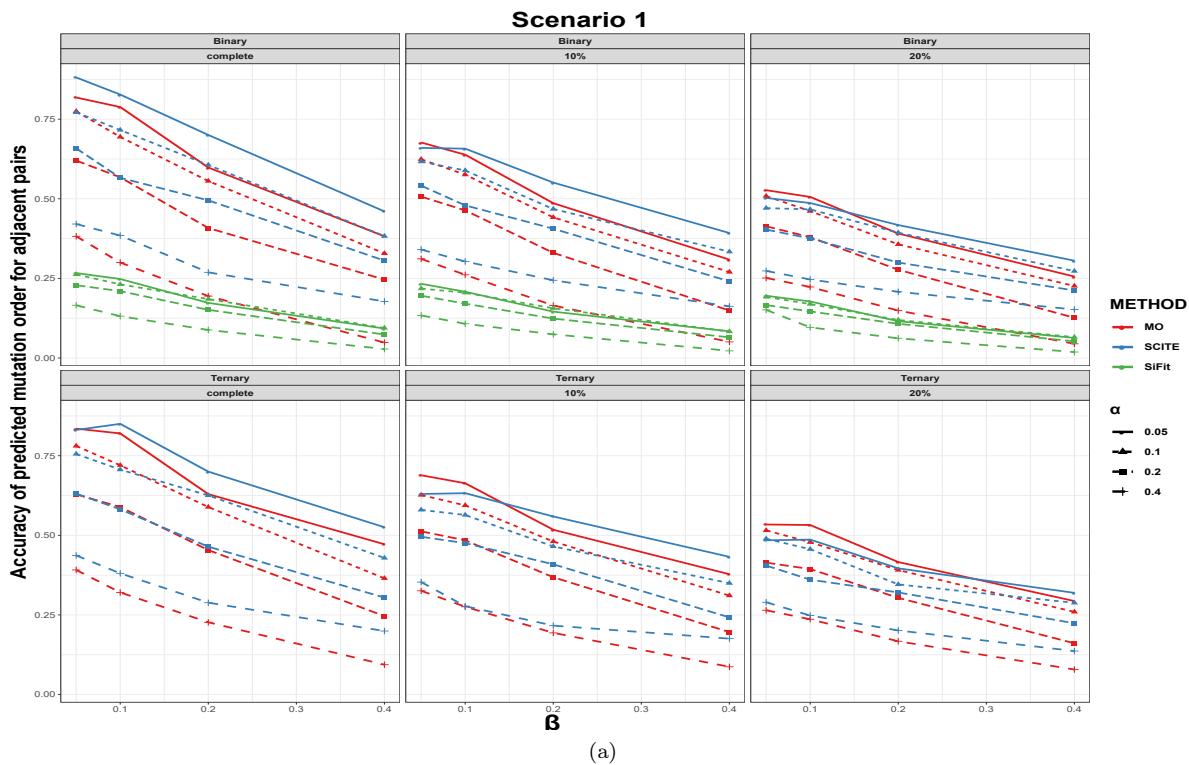


Fig. 6: Adjacent order accuracy in scenarios 1 and 2 for MO, SCITE and SiFit. Each panel includes the results from the specific type of genotype and missing data percentage. In each panel, red, blue and green colors correspond to MO, SCITE and SiFit, respectively. Each plotting symbol on the line represents a different probability of a false positive error, α . The x-axis is the probability of a false negative error, β , and the y-axis is adjacent order accuracy.

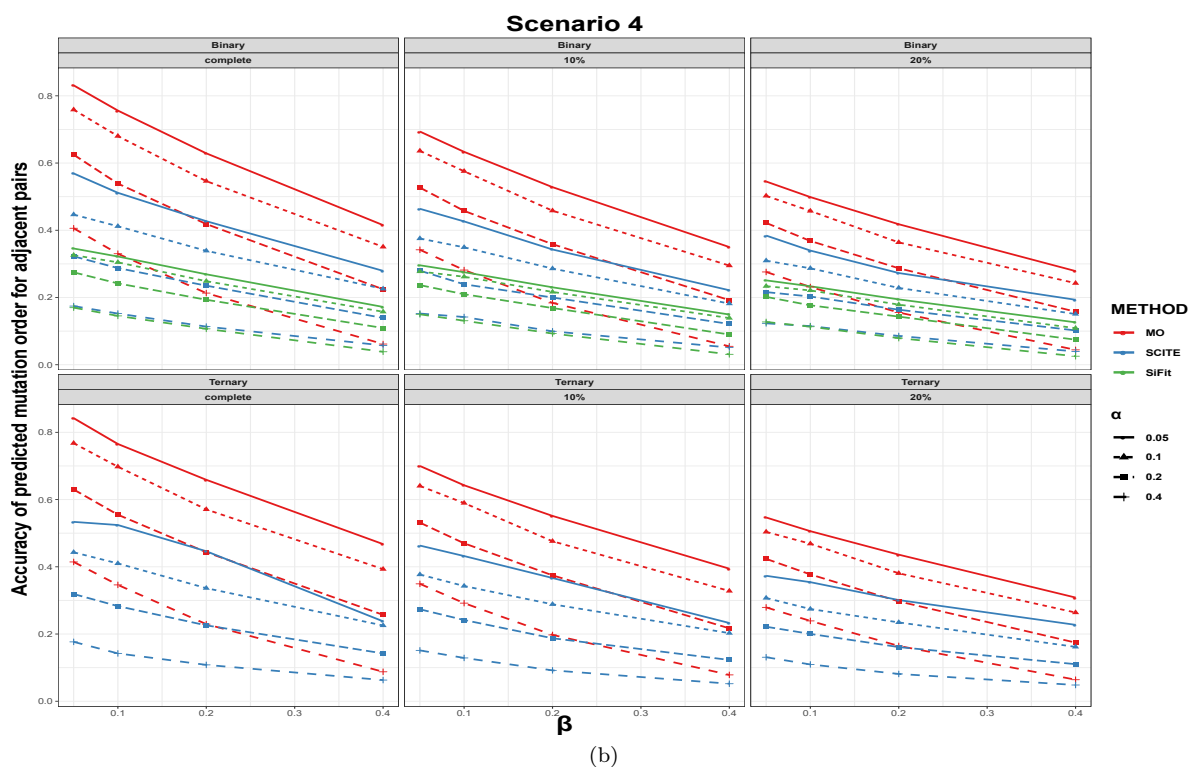
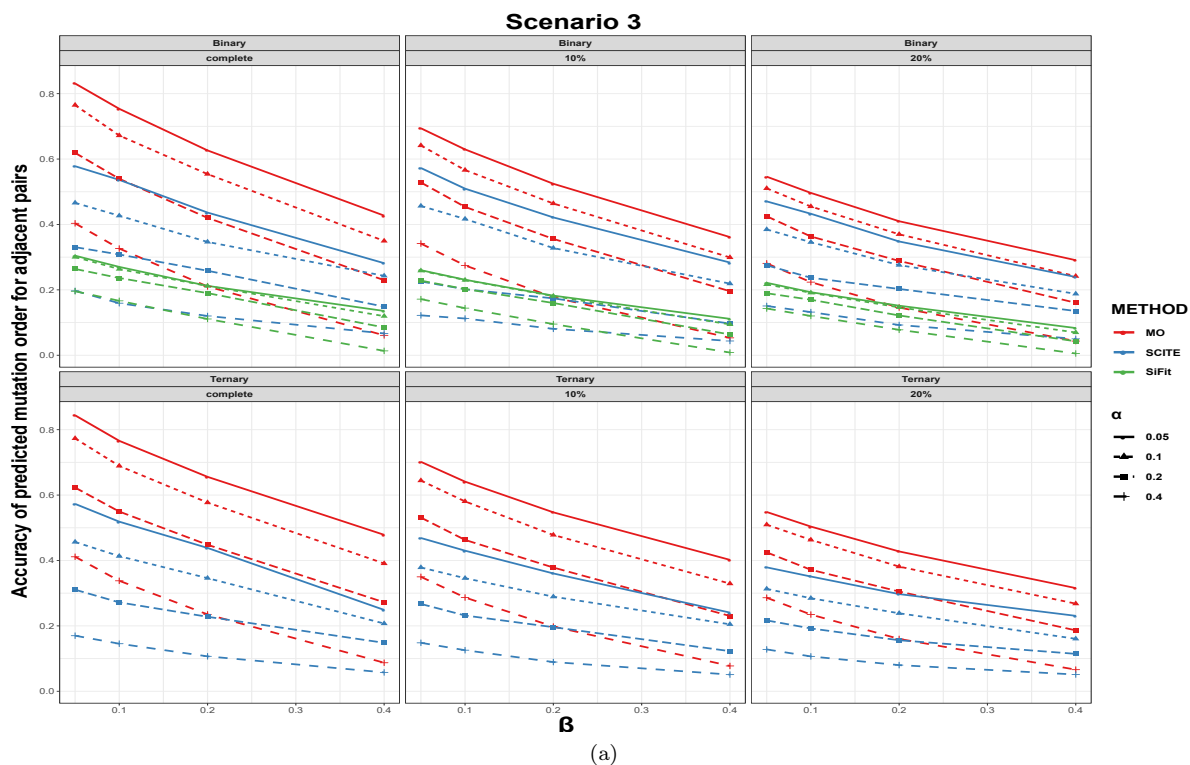


Fig. 7: Adjacent order accuracy in scenarios 3 and 4 for MO, SCITE and SiFit. Each panel includes the results from the specific type of genotype and missing data percentage. In each panel, red, blue and green colors correspond to MO, SCITE and SiFit, respectively. Each plotting symbol on the line represents a different probability of a false positive error, α . The x-axis is the probability of a false negative error, β , and the y-axis is adjacent order accuracy.

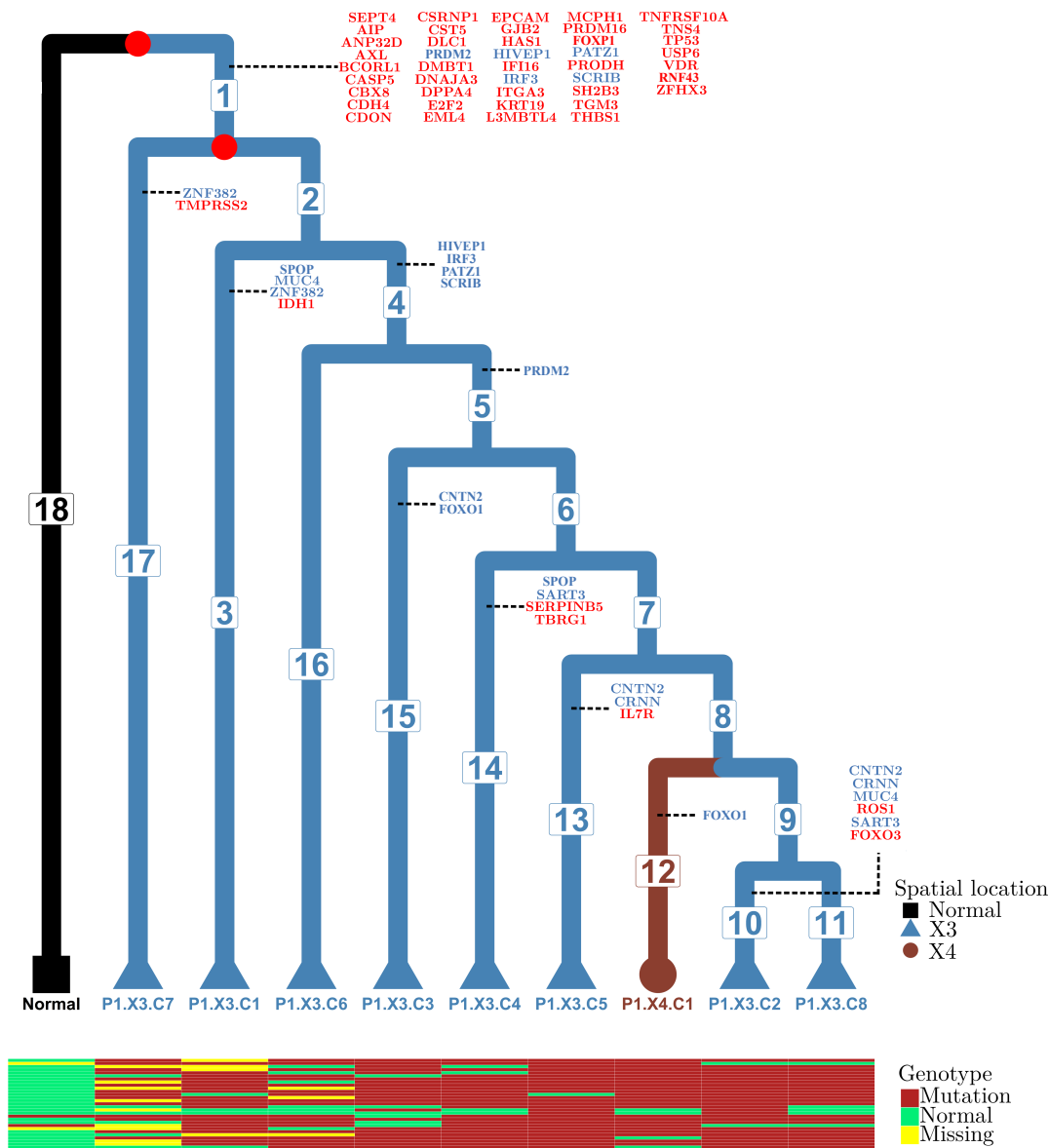


Fig. 8: P1 tumor phylogenetic tree and inferred temporal order of the mutations. The normal cell is set as the outgroup. There are 18 branches in this tree. We do not assume the molecular clock when estimating the branch lengths. Branch lengths in this figure are not drawn to scale. The color and tip shape represent the spatial locations of the samples (normal tissue, location X3 or location X4; see Su *and others* (2018)). The temporal order of the mutations is annotated on the branches of the tree. The uncertainty of mutation locations is highlighted with colors. Mutations with very strong signals (probability of occurring on one branch is greater than 0.7) are marked in red, while mutations with moderate signals (probabilities that sum to more than 0.7 on two or three branches) are marked in blue. Mutation data for 30 genes corresponding to the first 30 rows in Fig. S11 and Fig. S12 for each tip are shown in the heatmap matrix at the bottom.

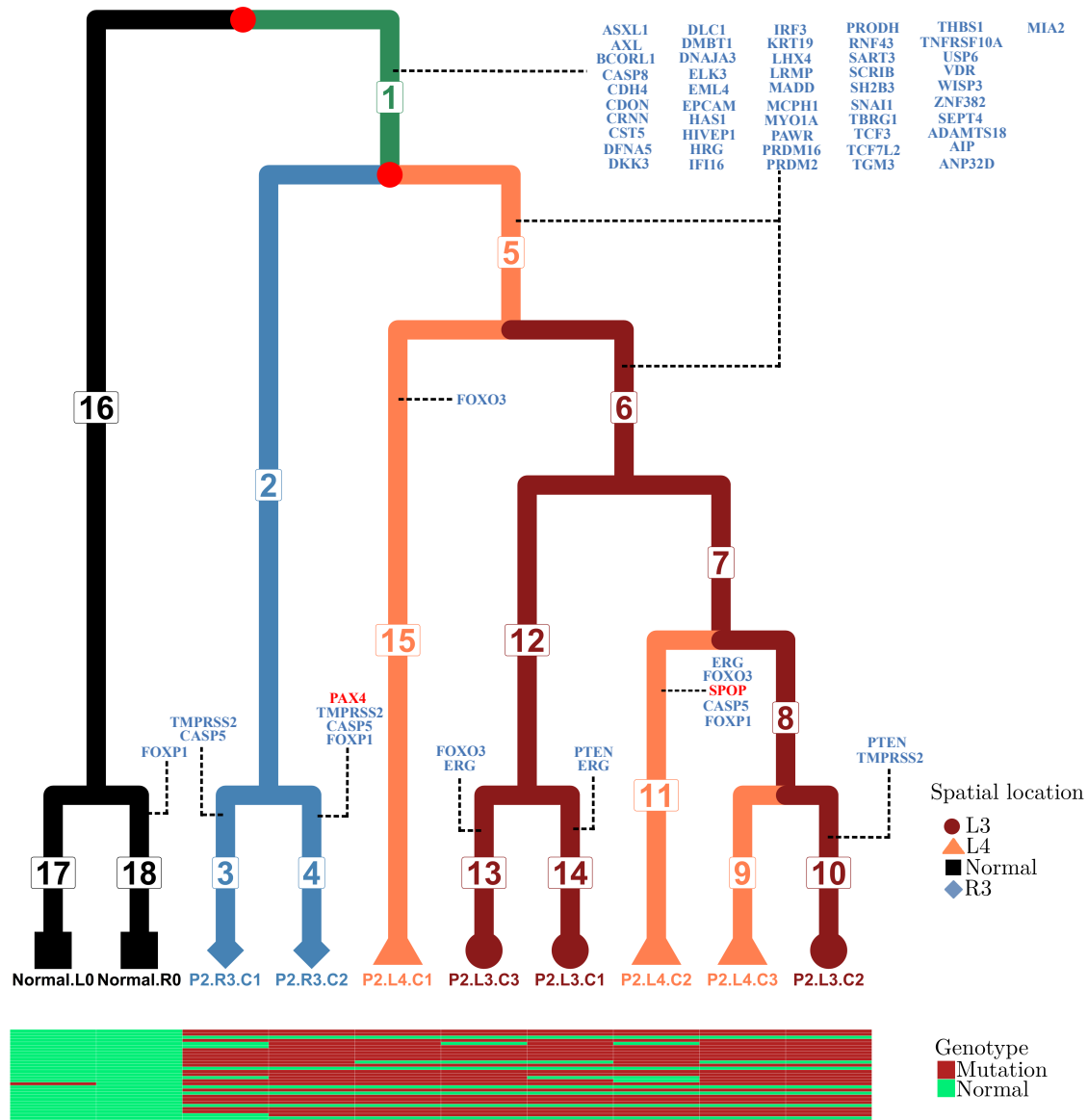


Fig. 9: P2 tumor phylogenetic tree and inferred temporal order of the mutations. Normal.R0 and Normal.L0 are normal cells from the right side and the left side of tissue, respectively, and are set as the outgroup. There are 18 branches in this tree. We do not assume the molecular clock when estimating the branch lengths, and branch lengths in this figure are not drawn to scale. The color and tip shape represent the spatial locations of the samples (normal tissue, left-side locations L3 or L4, or right-side location R3; see Su *and others* (2018)). The temporal order of the mutations is annotated on the branches of the tree. The uncertainty of mutation locations is highlighted with colors. Mutations with very strong signals (probability of occurring on one branch is greater than 0.7) are marked in red, while mutations with moderate signals (probabilities that sum to more than 0.7 on two or three branches) are marked in blue. Mutation data for 30 genes corresponding to the first 30 rows in Fig. S13 and Fig. S14 for each tip are shown in the heatmap matrix at the bottom.

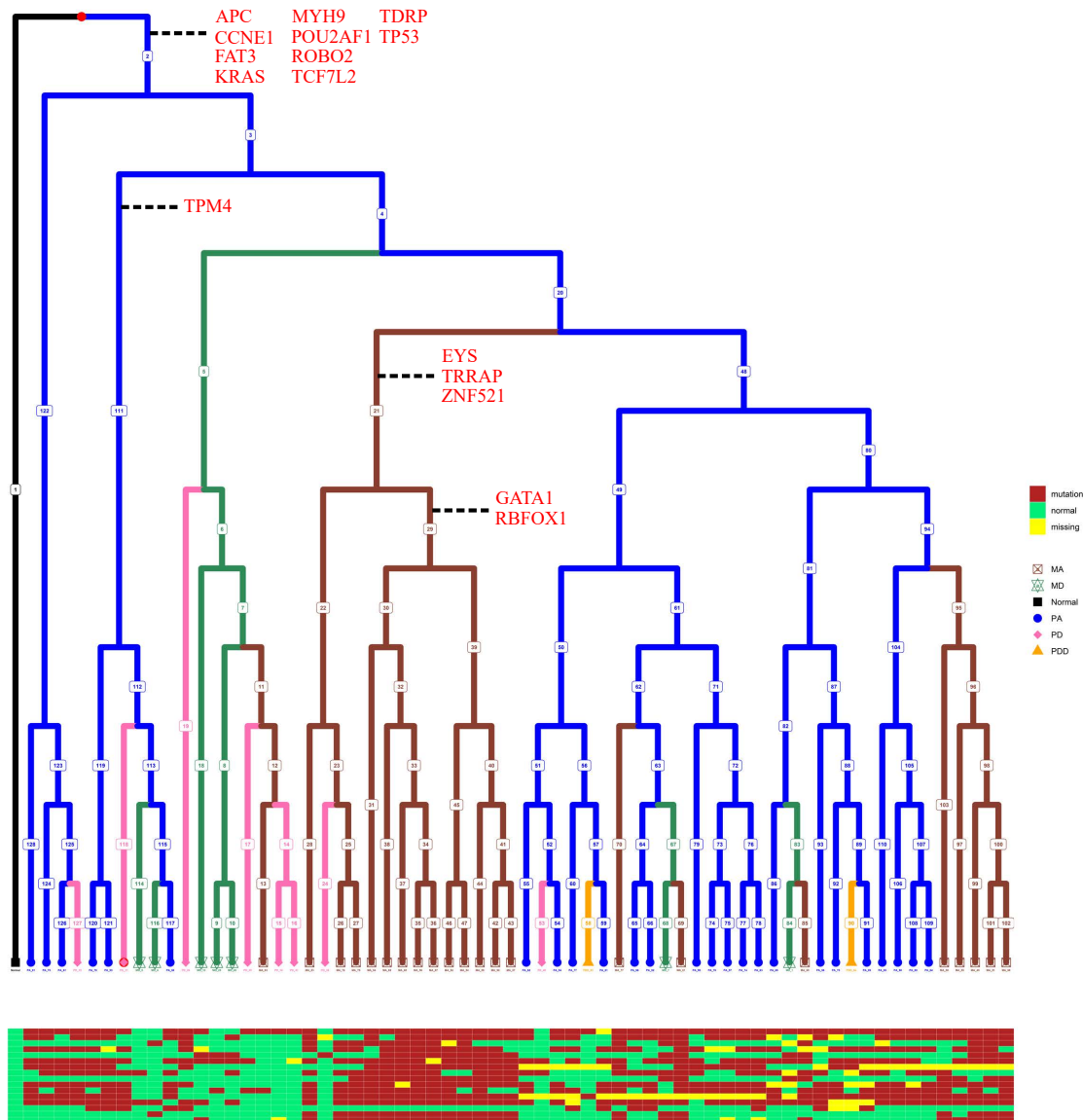


Fig. 10: CRC1 tumor phylogenetic tree and inferred temporal order of the mutations. The color and tip shape represent the spatial locations of the samples (Normal - normal tissue; PA - primary aneuploid; PD - primary diploid; MA - metastatic aneuploid; MD - metastatic diploid; see Leung *and others* (2017)). The temporal order of the mutations is annotated on the branches of the tree. The uncertainty of mutation locations is highlighted with colors. Mutations with very strong signals (probability of occurring on one branch is greater than 0.7) are marked in red, while genes with moderate signals (probabilities that sum to more than 0.7 on two or three branches) are marked in blue. The branch lengths are not scaled. Mutation data for the 16 genes corresponding to each tip are shown in the heatmap matrix at the bottom.

REFERENCES

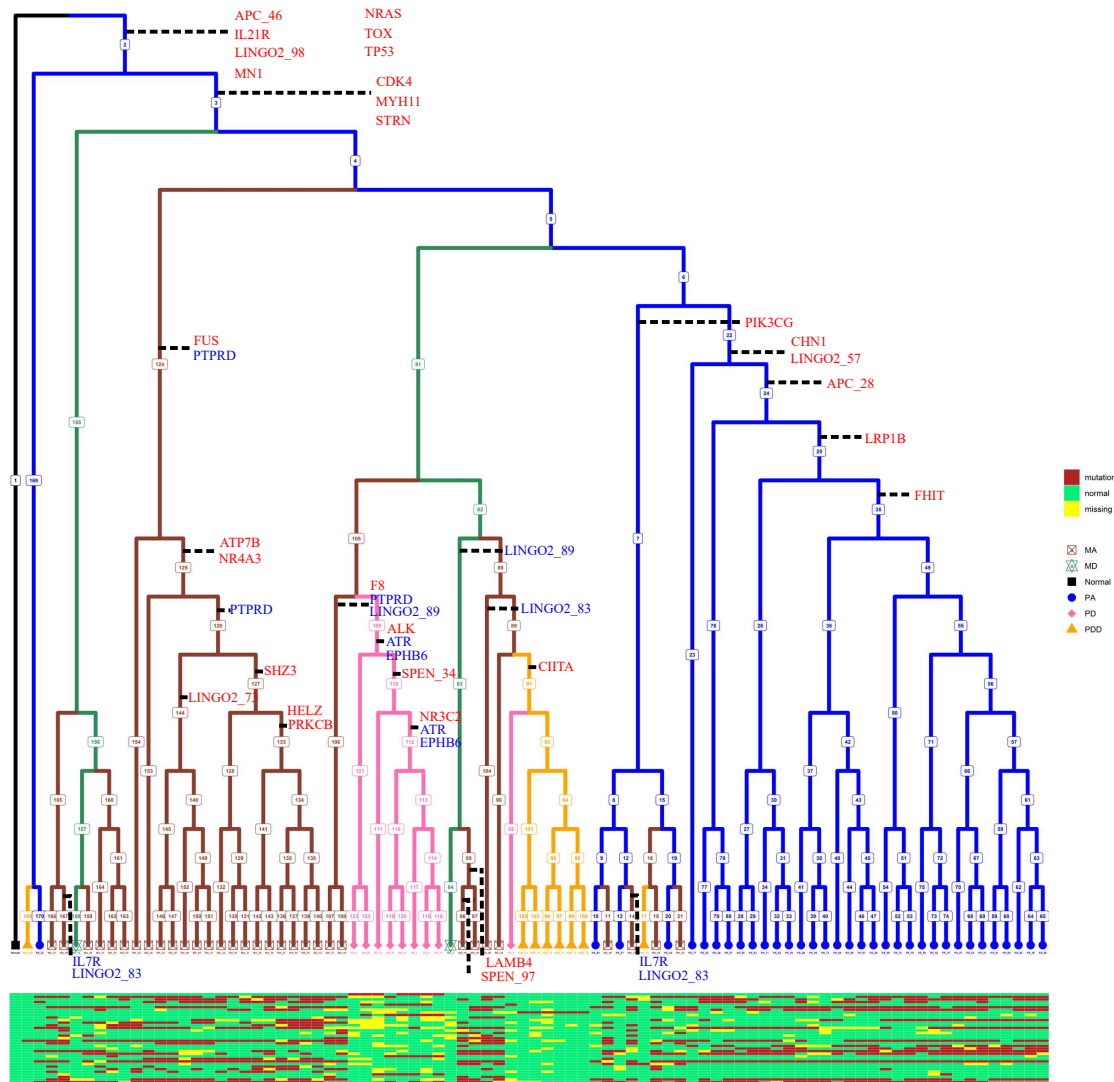


Fig. 11: CRC2 tumor phylogenetic tree and inferred temporal order of the mutations. The color and tip shape represent the spatial locations of the samples (Normal - normal tissue; PA - primary aneuploid; PD - primary diploid; MA - metastatic aneuploid; MD - metastatic diploid; see Leung *and others* (2017)). The temporal order of the mutations is annotated on the branches of the tree. The uncertainty of mutation locations is highlighted with colors. Mutations with very strong signals (probability of occurring on one branch is greater than 0.7) are marked in red, while mutations with moderate signals (probabilities that sum to more than 0.7 on two or three branches) are marked in blue. The branch lengths are not scaled. Mutation data for the 36 genomic sites corresponding to each tip are shown in the heatmap matrix at the bottom.

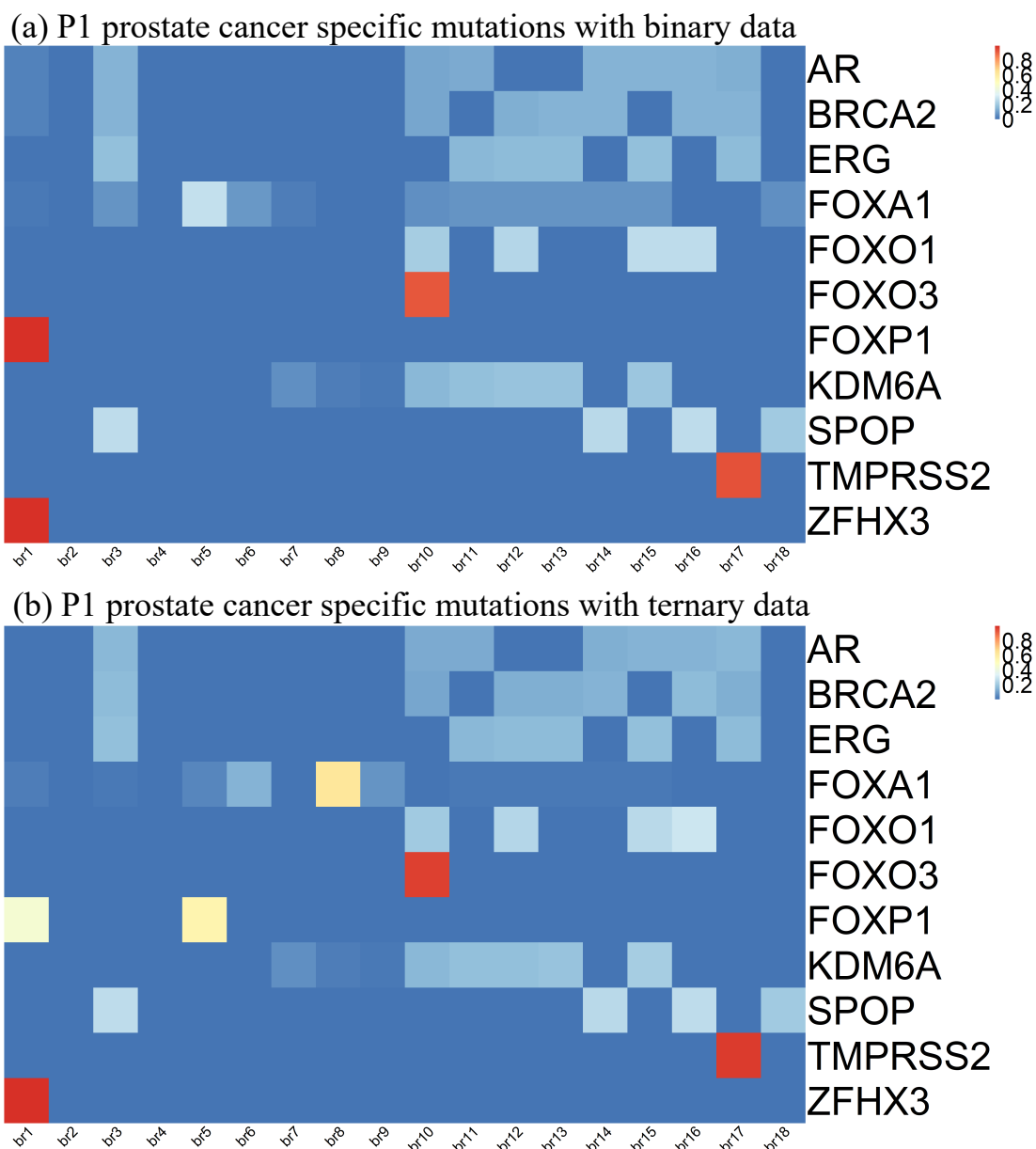


Fig. 12: Heatmap of posterior probabilities of mutation branch for P1 using (a) binary or (b) ternary data. This heatmap is for the prostate cancer-specific genes. Color indicates the magnitude of the probability, with red indicating probability close to 1 and blue indicating probability close to 0. For P1, the prior for α is set to $\alpha|\mathbf{S}_i \sim \text{Beta}(0.29, 0.71)$ (larger variance). The prior for β is set to $\beta|\mathbf{S}_i \sim \text{Beta}(0.02, 0.98)$ (larger variance). The distribution of transition rate λ_1 is set to $\lambda_1|\mathbf{S}_i \sim \text{Gamma}(2, 5.0 \times 10^{-8})$ (larger variance). The distribution of transition rate λ_2 is set to $\lambda_2|\mathbf{S}_i \sim \text{Gamma}(2, 5.0 \times 10^{-3})$ (larger variance) (cont.).

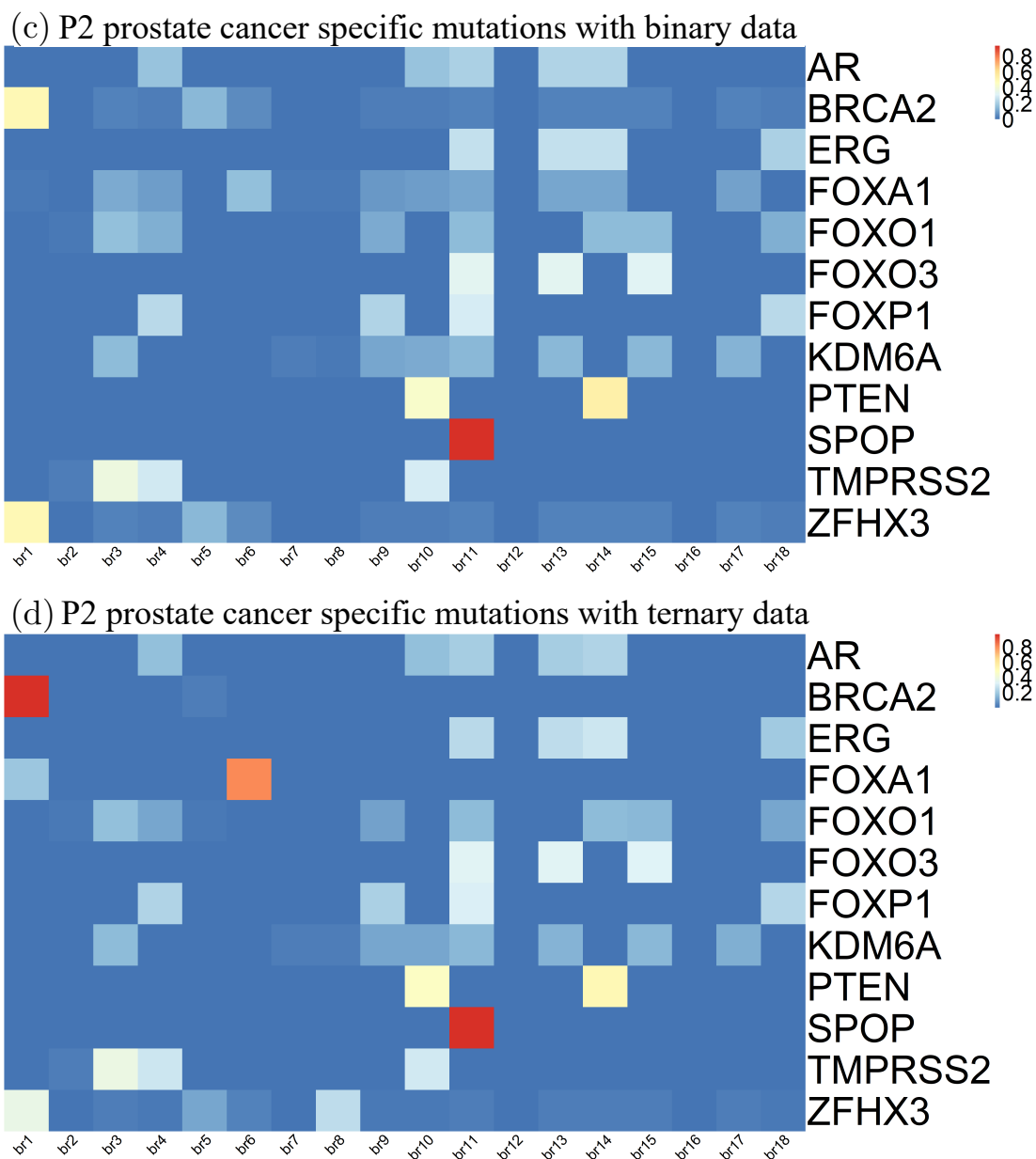


Fig. 12: Heatmap of posterior probabilities of mutation branch for P2 using (c) binary or (d) ternary data. This heatmap is for the prostate cancer-specific genes. Color indicates the magnitude of the probability, with red indicating probability close to 1 and blue indicating probability close to 0. For P2, the prior distribution for α is set to $\alpha|\mathbf{S}_i \sim \text{Beta}(0.31, 0.69)$ (larger variance). The prior distribution for β is set to $\beta|\mathbf{S}_i \sim \text{Beta}(0.02, 0.98)$ (larger variance). The distribution of transition rate λ_1 is set to $\lambda_1|\mathbf{S}_i \sim \text{Gamma}(2, 5.0 \times 10^{-8})$ (larger variance). The distribution of transition rate λ_2 is set to $\lambda_2|\mathbf{S}_i \sim \text{Gamma}(2, 5.0 \times 10^{-3})$ (larger variance).