

# 1 Cancer phylogenetic tree inference at scale from 1000s 2 of single cell genomes

3 Sohrab Salehi <sup>\*1</sup>, Fatemeh Dorri <sup>\*2</sup>, Kevin Chern<sup>1</sup>, Farhia Kabeer<sup>3</sup>, Nicole Rusk<sup>4</sup>,  
4 Tyler Funnell<sup>4</sup>, Marc J Williams<sup>4</sup>, Daniel Lai<sup>3,5</sup>, Mirela Andronescu<sup>3,5</sup>, Kieran R.  
5 Campbell<sup>6,7,8</sup>, Andrew McPherson<sup>4</sup>, Samuel Aparicio<sup>3,5</sup>, Andrew Roth<sup>2,3,5</sup>, Sohrab  
6 Shah<sup>4</sup>, and Alexandre Bouchard-Côté<sup>1,\*</sup>

7 <sup>1</sup> Department of Statistics, University of British Columbia

8 <sup>2</sup> Department of Computer Science, University of British Columbia

9 <sup>3</sup> Department of Pathology and Laboratory Medicine, University of British Columbia

10 <sup>4</sup> Computational Oncology, Department of Epidemiology and Biostatistics, Memorial Sloan Kettering Cancer  
11 Center

12 <sup>5</sup> Department of Molecular Oncology, BC Cancer Research Centre

13 <sup>6</sup> Lunenfeld-Tanenbaum Research Institute, Mount Sinai Hospital

14 <sup>7</sup> Department of Molecular Genetics, University of Toronto

15 <sup>8</sup> Department of Statistical Sciences, University of Toronto

16 <sup>\*</sup>Correspondence: [bouchard@stat.ubc.ca](mailto:bouchard@stat.ubc.ca)

## 17 **Abstract**

18 A new generation of scalable single cell whole genome sequencing (scWGS) meth-  
19 ods allows unprecedented high resolution measurement of the evolutionary dynamics  
20 of cancer cell populations. Phylogenetic reconstruction is central to identifying sub-  
21 populations and distinguishing the mutational processes that gave rise to them. Ex-  
22 isting phylogenetic tree building models do not scale to the tens of thousands of high  
23 resolution genomes achievable with current scWGS methods. We constructed a phy-  
24 logenetic model and associated Bayesian inference procedure, *sitka*, specifically for  
25 scWGS data. The method is based on a novel phylogenetic encoding of copy num-  
26 ber (CN) data, the *sitka* transformation, that simplifies the site dependencies induced  
27 by rearrangements while still forming a sound foundation to phylogenetic inference.  
28 The *sitka* transformation allows us to design novel scalable Markov chain Monte Carlo  
29 (MCMC) algorithms. Moreover, we introduce a novel point mutation calling method  
30 that incorporates the CN data and the underlying phylogenetic tree to overcome the  
31 low per-cell coverage of scWGS. We demonstrate our method on three single cell  
32 datasets, including a novel PDX series, and analyse the topological properties of the  
33 inferred trees. *Sitka* is freely available at <https://github.com/UBC-Stat-ML/sitkatree.git>.

---

\*Equal contribution

## 34 1 Introduction

35 A main challenge in investigating cancer evolution is the need to resolve the subpopulation  
36 structure of a heterogeneous tumour sample. Advances in next generation scWGS have  
37 enabled more accurate, quantitative measurements of tumours as they evolve [1, 2, 3, 4].  
38 Phylogenetic reconstruction is central to identifying clones in longitudinal xenoengraftment  
39 [5, 6] as well as patients [7], and has been used to approximate the rate and timing of  
40 mutation [8] to determine the origins and clonality of metastasis [9, 10]. Single cell cancer  
41 phylogenetics is an evolving field. Multiple approaches, spanning different study designs  
42 and data sources are reviewed in [11]. Many phylogenetic inference methods assume point  
43 mutations as input or a small number of leaf nodes [12, 13, 14, 15]. However, emerging  
44 single cell platforms produce up to thousands of single cell genomes and are suitable  
45 for determining copy number aberrations (CNA) [16, 1]. The method of [17] assumes  
46 a tree inferred from CNA exists and incorporates it in inference of point mutation based  
47 phylogenies. Distance based and agglomerative clustering methods such as neighbour  
48 joining are scalable and are used to elucidate hierarchical structures over cells [18, 19].  
49 While these are useful heuristics, they are statistically sub-optimal relative to likelihood  
50 based methods [20].

51 We describe *sitka*, a phylogenetic model and the associated Bayesian inference proce-  
52 dure designed specifically for inference based on CN information extracted from scWGS  
53 data. Our method addresses two key challenges: first, each CN event typically affects a  
54 large number of genomic sites, breaking the independence assumptions required by exist-  
55 ing phylogenetic methods [21, 15, 13, 22]; second, while detailed modelling of dependent  
56 evolutionary processes is in principle possible, they entail computational requirements in-  
57 compatible with the scale of modern scWGS data [23]. To confront these two difficulties,  
58 *sitka* uses a novel phylogenetic encoding of CN data, providing a statistical-computational  
59 trade-off by simplifying the site dependencies induced by rearrangements, while still form-  
60 ing a sound foundation to phylogenetic inference. Based on this encoding, we propose an  
61 innovative phylogenetic tree exploration move which makes the cost of Markov chain Monte  
62 Carlo (MCMC) iterations bounded by  $O(|C| + |L|)$ , where  $|C|$  is the number of cells and  $|L|$   
63 is the number of loci. In contrast, existing off-the-shelf likelihood-based methods incur an  
64 iteration cost of  $O(|C| |L|)$  [24, 13, 15]. Moreover, the novel move considers an exponential  
65 number of neighbouring trees whereas off-the-shelf moves consider a polynomial size set  
66 of neighbours.

67 We compare *sitka* with other tree-inference methods on three real-world datasets, includ-  
68 ing triple negative breast cancer patient derived xenograft samples, high grade serous  
69 ovarian primary and matched relapse samples. Since the true phylogeny is unknown, we  
70 design a phylogenetic goodness-of-fit framework to quantitatively assess the performance  
71 of our method and to visualize reconstruction confidence as well as violations of our as-  
72 sumptions.

73 We use the *sitka* inferred trees to analyse the topological properties of the real-world  
74 datasets. Finally, we introduce a model extension that enables the placement of single  
75 nucleotide variants (SNV) with high levels of missingness on a tree inferred from the CN  
76 data.

## 77 2 Results

### 78 2.1 Sitka: scalable single cell phylogenetic tree inference

79 **Fig. 1** shows the workflow of the sitka method. Sitka is based on a transformation of single  
80 cell copy number matrices retaining only presence or absence of changes in copy number  
81 profiles between contiguous genomic bins. This transformation allows us to approximate  
82 a complex evolutionary process (integer-valued copy numbers, prone to a high degree of  
83 homoplasy and dense dependence structure across sites) using a probabilistic version  
84 of a perfect phylogeny (see **Supplementary Fig. 1**). We leverage the special structure  
85 created by the change point transformation to build a special purpose MCMC kernel, which  
86 has better computational scalability per move compared to classical phylogenetic kernels  
87 (Methods section 9.4.3).

88 We visualise the input data to sitka in a colour-coded matrix exemplified in **Supplementary**  
89 **Fig. 1-a**. Each row in the matrix corresponds to an individual cell that has been sequenced  
90 in a single-cell platform. Each column in the matrix is a locus that is represented by a bin  
91 (a contiguous set of genomic positions). We assume that the integer copy number of each  
92 bin has been estimated as a preprocessing step, e.g., using a hidden Markov model [16].  
93 In **Supplementary Fig. 1-a** the copy number state is encoded by the colour of each entry  
94 in the matrix.

95 The output of sitka includes two types of directed rooted trees. Type I is the tree used for  
96 MCMC sampling in the inference procedure, and type II, which is derived from type I, is  
97 used in visualisation (**Fig. 2-a-c**). The set of nodes in a type I tree is given by the union of  
98 the cells, the CN change points (markers) under study, and a root node  $v^*$ . The topology  
99 of a type I tree bears the following phylogenetic interpretation: given a cell  $c$  in the tree,  $c$   
100 is hypothesized to harbour the markers in the shortest path between  $c$  and the root node  
101  $v^*$ , and only those markers. We enforce the constraint that all cells are leaf nodes, while  
102 markers can be either internal or leaf nodes. Markers placed at the leaves are interpreted  
103 as outliers, for example measured CN change points that are false positives.

104 We remove from the type I tree all marker nodes that are leaf nodes, i.e., markers that are  
105 not present in any cells. We also collapse into a single node, the list of connected marker  
106 nodes that have exactly one descendent (i.e., chains). **Supplementary Fig. 2** shows a  
107 small *type I tree*, its transformation to a *type II tree* and the respective marker matrix. We  
108 visualise the input matrix and the estimated tree simultaneously by sorting the individual  
109 cells (rows of the matrix) such that they line up with the position of the corresponding leaves  
110 of the tree.

111 Sitka uses change points as phylogenetic traits modelled using a relaxation of the perfect  
112 phylogeny assumption. Change points arising from non-overlapping CNA events do not  
113 break the perfect phylogeny assumption. **Supplementary Fig. 3** shows examples of over-  
114 lapping CNA events and their effect on markers. The two scenarios that can lead to the  
115 violation of the perfect phylogeny assumption are (i) when a CNA gain event is followed  
116 by an overlapping loss event or (ii) when a loss event is followed by an overlapping loss  
117 event, and the second event removes either end-point of the first event. For both (i) and  
118 (ii), a violation occurs only when the second overlapping event hits the same copy as the  
119 first event.

120 Imposing a perfect phylogeny on the *observed* change points is restrictive, as we expect  
121 both violations of the assumptions (e.g., due to homoplasy), and measurement noise. To  
122 address this we use an observation model (Methods section 9.4.1) which assigns positive  
123 probability to arbitrary deviations from the perfect phylogeny assumption, while encourag-  
124 ing configurations where few loci and cells are involved in violations. Subsequently we  
125 impose the perfect phylogeny assumption on a *latent* marker matrix defined as follows.  
126 Given a type I tree  $t$ , the latent marker matrix  $x$  is a deterministic function  $x = x(t)$ . We  
127 compute  $x : t \rightarrow \{0, 1\}^{C \times L}$  by setting  $x_{c,l} = 1$  if the single-cell  $c$  is a descendent of the  
128 marker node  $l$  in tree  $t$ , and otherwise  $x_{c,l} = 0$ . We use  $y_{c,l}$  to refer to the observed change  
129 point  $l$  in individual cell  $c$  (Methods section 9.4.1).

130 Synthetic experiments show that sitka's performance degrades gracefully in the face of  
131 some of the key types of expected violation of the perfect phylogeny assumption (**Fig. 3-**  
132 **a,b**, Methods section 9.5).

## 133 2.2 Performance of sitka relative to alternative approaches

134 We compare the performance of sitka to alternative approaches on three scWGS datasets  
135 introduced here (**Fig. 2-a-c**). The first dataset, *SA535*, is generated for this project and  
136 contains 679 cells from three passages of a triple negative breast cancer (TNBC) patient  
137 derived xenograft sample. Passages X1, X5, and X8 had 62, 369, and 231 cells post quality  
138 filtering (Methods section 9.1) respectively. We also include 17 mostly diploid control cells.  
139 These cells are combined to generate the input to the analysis pipeline (**Supplementary**  
140 **Fig. 6**). The second dataset, labelled *OVA*, consists of cells from three samples taken  
141 from a patient with high grade serous (HGS) ovarian cancer. The first sample, *SA1090*,  
142 was from an ascites pre-treatment, while *SA922* was from an ascites post-treatment. The  
143 third sample, *SA921*, was taken from the ovary. See **Supplementary Fig. 7** for the tree and  
144 the CNA profile heatmap for this dataset. The final dataset, *SA501* [25], is another TNBC  
145 xenograft tumour from 6 untreated passages, namely X2, X5, X6, X8, X11, and X15. After  
146 filtering, 515, 236, 328, 189, 836, and 308 cells remain in each passage respectively (for  
147 a total of 2,412 cells, see **Supplementary Fig. 8**). Table 1 shows the attrition after each  
148 step of filtering cells per passage in each dataset.

149 To evaluate inferred trees from sitka and other tree reconstruction methods, we use a good-  
150 ness of fit performance metric, which compares the compatibility of observed CN change  
151 points with a given phylogeny using Youden's J index (Methods section 9.6, **Fig. 2-d**). Sitka  
152 has the highest Youden's index across all three datasets. UPGMA and WPGMA perform simi-  
153 larly on *SA501* and *SA535*. UPGMA performs slightly better than WPGMA on the *OVA* dataset.  
154 HDBSCAN has a close but slightly smaller Youden's index than UPGMA over the *SA535* and  
155 *OVA* datasets, but performs marginally better on *SA501*. NJ trails WPGMA on *SA501* and  
156 the *OVA* datasets, and has the lowest Youden's index on *SA535*. MrBayes performs well  
157 on the smallest dataset, *SA535*, with MrBayes-np2 and MrBayes-np8 performing similar to  
158 WPGMA, and MrBayesWithBinaryInput having achieved the second highest Youden's in-  
159 dex. On the *OVA* data, MrBayesWithBinaryInput and MrBayes-np2 trail behind NJ, while  
160 MrBayes-np2 has the lowest Youden's index among all methods on all datasets. Similar  
161 to the *OVA* case, MrBayesWithBinaryInput and MrBayes-np2 trail behind NJ over the  
162 *SA501* dataset. Following [25], we run MrBayes for 10,000,000 generations. MrBayes-np8  
163 had completed only 278,000 iterations running on *SA501* after several days. The results  
164 in this comparison suggest that sitka performs better than the baseline methods. Running

165 sitka on the real-world datasets took on average 22.3, 46.6, and 12.9 hours for the *OVA*,  
166 *SA501*, and *SA535* datasets respectively, on a Linux workstation with 72 Intel Xeon Plat-  
167 inum 8272CL 2.60GHz CPU processors and 144 GB of memory. We complement these  
168 benchmarking results with experiments on synthetic data where sitka is also the highest  
169 performing method based on metrics measuring phylogenetic tree distance.

### 170 2.3 Single cell resolution phylogenetic inference in PDX

171 Here we analyse the foregoing three multi-sample datasets. To visualise the tree inference  
172 results we arrange the inferred consensus tree  $t$  (Methods section 9.4.5) and the cell-by-  
173 locus CN matrix side by side where the rows of the matrix correspond to the position of  
174 individual cells on the tree and the markers are arranged by their genomic position (**Fig. 1-**  
175 **h**). **Fig. 2-a-c** shows examples of the multi-channel visualisation where each marker is  
176 represented by a tuple of three different data-types or *channels*, namely: (i) the latent  
177 markers induced by the consensus tree,  $x(t)$ ; (ii) the matrix of marginal posterior probability  
178 that cell  $c$  is a descendent of marker  $l$ , computed via the average  $\bar{m}$  (**Fig. 1-g**, Methods  
179 section 9.4.5); and (iii) the sitka transformed input data  $y_{c,l}$ .

180 We use this view to assess potential discrepancies between the input data and the inferred  
181 tree. In most cells and loci (as quantified in **Supplementary Fig. 9.6**), the observed data  
182 is in close agreement with the inferred tree. In the following we provide some examples  
183 of disagreements. Consider first the ChrX in the *OV2295* dataset (**Fig. 2-a**). ChrX has a  
184 long orange band (inferred marker in channel (i)) not matched by a black band (observed  
185 marker in channel (iii)) suggesting that a perfect phylogeny violation has occurred. The  
186 pattern in this marker is consistent with the presence of an ancestral event followed by a  
187 deletion. In **Fig. 2-b**, a set of diploid cells are attached to the root of the tree. These are  
188 control cells included in the experiment and correspond to a region in the bottom of the  
189 matrix with no inferred markers (orange bands) and almost no observed markers (black  
190 bands). In this dataset, there are change points where the observed marker has a high  
191 density (black band), but the tree is reconstructed with the marker absent (no matching  
192 orange band). Examples can be found in Chr1, Chr7 and Chr16. One possible explanation  
193 could be that the end-points of each event were detected as slightly shifted across cells.  
194 For instance, in **Supplementary Fig. 8** there are two loci with an amplification (CN state  
195 equal to three) in Chr1p where cells that harbour a mutation in the first locus appear not to  
196 have a mutation in the second locus, suggesting that the same event was called in the first  
197 locus in some cells, and in the second locus in others. An alternative hypothesis is that  
198 the cells in this dataset have a mutator phenotype that promotes CN mutations in these  
199 loci.

200 **Supplementary Fig. 9** shows the distribution of mismatch rates for each dataset, de-  
201 fined as the fraction of times that the observed and inferred markers do not match,  
202 i.e.,  $\frac{1}{C} \sum_{c \in C} \mathbf{1}[y_{c,l} \neq x_{c,l}]$  for  $l \in L$  (corresponding to the black and orange bands in  
203 **Fig. 2-a**). In *OV2295*, 41 markers (11%) have a mismatch rate of over 50%, where  
204 marker *chr15\_67000001\_67500000* has the highest mismatch rate at 70%. In *SA501*,  
205 30 markers (11%) have a mismatch rate of over 50%, 13 of which (5%) have a mis-  
206 match rate of over 75%. *SA535* has the lowest maximum mismatch rate at 49% (marker  
207 *15\_72000001\_72500000*).

## 208 2.4 Placement of SNVs using the CNA inferred tree

209 To determine the presence or absence of SNVs in cells using data with high levels of  
210 missingness, we develop an extension of sitka, the sitka-snv model. Given single cell level  
211 variant read counts, the model incorporates CN data to place SNVs on the sitka-inferred  
212 phylogenetic tree. This *backbone* CN tree, provides a principled way to pool statistical  
213 strength across groups of single cells sequenced at low coverage, including data from the  
214 DLP+ platform [16]. The output of the sitka-snv model is an *extended* tree that has marker  
215 nodes that comprise SNVs in addition to the original CNAs.

216 The SNVs are added to the existing CNA-based tree with the computational complexity  
217 of  $O(|C| + |L|)$  per SNV. **Fig. 3-c** shows the result of SNV placement with the number  
218 of variant reads in *SA535*, corresponding to the tree shown in **Fig. 2-c**. **Supplementary**  
219 **Figs. 17, 18, and 19** show the number of variant reads and the matching SNV call proba-  
220 bilities for the *SA535*, *OVA* and *SA501* datasets respectively. Sitka and sitka-snv provide  
221 a comprehensive genomic analysis tool for large scale low-coverage scWGS.

## 222 3 Discussion

223 In this work we use data in which the genome of the single cells CNA profiles are partitioned  
224 into bins of a fixed size (500Kb), each assigned a constant integer CN state. The relatively  
225 large size is due to the low coverage inherent to the scWGS platform, but it implies that the  
226 same bin may harbour multiple CNA events. Biological processes that result in complex  
227 DNA rearrangements could further increase the probability of having two hits in one bin [26,  
228 27]. Such multiple hits can violate the perfect phylogeny assumptions. This highlights  
229 the importance of our goodness-of-fit and visualisation methods as they can detect such  
230 violations.

231 Structural variations such as chromothripsis, that affect multiple segments of the genome  
232 at the same time, make it difficult to determine the rate of CNA events and suggest that  
233 CNA events may not be suitable molecular clocks to estimate branch lengths. One possi-  
234 ble remedy is to first infer the tree topology via markers based on CNA events and then  
235 conditioned on this topology, add SNVs to the tree. The number of SNVs on each edge of  
236 the tree may be used to inform branch lengths.

237 Our preprocessing pipeline excludes multiple cells from the analysis (see Table 1). We filter  
238 out a fraction of cells to remove contaminated cells, either doublets or mouse cells, cells  
239 with too many erroneous sequencing artefacts, and cycling cells. Removing a portion of the  
240 sequenced cells will decrease the statistical power to determine the subclonal structure of  
241 the population—an important application of this work—, and may bias the sampling against  
242 clones that have a higher division rate. We expect this will be an intrinsic limitation to any  
243 scWGS phylogenetic methods and this motivates the design of improved classification  
244 methods detecting cell cycling from genomic and imaging data.

245 Evaluating the performance of a phylogenetic reconstruction method on real-world  
246 datasets is difficult, mainly due to a lack of ground truth. One promising area of research is  
247 the use of CRISPR-Cas9 based lineage tracing [6]. In absence of ground truth data, we de-  
248 veloped a goodness-of-fit framework that to our knowledge enables a first of a kind bench-  
249 marking of phylogenetic inference methods over real-world scWGS CNA datasets.

250 Phylogenetic tree reconstruction is a principled way to identify subpopulations in a hetero-  
251 geneous single-cell population. This in turn enables the use of population genetics models  
252 that track the abundance of subpopulations over multiple timepoints [5] and to make infer-  
253 ences about the evolutionary forces acting on each clone. Further study with timeseries  
254 modelling will provide insight into therapeutic strategies promoting early intervention, drug  
255 combinations and evolution-aware approaches to clinical management.

## 256 **4 Acknowledgements**

257 This project was generously supported by the BC Cancer Foundation at BC Cancer  
258 and Cycle for Survival supporting Memorial Sloan Kettering Cancer Center. SPS holds  
259 the Nicholls Biondi Chair in Computational Oncology and is a Susan G. Komen Scholar  
260 (#GC233085). SA holds the Nan and Lorraine Robertson Chair in Breast Cancer and is a  
261 Canada Research Chair in Molecular Oncology (950-230610). Additional funding provided  
262 by the Terry Fox Research Institute Grant 1082, Canadian Cancer Society Research In-  
263 stitute Impact program Grant 705617, CIHR Grant FDN-148429, Breast Cancer Research  
264 Foundation award (BCRF-18-180, BCRF-19-180 and BCRF-20-180), MSK Cancer Cen-  
265 ter Support Grant/Core Grant (P30 CA008748), National Institutes of Health Grant (1RM1  
266 HG011014-01), CCSRI Grant (#705636), the Cancer Research UK Grand Challenge Pro-  
267 gram, Canada Foundation for Innovation (40044) to SA, SPS and ABC. We extend our  
268 gratitude to Sarah P. Otto for her helpful comments on a draft of this manuscript.

## 269 **5 Funding**

## 270 **6 Author Contributions**

271 SS, FD: computational method development, data analysis, manuscript writing; KC: data  
272 analysis, manuscript writing; KRC, AR: method development; FK, data generation; DL,  
273 MA, AM, MW, TF: computational biology, data analysis; NR: manuscript editing; SA: data  
274 generation and oversight; SPS: method development and oversight; ABC: project concep-  
275 tion and oversight, statistical inference method development, manuscript writing, senior  
276 responsible author;

## 277 **7 Competing Interests**

278 S.P.S. and S.A. are founders, shareholders, and consultants of Canexia Health Inc.

## 279 **8 Code availability**

280 Sitka is available at <https://github.com/UBC-Stat-ML/sitkatree.git>.



## 281 9 Methods

### 282 9.1 Pre-processing

283 The raw data contain cells that are either contaminated (e.g., contains biological material  
284 from mice) or have undesired sequencing artefacts. These include cells that were captured  
285 for DNA sequencing when undergoing mitosis. Since the sitka model does not account for  
286 such phenomena, the filtering is an important step. **Supplementary Fig. 15** shows the  
287 steps taken from pulling the raw data to the CNA integer matrix ready for sitka transfor-  
288 mation (details in the Supplementary Information). Briefly, we remove control cells, cells  
289 with highly-noisy CN calls, and cells that have very few mapped reads. We also remove  
290 copy number bins that lie in difficult to sequence regions of the genome (bins with low-  
291 mappability). Finally, we drop cells that, based on their CNA profile, are suspected to be  
292 cycling cells.

### 293 9.2 The sitka transformation

294 To obtain the  $C \times L_{\text{Markers}}$  phylogenetic markers matrix  $y$  that comprises the input to the  
295 sitka model, we apply a lossy transformation to the  $C \times L_{\text{Bins}}$  CNA matrix  $a$  that involves  
296 computing the change in copy number state between two consecutive bins. **Supplemen-**  
297 **tary Fig. 1** shows a small CNA matrix and its corresponding transformation into the marker  
298 matrix. For brevity, in what follows we assume that only one chromosome is used, so that  
299  $L_{\text{Bins}} = L$  and  $L_{\text{Markers}} = L_{\text{Bins}} - 1$ . In practice, we use all available chromosomes, and  
300  $L_{\text{Markers}} = L_{\text{Bins}} - N_{\text{Chr}}$  where  $N_{\text{Chr}}$  denotes the total number of chromosomes used.

301 Given a filtered cell-by-locus matrix  $a$ , we sort bins by their genomic position. Then in each  
302 chromosome, we compute markers as the binarised difference between consecutive bins.  
303 In other words,  $y = (y_{c,l'})$  and  $l' \in \{1, \dots, L - 1\}$ , and

$$y_{c,l'} := \mathbf{1}(|a_{c,l'} - a_{c,l'+1}| > 0), \quad (1)$$

304 where  $\mathbf{1}(x)$  is the indicator function.

### 305 9.3 Fixing jitter and selection of phylogenetic markers

306 The copy numbers available to us in this work are estimated independently for each cell.  
307 This is one reason why the start position (bin) of the same CN change event may be  
308 slightly different across cells, generating some *jitter*. We address this by enumerating  
309 each change point column in order of decreasing density (where the density of column  $l$  is  
310 given by  $\sum_{c \in C} y_{c,l} / |C|$ ) and merging the column with its  $k = 2$  immediate neighbours (see  
311 Algorithm 1 for details). An example of the result of the jitter correction heuristic is shown  
312 in **Fig. 1** panel **c**. To speed-up computation, only a subset of markers present in at least  
313 a minimum number of cells are chosen for phylogenetic inference. That is, we removed  
314 columns  $l$  in  $y$  with relative density  $\sum_{c \in C} y_{c,l} / |C|$  less than a threshold, set to 5%. Larger  
315 values of this threshold may lead to less resolved clades in the inferred tree.

---

### Algorithm 1 JitterFix

---

```

1: procedure JITTER-FIX( $y, k$ )
2:   column-queue  $\leftarrow$  OrderByDensityDecreasing( $y$ )
3:   columns-visited  $\leftarrow$  {}
4:   for column-index  $c$  in column-queue do
5:     neighbours  $\leftarrow$  neighbours( $c, y, k$ )
6:     for column-index  $n$  in neighbours do  $\triangleright$  The function neighbours is defined as the  $k$  columns to the
       left and  $k$  to the right of  $c$  (when applicable)
7:       if  $n \notin$  columns-visited then
8:          $y_{1:C,c} \leftarrow y_{1:C,c} \vee y_{1:C,n}$ 
9:          $y_{1:C,n} \leftarrow 0$ 
10:        columns-visited  $\leftarrow$  columns-visited  $\cup n$ 
11:   return  $y$ 

```

---

## 316 9.4 The sitka model

### 317 9.4.1 Model description

318 The sitka model starts with the perfect phylogeny assumption for the latent variables  $x_{c,l}$   
319 but allows deviation from it via allowing noisy observations  $y_{c,l}$ . In a perfect phylogeny  
320 model, each phylogenetic trait arises only once on the rooted tree topology and all cells  
321 descending from that position will inherit that trait and no deletions are allowed.

322 Let  $C$  and  $L$  denote the disjoint sets of cells and loci respectively.

We posit an observation probability model  $p(y|x, \theta)$ , where  $\theta$  are model parameters described shortly, and both  $x$  and  $y$  are cell by locus matrices, the former being latent (derived from the unobserved tree via  $x = x(t)$ ), while the latter is the matrix obtained from the sitka transformation. To model errors in copy number calls as well as perfect phylogeny violations, we introduce false positive and negative rate parameters  $r^{\text{FP}} \in (0, 1)$  and  $r^{\text{FN}} \in (0, 1)$  respectively, and an error matrix

$$e^{r^{\text{FP}}, r^{\text{FN}}} = \begin{bmatrix} 1 - r^{\text{FP}} & r^{\text{FP}} \\ r^{\text{FN}} & 1 - r^{\text{FN}} \end{bmatrix},$$

$$p(y_{c,l}|x_{c,l}, r^{\text{FP}}, r^{\text{FN}}) = e^{r^{\text{FP}}, r^{\text{FN}}}_{x_{c,l}, y_{c,l}},$$

323 from which we set:

$$p(y|x, \theta) = \prod_{l \in L} \prod_{c \in C} p(y_{c,l}|x_{c,l}, r_{c,l}^{\text{FP}}(\theta), r_{c,l}^{\text{FN}}(\theta)).$$

324 We define two type of models, differing in the choice of functions  $r_{c,l}(\cdot)$  and dimension-  
325 ality of  $\theta$ : one based on global error parameters, and one based on locus-specific error  
326 parameters.

327 For the global parameterization,  $\theta = \theta_{\text{global}} = (r_{\text{global}}^{\text{FN}}, r_{\text{global}}^{\text{FP}})$ , and the false positive and  
328 false negative functions are given by  $r_{c,l}^{\text{FP}}(\theta_{\text{global}}) = r_{\text{global}}^{\text{FP}}$  and  $r_{c,l}^{\text{FN}}(\theta_{\text{global}}) = r_{\text{global}}^{\text{FN}}$ .

329 For the locus-specific error model, we set the error rates to be locus-dependent:  $\theta =$   
330  $(r_1^{\text{FP}}, r_2^{\text{FP}}, \dots, r_{|L|}^{\text{FP}}, r_1^{\text{FN}}, r_2^{\text{FN}}, \dots, r_{|L|}^{\text{FN}})$ ,  $r_{c,l}^{\text{FP}}(\theta) = r_l^{\text{FP}}$  and  $r_{c,l}^{\text{FN}}(\theta) = r_l^{\text{FN}}$ . With this extra flexibil-  
331 ity, the model can discount the effect of a trait violating the perfect phylogeny assumption,  
332 by setting high error rates for the trait's locus.

333 The two parameterizations are compared in the Supplementary Information. We use the  
 334 global parameterization by default unless mentioned otherwise.

335 In both the global and locus-specific parameterizations, we need to construct a prior dis-  
 336 tribution  $p(\theta)$  over the error parameters. Using a uniform prior distribution with support on  
 337  $[0, 1]$  can lead to pathological cases as shown in **Supplementary Fig. 4**. To avoid that, we  
 338 use the following prior distributions on the two types of error:

$$r^{\text{FP}} \sim \text{Uniform}\left(0, \overline{r^{\text{FP}}}\right),$$

$$r^{\text{FN}} \sim \text{Uniform}\left(0, \overline{r^{\text{FN}}}\right).$$

339 We use  $\overline{r^{\text{FP}}} = 1/10$  and  $\overline{r^{\text{FN}}} = 1/2$  as default in our experiments.

340 Next, we describe the prior  $p(t)$  on phylogenies using a two-step generative process:

341 **Sampling a mutation tree:** let  $\mathcal{V}^m = L \cup \{v^*\}$  denote a vertex set composed of one vertex  
 342 for each of the  $|L|$  loci plus one artificial root node  $v^*$ . The artificial root node induces  
 343 an implicit notion of direction on the edges, viewing them as pointing away from  $v^*$ .  
 344 Let  $\mathcal{T}^m$  denote the set of trees  $t^m$  spanning  $\mathcal{V}^m$ . The interpretation of  $t^m$  is as follows:  
 345 there is a directed path from vertex/locus  $l$  to  $l'$  in  $t^m$  if and only if the trait indexed by  
 346  $l$  is hypothesized to have emerged in a cell which is ancestral to the cell in which  $l'$   
 347 emerged. Pick one element  $t^m \in \mathcal{T}^m$ .

348 **Sampling cell assignments:** assign each cell to a vertex in  $t^m$ . The interpretation of  
 349 assigning cell  $c$  to locus  $l$  is that among the traits under study,  $c$  is hypothesized to  
 350 possess only the traits visited by the shortest path from  $v^*$  to  $l$  in  $t_m$ . If a cell  $c$  is  
 351 assigned to  $v^*$ , the interpretation is that  $c$  is hypothesized to possess none of the  
 352 traits under study.

The number of possible trees obtained from this two-step sampling process is:

$$\begin{aligned} |\mathcal{T}| &= |\mathcal{T}^m| |\{f : C \rightarrow L \cup \{v^*\}\}| \\ &= (|L| + 1)^{(|L|+1)-2} (|L| + 1)^{|C|} \\ &= (|L| + 1)^{|L|+|C|-1}, \end{aligned}$$

353 where we use Cayley's formula to compute  $|\mathcal{T}^m|$ . Hence the uniform prior probability mass  
 354 function over the possible outputs of this two-step sampling process is given by:

$$p(t) = \frac{\mathbf{1}[t \in \mathcal{T}]}{(|L| + 1)^{|L|+|C|-1}},$$

355 where  $\mathcal{T}$  is the set of all perfect phylogenetic trees that result from the two step generative  
 356 process described above. This simple prior has a useful property: if a collection of say  
 357 two splits are supported by  $m_1$  and  $m_2$  traits, then the prior probability for an additional  
 358 trait to support the first versus second split is proportional to  $(m_1 + 1, m_2 + 1)$ . Therefore,  
 359 there is a "rich gets richer" behaviour built-in into the prior, which is viewed as useful in  
 360 many Bayesian non-parametric models. Of course, more complicated priors over  $\mathcal{T}$  could  
 361 be easily incorporated as the complexity of inference typically comes from the likelihood  
 362 rather than the prior. Simulation from the prior can be performed using Wilson's algorithm  
 363 [28], followed by independent categorical sampling to simulate the cell assignments.

## 364 9.4.2 Inference

365 The posterior distribution,

$$\pi(t, \theta) \propto p(t)p(\theta)p(y|x(t), \theta),$$

366 is approximated using MCMC. Two MCMC moves are used, described in the next two  
367 sections. The posterior distribution is summarized using a Bayes estimator described in  
368 Section 9.4.5. The model is implemented in the Blang probabilistic programming language  
369 [29].

## 370 9.4.3 MCMC tree exploration move

371 Sitka uses a tree sampling move to efficiently explore, at each MCMC iteration, the pos-  
372 terior distribution in a large neighbourhood of a given tree. Given a tree  $t$  and locus  $l$ ,  
373 we define a neighbourhood  $N^l(t) \subset \mathcal{T}$  by removing  $l$  from  $t$ , and considering all possible  
374 ways to reattach  $l$  and hence defining a neighbourhood of phylogenetic trees (we also im-  
375 plemented a separate move reattaching cell nodes instead of locus nodes, its derivation  
376 follows similar lines as the move described in this section). The process of removing  $l$  is  
377 called an *edge-contraction* (removing an edge after connecting its two end-points) while  
378 the process of adding back a locus is called an *edge-insertion*. An edge insertion can be  
379 described as follows:

- 380 1. Pick a non-cell vertex  $v$ , i.e. an element from the set  $R = \{v^*\} \cup L \setminus \{l\}$  where  $v^*$  is  
381 the root node.
- 382 2. Pick any subset of  $v$ 's descendent subtrees and disconnect them from  $v$ .
- 383 3. Add a new node  $l$  under  $v$  and move the selected nodes from step 2 above and attach  
384 them to  $l$ .

385 **Fig. 1-f** (right) shows an example of an edge-insertion. A locus named *chr15\_5950*  
386 coloured red, has three children at MCMC iteration 100. This corresponds to node  $v$  in the  
387 above description. In step 2 of the edge insertion process, two of its children, namely cells  
388 *RC07C* and *RC05C4* are chosen and disconnected from  $v$ . They are then inserted under  
389 locus *chr1\_4900*, corresponding to  $l$ , which becomes a child of locus *chr15\_5950*.

390 In the following, we derive the probability distributions to be used in steps 1 and 2 above  
391 that lead to a Gibbs sampling algorithm (i.e. an MCMC move with no rejection step). The  
392 Gibbs sampler first selects a locus  $l$  from a fixed distribution (a tuning parameter), which  
393 we take for simplicity as being uniform over the  $|L|$  loci.

394 After having sampled  $l$ , we partition  $N^l(t_{\setminus l})$  into blocks corresponding to the choice of node  
395  $v$  made in Step 1,  $N^l(t_{\setminus l}) = \cup_v N_v^l(t_{\setminus l})$ . The Gibbs conditional probabilities required in step  
396 1 above are of the form:

$$\bar{\rho}_v = \frac{\rho_v}{\sum_{\tilde{v} \in R} \rho_{\tilde{v}}},$$

where:

$$\rho_v = \sum_{t \in N_v^l(t_{\setminus l})} p(t)p(y|x(t), \theta), \quad (2)$$

and  $t_{\setminus l}$  denotes the tree obtained after performing an edge contraction, where the contracted edge is between  $l$  and the parent node of  $l$ . To compute  $\rho_v$  efficiently, we start with the following likelihood recursion for all vertex  $v$  in  $t_{\setminus l}$ . First, for all vertices  $c$  corresponding to a cell and  $b \in \{0, 1\}$ , define:

$$p_c^b = p(y_{c,l} | b, \theta).$$

397 Next, we perform the following bottom-up recursion for all subtrees of  $t_{\setminus l}$ : for all  $v \in R$ ,  
 398  $b \in \{0, 1\}$ ,

$$p_v^b = \prod_{v'' \in \text{children}(v)} p_{v''}^b,$$

399 where  $\text{children}(v)$  denotes the list of children of vertex  $v$ .

400 We can now return to the problem of computing  $\bar{\rho}_v$ . First, observe that the sum in Equa-  
 401 tion (2) can be re-indexed by a bit vector  $\mathbf{b} = (b_1, b_2, \dots, b_k)$ ,  $b_{v''} \in \{0, 1\}$  of length equal  
 402 to  $k = |\text{children}(v)|$ . Each bit  $b_{v''}$  is equal to one if children  $v''$  is to be moved into a child  
 403 of  $v'$  (refer to **Supplementary Fig. 5**), and zero if it is to stay as a child of  $v$ . For each  
 404 possible assignment, we obtain a tree  $t \in N_v^l(t_{\setminus l})$ , and its probability can be decomposed  
 405 into factors corresponding to cells that are descendant of  $v$  (denoted  $C_v$ , solid red thick line  
 406 under the tree of **Supplementary Fig. 5-B**) and those that are not (denoted  $C_{\setminus v}$ , dashed  
 407 green thick line under the tree of **Supplementary Fig. 5-B**).

408 The product of the likelihood factors corresponding to cells that are not descendants of  $v$   
 409 (“outside product”) does not depend on the choice of the bit vector. This outside product  
 410 can be obtained as follows:

$$\prod_{c \in C_{\setminus v}} p_c^0 = \frac{p_{v^*}^0}{p_v^0}.$$

411 Note that this assumes  $p_v^0 > 0$ . As a workaround to cases where there are structural  
 412 zeros, we recommend injecting small numerical values if  $p_v^0 = 0$  (we used  $10^{-6}$  in our  
 413 implementation).

For the cells under  $v$ , we now have to take into account whether they are selected under the newly introduced locus or not. More precisely, for each of the children  $v_1, v_2, \dots, v_k$ , we have to take into account the value of the bit vector  $\mathbf{b} = (b_1, b_2, \dots, b_k)$ . The sum over possible assignments written naively has a number of terms which is exponential in  $k$ , but can be rewritten into a product over  $k$  factors:

$$\sum_{t \in N_v^l(t_{\setminus l})} \prod_{c \in C_v} p_c^{x_{c,l}(t)} = \sum_{b_1=0}^1 \dots \sum_{b_k=0}^1 \prod_{i=1}^k p_{v_i}^{b_i} = \prod_{i=1}^k (p_{v_i}^0 + p_{v_i}^1).$$

Putting it all together, we obtain for some constants  $K_i$  independent of  $v$ :

$$\begin{aligned} \rho_v &= K_1 \sum_{t \in N_v^l(t_{\setminus l})} p(y|x(t), \theta) \\ &= K_1 \sum_{t \in N_v^l(t_{\setminus l})} \prod_{l' \in L} \prod_{c \in C} p(y_{c,l'} | x_{c,l'}(t), r_{c,l'}^{\text{FP}}(\theta), r_{c,l'}^{\text{FN}}(\theta)) \end{aligned}$$

$$\begin{aligned}
&= K_1 \left( \prod_{l' \in L, l' \neq l} \prod_{c \in C} p \left( y_{c,l'} | x_{c,l'}(t), r_{c,l'}^{\text{FP}}(\theta), r_{c,l'}^{\text{FN}}(\theta) \right) \right) \sum_{t \in N_v^l(t_{\setminus l})} \prod_{c \in C} p \left( y_{c,l} | x_{c,l}(t), r_{c,l}^{\text{FP}}(\theta), r_{c,l}^{\text{FN}}(\theta) \right) \\
&= K_1 K_2 \sum_{t \in N_v^l(t_{\setminus l})} \prod_{c \in C} p \left( y_{c,l} | x_{c,l}(t), r_{c,l}^{\text{FP}}(\theta), r_{c,l}^{\text{FN}}(\theta) \right) \\
&= K_1 K_2 \sum_{t \in N_v^l(t_{\setminus l})} \prod_{c \in C} p_c^{x_{c,l}(t)} \\
&= K_1 K_2 \sum_{t \in N_v^l(t_{\setminus l})} \left( \prod_{c \in C_v} p_c^{x_{c,l}(t)} \right) \left( \prod_{c \in C \setminus v} p_c^{x_{c,l}(t)} \right) \\
&= K_1 K_2 \left( \prod_{c \in C \setminus v} p_c^{x_{c,l}(t)} \right) \sum_{t \in N_v^l(t_{\setminus l})} \prod_{c \in C_v} p_c^{x_{c,l}(t)} \\
&= K_1 K_2 \left( \frac{p_v^0}{p_v^0} \right) \sum_{t \in N_v^l(t_{\setminus l})} \prod_{c \in C_v} p_c^{x_{c,l}(t)} \\
&= K_1 K_2 \left( \frac{p_v^0}{p_v^0} \right) \prod_{i=1}^k (p_{v_i}^0 + p_{v_i}^1) \\
&= K_1 K_2 K_3 \frac{\prod_{i=1}^k (p_{v_i}^0 + p_{v_i}^1)}{p_v^0}.
\end{aligned}$$

414 Putting these together we can compute the probabilities required in step 1 above:

$$\bar{\rho}_v = \frac{\rho_v}{\sum_{\tilde{v} \in R} \rho_{\tilde{v}}} \quad (3)$$

$$\begin{aligned}
&= \frac{\left( \frac{\prod_{v_i \in \text{children}(v)} (p_{v_i}^0 + p_{v_i}^1)}{p_v^0} \right)}{\sum_{\tilde{v} \in R} \left( \frac{\prod_{v'_i \in \text{children}(\tilde{v})} (p_{v'_i}^0 + p_{v'_i}^1)}{p_{\tilde{v}}^0} \right)}. \quad (4)
\end{aligned}$$

415 Once  $v$  is sampled, we choose a subset of its children to move to  $v'$  by sampling  $k$  inde-  
416 pendent Bernoulli random variables with the  $i$ -th one having bias

$$\frac{p_{v_i}^1}{p_{v_i}^0 + p_{v_i}^1},$$

417 and selecting children with corresponding Bernoulli realisations of 1.

#### 418 9.4.4 MCMC parameter exploration move

419 To resample the parameters  $\theta$  we condition on the tree  $t$ , and hence on the hidden state  
420 matrix  $x$ , and update  $\theta$  in a Metropolis-within-Gibbs framework. There are two different

421 samplers depending on whether the global or locus-specific parameterization is used. We  
 422 start with describing the former.

423 We compute two sufficient statistics from the matrix  $x$  (i) the number of false positive in-  
 424 stances,  $n^{\text{FP}}$ , and (ii) the number of false negative instances,  $n^{\text{FN}}$ ,

$$n^{\text{FP}} = n^{\text{FP}}(x) = \sum_{c \in C} \sum_{l \in L} \mathbf{1}[x_{c,l} = 0, y_{c,l} = 1]$$

$$n^{\text{FN}} = n^{\text{FN}}(x) = \sum_{c \in C} \sum_{l \in L} \mathbf{1}[x_{c,l} = 1, y_{c,l} = 0].$$

Based on these cached statistics, we obtain:

$$p(y|x, \theta_{\text{global}}) \propto (r^{\text{FP}})^{n^{\text{FP}}} (r^{\text{FN}})^{n^{\text{FN}}} (1 - r^{\text{FP}})^{n^{\text{N}} - n^{\text{FN}}} (1 - r^{\text{FN}})^{n^{\text{P}} - n^{\text{FP}}}, \quad (5)$$

where the the number of positive  $n^{\text{P}}$  and negative  $n^{\text{N}}$  instances in the data can be pre-computed,

$$n^{\text{P}} = \sum_{c \in C} \sum_{l \in L} \mathbf{1}[y_{c,l} = 1]$$

$$n^{\text{N}} = |C||L| - n^{\text{P}}.$$

425 Based on the above expression, which can be evaluated in  $O(1)$  once the statistics are  
 426 computed, we then use a slice sampling algorithm to update the parameters [30].

The sampler for the locus-specific parameterization is very similar. The main difference is that we compute the statistics for each locus  $l$ :

$$n_l^{\text{FP}} = n_l^{\text{FP}}(x) = \sum_{c \in C} \mathbf{1}[x_{c,l} = 0, y_{c,l} = 1]$$

$$n_l^{\text{FN}} = n_l^{\text{FN}}(x) = \sum_{c \in C} \mathbf{1}[x_{c,l} = 1, y_{c,l} = 0]$$

$$n_l^{\text{P}} = \sum_{c \in C} \mathbf{1}[y_{c,l} = 1]$$

$$n_l^{\text{N}} = |C| - n_l^{\text{P}}$$

$$p(y|x, \theta) = \prod_l (r_l^{\text{FP}})^{n_l^{\text{FP}}} (r_l^{\text{FN}})^{n_l^{\text{FN}}} (1 - r_l^{\text{FP}})^{n_l^{\text{N}} - n_l^{\text{FN}}} (1 - r_l^{\text{FN}})^{n_l^{\text{P}} - n_l^{\text{FP}}}.$$

427 Then a slice sampling move is applied to each locus-specific parameter.

#### 428 9.4.5 Posterior summarization

429 Here we approximate the Bayes estimator by minimising the Bayes risk:

$$\operatorname{argmin}_{t \in \mathcal{T}} \sum_{t' \in \mathcal{T}} \int L(t, t') \pi(t, d\theta), \quad (6)$$

430 using the L1 metric on the matrices of induced indicators  $x(t)$  as the loss function:

$$L(t, t') = \sum_{l \in L} \sum_{c \in C} |x_{c,l}(t) - x_{c,l}(t')|.$$

431 It is useful to define the marginal indicators  $m_{c,l}$  that can be conceptualised as the posterior  
432 probability of cell  $c$  to have trait  $l$ :

$$m_{c,l} = \sum_{t \in \mathcal{T}} \int \mathbf{1}[x_{c,l}(t) = 1] \pi(t, d\theta).$$

433 Using the MCMC samples  $t^1, t^2, \dots, t^N$ , we obtain a Monte Carlo approximation:

$$\bar{m}_{c,l} = \frac{1}{N} \sum_{i=1}^N x_{c,l}(t^i) \rightarrow m_{c,l},$$

434 with probability one.

**Fig. 1-g** shows an example of the matrix  $m$  each element of which is one of the approximated  $\bar{m}_{c,l}$ . We can now write the objective function of Equation (6) via the above marginal indicators:

$$\begin{aligned} \sum_{t' \in \mathcal{T}} \int L(t, t') \pi(t, d\theta) &= \sum_{t' \in \mathcal{T}} \int \sum_{l \in L} \sum_{c \in C} |x_{c,l}(t) - x_{c,l}(t')| \pi(t, d\theta) \\ &= \sum_{l \in L} \sum_{c \in C} \sum_{t' \in \mathcal{T}} \int |x_{c,l}(t) - x_{c,l}(t')| \pi(t, d\theta) \\ &= \sum_{l \in L} \sum_{c \in C} \{m_{c,l}(1 - x_{c,l}(t)) + (1 - m_{c,l})x_{c,l}(t)\} \\ &= \sum_{l \in L} \sum_{c \in C} \{x_{c,l}(t) - 2m_{c,l}x_{c,l}(t)\} + \text{constant}. \end{aligned} \quad (7)$$

435 We use a greedy algorithm to approximately minimize Equation (7). We start with a star  
436 tree with leaves  $C$  rooted at  $v^*$  and add loci from  $L$  one by one from a locus queue sorted  
437 by priority score. The priority score of each locus  $l$  is computed as

$$\text{priority}(l) = \max_{t' \in N^l(t)} \frac{q(t')}{\sum_{t'' \in N^l(t)} q(t'')},$$

438 where

$$\begin{aligned} q(x) &= \prod_{c \in C} \prod_{l \in L(x)} q_{c,l}(x_{c,l}) \\ q_{c,l}(x_{c,l}) &= 2m_{c,l}x_{c,l} - x_{c,l}. \end{aligned}$$

439 The quantities in the priority queue can be computed as in Section 9.4.3. We take the  
440 result of the minimization of the Bayes risk as the consensus tree.



#### 441 9.4.6 Consensus tree and CNA heatmap visualisation

442 To visualize the consensus tree, we collapse the chains (sequence of loci having only one  
443 child) as well as remove the subtrees containing no cells. We align the leaves of the tree  
444 which correspond to cells after collapsing to the rows of a cell-locus matrix.

### 445 9.5 Synthetic experiments

#### 446 9.5.1 Benchmarking

447 To assess the performance of sitka against alternative approaches, we ran inference on 72  
448 simulated datasets of varying characteristics. We will refer to this set of datasets as  $S72$ ;  
449 its simulation procedure is described in Section 9.5.3. For each dataset in  $S72$ , we scored  
450 each method by computing the Robinson-Foulds (RF) [31] distance between the simulated  
451 tree and the inferred tree. The scores were normalized within each dataset by dividing  
452 each method's score by the worst performing method's score.

453 We compared sitka against the following baseline methods: UPGMA, WPGMA, NJ, HDBSCAN,  
454 and balanced and ordinary least-squares minimum-evolution methods (BME, OME respec-  
455 tively) of [32]. We also report the score of a uniformly random bifurcating tree, *Uniform*,  
456 to help interpret the absolute scores. Each method was given raw data from  $S72$ , as well  
457 as input identical to that of sitka, i.e., filtered binary marker data. Sitka's inference settings  
458 are summarized in **Supplementary Table 2**.

459 Baseline methods performed significantly worse with sitka's input and are thus omitted from  
460 the following summary. Sitka's normalized RF score ( $0.62 \pm 0.06$ ) dominated all baseline  
461 methods, the next best performer was BME ( $0.90 \pm 0.08$ ). Sitka ranked first in all 72 but one  
462 set of data, where it ranked 6 for one dataset of size  $500 \times 800$ . Summing each method's  
463 rank over all datasets, sitka scored a total rank of 77, while BME scored 193.5 (lower is  
464 better). These results are summarized in **Supplementary Fig. 12**.

#### 465 9.5.2 Exploratory experiments within sitka

466 To explore the effectiveness of global versus *local* (locus-specific) parameterization (Sec-  
467 tion 9.4.1), and the posterior summarization method (Section 9.4.5), we ran inference on 10  
468 datasets. We will refer to this set of datasets as  $S10$ ; its simulation procedure is described  
469 in Section 9.5.3. Inference settings are summarized in **Supplementary Table 2**.

470 RF distances from the *best-possible tree* were computed as a metric. The best-possible  
471 tree is defined as the perfect phylogenetic tree constructed from the noiseless synthetic,  
472 unviolated cell-locus matrix data. For a baseline to compare the greedy estimator (GE) of  
473 Section 9.4.5 with, consider the *trace search estimator* (TSE). The TSE is defined as a  
474 tree in the sampler trace that minimizes the sample L1 distance (Section 9.4.5).

475 The GE outperformed the TSE under both models. This suggests the proposed GE can,  
476 informally, harness more information from the posterior and more accurately summarize a  
477 posterior to arrive at a consensus tree than, say, a search over the posterior under some  
478 criterion. Under the TSE, the global model ( $0.44 \pm 0.09$ ) outperformed the local model  
479 ( $0.71 \pm 0.06$ ). This observation suggests that the local parameterization has a strong in-  
480 fluence on the trace (in tree space) of our sampler, as the TSE is essentially a search

481 over the posterior sample. Under the GE, the global model ( $0.31 \pm 0.07$ ) and local model  
482 ( $0.30 \pm 0.07$ ) performed evenly well. This observation suggests that the choice of param-  
483 eterization does not heavily influence the information contained in the marginal posterior  
484 over trees. Ultimately this experiment suggests that the GE summarizes the marginal pos-  
485 terior sufficiently well such that the global model, the simpler model of the two, suffices  
486 for reconstructing phylogenies and should be the preferred model. A summarizing plot is  
487 shown in **Supplementary Fig. 13**.

488 In our final synthetic experiment, we aimed to study the effects of perfect phylogeny as-  
489 sumption violations on the reconstruction of trees, and attempted to draw connections to  
490 real world data. The two violations considered are infinite sites and loss violations, de-  
491 scribed in Section 9.5.3. Inference was performed on 130 datasets ( $S_{130}$ ). Inference set-  
492 tings are summarized in **Supplementary Table 2**, and the simulation procedure for  $S_{130}$   
493 is described in Section 9.5.3.

494 The experiment results are summarized in **Fig. 3-a**. Holding one violation rate fixed at  
495 zero and varying the other, we observed linear effects for both types of violations. The  
496 results suggest sitka is more robust to infinite sites violations, with estimated effects to be  
497  $0.31 \pm 0.07$ , which is much less than loss violations ( $0.47 \pm 0.07$ ). When varied together,  
498 the linear effects were estimated to be  $0.25 \pm 0.04$ ,  $0.38 \pm 0.04$  respectively. In an attempt  
499 to draw connections to real datasets, we estimated both violation rates of real data to be  
500 less than 0.25 (the estimation procedure is described below; **Fig. 3-b**). These observations  
501 suggest sitka should perform reasonably well for the real datasets considered in this study,  
502 with RF distances in the vicinity of (0.2, 0.3).

503 The violation rate estimation procedure was performed post-inference, and can be de-  
504 scribed as follows. Given the inferred tree and its corresponding marker matrix  $x$  (as in  
505 Section 9.4.1), and the sitka-transformed marker matrix  $y$  (as in Section 9.2), define the  
506 difference matrix  $z := x - y$ , i.e.,  $z$  has entries  $z_{i,j} = x_{i,j} - y_{i,j}$ . Next, define  $z_{\text{Loss}}$  with entries  
507  $z_{i,j}^{\text{Loss}} := \mathbf{1}(z_{i,j} > 0)$ , and similarly  $z_{\text{IS}}$  with entries  $z_{i,j}^{\text{IS}} := \mathbf{1}(z_{i,j} < 0)$ . Given an integer-valued  
508 threshold  $\epsilon_v > 0$ , we say a column or trait  $l$  in  $z_v$  (for  $v \in \{\text{Loss, IS}\}$ ) has a violation if there  
509 exists an *island* of size at least as large as  $\epsilon_v$ . An island in column  $l$  is defined to be  
510 any sequence of row indices  $i, i + 1, \dots, i + s$  such that  $z_{i,l}^v = z_{i+1,l}^v = \dots = z_{i+s,l}^v = 1$   
511 and  $z_{i-1,l}^v, z_{i+s+1,l}^v$  are, not necessarily the same, 0 or undefined. Finally, the proportion of  
512 columns with a given type of violation, loss or infinite sites, is taken to be the violation rate  
513 estimate. The intuition behind this estimation procedure is to identify the proportion of loci  
514 where the inferred tree (or its marker matrix) is in contradiction with observations.

### 515 9.5.3 Data simulation

516 Datasets in  $S_{72}$  were generated in two steps: (i) simulate a cell tree and its corresponding  
517 CNA data, and (ii) inject noise into the CNA data from step one.

518 In the first step we used the simulator of [33] to generate trees along with CNAs, where  
519 leaf nodes represent observed cells and internal nodes represent latent ancestral cells,  
520 i.e., unobserved cells. An edge in the tree represents an ancestral relationship between  
521 the respective cells.

522 The simulator of [33] itself consists of two parts, which we briefly describe as follows.  
523 First, the simulator samples a tree based on a generalization of the Blum-François Beta-

524 splitting model [34, 35], which is inspired by the Beta-splitting model of [36]. The Beta-  
525 splitting model is particularly well-suited for generating a wide range of topologies, varying  
526 from balanced to imbalanced tree structures. Second, given a tree, CNAs are simulated  
527 on the edges of the tree where the number and size of CNAs are drawn from Poisson  
528 and exponential distributions respectively. The simulator also accounts for clonal whole  
529 chromosome amplification events, motivated by punctuated evolution models [37].

530 The second step of our synthetic data simulation process, independent of [33], injects  
531 noise into a cell by locus input CNA matrix  $y$ , and outputs a noisy matrix of the same size.  
532 Three types of noise were employed, namely, uniform noise, jitter noise, and a doubling  
533 noise.

534 The uniform noise is parameterized by false positive (FPR) and false negative (FNR)  
535 rate parameters. For each element of the input matrix  $y_{ij}$ , add an integer  $N_{ij} \sim$   
536  $\text{Binomial}(y_{ij}, \text{FNR})$  or subtract an integer  $M_{ij} \sim \text{Binomial}(1, \text{FPR})$ .

537 The doubling noise is parameterized by a probability  $p_d$ : for each row of the CNA matrix  
538  $y$ , draw a factor  $K$  where  $K - 1 \sim \text{Binomial}(1, p_d)$ , which is then multiplied to the row of  
539 the CNA matrix as noise. This procedure effectively, on average, doubles the copy number  
540 values for  $p_d$  proportion of cells in the sample.

541 The jitter noise is parameterized by a probability  $p_j$ . First, map the CNA matrix to its marker  
542 matrix. Then for each marker, the locus corresponding to the marker is randomly duplicated  
543 to the previous bin(s), or the next bin(s). The number of bins  $J$  to be overwritten — zero,  
544 one, or two — is drawn from a  $\text{Binomial}(2, p_j)$  distribution.

545 Datasets in  $S72$  were of sizes  $\{500, 1000, 1500, 2000, 2500, 3000\}$  cells by (approximately)  
546  $\{400, 600, 800\}$  markers. For each combination of sizes, we generated four datasets based  
547 on different random seeds to make a total of  $6 \times 3 \times 4 = 72$  datasets. The approximate  
548 number of markers is the target number of markers after correcting for jitter and filtering.  
549 **Supplementary Fig. 11** shows the CNA profiles of a subset of simulated data.

550 To describe the simulation parameters used for  $S72$ , we follow the terminologies and nota-  
551 tion used in [33]. For generating trees, the  $\alpha$  and  $\beta$  values parameterize the generalized  
552 Beta-splitting model. We drew  $\alpha, \beta$  from a uniform distribution on the interval  $(-1, 10)$ . For  
553 generating CNA data, the mean number of CNA to be added to a branch in the tree was  
554 chosen to generate data with approximately the number of desired markers post filtering  
555 and jitter-fixing. The multiplier of the mean CNA on the root was set to 8, the whole am-  
556 plification rate (rate of an allele chosen to be amplified) was set to 0.5. The remaining  
557 parameters used default settings. See [33] for a more thorough description of parame-  
558 ters.

559 For injecting noise, we drew the uniform noise parameters FPR and FNR from uniform  
560 distributions on the intervals  $(0.001, 0.01)$ ,  $(0.01, 0.03)$  respectively. The doubling noise pa-  
561 rameter  $p_d$  was drawn from a  $\text{Uniform}(0.03, 0.07)$  distribution. The jitter noise parameter  $p_j$   
562 was drawn from a  $\text{Uniform}(0.3, 0.7)$  distribution.

563 Datasets in  $S10$  and  $S130$  were also generated in two steps: (i) simulate a cell tree and its  
564 corresponding binary marker data satisfying perfect phylogeny assumptions, and (ii) inject  
565 noise and/or violations into the the binary marker data from step one.

566 In the first step, a tree is generated via Kingman's coalescent [38].<sup>1</sup> Briefly, we sample  
567 a coalescent tree for the set of cells  $C$  by uniformly selecting pairs of cells  $c_i, c_j \in C$  to  
568 coalesce backwards in time. The waiting time, or the branch length, between each event  
569 is exponentially distributed. Conditionally on the coalescent tree and given a set of loci  $L$ ,  
570 we simulate a  $|C| \times |L|$  marker matrix  $y$ . Every entry  $y_{i,j}$  is initialized to 0. Then for each  
571 column  $l$ , we select a subset of cells  $C'$  from  $C$  to set  $y_{i,l}$  to 1, for all  $i \in C'$ . The subset  
572 of cells is sampled by choosing a branch on the tree with probability proportional to the  
573 branch length, and selecting all cells descendant from the selected branch. In essence,  
574 we are simulating the number of events via a Poisson process, and directly mapping these  
575 events to the cell-locus marker matrix. The above concludes the data generation procedure  
576 satisfying perfect phylogeny assumptions.

577 In the second step of  $S10$ 's simulator, we injected artificial noise by introducing standard  
578 false positive and negative values into  $y$ . This concludes  $S10$ 's simulator. The simulator  
579 for  $S130$  has an additional sampling step for controlling the degree of perfect phylogeny  
580 violations. We considered two types of violations: (i) the loss of markers along a tree's  
581 branches, and (ii) the violation of the infinite sites (IS) assumption, that is, the occurrence  
582 of multiple distinct events in the same locus.

583 The procedure for simulating loss of marker events can be described as follows. First,  
584 randomly select a locus  $l$ , then identify the most recent common ancestor  $a$  for the set of  
585 cells  $\{i : y_{i,l} = 1\}$ . Given  $a$ , sample a cell  $d$  descendant of  $a$  (including  $a$ ). Finally, the loss  
586 event is simulated by reverting  $y_{i,l}$  to 0, for all  $i$  descendant of, and including,  $d$ .

587 IS model violations were simulated as follows. Uniformly sample a pair of loci  $(j, k)$ , and  
588 merge  $y_{.,j}, y_{.,k}$  into one column, yielding a cell-locus matrix of size one less than the original  
589 size. However, to maintain control over  $|L|$ , datasets in  $S130$  were simulated with  $|L| + N_{\text{IS}}$   
590 loci such that after simulating IS violations, we recover a matrix of size  $|C| \times |L|$ , where  $N_{\text{IS}}$   
591 is the number of IS violations.

592 The total number of loss and infinite sites violation events ( $N_{\text{Loss}}, N_{\text{IS}}$ ) were drawn from  
593 binomial distributions with probability  $p_{\text{Loss}}, p_{\text{IS}}$  respectively (and size  $|L|$ ). As a final step,  
594 false positives and negatives were artificially injected.

595 For both  $S10$  and  $S130$ , datasets of size  $|C| \times |L| = 500 \times 100$  with FNR and FPR both  
596 set to 0.002 were generated. For  $S130$ , the unordered pair  $(p_{\text{Loss}}, p_{\text{IS}})$  were set to values in  
597  $\{(0, 0), (0.1, 0.1), \dots, (0.4, 0.4)\} \cup \{(0, 0.1), (0, 0.2), (0, 0.3), (0, 0.4)\}$ . For each configuration  
598 of simulation parameters, 10 different seeds were used to generate a total of 10 and 130  
599 datasets for  $S10$  and  $S130$  respectively.

## 600 9.6 Goodness-of-fit

601 To evaluate the goodness-of-fit of inferred trees on real data, we suggest a test comparing  
602 the posterior distribution over entries of the matrix  $x$  with the data  $y$ .

603 Consider an inferred tree,  $\tau$  and the corresponding genotype matrix  $g = g(\tau)$ . We set  
604  $g(\tau) = x(\tau)$  for trees inferred from *sitka*. For trees inferred from the baseline methods, we  
605 define  $g(\tau)$  as  $x(\tau)$  except that  $g : \tau \rightarrow \{0, 1\}^{C \times U}$  where  $U$  the set of internal nodes of  $\tau$   
606 (Methods section 9.4.1). In general the inferred trees from the baseline methods do not

---

<sup>1</sup>We used the R packages [39, 40] for simulation.

607 have named internal nodes, nor do they have the same number of internal nodes as the  
608 number of loci  $L$ . Therefore we do not know which locus in the inferred tree  $\tau$  corresponds  
609 to which locus in the matrix  $y$ . We note that this is not the case with trees inferred from sitka  
610 where the internal nodes of the tree correspond to the columns of the induced genotype  
611 matrix  $g$ . As a result, for methods other than sitka, for each column in the input data  
612 matrix, we pick a clade in  $\tau$  that has the highest prediction accuracy for the entries in that  
613 column.

614 For each method, we report Youden's  $J$  index [41] which is equal to the sum of the sensitiv-  
615 ity and specificity minus 1. We now define a binary classification counts matrix function  $h$ ,  
616 i.e., a function which, for two vectors  $w$  and  $z$  of length  $C$ , forms the confusion matrix:

$$h_{i,j}(w, z) = \sum_{c \in C} \mathbf{1}(w_c = i) \mathbf{1}(z_c = j).$$

617 For example  $h_{0,0}(w, z)$  would count the number of times both elements of  $w$  and  $z$  were  
618 equal to zero (or *true negative*). We define accuracy for a given confusion matrix  $o$  com-  
619 puted from the  $h$  map above as:

$$\text{acc}(o) := \frac{o_{0,0} + o_{1,1}}{\sum_{i,j} o_{i,j}}.$$

We further define sensitivity and specificity as

$$\text{sensitivity}(o) := \frac{o_{1,1}}{o_{1,1} + o_{1,0}},$$

$$\text{specificity}(o) := \frac{o_{0,0}}{o_{0,0} + o_{0,1}},$$

$$\text{youden}(o) := \text{sensitivity}(o) + \text{specificity}(o) - 1.$$

620 For a given tree  $\tau$  and its corresponding matrix  $g$  we compute the Youden's score as fol-  
621 lows:

622 1. for all locus  $l$  in  $y$ ,  $o_l = \text{argmax}_{o'_l, l' \in \text{columns}(g)} \text{acc}(o_{l'})$ ,

623 2.  $o_\tau = \sum_{l' \in \text{columns}(g)} o_{l'}$

624 3.  $\text{youden}_\tau := \text{youden}(o_\tau)$ .

625 That is for each locus in  $y$ , we take the clade that among all possible clades in  $\tau$  maximizes  
626 the accuracy in predicting which cells are present in the  $l$ -th column of  $y$ . We then sum over  
627 all these scores to compute a confusion matrix for  $\tau$  and use this agglomerative matrix to  
628 compute the Youden's score for the tree. We use the delta method to calculate confidence  
629 intervals. **Fig. 2-d** shows the Youden's score and its 95% confidence interval for sitka and  
630 6 baseline methods on 3 different real-world datasets. Sitka has a higher score than all  
631 competing methods.

## 632 9.7 Application: assignment of single nucleotide variants

633 Here we posit an observation probability model for adding single nucleotide variant (SNV)  
634 data to an existing phylogenetic tree.

635 For locus  $l$  in cell  $c$ , let  $y_{c,l}^{SNV} = (d_{c,l}, \nu_{c,l}, c_{c,l})$  denote the observed SNV data where the  
 636 total number of reads, the number of reads with a variant allele, and the corresponding  
 637 copy number are indicated by  $d_{c,l}$ ,  $\nu_{c,l}$ , and  $c_{c,l}$  respectively.

638 We use  $x_{c,l}^{SNV}$  to denote an indicator variable taking the value one if and only if an an-  
 639 cestor of cell  $c$  harboured a single nucleotide alteration event at locus  $l$ . This variable is  
 640 unobserved and the focus of inference in this section. As in the sitka model, we assume a  
 641 perfect phylogeny structure on these indicator variables, and add an error model to relate  
 642  $x_{c,l}^{SNV}$  to the observed data while allowing violations of the perfect phylogeny assumption  
 643 and measurement noise. In the context of single nucleotide data, this is similar to [12]. The  
 644 parameters of the error model are denoted  $\theta^{SNV} = (\epsilon_{FP}, \epsilon_{FN})$ , where  $\epsilon_{FP}$  and  $\epsilon_{FN}$  are  
 645 false positive rate and false negative rates, respectively. Define:

$$q_{c,l}^b = p(y_{c,l}^{SNV} | x_{c,l}^{SNV}, \theta^{SNV}) = p(\nu_{c,l} | d_{c,l}, c_{c,l}, x_{c,l}^{SNV} = b, \theta^{SNV}), \quad (8)$$

646 where  $d_{c,l}$  and  $c_{c,l}$  are given inputs. The likelihood probability of cell node  $c$  is denoted by  
 647  $q_{c,l}^b$ , where  $b \in \{0, 1\}$ . For  $b = 1$ ,  $q_{c,l}^b$  reflects the likelihood of cell  $c$  being mutated at locus  $l$ ;  
 648 and for  $b = 0$ ,  $q_{c,l}^b$  reflects the likelihood of cell  $c$  not being mutated at locus  $l$ . For  $d_{c,l} = 0$ ,  
 649 we set  $q_{c,l}^b = 0.5$ .

650 The probability  $q_{c,l}^b$  is obtained by marginalizing a mixture of binomial distributions depend-  
 651 ing on all possible genotype states of locus  $l$  at cell  $c$ . Given the copy number  $c_{c,l}$ , the  
 652 possible genotype states are  $\mathcal{G} = \{A \dots A, AA \dots B, A \dots BB, \dots, B \dots B\}$ , where each  
 653 element has a length equal to  $c_{c,l}$ . For example, the genotype  $AAB$  refers to a genotype  
 654 with one variant allele  $B$  and two reference alleles  $A$ . For each genotype state  $g_i$ , where  $i$   
 655 indexes the elements of  $\mathcal{G}$ , the mean parameter of the corresponding binomial distribution  
 656 is denoted by  $\xi_{c,l}^i$ :

$$\xi_{c,l}^i = \begin{cases} \frac{\mathcal{B}(g_i)}{c_{c,l}}, & 1 \leq \mathcal{B}(g_i) < c_{c,l}, \\ 1 - \epsilon_{FP}, & \mathcal{B}(g_i) = c_{c,l}, \\ \epsilon_{FP}, & \text{otherwise,} \end{cases} \quad (9)$$

where  $\mathcal{B}(g_i)$  represents the number of variant alleles of genotype  $g_i$ . Therefore, for  $b =$   
 1,

$$q_{c,l}^1 = p(\nu_{c,l} | d_{c,l}, c_{c,l}, x_{c,l}^{SNV} = 1, \theta^{SNV}) \quad (10)$$

$$= \sum_{i=1}^{c_{c,l}} p(g_i) [\xi_{c,l}^{\nu_{c,l}} (1 - \xi_{c,l})^{d_{c,l} - \nu_{c,l}}] \quad (11)$$

$$+ \epsilon_{FN} [\epsilon_{FP}^{\nu_{c,l}} (1 - \epsilon_{FP})^{d_{c,l} - \nu_{c,l}}].$$

657 The value of  $p(g_i)$  equals  $\frac{1 - \epsilon_{FN}}{c_{c,l}}$ , and  $\epsilon_{FN}$  represents the error due to mutation loss or tree  
 658 errors.

659 If the mutation status of cell  $c$  at locus  $l$  is a wildtype (i.e., mutation is not present), then  
 660 the possible genotype states should not have any variant allele. The only possible geno-  
 661 type state is  $\{A \dots A\}$ . The mean parameter of the binomial distribution equals  $\epsilon_{FP}$  (false  
 662 positive rate). Therefore,

$$q_{c,l}^0 = p(\nu_{c,l} | d_{c,l}, c_{c,l}, x_{c,l}^{SNV} = 0, \epsilon_{FP}). \quad (12)$$

663 With the proposed probability model for SNVs, we can incorporate both SNV data and  
 664 CNA data to infer the underlying tree phylogeny in the sitka model. Therefore,

$$p(y|x, \theta) = \prod_{c \in C} \prod_{l \in L_{CNA}} p(y_{c,l}^{CNA} | x_{c,l}^{CNA}, \theta^{CNA}) \prod_{l \in L_{SNV}} p(y_{c,l}^{SNV} | x_{c,l}^{SNV}, \theta^{SNV}), \quad (13)$$

665 where  $C$  and  $L$  are the disjoint set of cells and loci, respectively. In this section, the loci set  
 666  $L$  includes both CNA and SNV traits.

667 Assume now that we seek to add one locus to an existing tree. We proceed similarly to  
 668 Section 9.4.3. Equation (4) can be rewritten in the following form:

$$\bar{\rho}_v = \frac{\left( \frac{\prod_{v_i \in \text{children}(v)} (\gamma_{v_i}^0 + \gamma_{v_i}^1)}{\gamma_v^0} \right)}{\sum_{\bar{v} \in R} \left( \frac{\prod_{\bar{v}_i \in \text{children}(\bar{v})} (\gamma_{\bar{v}_i}^0 + \gamma_{\bar{v}_i}^1)}{\gamma_{\bar{v}}^0} \right)}, \quad (14)$$

where  $\gamma_v^b$ , for  $b \in \{0, 1\}$  is:

$$\gamma_v^b = \begin{cases} p_v^b, & \text{if } l \text{ represents a CNA loci,} \\ q_v^b, & \text{if } l \text{ represents a SNV loci.} \end{cases}$$

669 For  $v \in R = \{v^*\} \cup L \setminus \{l\}$ , and  $b \in \{0, 1\}$ , the value of  $q_v^b$  is

$$q_v^b = \prod_{v'' \in \text{children}(v)} q_{v''}^b. \quad (15)$$

670 For the cell nodes that are the leaves of the tree  $q_v^b = q_{c,l}^b$ .

### 671 9.7.1 Detection of SNVs for individual cells

672 Given a fixed CNA tree (denoted by  $t$ ) and the read counts data ( $y^{SNV}$  denoted by  $y$  for  
 673 simplicity), here the goal is to calculate the posterior distribution of  $x_{c,l}^{SNV}$ , the mutation  
 674 status of locus  $l$  at cell  $c$ , which we denote by  $x_{c,l}$  for simplicity.

675 The joint probability distribution of  $x_{c,l}$ ,  $y$  and  $t$  can be written as:

$$p(x_{c,l}, y, t) = \sum_{v \in R} \sum_{t' \in \mathcal{N}_v^l(t \setminus l)} p(x_{c,l}, t', y) \quad (16)$$

$$= \sum_{v \in R} \sum_{t' \in \mathcal{N}_v^l(t \setminus l)} p(x_{c,l} | t') p(y | t') p(t'), \quad (17)$$

676 where  $R$  is the set of all loci nodes in the tree (including the root) excluding locus  $l$ . The  
 677 joint probability distribution is calculated as

$$p(x_{c,l} = 1, y, t) = \sum_{v \in \mathcal{P}(c,t)} \sum_{t' \in \mathcal{N}_v^l(t \setminus l)} p(y | t') p(t'). \quad (18)$$

678 The set  $\mathcal{P}(c, t)$  denotes all nodes on the shortest path from cell  $c$  to the root of the tree  
 679 (including the root and excluding the cell  $c$  node). An example of the path on an imaginary  
 680 tree is depicted in **Supplementary Fig. 14**. The nodes coloured in green belong to  $\mathcal{P}(c, t)$ .  
 681 Therefore, the posterior probability distribution of  $x_{c,l} = 1$  yields

$$p(x_{c,l} = 1|y, t) = \frac{p(x_{c,l} = 1, y, t)}{p(y, t)} = \frac{\sum_{v \in \mathcal{P}(c,t)} \sum_{t' \in \mathcal{N}_v^l(t \setminus l)} p(y|t') p(t')}{p(y, t)}. \quad (19)$$

Rewriting Equation (19) assuming uniform probability distribution for  $p(t')$  yields:

$$\begin{aligned} p(x_{c,l} = 1|y, t) &\propto \sum_{v \in \mathcal{P}(c,t)} \sum_{t' \in \mathcal{N}_v^l(t \setminus l)} p(y|t'), \\ &= \sum_{v \in \mathcal{P}(c,t)} \sum_{t' \in \mathcal{N}_v^l(t \setminus l)} \prod_{l' \in L} \prod_{c' \in C} p(y_{c',l'}|t'), \\ &= \sum_{v \in \mathcal{P}(c,t)} \sum_{t' \in \mathcal{N}_v^l(t \setminus l)} \prod_{\substack{l' \in L \\ l' \neq l}} \prod_{c' \in C} p(y_{c',l'}|t') \prod_{m' \in C} p(y_{c',l}|t'), \\ &= K_1 \sum_{v \in \mathcal{P}(c,t)} \sum_{t' \in \mathcal{N}_v^l(t \setminus l)} \prod_{c' \in C} p(y_{c',l}|t'), \\ &= K_1 \sum_{v \in \mathcal{P}(c,t)} \sum_{t' \in \mathcal{N}_v^l(t \setminus l)} \prod_{c' \in C \setminus v} p(y_{c',l}|t') \prod_{c' \in L_v} p(y_{c',l}|t'), \end{aligned}$$

where  $N$  denotes the set of all trait nodes,  $C$  denotes the set of all cell nodes,  $C_v$  denotes the cells that are a descendant of node  $v$ , and  $C \setminus v$  denotes the cells that are not descendant of node  $v$ . The product of the likelihood contributions for non-descendant nodes can be calculated by taking the product of  $q_c^0$  for all cells, divided by the ones that are descendant of  $v$ :

$$\prod_{c' \in C \setminus v} q_{c'}^0 = \frac{q_v^0}{q_v^*}.$$

Therefore:

$$p(x_{c,l} = 1|y, t) \propto K_1 \sum_{v \in \mathcal{P}(c,t)} \frac{q_v^0}{q_v^*} \sum_{t' \in \mathcal{N}_v^l(t \setminus l)} \prod_{c' \in C_v} p(y_{c',l}|t'). \quad (20)$$

682 The likelihood contribution of descendant cells can be re-indexed by a binary vector  $\mathbf{b} =$   
 683  $(b_1, b_2, \dots, b_k)$ , where  $b_i \in \{0, 1\}$ , and  $b_i = 1$  if the child  $v$  is to be moved into a child of the  
 684 node  $l$ . The value of  $k$  denotes the number of children of  $v$ . The  $i^*$ th child of  $v$  which is  
 685 on the path from node  $v$  to cell  $c$  is called  $v_i^*$ . This implies  $b_{i^*} = 1$  (See **Supplementary**  
 686 **Fig. 14**). Therefore:

$$\sum_{t' \in \mathcal{N}_v^l(t \setminus l)} \prod_{c' \in C_v} p(y_{c',l}|t') = q_{v_i^*}^1 \sum_{b_1=0}^1 \sum_{b_2=0}^1 \dots \sum_{b_{i-1}=0}^1 \sum_{b_{i+1}=0}^1 \dots \sum_{b_k=0}^1 \prod_{\substack{i=1 \\ i \neq i^*}}^k q_{v_i}^{b_i}. \quad (21)$$



Rewriting Equation (20) using Equation (21) yields:

$$\begin{aligned}
 p(x_{c,l} = 1|y, t) &\propto K_1 \sum_{v \in \mathcal{P}(c,t)} \frac{q_{v^*}^0}{q_v^0} q_{v^*}^1 \sum_{b_1=0}^1 \sum_{b_2=0}^1 \cdots \sum_{b_{i-1}=0}^1 \sum_{b_{i+1}=0}^1 \cdots \sum_{b_k=0}^1 \prod_{\substack{i=1 \\ i \neq i^*}}^k q_{v_i}^{b_i}, \\
 &= K_1 \sum_{v \in \mathcal{P}(c,t)} \frac{q_{v^*}^0}{q_v^0} q_{v^*}^1 \prod_{\substack{i=1 \\ i \neq i^*}}^k (q_{v_i}^0 + q_{v_i}^1), \\
 &= K_1 \sum_{v \in \mathcal{P}(c,t)} \frac{q_{v^*}^0}{q_v^0} \frac{\prod_{i=1}^k (q_{v_i}^0 + q_{v_i}^1)}{(q_{v_{i^*}}^0 + q_{v_{i^*}}^1)} q_{v_{i^*}}^1, \\
 &= K_1 q_{v^*}^0 \sum_{v \in \mathcal{P}(c,t)} \frac{q_{v_{i^*}}^1}{q_v^0 (q_{v_{i^*}}^0 + q_{v_{i^*}}^1)} \prod_{i=1}^k (q_{v_i}^0 + q_{v_i}^1). \tag{22}
 \end{aligned}$$

## 687 9.8 Computational complexity of the SNV calling algorithm

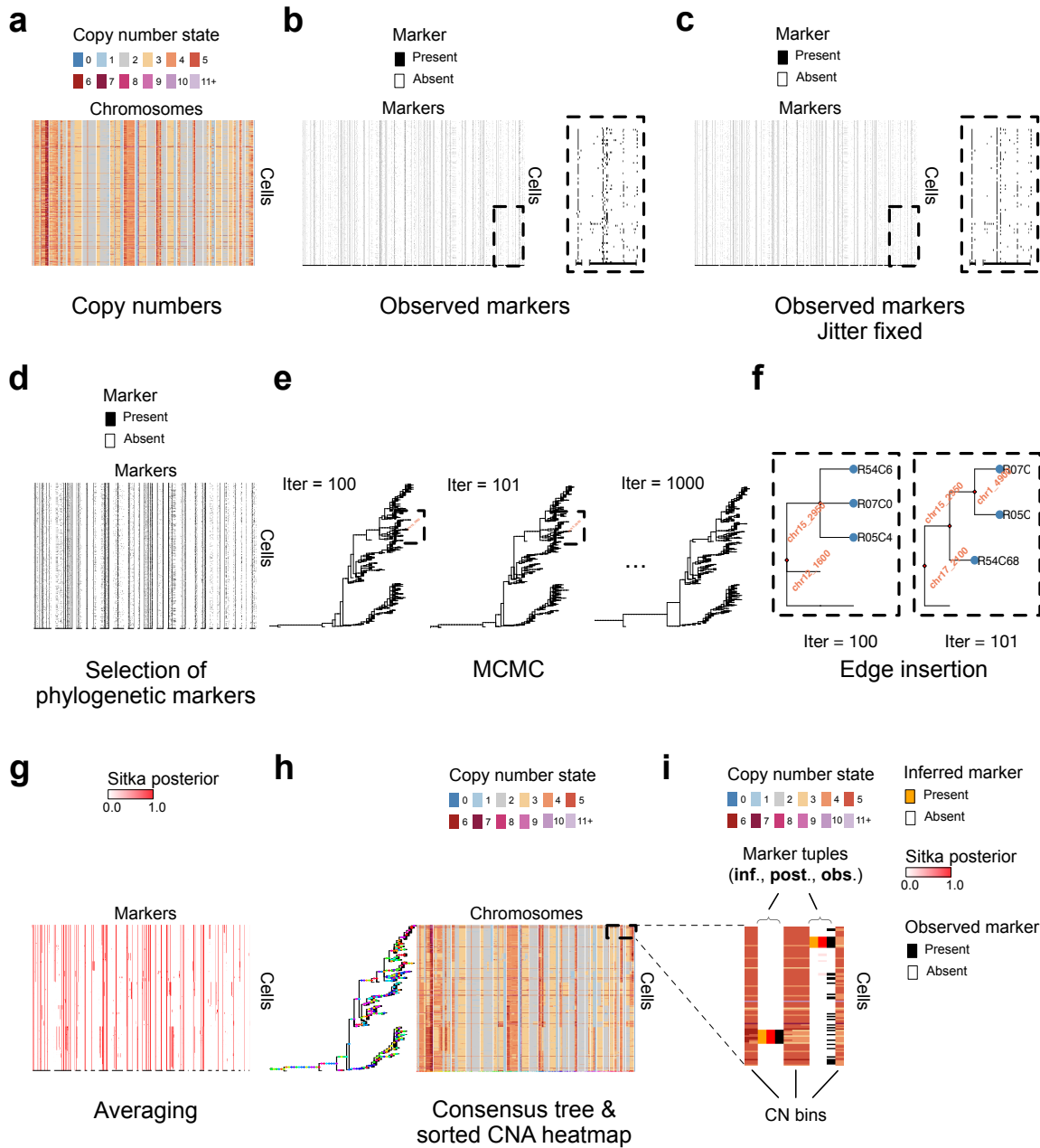
688 The computational complexity of Equation (22) is  $O(|C| \cdot |L|)$  with  $|C|$  the number of cells  
 689 and  $|L|$  the number of loci. In order to reduce the complexity of calculating  $p(x_{c,l} = 1|y, t)$   
 690 for each locus and cell,  $\mathcal{P}'(c, t)$  is defined to denote the nodes sitting on the path from root  
 691 to cell  $c$ , excluding the root node and including the cell  $c$  node. Then,

$$q_v^* = \prod_{i=1}^k (q_{v_i}^0 + q_{v_i}^1). \tag{23}$$

Therefore,

$$K_1 q_{v^*}^0 \sum_{v \in \mathcal{P}(c,t)} \frac{q_{v_{i^*}}^1}{q_v^0 (q_{v_{i^*}}^0 + q_{v_{i^*}}^1)} \prod_{i=1}^k (q_{v_i}^0 + q_{v_i}^1) = K_1 q_{v^*}^0 \sum_{v \in \mathcal{P}'(c,t)} \frac{q_v^1}{(q_v^0 + q_v^1)} \frac{q_{\text{parent}(v)}^*}{q_{\text{parent}(v)}^0}.$$

692 Calculating  $p(x_{c,l} = 1|y, t)$  with a recursive approach reduces the complexity from  $O(|C||L|)$   
 693 to  $O(|C| + |L|)$ , where as in the last section  $L$  is the union of SNV and CNA loci.



**Figure 1**

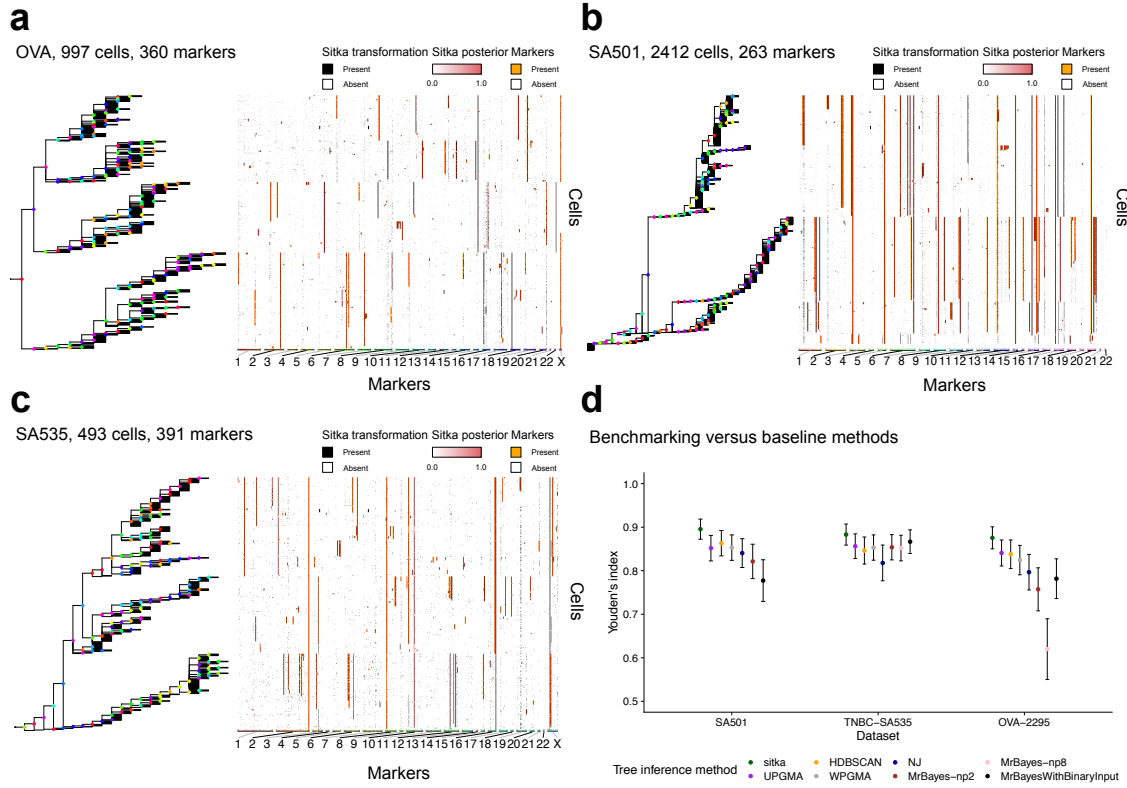


Figure 2.

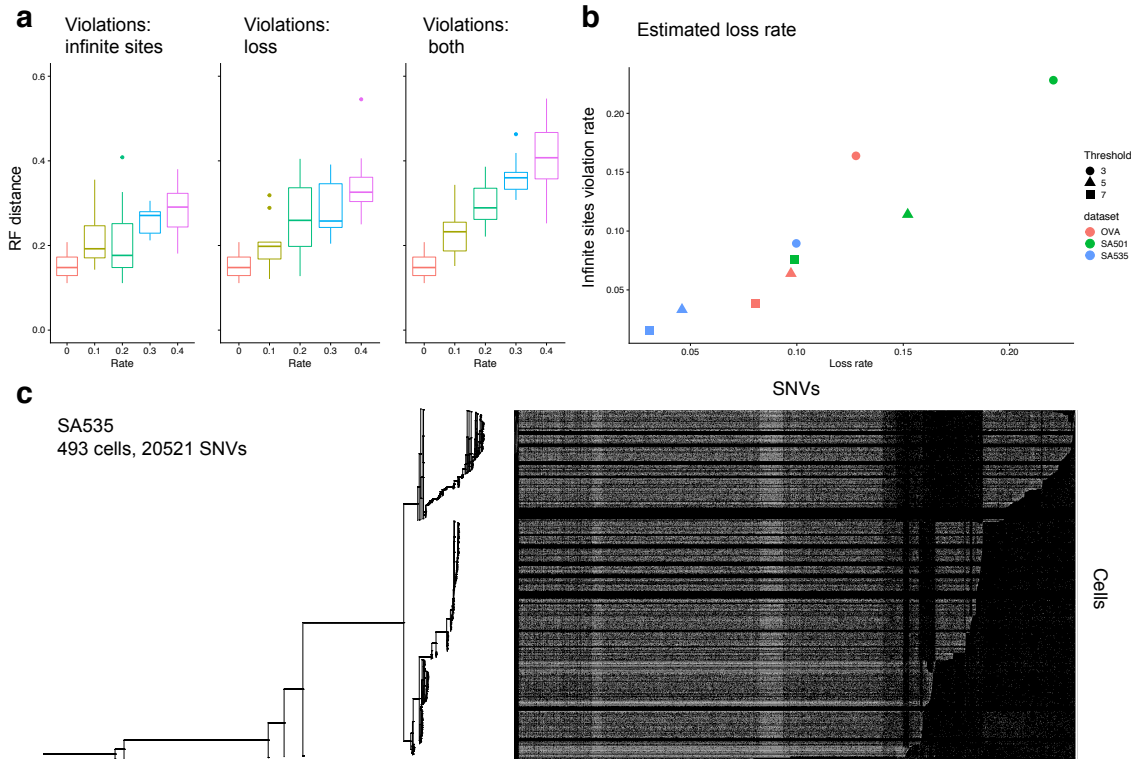


Figure 3

## 694 List of Figures

695 **Figure 1.** Workflow of *sitka*. **(a)** *Sitka* takes copy number calls data from a heterogeneous  
696 single-cell population. The cells (rows of the copy number matrix) are randomly sorted. **(b)**  
697 A lossy binary transformation is applied to obtain markers data. (Methods section 9.2 and  
698 **Supplementary Fig. 1**). Note that each single-cell is now represented by the presence  
699 or absence of CN changes between consecutive bins. **(c)** The boundary conditions are  
700 smoothed to account for cell-specific marker miss-alignment. (Methods section 9.3) to  
701 correct for this marker misalignment. Note how the columns in the inset in panel-**c** are  
702 less noisy than their counterpart in panel-**b**. **(d)** A subset of markers present in at least  
703 5 percent of the cells are chosen for input to the tree inference algorithm. **(e)** An MCMC  
704 algorithm efficiently explores the tree space. **(f)** An example of an edge-insertion. **(g)** The  
705 indicator matrix of all post-burn-in MCMC trees are averaged to generate a matrix indicating  
706 the posterior probability of a cell being attached to a marker (Methods section 9.4.5). **(h)**  
707 The copy number data in **(a)** is sorted according to the inferred consensus tree, shown on  
708 the left of the matrix. **(i)** The inset shows the tuple of marker columns in the context of  
709 the copy number calls, namely **inf.** (inferred markers, i.e., latent state  $x_{c,l}$ ), **post.** (posterior  
710 probability of the latent state  $x_{c,l}$ ), and **obs.** (observed markers), interlaced with the CN  
711 columns (similar to **Supplementary Fig. 1**). The results are from the *SA535* dataset, a  
712 triple negative breast cancer patient derived xenograft sample (Methods section 2.2).

713 **Figure 2.** Results over real-datasets and benchmarking against baseline methods. **(a)**,  
714 **(b)**, and **(c)** show the consensus tree and marker-space matrix for the *OVA*, *SA501*, and  
715 *SA535* datasets respectively. **(d)** Comparison to baseline methods.

716 **Figure 3.** Synthetic experiments and an application to point mutation placement. **(a)** RF  
717 distance of Bayes tree estimate to the best-possible tree. The first plot holds  $p_{is}$  constant  
718 at zero. The second plot holds  $p_{loss}$  constant at 0. The third plot varies  $p_{is} = p_{loss}$  jointly. **(b)**  
719 Estimation of violation rates in real data and a set of synthetic data. **(c)** Over 20,000 SNV's  
720 with high levels of missingness are placed on a backbone tree inferred from the CNA data  
721 for *SA535*.

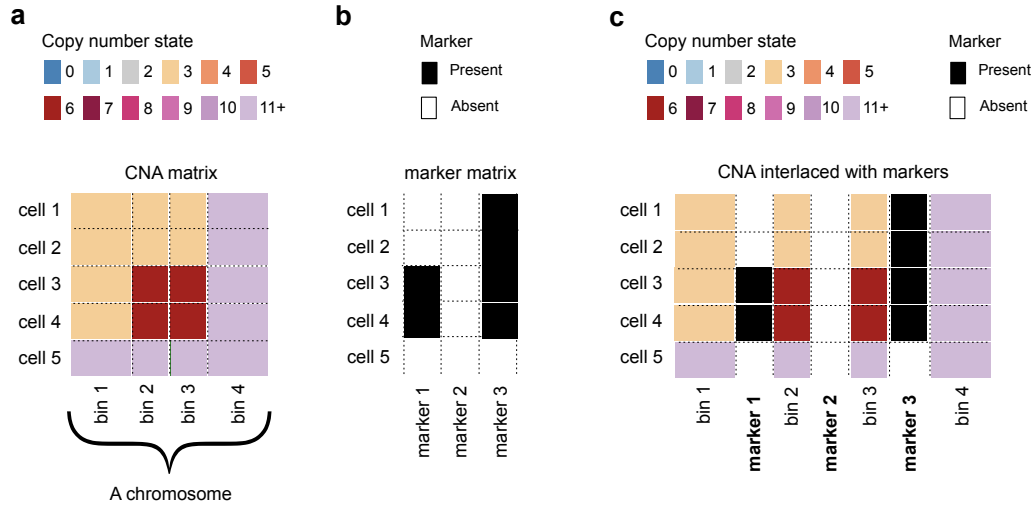
## 722 References

- 723 [1] E. Laks et al. “Clonal Decomposition and DNA Replication States Defined by Scaled  
724 Single-Cell Genome Sequencing”. In: *Cell* 179.5 (2019), 1207–1221.e22.
- 725 [2] M. Pellegrino et al. “High-throughput single-cell DNA sequencing of acute myeloid  
726 leukemia tumors with droplet microfluidics”. In: *Genome research* 28.9 (2018),  
727 pp. 1345–1352.
- 728 [3] T. Baslan et al. “Genome-wide copy number analysis of single cells”. In: *Nature pro-  
729 tocols* 7.6 (2012), pp. 1024–1041.
- 730 [4] C. Gawad, W. Koh, and S. R. Quake. “Single-cell genome sequencing: current state  
731 of the science”. In: *Nature Reviews Genetics* 17.3 (2016), p. 175.
- 732 [5] S. Salehi et al. “Clonal fitness inferred from time-series modelling of single-cell can-  
733 cer genomes”. In: *Nature* (2021). DOI: 10.1038/s41586-021-03648-3. URL:  
734 <https://doi.org/10.1038/s41586-021-03648-3>.
- 735 [6] J. J. Quinn et al. “Single-cell lineages reveal the rates, routes, and drivers of metas-  
736 tasis in cancer xenografts”. In: *Science* (2021). ISSN: 0036-8075. DOI: 10.1126/  
737 science.abc1944. eprint: [https://science.sciencemag.org/content/early/  
738 2021/01/21/science.abc1944.full.pdf](https://science.sciencemag.org/content/early/2021/01/21/science.abc1944.full.pdf). URL: [https://science.sciencemag.  
739 org/content/early/2021/01/21/science.abc1944](https://science.sciencemag.org/content/early/2021/01/21/science.abc1944).
- 740 [7] C. Abbosh et al. “Phylogenetic ctDNA analysis depicts early-stage lung cancer evo-  
741 lution”. In: *Nature* 545.7655 (2017), pp. 446–451.
- 742 [8] Y. Wang et al. “Clonal evolution in breast cancer revealed by single nucleus genome  
743 sequencing”. In: *Nature* 512.7513 (2014), pp. 155–160.
- 744 [9] M. L. Leung et al. “Single-cell DNA sequencing reveals a late-dissemination model  
745 in metastatic colorectal cancer”. In: *Genome research* 27.8 (2017), pp. 1287–1299.
- 746 [10] C. Yu et al. “Discovery of biclonal origin and a novel oncogene SLC12A5 in colon  
747 cancer by single-cell sequencing”. In: *Cell research* 24.6 (2014), pp. 701–712.
- 748 [11] R. Schwartz and A. A. Schäffer. “The evolution of tumour phylogenetics: principles  
749 and practice”. In: *Nature Reviews Genetics* 18.4 (Apr. 2017), pp. 213–229. ISSN:  
750 1471-0064.
- 751 [12] K. Jahn, J. Kuipers, and N. Beerenwinkel. “Tree inference for single-cell data”. In:  
752 *Genome Biology* 17 (2016), p. 86. ISSN: 1474-760X.
- 753 [13] E. M. Ross and F. Markowetz. “OncoNEM: inferring tumor evolution from single-cell  
754 sequencing data”. In: *Genome Biology* 17 (Apr. 2016), p. 69. ISSN: 1474-760X.
- 755 [14] C. A. Miller et al. “SciClone: inferring clonal architecture and tracking the spatial and  
756 temporal patterns of tumor evolution”. In: *PLoS Comput Biol* 10.8 (2014), e1003665.
- 757 [15] Jochen Singer et al. “Single-cell mutation identification via phylogenetic inference”.  
758 In: *Nature communications* 9.1 (2018), pp. 1–8.
- 759 [16] H. Zahn et al. “Scalable whole-genome single-cell library preparation without pream-  
760 plification”. In: *Nature Methods* 14.2 (2017), pp. 167–173. ISSN: 1548-7105.
- 761 [17] Gryte Satas et al. “Scarlet: Single-cell tumor phylogeny inference with copy-number  
762 constrained mutation losses”. In: *Cell Systems* 10.4 (2020), pp. 323–332.
- 763 [18] F. Wang et al. “Single-cell copy number lineage tracing enabling gene discovery”.  
764 In: *bioRxiv* (2020). DOI: 10.1101/2020.04.12.038281. eprint: [https://www.  
765 biorxiv.org/content/early/2020/04/13/2020.04.12.038281.full.pdf](https://www.biorxiv.org/content/early/2020/04/13/2020.04.12.038281.full.pdf).  
766 URL: [https://www.biorxiv.org/content/early/2020/04/13/2020.04.12.  
767 038281](https://www.biorxiv.org/content/early/2020/04/13/2020.04.12.038281).

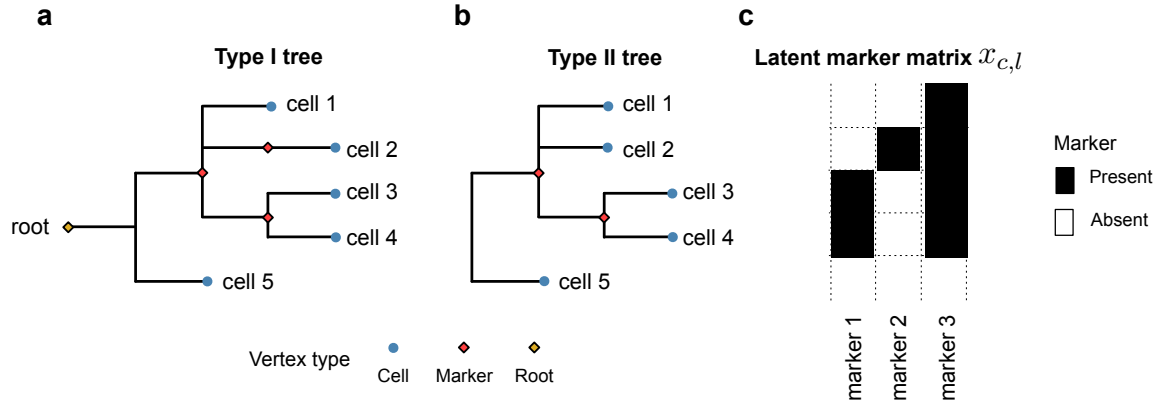
- 768 [19] X. Xu et al. “Single-Cell Exome Sequencing Reveals Single-Nucleotide Mutation  
769 Characteristics of a Kidney Tumor”. In: *Cell* 148.5 (Mar. 2012), pp. 886–895. ISSN:  
770 0092-8674.
- 771 [20] T. L. Williams and B. M. E. Moret. “An investigation of phylogenetic likelihood meth-  
772 ods”. In: *Third IEEE Symposium on Bioinformatics and Bioengineering, 2003. Pro-  
773 ceedings*. Mar. 2003, pp. 79–86.
- 774 [21] S. Malikic et al. “Integrative inference of subclonal tumour evolution from single-cell  
775 and bulk sequencing data”. In: *Nature communications* 10.1 (2019), pp. 1–12.
- 776 [22] J. Ma et al. “The infinite sites model of genome evolution”. In: *Proceedings of the  
777 National Academy of Sciences* 105.38 (2008), pp. 14254–14261.
- 778 [23] C. D. Greenman et al. “Estimation of rearrangement phylogeny for cancer genomes”.  
779 In: *Genome Research* 22.2 (Feb. 2012), pp. 346–361. ISSN: 1088-9051, 1549-5469.
- 780 [24] H. Zafar et al. “SiFit: A Method for Inferring Tumor Trees from Single-Cell Sequencing  
781 Data under Finite-site Models”. In: *bioRxiv* (Dec. 2016), p. 091595.
- 782 [25] P. Eirew et al. “Dynamics of genomic clones in breast cancer patient xenografts at  
783 single-cell resolution”. In: *Nature* 518.7539 (2015), pp. 422–426.
- 784 [26] K. Yi and Y. Seok Ju. “Patterns and mechanisms of structural variations in human  
785 cancer”. In: *Experimental & Molecular Medicine* 50.8 (2018), p. 98. DOI: 10.1038/  
786 s12276-018-0112-3. URL: <https://doi.org/10.1038/s12276-018-0112-3>.
- 787 [27] S. Mishra and J. R. Whetstone. “Different Facets of Copy Number Changes: Per-  
788 manent, Transient, and Adaptive”. In: *Molecular and Cellular Biology* 36.7 (2016),  
789 pp. 1050–1063. ISSN: 0270-7306. DOI: 10.1128/MCB.00652-15. eprint: <https://mcb.asm.org/content/36/7/1050.full.pdf>. URL: <https://mcb.asm.org/content/36/7/1050>.
- 792 [28] D.B. Wilson. “Generating Random Spanning Trees More Quickly Than the Cover  
793 Time”. In: *Proceedings of the Twenty-eighth Annual ACM Symposium on Theory of  
794 Computing*. STOC '96. New York, NY, USA: ACM, 1996, pp. 296–303. ISBN: 978-0-  
795 89791-785-8.
- 796 [29] A. Bouchard-Côté et al. “Blang: Bayesian declarative modelling of general data  
797 structures and inference via algorithms based on distribution continua”. In:  
798 *arXiv:1912.10396 [stat]* (2021). arXiv: 1912.10396 [stat.CO].
- 799 [30] R.M. Neal. “Slice sampling”. In: *The Annals of Statistics* 31.3 (June 2003), pp. 705–  
800 767. ISSN: 0090-5364, 2168-8966.
- 801 [31] D. F. Robinson and L. R. Foulds. “Comparison of phylogenetic trees”. In: *Mathemat-  
802 ical biosciences* 53.1-2 (1981), pp. 131–147.
- 803 [32] R. Desper and O. Gascuel. “Fast and accurate phylogeny reconstruction algorithms  
804 based on the minimum-evolution principle”. In: *International Workshop on Algorithms  
805 in Bioinformatics*. Springer. 2002, pp. 357–374.
- 806 [33] X. F. Mallory et al. “Methods for copy number aberration detection from single-cell  
807 DNA-sequencing data”. In: *Genome biology* 21.1 (2020), pp. 1–22.
- 808 [34] R. Sainudiin and A. Véber. “A Beta-splitting model for evolutionary trees”. In: *Royal  
809 Society open science* 3.5 (2016), p. 160016.
- 810 [35] M. G. B. Blum and O. François. “Which random processes describe the tree of life?  
811 A large-scale study of phylogenetic tree imbalance”. In: *Systematic Biology* 55.4  
812 (2006), pp. 685–691.
- 813 [36] D. Aldous. “Probability distributions on cladograms”. In: *Random discrete structures*.  
814 Springer, 1996, pp. 1–18.

- 815 [37] R. Gao et al. “Punctuated copy number evolution and clonal stasis in triple-negative  
816 breast cancer”. In: *Nature genetics* 48.10 (2016), p. 1119.
- 817 [38] J. F. C. Kingman. “The coalescent”. In: *Stochastic processes and their applications*  
818 13.3 (1982), pp. 235–248.
- 819 [39] K.P. Schliep. “phangorn: phylogenetic analysis in R”. In: *Bioinformatics* 27.4 (2011),  
820 pp. 592–593.
- 821 [40] P. R. Staab and D. Metzler. “Coala: an R framework for coalescent simulation”. In:  
822 *Bioinformatics* (2016). DOI: 10.1093/bioinformatics/btw098.
- 823 [41] W. J. Youden. “Index for rating diagnostic tests”. In: *Cancer* 3.1 (1950), pp. 32–35.

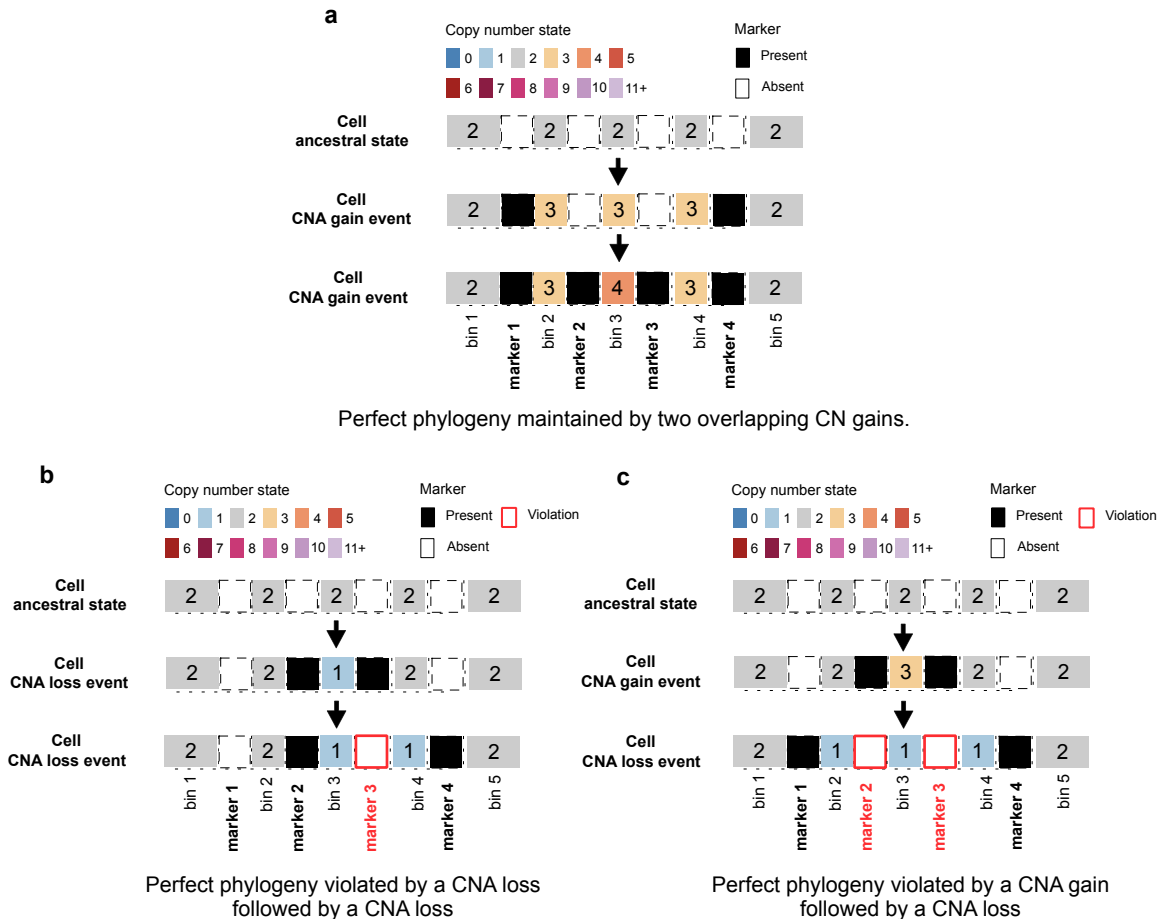




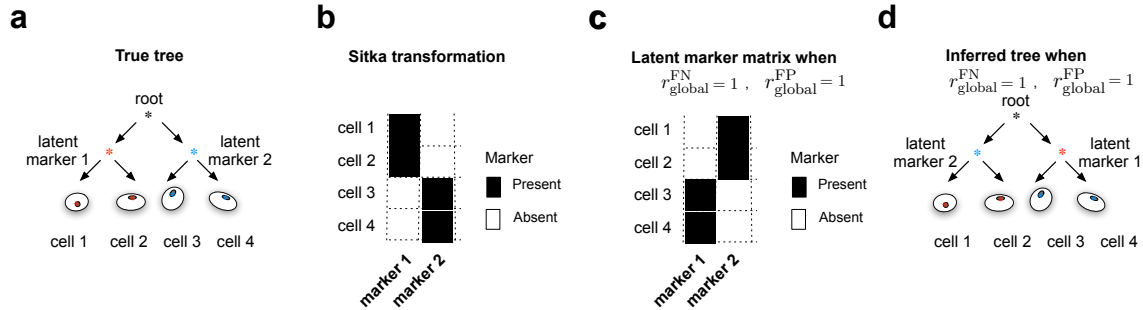
**Supplemental Figure 1.** Description of the process involved in the construction of *markers*, the input to the *sitka* model. A *bin* is a contiguous set of genomic positions. Each pair of consecutive bins (e.g., bins 1 and 2 in (a)) is associated with a *marker* (e.g., marker 1) that measures for each individual cell, whether there is a difference between the CNA states of the two bins. (a) The observed CNA matrix for a subset of bins on a chromosome. The rows are sequenced single cells, and the columns are bins. The CN states are colour-coded. (b) The three markers shown are associated with the four bins. Each marker records the presence (black) or absence (white) of a CN state change between a pair of consecutive bins. Note that in the CNA matrix, there is a CN change at row 3 from bin 1 to bin 2 (CN state 3 to 6). This is reflected in the marker matrix, at row 3 of marker 1 with a black square. There are no changes between bins 2 and 3 across any rows in the CNA matrix. This is reflected in marker 2 comprising all white squares. (c) For visualisation purposes, the CNA matrix can be interlaced with the marker matrix to more clearly show where the CNA changes occur. Each column of the marker matrix is inserted between the associated pair of columns in the CNA matrix. The resulting matrix is an example of an *augmented* view that combines data from two or more sources (here the CNA matrix and the marker matrix). In an augmented view, we call columns from each source a *channel*.



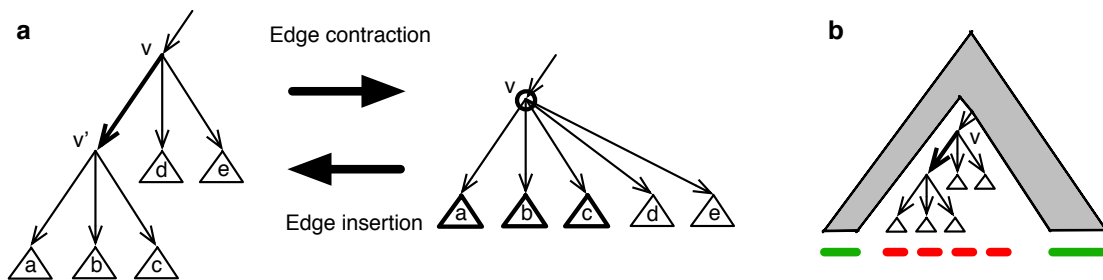
**Supplemental Figure 2.** Visualisation of a small type I tree  $t$  (a), its transformation into a type II tree (b), and the corresponding marker matrix  $x = (x_{c,l})$  (c). Given a tree  $t$ , the latent marker matrix  $x$  is a deterministic function  $x = x(t)$ . Note that the clade comprising single-cells 3 and 4 has support in both markers 1 and 3. For clarity, we do not visualise type I trees, but plot their transformation, i.e., type II trees as follows. We remove from the type I tree all marker nodes that have  $x_{c,l} = 0$  for all single-cells  $c$ . Lists of connected edges that have exactly one descendent (i.e., chains) are also collapsed into a single edge, e.g., the edge corresponding to markers 2 and 3 are collapsed into one edge (since marker 2 has only one descendent, namely single-cell 2).



**Supplemental Figure 3.** The effects of overlapping CNA events on the perfect phylogeny assumption. A segment of a chromosome with five consecutive bins and their four corresponding markers are shown. Each panel follows the CN states interlaced with markers for a cell at the ancestral state (top), after a CNA event (middle), and after a second overlapping CNA event (bottom). The numbers in the CNA squares show the integer CN state (e.g., the ancestral state has two copies of the 5-bins long segment). **(a)** Two overlapping CNA gains maintain the perfect phylogeny assumption. By the infinite site argument, it is unlikely for the end-points of the two gain events to exactly match. The same argument holds for a CNA loss followed by a CNA gain event. Note that in these cases, once a change point is acquired, it is not lost. **(b)** If a loss event is followed by another loss event in which either end-points of the first event is removed, the perfect phylogeny assumption will be violated (e.g., marker 3 is lost after the second loss event). Note that a violation does not occur if the loss events hit different copies of a segment. **(c)** Similarly, if a gain event is followed by a loss event, only if the latter erases the end-points of the former is the perfect phylogeny violated. Note how marker 2 and marker 3 are lost after the second CNA event.

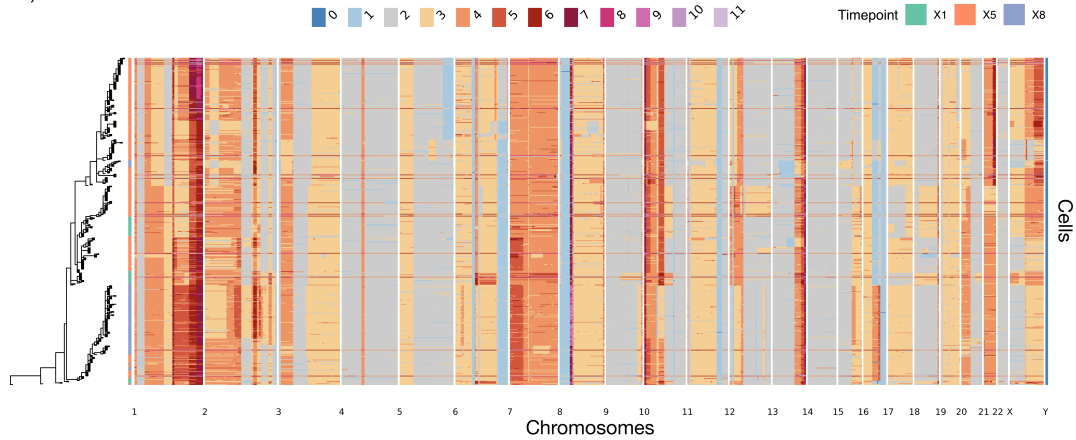


**Supplemental Figure 4.** Pathological tree reconstruction under default observation prior. (a) The true tree reconstruction in a simple example with a balanced phylogeny with two clades of size two, and two unique markers, coloured red and blue, that distinguish the left and right clades respectively. (b) The binarised input matrix corresponding to the four cells at the two markers. The desired observation error rates should be zero and the latent and observed marker matrices should match exactly, as the perfect phylogeny assumption holds. If the observation error parameters are set to one, that is  $r_{\text{global}}^{\text{FP}} = 1$  and  $r_{\text{global}}^{\text{FN}} = 1$ , then the latent marker matrix with all entries flipped as shown in (c) will have an equal likelihood under this setting as the desired latent matrix has when error rates are set to zero. (d) The incorrect tree reconstruction where the left and right clades are erroneously assigned to the blue and red markers.



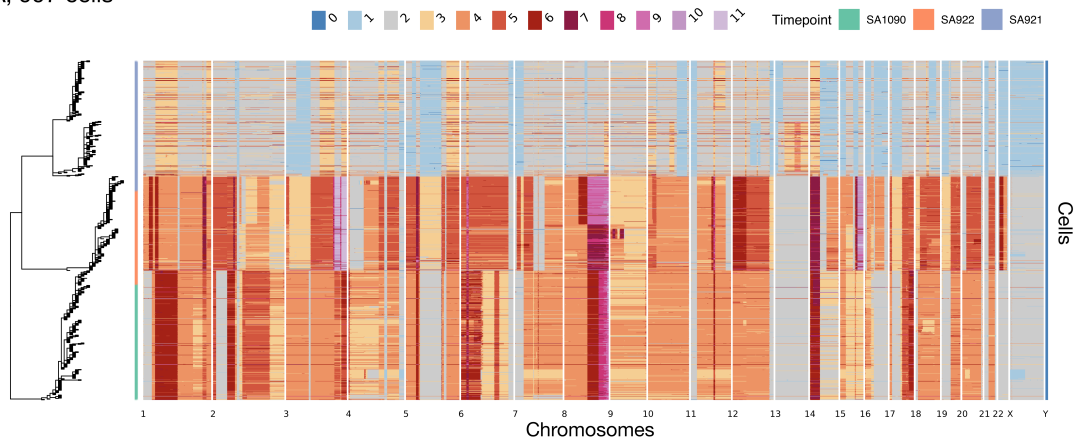
**Supplemental Figure 5.** (a) Reading from left to right: the interpretation of removing a column in the matrix  $x$  is to perform contraction of an edge corresponding to a locus shown in bold. Reading from right to left: the interpretation of inserting back a column while assigning new binary values is an edge insertion. The circled node  $v$  refers to Step 1. The subtrees in bold refer to those selected in Step 2. The edge in bold, the one introduced in Step 3. (b) Decomposition used for the recursion of Section 9.4.3.

SA535, 493 cells



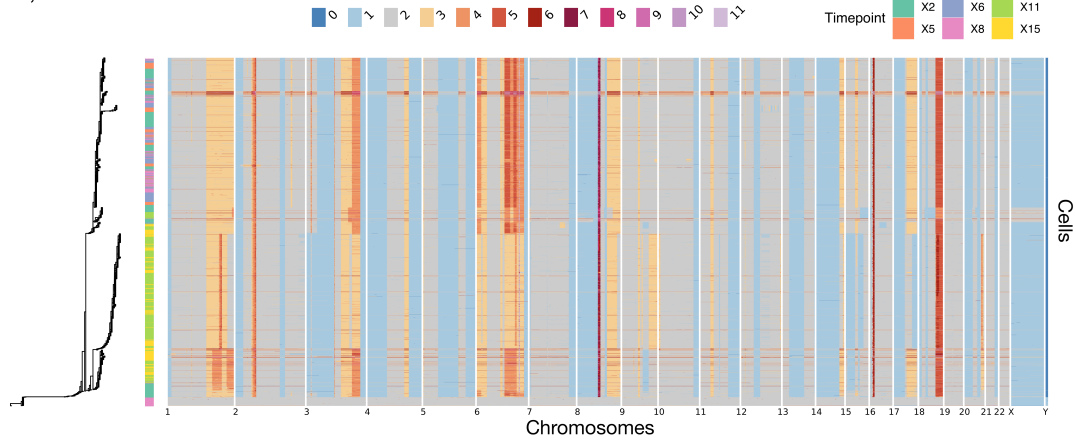
**Supplemental Figure 6.** Phylogenetic tree and CNA profile heatmap for the SA535 dataset. The rows of the heatmap are sorted according to the placement of cells on the phylogenetic tree. The columns of the heatmap are sorted by their genomic position.

OVA, 997 cells

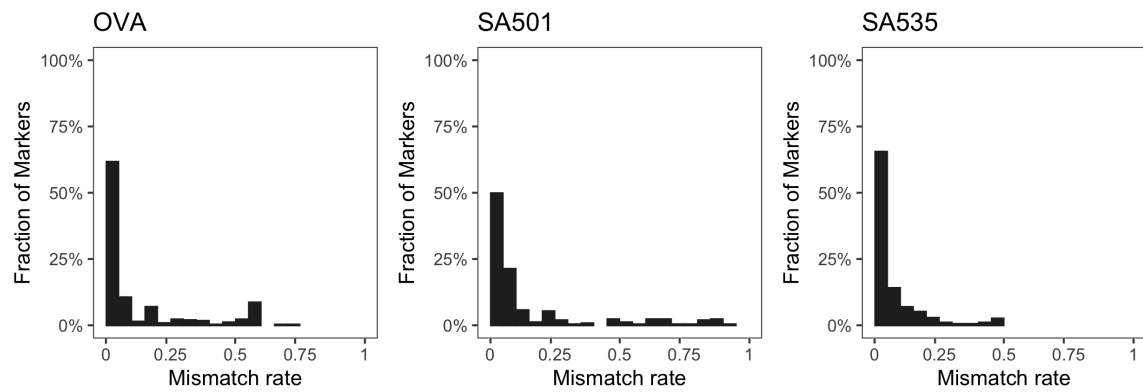


**Supplemental Figure 7.** Phylogenetic tree and CNA profile heatmap for the *OVA* dataset. The nearly diploid cells with the loss of heterozygosity on chromosome X are from SA1090. The cells with an amplification on chromosome 22 are from SA922. The rest belong to SA921.

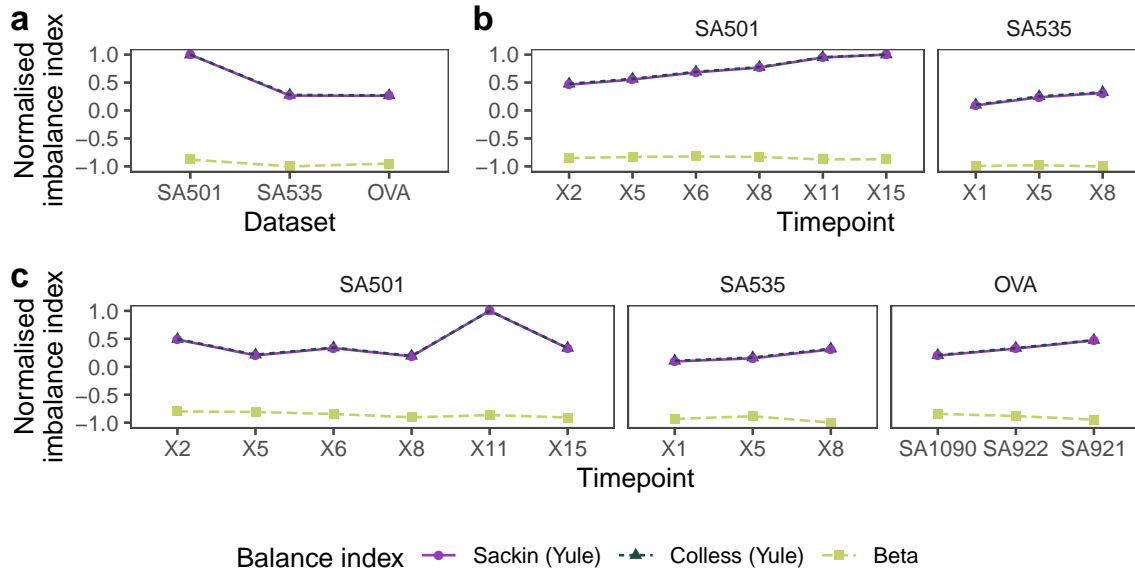
SA501, 2412 cells



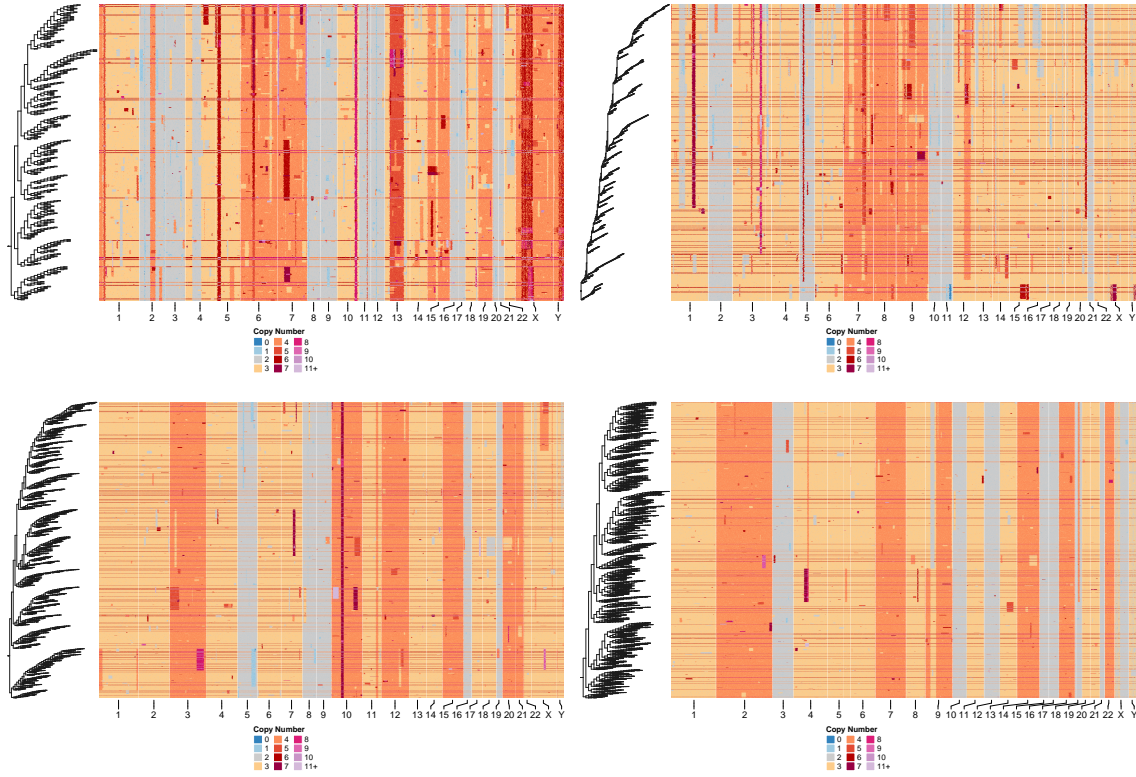
**Supplemental Figure 8.** Phylogenetic tree and CNA profile heatmap for the SA501 dataset. Note that the diploid cells at the bottom of the heatmap are control cells that were included in the experiment.



**Supplemental Figure 9.** the distribution of mismatch rate defined as the fraction of cells that have a mismatch between the inferred and jitter-fixed value of a marker.

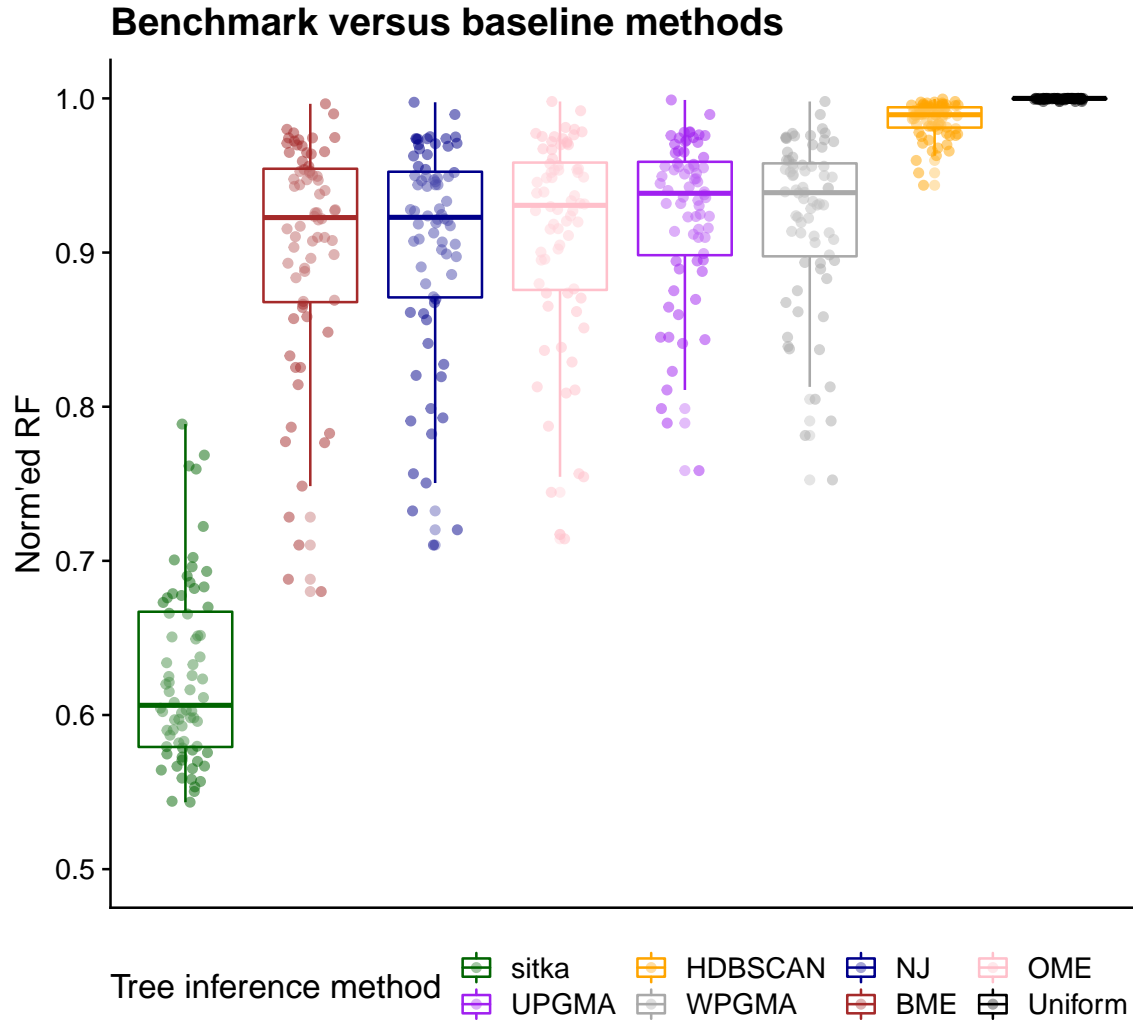


**Supplemental Figure 10.** (a) Tree imbalance index where zero indicates that the tree is consistent with one simulated from a Yule model (completely balanced) and positive values indicate deviation from the Yule model (more imbalanced). For ease of plotting, each balance index is normalised by the absolute value of the maximum estimated statistic among all samples. Cumulatively adding more timepoints (b), or for the maximal subtree comprising cells of a specific timepoint (c).

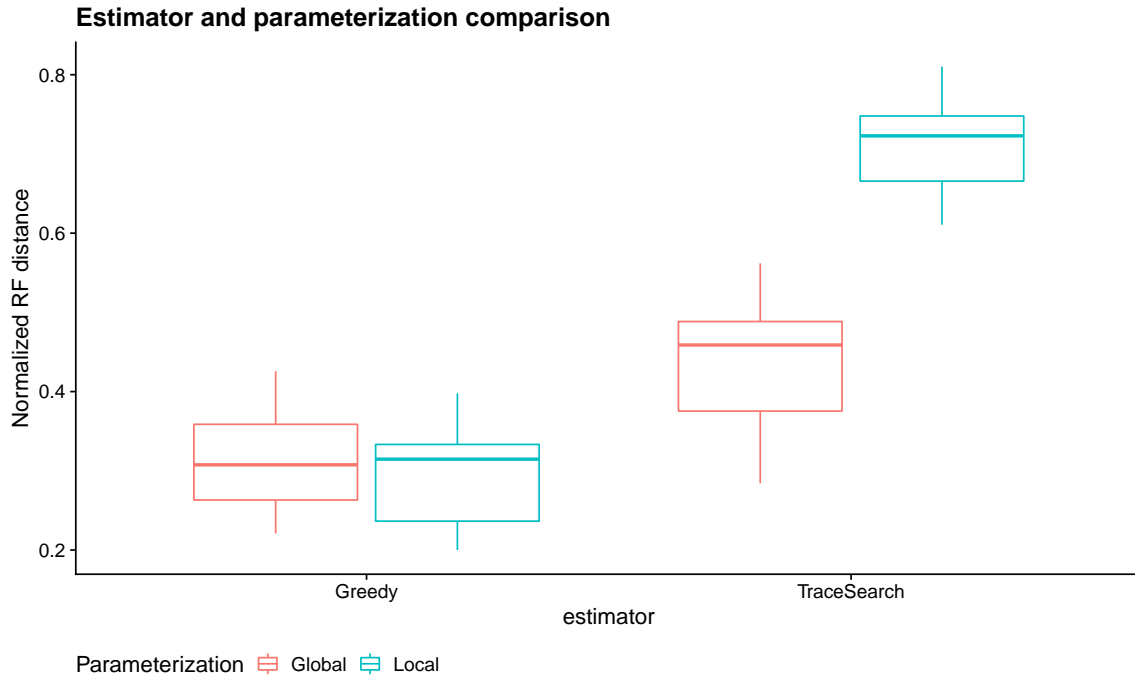


**Supplemental Figure 11.** Synthetic datasets simulated from Beta-splitting processes.

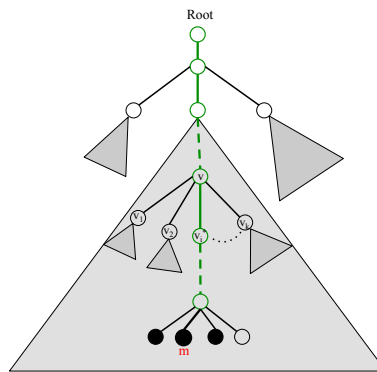




**Supplemental Figure 12.** Tree reconstruction evaluation using a normalized Robinson–Foulds metric on synthetic datasets from  $S72$ , simulated from Beta-splitting processes. Here normalization is done by dividing the RF distance of each inference method by the worst performer per dataset.



**Supplemental Figure 13.** A model and estimator comparison based on tree reconstruction accuracy for datasets from  $S_{10}$ . For each dataset, inference was performed on both the globally- and locally-parameterized model. Both the greedy and trace search estimates were computed for each inference result.



**Supplemental Figure 14.** A schematic view of the underlying tree inferred from CNA and SNV loci across multiple cells. Black and white nodes represent cells and loci, respectively. The grey triangle represents a subtree rooted at a node. It includes all of the nodes and edges in the subtree.

Copy number and cell meta-data

cn.csv

drop low-mapability bins

cn\_bin\_filtered.csv

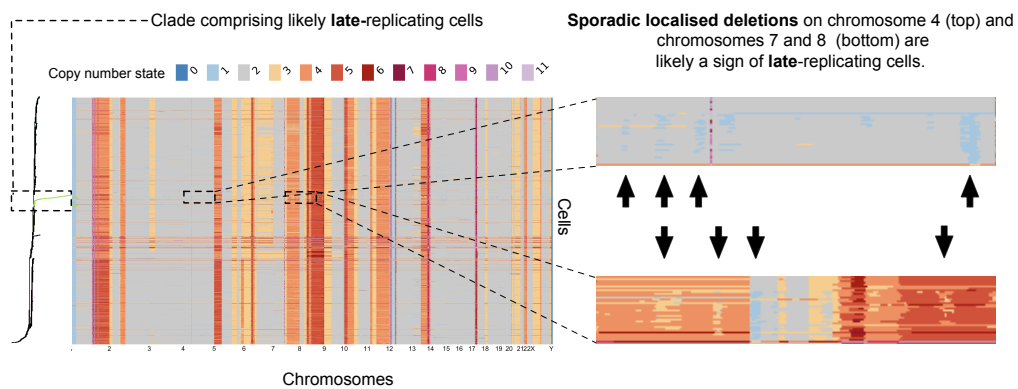
drop low-quality,  
contaminated,  
and cycling cells

cn\_bin\_cell\_filtered.csv

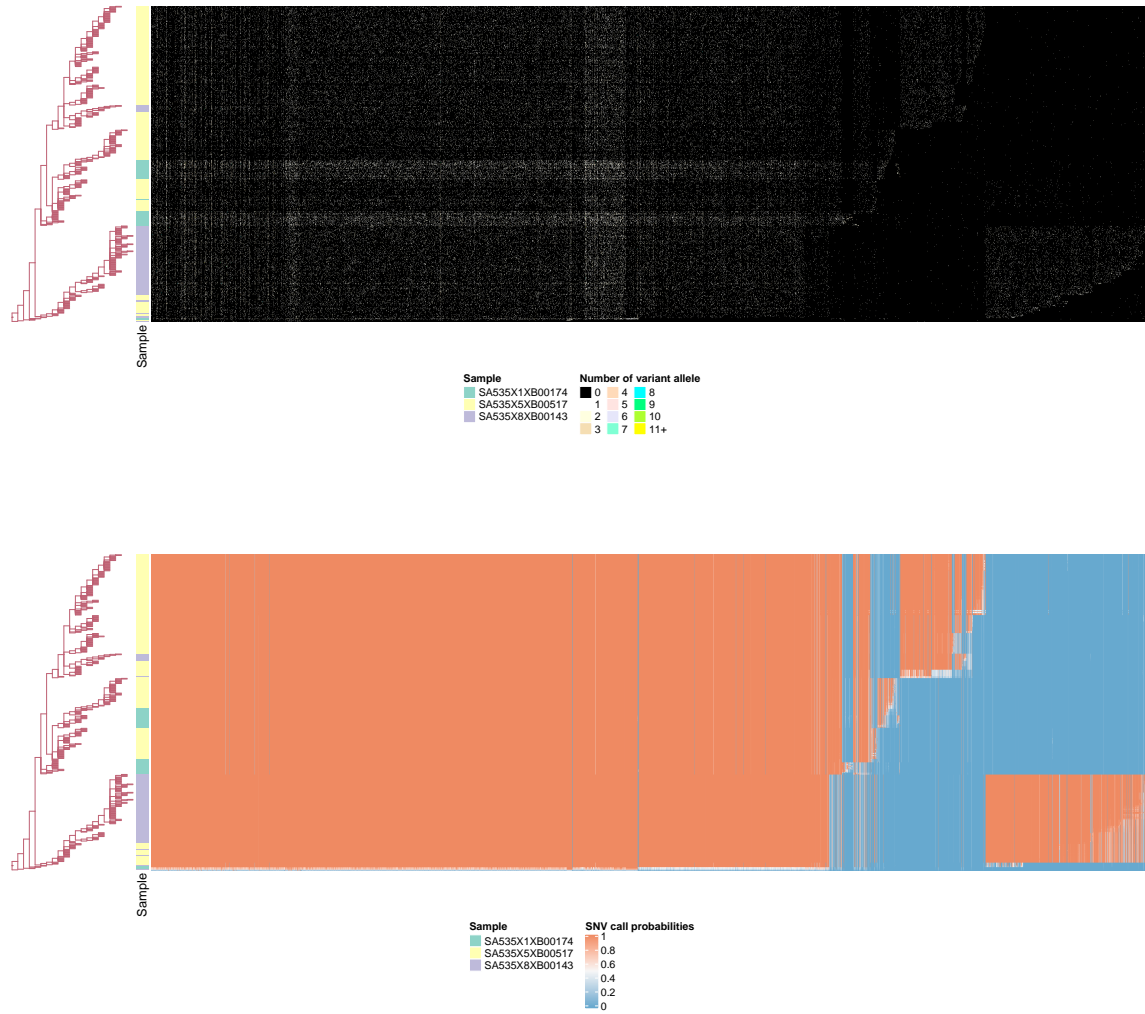
drop suspect cycling cells

cn\_bin\_cell\_filtered\_no\_jump.csv

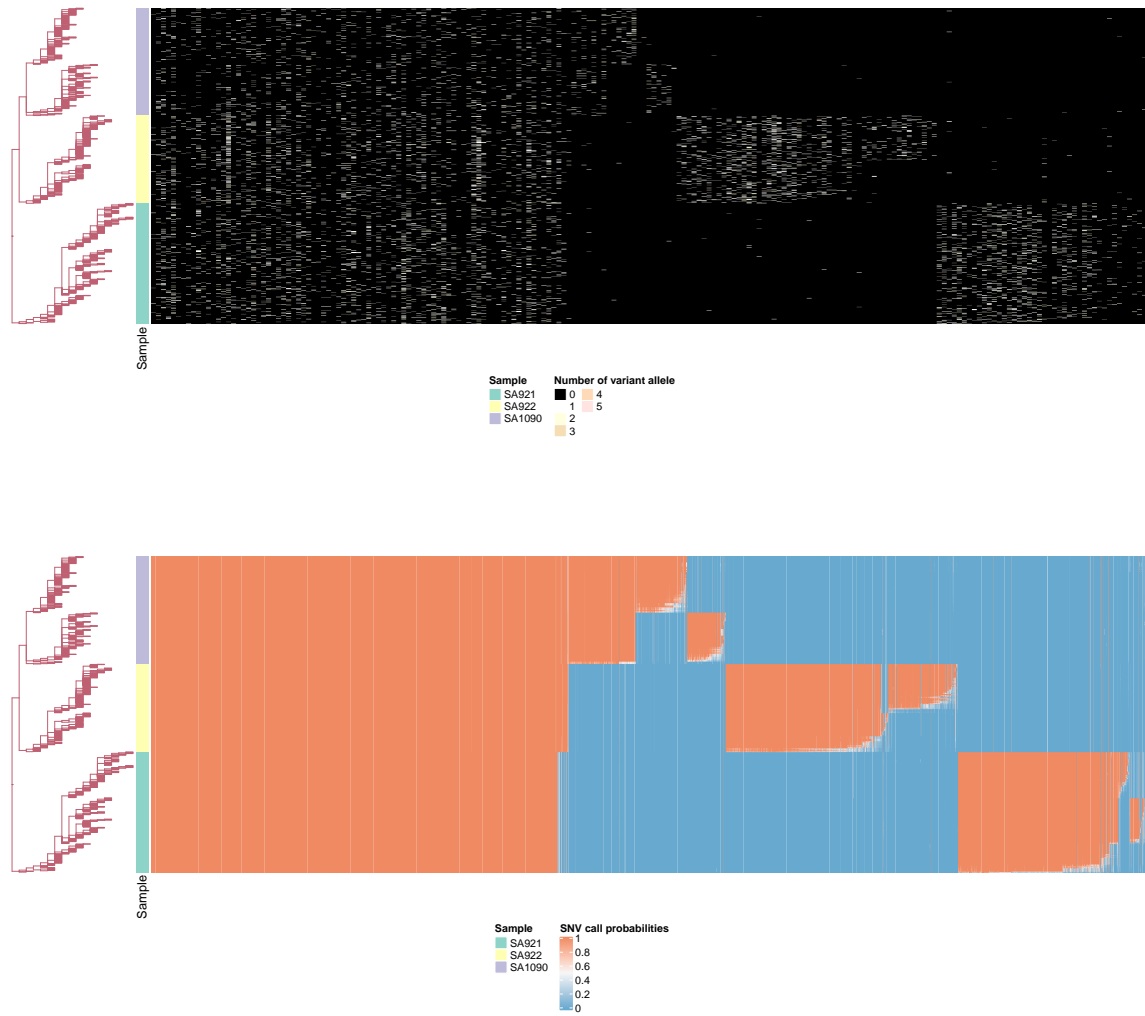
**Supplemental Figure 15.** Filtering the CNA data for tree inference.



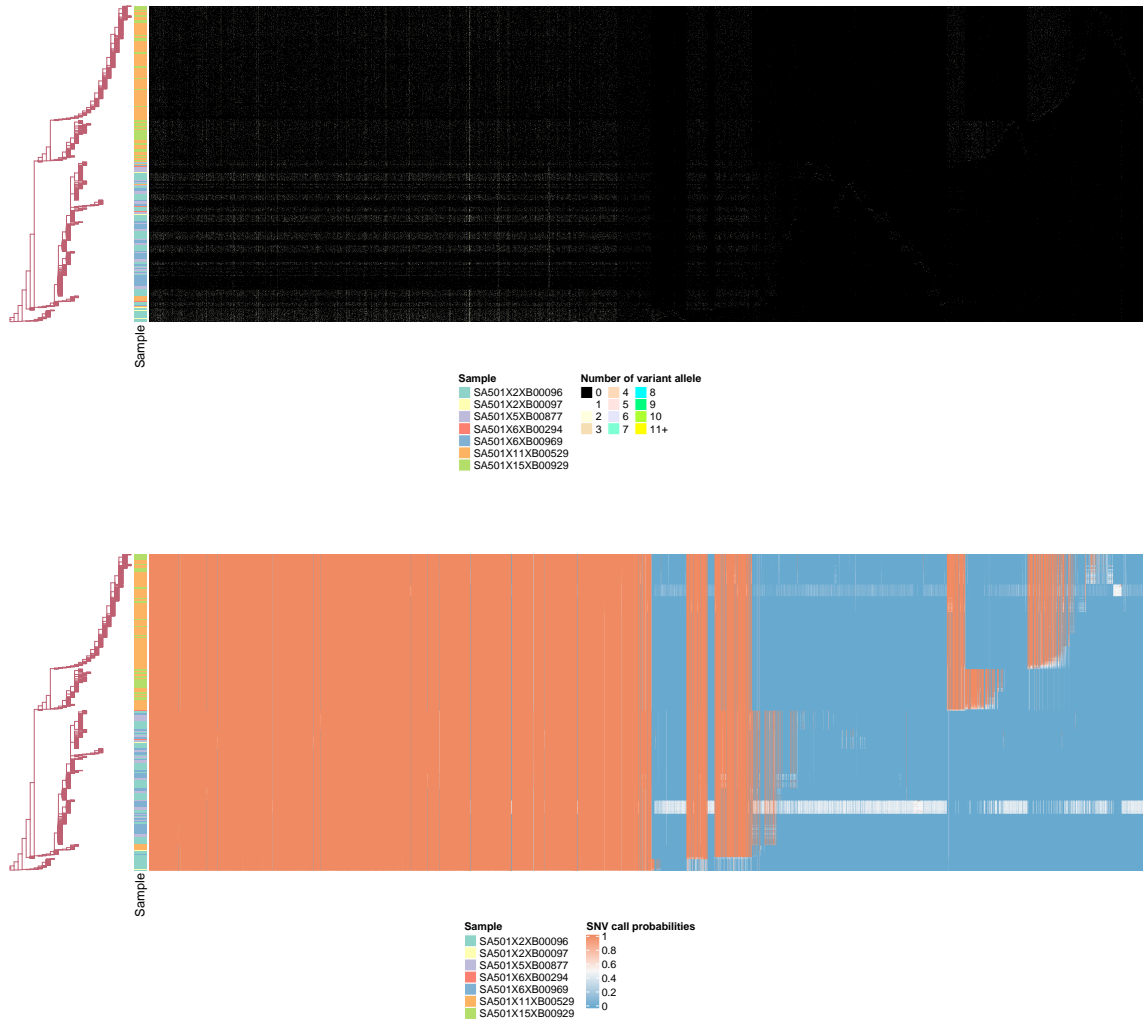
**Supplemental Figure 16.** An example of replicating cells. Note the scattered localised deletions. This heatmap is from a HER2+ PDX line. These late replicating cells form a *finger* like clade in the tree. The top inset shows chromosome 4 while the bottom inset spans chromosomes 7 and 8.



**Supplemental Figure 17.** SNV variant reads data and SNV call probabilities for SA535 dataset beside the underlying phylogenetic tree.



**Supplemental Figure 18.** SNV variant reads data and SNV call probabilities for OVA dataset beside the underlying phylogenetic tree..



**Supplemental Figure 19.** SNV variant reads data and SNV call probabilities for SA501 dataset beside the underlying phylogenetic tree.

## 824 List of Supplementary Figures

825 **Supplemental Figure 1.** Description of the process involved in the construction of *mark-*  
826 *ers*, the input to the sitka model

827 **Supplemental Figure 2.** Visualisation of a small tree and its corresponding marker ma-  
828 trix

829 **Supplemental Figure 3.** Perfect phylogeny and effects of overlapping CNA events.

830 **Supplemental Figure 4.** Pathological tree reconstruction under default observation  
831 prior

832 **Supplemental Figure 5.** (a) Reading from left to right: the interpretation of removing a  
833 column in the matrix  $x$  is to perform contraction of an edge corresponding to a locus shown  
834 in bold. Reading from right to left: the interpretation of inserting back a column while  
835 assigning new binary values is an edge insertion. The circled node  $v$  refers to Step 1. The  
836 subtrees in bold refer to those selected in Step 2. The edge in bold, the one introduced in  
837 Step 3. (b) Decomposition used for the recursion of Section 9.4.3.

838 **Supplemental Figure 6.** Phylogenetic tree and CNA profile heatmap for the SA535  
839 dataset

840 **Supplemental Figure 7.** Phylogenetic tree and CNA profile heatmap for the *OVA*  
841 dataset

842 **Supplemental Figure 8.** Phylogenetic tree and CNA profile heatmap for the SA501  
843 dataset

844 **Supplemental Figure 9.** The distribution of mismatch rate

845 **Supplemental Figure 10.** The tree topology balance

846 **Supplemental Figure 11.** A subset of synthetic datasets used for benchmarking

847 **Supplemental Figure 12.** Tree reconstruction evaluation on synthetic datasets.

848 **Supplemental Figure 13.** Synthetic experiment comparing global versus local parameter-  
849 izations.

850 **Supplemental Figure 14.** A schematic representation of CNA and SNV tree

851 **Supplemental Figure 15.** Filtering the CNA data for tree inference.

852 **Supplemental Figure 16.** An example of replicating cells

853 **Supplemental Figure 17.** SNV variant reads data and SNV call probabilities for SA535  
854 dataset beside the underlying phylogenetic tree.

855 **Supplemental Figure 18.** SNV variant reads data and SNV call probabilities for OVA  
856 dataset beside the underlying phylogenetic tree..

857 **Supplemental Figure 19.** SNV variant reads data and SNV call probabilities for SA501  
858 dataset beside the underlying phylogenetic tree.

**Supplemental Table 1.** Summary of real-world datasets used. *final* is the final number of cells after all filters except for *!lmr* are applied. *final* additionally filters out *lmr* cells, those that have total mapped reads fewer than 500,000. Abbreviations used are *tp*: time point; *qual.* : quality; *!sphase*: not *sphase*; *!lmr*: not low mapped reads.

Dataset	parameter	value
Real datasets	engine	PT
Real datasets	globalParameterization	true
Real datasets	fprBound	0.1
Real datasets	fnrBound	0.5
Real datasets	nChains	1
Real datasets	nScans	1000
Real datasets	nPassesPerScan	1
Real datasets	thinning	1
Real datasets	burnin fraction	0.5
<i>S72</i>	engine	PT
<i>S72</i>	globalParameterization	true
<i>S72</i>	fprBound	0.1
<i>S72</i>	fnrBound	0.5
<i>S72</i>	nChains	1
<i>S72</i>	nScans	20000
<i>S72</i>	nPassesPerScan	1
<i>S72</i>	thinning	1
<i>S72</i>	burnin fraction	0.5
<i>S10</i>	globalParameterization	true, false
<i>S130</i>	globalParameterization	true
<i>S10,S130</i>	engine	PT
<i>S10,S130</i>	fprBound	0.1
<i>S10,S130</i>	fnrBound	0.5
<i>S10,S130</i>	nChains	8
<i>S10,S130</i>	nScans	5000
<i>S10,S130</i>	nPassesPerScan	10
<i>S10,S130</i>	thinning	1
<i>S10,S130</i>	burnin fraction	0.5

**Supplemental Table 2.** Inference settings used for each dataset.