

A method for phylogenetic reconstruction of aneuploid cancers based on multiregional genotyping data

Natalie Andersson (✉ natalie.andersson@med.lu.se)

Lund University <https://orcid.org/0000-0002-3643-4404>

Subhayan Chattopadhyay

Lund University

Anders Valind

Lund University

Jenny Karlsson

Lund University

David Gisselsson

Lund University <https://orcid.org/0000-0002-0301-426X>

Article

Keywords: tumor heterogeneity, pediatric cancer, phylogenetic reconstruction

Posted Date: January 22nd, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-140537/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Version of Record: A version of this preprint was published at Communications Biology on September 20th, 2021. See the published version at <https://doi.org/10.1038/s42003-021-02637-6>.

A method for phylogenetic reconstruction of aneuploid cancers based on multiregional genotyping data

Natalie Andersson¹, Subhayan Chattopadhyay¹, Anders Valind^{1,2}, Jenny Karlsson¹ & David Gisselsson^{1,3,4}

¹Division of Clinical Genetics, Department of Laboratory Medicine, Lund University, Lund, Sweden.

²Department of Pediatrics, Skåne University Hospital, Lund, Sweden. ³Division of Oncology-Pathology, Department of Clinical Sciences, Lund University, Lund, Sweden. ⁴Clinical Genetics and Pathology, Laboratory Medicine, Lund University Hospital, Skåne Healthcare Region.

Correspondence and request for materials should be addressed to N.A. (email: natalie.andersson@med.lu.se).

Abstract word count: 129

Text word count: 6563

References: 37

Figures: 4

Tables: 1

Supplementary Figures: 7

Supplementary Datasets: 17

1 **Abstract**

2 Phylogenetic reconstruction of cancer cell populations remains challenging. There is a particular lack of
3 tools that deconvolve clones based on copy number aberration analyses of multiple tumor biopsies
4 separated in time and space from the same patient. This has hampered investigations of tumors rich in
5 aneuploidy but few point mutations, as in many childhood cancers. Here, we present DEVOLUTION,
6 an algorithm for subclonal deconvolution followed by phylogenetic reconstruction from bulk
7 genotyping data. It integrates copy number and sequencing information across multiple tumor regions
8 throughout the inference process, provided that the mutated clone fraction for each mutation is known.
9 We validate DEVOLUTION on SNP-array data from 56 pediatric tumors comprising 253 tumor
10 biopsies confirming concordance to biological mechanisms and show a robust performance on extensive
11 simulations of bulk genotyping data.

12

13

14

15

16

17

18

19

20

21

22

23

24 **Introduction**

25 Neoplasms are a heterogenous group of diseases driven by Darwinian selection. Most cancers are
26 presumed to originate from a single mutated cell, from which each mutation is conveyed to its daughter
27 cells, that in turn can acquire additional aberrations, establishing subpopulations (subclones) of cells
28 with diverse genetic compositions within the tumor¹. This evolution of cancer cells is further shaped by
29 genetic drift and selection pressure from the tumor microenvironment and oncological treatment²⁻⁵. Due
30 to this process, many cancers exhibit vast intratumor heterogeneity (ITH) as well as intertumor
31 heterogeneity between the primary tumor and its metastases⁶⁻¹⁰. Knowledge about how ITH emerges
32 over time remains limited and multiple models have been proposed to explain it such as punctuated,
33 neutral, linear, and branched evolution as well as a big bang model of tumor growth followed by neutral
34 evolution¹¹⁻¹³. By analyzing the genetic variation of the tumor spatially as well as temporally,
35 mathematical methods can be employed in order to reconstruct its evolution, commonly in the form of
36 a phylogenetic tree that links together distinct cancer cell subpopulations in an inferred temporal order.
37 Such phylogenetic reconstructions can improve the understanding of tumorigenesis, progression to
38 metastatic disease, and aid the development of novel therapeutic strategies^{7,8,14}.

39 One of the biggest challenges in phylogenetic analysis of bulk sample data from tumors is that the
40 genetic analysis of each sample is conducted on millions of cells at once, usually constituting multiple
41 subclones. The relative proportions of the subclones within each biopsy may also vary across the
42 biopsied regions, stressing the need to integrate information from multiple biopsies separated in space
43 to thoroughly assess the genetic profile of the tumor. Not addressing this may result in the prediction of
44 illicit biological trajectories and so-called biopsy trees, not constituting true phylogenies¹⁵. Phylogenetic
45 relationships should thus ideally be constructed based on the deconvolved clonal structure, i.e. one
46 should infer which subclones are characterized by which alterations and, in addition, the ancestral order
47 of these.

48 Even though single cell sequencing (SCS) has emerged as an important tool for temporal reconstruction
49 that circumvents the issue of clonal deconvolution it is very costly to implement on the scales needed in

50 the clinic and usually provides limited sequence coverage¹⁶. A more cost-effective alternative is to
51 perform computerized deconvolution of bulk genotyping data derived from single nucleotide
52 polymorphism array (SNP-array), targeted deep sequencing (TDS), whole exome sequencing (WES) or
53 whole genome sequencing (WGS). Bulk genotyping yields a set of genetic alterations present in each
54 biopsy along with information that can be utilized to estimate the proportion of cells harboring each
55 aberration in each biopsy, denoted the mutated clone fraction (MCF)¹⁷. However, most tools developed
56 for computerized deconvolution of bulk genotyping data focuses solely on somatic point mutations,
57 presumes a diploid background, and lacks specific pipelines to handle intratumoral heterogeneity of
58 copy number alterations (CNAs) of chromosomal segments or whole chromosomes. In addition, they
59 do not provide the possibility to infer phylogenetic trees solely on copy number aberration data from
60 multiple biopsies separated in time and space¹⁸⁻²¹. Since most cancers are aneuploid to some degree²²,
61 this is a serious shortcoming, especially for cancer types where aneuploidy is a common feature such as
62 high-grade adult carcinomas and many childhood cancers^{8,23}. Consequently, there is a particular need
63 for tools capable to infer phylogenetic trees based on multiregional copy number data.

64 To fill this methodological gap, we introduce DEVOLUTION, an algorithm for subclonal deconvolution
65 followed by phylogenetic reconstruction from bulk genome profiles including high-resolution copy
66 number data (e.g. from SNP-array, WES or WGS) and sequencing information (e.g. from WES, WGS
67 or TDS) separately or in unison. The deconvolution is based on à priori MCF-estimation of the
68 individual aberrations in each sample and the algorithm systematically combines information from all
69 available biopsies throughout the inference process to reconcile the most probable temporal evolution
70 of the tumor by inferring an event matrix that is used to reconstruct phylogenetic trees. Importantly it
71 can deduce evolutionary trajectories based on copy number data alone. In addition, predictions of the
72 subclonal size and compositions across biopsies are visualized directly in the phylogenetic tree.
73 DEVOLUTION provides an objective framework for creating event matrices and phylogenetic trees
74 from bulk genotyping data, avoiding subjective bias compromising the validity of tree-to-tree
75 comparisons (S. Table 1).

76 To demonstrate DEVOLUTION's utility, the algorithm was evaluated using SNP-array data from 253
77 tumor regions from 56 pediatric cancers including neuroblastoma (NB), Wilms tumor (WT), and
78 rhabdomyosarcoma (RMS), comprising the most common extracranial, solid tumors in children.
79 Additionally, the algorithm showed a robust performance on simulated multiregional bulk genotyping
80 data. DEVOLUTION holds the potential to facilitate further insights into the development, progression,
81 and response to treatment, particularly in tumors with high burden of chromosomal copy number
82 alterations.

83

84

85 **Results**

86 **Overview of the algorithm workflow**

87 The algorithm operates on multiregional sampling data analyzed using whole genome profiling followed
88 by MCF-computation (**Figure 1a-c**). The input file is an $u \times v$ dimensional matrix, containing
89 information about the u genetic alterations detected in a tumor. For each alteration there are v columns
90 indicating the genetic position, alteration type as well as the proportion of cells in the biopsy harboring
91 the alteration (**Supplementary Figure 1**). For copy number aberrations, the matrix is subjected to an
92 algorithm identifying all unique events across the samples while considering the uncertainty in the
93 aberration breakpoint measurement (**Supplementary Figure 2**). The DBSCAN (density-based spatial
94 clustering of applications with noise) algorithm is then used to identify clusters of genetic alterations
95 having similar cellular proportions across multiple samples, indicating that they might be reflecting a
96 group of cells having an identical genetic profile (**Figure 1d**). By identifying these clusters of genetic
97 alterations, the computational load can be decreased and the unfolding of the subclonal composition
98 aided.

99 DBSCAN prepares the data set to be subjected to the subclonal deconvolution algorithm, which is
100 employed to elucidate the temporal order of the clusters of mutations. Information from multiple

101 biopsies is integrated throughout this process to minimize the occurrence of parallel evolution (PLC)
102 and back mutation contradictions (BMC). PLC in copy number data means that the same type of genetic
103 alteration with the same genomic start- and endpoints in the chromosome appears independently in
104 different cells within the tumor. This is, in most cases, unlikely from a biological standpoint, unless
105 there is a copy number alteration including an entire chromosome. BMC incorporates genetic alterations
106 that are gained and then lost further down the evolutionary history in the tree, which may be feasible
107 scenarios for some types of genetic alterations such as a gain of a whole chromosome that is later lost
108 but are less likely to occur for structural chromosomal aberrations and point mutations, and should never
109 be occur for loss of heterozygosity events²⁴.

110 In addition, the user can provide a matrix containing information about illicit orders of genetic events,
111 that can be taken into consideration during the deconvolution (**Figure 1e**). The deconvolution
112 culminates in a suggestion of the most likely temporal order of all genetic aberrations, constituting the
113 basis for the creation of an event matrix, illustrating the distribution of genetic alterations across
114 subclones (**Figure 1f**). Using this event matrix, the biological distance between the subclones is
115 calculated using the Hamming distance²⁵ and phylogenetic trees are reconstructed using the maximum
116 likelihood and parsimony methods. In addition, the algorithm provides the distribution and size of the
117 clusters across the samples, resulting in an overview of the dynamics, spatial distribution, and
118 dissemination of the tumor (**Figure 1g**).

119

120 **Validation on pediatric tumors confirms concordance to biological mechanisms**

121 DEVOLUTION was applied to a previously reported dataset of 56 pediatric cancers and phylogenetic
122 trees were generated based on copy number data for 22 neuroblastomas, 20 Wilms tumors and 8
123 rhabdomyosarcomas comprising a total of 253 biopsies (**Figure 2, Supplementary Figure 3,**
124 **Supplementary File 1**). Event matrices but not phylogenetic trees could be reconstructed from six
125 patients (NB1, NB24, WT1, WT2, WT3 and WT5) in which all samples from the same patient had
126 identical genomic profiles.

127 The phylogenetic trees of the remaining 50 tumors all represented plausible biological scenarios and
128 often illustrated key events in tumor evolution. For example, in NB5 (**Figure 2a**), where a primary tumor
129 and a metastasis presented at the same time, the metastasis was demonstrated to originate from a
130 population of cells having the genetic alterations of the stem as well as one group of cells encompassing
131 a subclone also present in the primary tumor. This indicates polyclonal seeding. A more complex pattern
132 of polyclonal seeding was observed in NB22 (**Figure 2b**), a patient with progressive tumor growth
133 across multiple metastatic locations. Here there was more subclonal variation among the lymph node
134 metastasis than among metastases to distant organs. The metastases to distant organs often presented as
135 solitary branches, exemplified by the same subclone colonizing both the lung and skull. This might
136 indicate that the threshold for tumor cells to escape the primary tumor and colonize the lymph nodes is
137 lower than for colonization of distant organs, displayed as a wider variety of different subclones across
138 lymph node locations compared to an extensive selection for a certain subclone in distant loci. WT11
139 and WT19 (**Figure 2c and d**) both showed subclones that were distributed across several locations
140 within the primary tumors, a phenomenon which has previously been demonstrated to be common in
141 Wilms tumors¹⁷. RMS8, an alveolar rhabdomyosarcoma (**Figure 2e**), displayed an intricate evolutionary
142 pattern with many genetic alterations. Here the primary tumour's subclones form a cluster at the root of
143 the tree, while the cell populations from a metastasis and a local relapse share a branch having a vast
144 amount of additional genetic alterations. Hence, phylogenetic trees produced from copy number profiles
145 by DEVOLUTION, can provide biological insights that might aid the understanding of how cancer
146 develops and progresses in individual patients.

147

148 **Contradictions are rarely seen in the phylogenetic trees**

149 In 80 % of all tumors analyzed, the maximum likelihood (ML) and parsimony (MP) methods resulted
150 in identical phylogenetic trees (19/22 NB, 17/20 WT, 4/8 RMS) (**Figure 3a-d**). When the ML and MP
151 tree for the same case did differ from one another, the differences in the branching structure were minor
152 (**Supplementary Figure 3**). We identified the positions of the genetic copy number alterations in each

153 tree to identify contradictions based on prior knowledge about how genetic aberrations occur in cancer
154 cells. More specifically, we analyzed instances of PLC and BMC.

155 When the ML and MP trees differed from one another, this was always due to a PLC/BMC in the
156 phylogeny, often with both contradictions together in the same tree. PLC and/or BMC were found in
157 14/50 ML-trees and 14/50 MP-trees (5 NB, 5 WT, 4 RMS), hence 28/100 trees in total. Of these, 3/5
158 NB, 3/5 WT and 1/4 RMS trees contained only one single contradiction located among the leaves of the
159 trees i.e. it did not have any significant impact on the tree structure. In addition, there did not seem to
160 be any apparent difference between the frequency of the types of contradictions between ML and MP
161 trees (**Figure 3e**). Excluding the cases with PLC and BMC of whole chromosomes, which are plausible
162 events, only 8/50 ML-trees and 8/50 MP-trees (0 NB, 4 WT and 4 RMS) exhibited contradictions in the
163 tree structure. In these eight cases, merely a few genetic alterations caused the PLC and/or BMC. They
164 were caused by aberrations altering clone size compared to another event across samples or similar
165 aberrations that still fell outside the breakpoint cutoff for similarity causing them to be considered
166 separate events. These situations may be resolvable by critically reviewing the original data
167 (**Supplementary Figure 4**). Alternating clone sizes was particularly common in the RMS trees
168 (**Supplementary Figure 5**). The RMS tumors had a significantly higher total branch length than NB
169 and WT (**Figure 3b**), indicating a more complex genomic profile. They also had a mean number of
170 genetic alterations per biopsy of 20, most of them present in > 50 % of cells in a single biopsy thus
171 allowing just one single solution of the temporal evolution. To prompt review of original data when
172 pertinent, the software will warn the user that there is a contradiction in the data set and the tree might
173 therefore not be entirely biologically accurate.

174

175 **Evaluation using simulated data**

176 The reliability of the algorithm was further evaluated using simulated bulk sampling data. In the patient
177 data set, the median number of unique subclonal alterations in total across all biopsies were 6 for NB, 5
178 for WT and 13 for RMS. To accommodate this large variation, the simulation was conducted for three

179 different mutation frequencies resulting in 15, 50 and 100 subclonal genetic alterations distributed across
180 40000 virtual tumor cells (**Figure 4a**). Virtual biopsies were sampled randomly from the set of cells
181 while varying the number of biopsies from 1 to 10, generating a segment file along with a list of true
182 unique subclones across biopsies for each mutational frequency. Hence, analysis could be performed
183 using DEVOLUTION while having the true subclonal composition at hand.

184 As expected, when increasing the number of biopsies or the mutation frequency more genetic aberrations
185 were identified. In addition, a higher number of subclones were correctly allocated. (**Figure 4b**). A
186 higher number of overall mutations will provide the software with more information, which is why the
187 number of correctly allocated genetic alterations increase with mutation frequency for the same number
188 of biopsies. Specifically, the number of positions in the temporal sequence at which genetic alterations
189 can be allocated decreases with the number of genetic alterations. The number of unique subclones
190 increases while the proportion of cells in each biopsy representing each subclone decreases. However,
191 when more alterations are correctly allocated in absolute numbers, the proportion of correctly allocated
192 alterations will not. Throughout the different mutational frequencies $93 \pm 5.5 \%$ of the genetic
193 alterations were correctly allocated in the event matrix (**Figure 4c**). Thus, when increasing the number
194 of biopsies, the absolute number of correctly allocated genetic alterations increases, but the proportion
195 of correctly allocated alterations does not change significantly. The reason for this is that sampling
196 additional locations will also increase the chance of finding an area with a late genetic alteration that is
197 hard to correctly place in the phylogeny because of its low spatial dissemination. We further dissected
198 why not all subclones were correctly allocated. It was hypothesized to occur due to low spatial
199 dissemination, resulting in the presence of certain genetic alterations in a subset of biopsies. Excluding
200 all genetic alterations found in only one single biopsy in fact resulted in the correct allocation if $99.5 \pm$
201 0.8% of the genetic alterations in this mixture of entities (**Figure 4d-e**).

202

203

204

205 **Relationship to existing work**

206 Much progress has been made in the field of clonal deconvolution. Many methods are although limited
207 to integrating information from a single biopsy, such as TITAN and THetA in addition to only accepting
208 sequence data^{26,27}. PyClone and SciClone are, however, employable on multiple biopsies, but assume
209 that all detected CNAs are clonal i.e., present in all cells, do not infer the evolutionary relationship
210 between the identified clusters and sequence data is required for them to operate^{19,20}. In addition,
211 SciClone focuses exclusively on SNVs in copy number neutral and loss of heterozygosity (LOH) free
212 portions of the genome. Both PyClone and SciClone mainly operate as clustering algorithms and do not
213 infer the order of genetic alterations, hence their output could be used as an input to DEVOLUTION.
214 Hence, they fulfill different purposes and is not meant to exclude one other. PhyloWGS on the other
215 hand can use both SNAs and CNAs to infer a phylogeny and is employable on multiple samples. The
216 algorithm does although not integrate information between samples during the inference procedure,
217 representing a loss valuable information, and is limited to WGS data¹⁸. Also, other methods such as
218 Clomial, LiCheE and SCHISM are specifically designed for SNAs and cannot solely include CNAs to
219 infer a phylogeny²⁸⁻³⁰. REVOLVER also requires sequencing data, cannot use CNA-data alone and is
220 specifically designed to integrate phylogenies from large cohorts of patients in order to infer common
221 trajectories of repeated evolution. SPRUCE requires input of sequencing data to infer phylogenetic trees.
222 Hence, most methods available focus exclusively on sequencing data and there is currently no available
223 tool that can reconstruct a phylogenetic tree based on multiregional SNP-array data alone, which is a
224 commonly used genotyping method in the clinic. Additionally, many methods focused on SNA only
225 outputs clusters of genetic alterations and their MCF. Event matrices are subsequently often constructed
226 manually from these MCF estimates, posing a risk for unintentional subjective bias in the deconvolution
227 process, especially when integrating information from multiple biopsies conjointly. These
228 methodological gaps are filled by DEVOLUTION, which provides a subjective method to infer
229 phylogenies based on a priori MCF estimations based on preset rules that are employed equally across
230 all patient data, hence providing a standardized framework for inferring phylogenies from bulk
231 genotyping data, thus allowing tree-to-tree comparison without the risk of bias from subjective curation.

232

233

234

235 **Discussion**

236 A single biopsy from a tumor can consist of multiple distinct subclones and their prevalence may vary
237 across biopsied areas. Not addressing this fact when studying cancer cell evolution can be deleterious
238 and result in incorrect phylogenies¹⁵, stressing the need for multispatial and temporal sampling to unfold
239 the genomic landscape. DEVOLUTION thoroughly assesses the problem by combining information
240 obtained across multiple biopsied regions throughout the entire subclonal deconvolution effectively
241 deconvolving subclones transversing clonal territories, with the potential to concomitantly include both
242 point mutations and copy number alteration data. In contrast to other methods DEVOLUTION allows
243 phylogenetic trees to be constructed using copy number information alone and integrating information
244 from multiple biopsied areas throughout the inference procedure. The algorithm can combine data from
245 SNP-array, TDS, WES and WGS, provided the MCF for each genetic alteration is known, which can be
246 computed based on the log₂ ratio for copy number alterations or variant allele frequencies (VAF) for
247 point mutations as described extensively elsewhere^{17,19,20}. Since aneuploidy is a common feature across
248 many adult carcinomas and a majority of childhood cancers, the integration of copy number aberrations
249 in the phylogeny holds the potential to increase the understanding of the evolution of these diseases. If
250 needed, there is also a possibility for user curation, for example if it is known that certain genetic
251 alterations cannot co-exist.

252 Evaluating the software using high-throughput SNP-array data from 253 multitemporal and spatial
253 regions from 56 pediatric tumors produced phylogenetic trees that were well in concordance with prior
254 knowledge of how chromosomal aberrations occur in cancer cells. Surprisingly, generating phylogenetic
255 trees using ML and MP predominantly yielded identical tree structures. When thoroughly examining
256 these trees, contradictions such as PLC and BMC were identified in the cases where the ML/MP trees
257 differed, which were found to always be due to disagreements in the original data set. The user is

258 therefore encouraged to reevaluate this genetic alteration in the input segment file, since it may be
259 particularly subjected to noise or keep the tree if the phylogenetic situation is considered biologically
260 plausible. The extensive evaluation of the two methods did not indicate that the choice of mathematical
261 method favors a certain type of error in the phylogenetic tree.

262 Clustering genetic alterations using DBSCANs were sufficient for analyzing the pediatric tumors
263 incorporated in this study and the simulated data sets. The size of ϵ can be changed by the user to
264 increase the flexibility in the clustering, thus optimizing it further for the data set at hand to account for
265 noise. Note that the purpose of the ad hoc clustering mainly is to reduce the computational complexity.
266 If the clustering is too strict, the DEVOLUTION algorithm compensates by thoroughly integrating the
267 information across samples. The clustering algorithm can easily be changed.

268 In summary, we have seen how DEVOLUTION can be used to analyze the intratumoral heterogeneity
269 as well as the intertumoral heterogeneity between the primary tumor and metastasis through evaluation
270 on a dataset of pediatric tumors harboring extensive aneuploidy. By analyzing a cancer's phylogenetic
271 tree, an overview of its heterogeneity and temporal order of genetic alterations can be assessed, which
272 can be used to follow the tumor's evolutionary response to treatment^{31,32} and may aid the identification
273 of subclones posing a risk of metastasizing, relapsing or being resistant to therapy. It may also make it
274 possible to identify genetic alterations that seems to appear early in the tumor's development, posing
275 attractive targets for therapy since they are present in a large proportion or all cells in the tumor.

276

277 **Conclusions**

278 DEVOLUTION produces phylogenetic trees from multiregional cancer genomics data, integrating copy
279 number aberrations and sequence mutations separately or in unison, provided the MCF is known. The
280 algorithm reliably deconvolves subclones and infers the evolutionary history in well in concordance to
281 prior knowledge of how genetic aberrations occur in cancer cells. It performs robustly across real data
282 from multiple aneuploid tumor types and in simulated data across a wide variation in mutation rates and
283 sample numbers. DEVOLUTION holds the potential to facilitate insights into the development,

284 progression, and response to treatment, particularly in tumors with high burden of chromosomal copy
285 number alterations.

286

287

288 **Methods**

289 **The software**

290 The major structure of the software can be divided into five steps

- 291 1. Preprocessing of the data.
- 292 2. Clustering of genetic alterations based on information from multiregional sampling from the
293 same patient.
- 294 3. Subclonal deconvolution based on information from multiregional sampling from the same
295 patient.
- 296 4. Construction of an event matrix.
- 297 5. Usage of a mathematical model to reconstruct the phylogenetic trees, in this case
 - 298 a. Maximum likelihood
 - 299 b. Maximum parsimony

300

301 **Preprocessing of the data**

302 The input data consist of an $u \times v$ dimensional matrix containing information about the u detected
303 genetic alterations present in each biopsy. The matrix should also specify the genetic location of each
304 alteration, its type (gain, loss, cni etc.) as well as the proportion of cells harboring the alteration in that
305 particular biopsy (the mutated clone fraction, MCF), represented by the v columns (**Supplementary**
306 **Figure 1**). How the MCF can be computed from log2 ratios as well as VAFs is described extensively
307 elsewhere¹⁷. Also, methods such as PyClone and SciClone could be used to infer clusters of somatic
308 mutations and their corresponding VAFs across samples, which could then be used to calculate the MCF.

309 If allelic copy number alterations are considered, the user is advised to choose a cutoff for the detected
310 genetic alterations in the segment file to be considered separate events, reflecting the measurement
311 uncertainty regarding the start and end positions of the genetic alterations. The default cutoff is 1 Mbp.
312 The user can also choose which data types to include in the analysis. In this way e.g. SNP array and
313 sequencing data can be analyzed separately for comparison or in unison without having to separate the
314 matrix manually.

315 The algorithm scans the MCFs for missing values, indicating that the MCF has not been able to be
316 determined. If the event is considered to belong to the stem, based on biological knowledge or additional
317 data, the missing value is replaced by 100 %. Amplicon accumulation is an example of such a case when
318 it is not possible to determine the fraction of cells harboring it since the number is hypervariable. A stem
319 event is defined as the presence of the alteration in ≥ 90 % of the cells in all samples. The alterations
320 containing missing values for MCF are removed entirely if part of a subclone to not overestimate genetic
321 variation within the tumor.

322 The clustering algorithm was constructed to localize all unique genetic alterations throughout the tumor
323 samples. The program loops through the rows of the data file representing the genetic aberrations. For
324 each row it compares the genetic alterations and their position on the chromosome to all the other rows,
325 representing other detected genetic aberrations throughout the samples. If the events' start or end
326 positions differ by a certain cutoff, set by the user based on the measurement uncertainty of the data set,
327 and/or they are different aberration types, they are considered as two separate events, else they are
328 considered as the same event (**Supplementary Figure 2**). Thus, *all* conditions stated below must be met
329 for the algorithm to consider two alterations detected in the same patient to be the same.

- 330 1. Alteration 1 and 2 are localized on the same chromosome.
- 331 2. Alteration 1 and 2 harbor the same type of alteration.
- 332 3. Neither alteration 1 nor 2 should belong to the stem.
 - 333 a. Alterations belonging to the stem are always considered as separate events.
- 334 4. $X_1 = \|A_{1_{\text{start}}} - A_{2_{\text{start}}}\| \leq \text{co}_{\text{ev}}$

335 $5. X_2 = \|A_{1_{\text{end}}} - A_{2_{\text{end}}}\| \leq \text{co}_{\text{ev}}$

336 In the present study, considering allelic copy number aberrations, the cutoff (co_{ev}) for measurement
 337 uncertainty in the start and end position of the events was set to 1 Mbp, also constituting the default for
 338 DEVOLUTION. Since the chromosome sizes ranges from 48-250 Mbp this cutoff constitutes a start and
 339 end point deviation of 0.4-2 % of the chromosome length.

340

341 **Clustering of genetic alterations incorporating information from multiregional and temporal**
 342 **sampling from the same patient**

343 In our model a tumor is proposed to consist of multiple subpopulations of cells that harbor different sets
 344 of genetic alterations. Each individual alteration is part of a mutation space $m_i \in \{m_1, m_2 \dots m_\theta\}$
 345 comprising all mutations present in the tumor where $i, \theta \in \mathbb{N}^+$ and θ is the total number of mutations.
 346 The mutational profile obtained from the biopsies thus represent a subset of the total mutation space and
 347 is the information at hand to describe the evolutionary trajectory of the tumor. For this purpose, all
 348 detected mutations are combined with their respective MCF-values into a matrix representing their
 349 distribution across samples (**Supplementary Figure 6**). For a particular tumor, this results in a matrix
 350 $T_{M \times B}$ with the dimensions $M \times B$, where M is the total number of unique genetic alterations and B is the
 351 total number of biopsies available. Hence m_δ indicates a certain genetic alteration δ , and b_ω represent
 352 a biopsy ω . The value $t_{\delta\omega}$ consequently corresponds to the MCF for an alteration, m_δ , in a sample, b_ω
 353 where $t_{\delta\omega} \in [0,100]$ i.e. it is bound between 0 and 100 %. This can be written as

354
$$T_{M \times B} = \begin{bmatrix} t_{1,1} & t_{1,2} & \dots & t_{1,B} \\ t_{2,1} & t_{2,2} & \dots & t_{2,B} \\ \vdots & \vdots & \ddots & \vdots \\ t_{M,1} & t_{M,2} & \dots & t_{M,B} \end{bmatrix}$$

355 where $t_{\delta\omega} \in [0,100]$, $\delta \in \{1, \dots, M\}$, $\omega \in \{1, \dots, B\}$ and $\delta, \omega \in \mathbb{N}^+$

356 In order to generate phylogenetic trees illustrating the relationship between the subclones present in the
 357 tumor, which aberrations reside in the same cells as well as which subpopulations of cells the tumor
 358 consist of must be determined. To solve this, the idea is that a true subclone of cells should form a cluster

359 of unique genetic alterations that persist. They should remain grouped irrespective of inclusion of new
360 data from an additional region of the primary tumor or metastasis. Alterations that seem to follow each
361 other are more likely to be in the same cells. The first step is thus to yield a clustering to identify groups
362 of genetic alterations, uniquely identifying a certain subclone. The subsequent step is to determine the
363 temporal order of the alterations in question, since each subclone will represent a linear combination of
364 the clusters identified. Note that for DEVOLUTION the clustering is only used to reduce the
365 computational complexity for the upcoming subclonal deconvolution algorithm. Mostly alterations are
366 clustered that show similar MCF in all available biopsies. This is not to be confused with the more
367 intricate clustering methods used in for example SciClone.

368 Density based clustering techniques such as DBSCAN³³ are superior at unsupervised clustering of non-
369 uniform clusters. Furthermore, the number of clusters does not have to be specified beforehand, which
370 you have to do with many other established clustering algorithms. In addition, it does only have two
371 hyperparameters named *minPts*, which is the minimal number of points that is allowed in a cluster, and
372 ϵ representing the radius in which points i.e. the genetic alterations' position in the B-dimensional space,
373 are included, where B is the total number of biopsies. If ϵ is chosen too small a large part of the data
374 might not be clustered and choosing it too big will put all alterations in the same cluster. The choice of
375 ϵ can be aided by using a k-distance-graph which illustrates the distance to the $\text{minPts}-1 = k$ nearest
376 neighbor. The value to choose is when this plot shows an elbow which can be obtained by visual
377 inspection. Another method would be to create a vector between the start- and endpoint of the graph.
378 Then create a vector perpendicular to this having its start position at this line and its end position at our
379 data curve. The elbow can be computed by finding the vector with the largest magnitude. The clustering
380 method is not integrated directly in the subclonal deconvolution algorithm described below and is
381 therefore easy to replace or alter if needed for noisier data³³.

382 The algorithm provides a matrix containing all clusters of genetic alterations. Let $C_{K \times N}$ be the matrix
383 representing the clusters of genetic alterations. It has the dimensions $K \times N$ where K is the number of
384 genetic alterations in the cluster and N is the cluster number. All matrix positions $c_{kn} \neq 0$ are unique
385 i.e. the same genetic alteration cannot belong to multiple clusters.

386
$$C_{K \times N} = \begin{bmatrix} c_{1,1} & c_{1,2} & \dots & c_{1,N} \\ c_{2,1} & c_{2,2} & \dots & c_{2,N} \\ \vdots & \vdots & \ddots & \vdots \\ c_{K,1} & c_{K,2} & \dots & c_{K,N} \end{bmatrix}$$

387 where $c_{kn} \in m_\delta \wedge c_{kn} \neq c_{ed}, (\forall k, e \in \{1, \dots, K\} \& n, d \in \{1, \dots, N\} \wedge c_{kn} \neq 0)$

388 A matrix representing the clusters present in each biopsy and their size determined by the mean of the
389 aberrations in the cluster is also constructed.

390
$$Z_{C \times B} = \begin{bmatrix} z_{1,1} & z_{1,2} & \dots & z_{1,B} \\ z_{2,1} & z_{2,2} & \dots & z_{2,B} \\ \vdots & \vdots & \ddots & \vdots \\ z_{C,1} & z_{C,2} & \dots & z_{C,B} \end{bmatrix} \text{ where } z_{cb} \in [0,100]$$

391 Where c is a specific cluster of aberrations, b the biopsy and z_{cb} the size of the cluster c in sample b .

392

393 **Subclonal deconvolution based on information from multiple samples from the same patient**

394 The space of a single biopsy is 100 % and the space of all biopsies can thus be represented by a matrix
395 where p is the partitioning of the available space in the biopsy and s_{pb} is the space available in a specific
396 partitioning p in biopsy b . Initially $s_{1,b} = 100 \wedge s_{p \neq 1,b} = 0$.

397
$$S_{P \times B} = \begin{bmatrix} s_{1,1} & s_{1,2} & \dots & s_{1,B} \\ s_{2,1} & s_{2,2} & \dots & s_{2,B} \\ \vdots & \vdots & \ddots & \vdots \\ s_{P,1} & s_{P,2} & \dots & s_{P,B} \end{bmatrix} \text{ where } s_{pb} \in [0,100] \text{ and } \sum_{p=1}^P s_{p,b} = 100 \wedge b \in \{1, \dots, B\} \in \mathbb{N}^+$$

398 The clusters of aberrations in each biopsy, as supplied by $Z_{C \times B}$, are allocated to the space in decreasing
399 order, altering the magnitude of the spaces in $S_{P \times B}$ based on the MCF of the clusters allocated to it.

400 The allocation iteration algorithm is initially conducted considering each sample individually, resulting
401 in a matrix encompassing all possible allocations of the clusters in every biopsy. Subsequently, all
402 possible allocations throughout the samples are addressed to minimize the occurrence of parallel
403 evolution. The algorithm hence tries to produce one uniform solution of the temporal order of events
404 that does not contradict any information provided in the biopsies. The solution should be in concordance
405 to every biopsy provided. If not possible, parallel evolution or back mutations will occur in the final

406 phylogenetic tree and the user is advised to reconsider the original data set, since it may be biologically
407 unlikely. It may be possible to allocate a cluster to multiple positions without producing contradicting
408 temporal orders in any of the samples, for which the largest available space assumption is employed to
409 make an objective decision, based on the presumption that the mutational frequency is equal in all cells
410 within the biopsy no matter how many mutations they have. The cluster will consequently be placed as
411 a descendant to the cluster constituting the largest proportion of the biopsy, in the absence of further
412 biological information steering it elsewhere. Clusters presenting with only one possible allocation in a
413 biopsy provides especially valuable information concerning the temporal order of events, e.g. having a
414 group of alterations that are all present in all cells in a biopsy clearly indicates that there exists a group
415 of cells in the tumor harboring all of these alterations, aiding the temporal allocation in other samples
416 where these alterations may present themselves as subclonal. DEVOLUTION does allow some overlap
417 in cellular frequency in the allocation algorithm taking into consideration the measurement uncertainty
418 of MCF. This iterative computation results in the subclones present in the biopsies along with an
419 estimation of their size and distribution across the samples.

420

421 **Incorporating user-controlled rules for avoiding imposition of illicit biological trajectories**

422 Some genetic aberrations present in the data set might be known to never occur in the same cell for some
423 well-known biological reason. Such constraints should optimally be supplied to the algorithm to ensure
424 biologically plausible solutions. The user can therefore provide the DEVOLUTION algorithm with a
425 matrix indicating which genetic aberrations in the data set that cannot be placed after one another. The
426 first column represents a mother genetic alteration that the daughter alteration specified in the second
427 column, cannot have (**Supplementary Figure 7**). The subclonal deconvolution algorithm extracts a list
428 for each genetic alteration containing information about in how many of the samples it can be allocated
429 after a certain cluster. There might be multiple possible solutions, equally prevalent. In this instance the
430 matrix containing information about illicit biological orders can aid the program in taking a decision
431 regarding which of these allocations are less likely, subsequently discarding them. These rules will thus
432 only be employed if the data set allows the genetic alterations to be placed in any other way. If the only

433 possible way for the events to be allocated is to place them as descendants, the user will be advised to
434 revise the original data set. No such rules were integrated in the analysis of the 56 pediatric tumors in
435 the present study.

436

437 **Construction of an event matrix**

438 Based on the estimated subclonal composition, an event matrix $E = [\hat{a}_1, \hat{a}_2 \dots \hat{a}_k]$ was constructed
439 illustrating the distribution of genetic alterations across the identified subclones. Each \hat{a}_i is binary vector
440 belonging to subclone i . Each row represents a genetic alteration indicated with a 1 if present or 0 if
441 absent in the subclone. The event matrix is used as the foundation for phylogenetic tree generation,
442 illustrating the relationship between the subclones within the tumor.

443

444 **Reconstruction of phylogenetic trees**

445 In order to generate the phylogenetic trees, the genetic distances between all the subclones must be
446 computed. Here, the Hamming distance was used to assess a distance matrix displaying the genetic
447 distance between each of the subclones. It computes the distance between two vectors by adding all
448 positions in which they differ from one another. Note that it does not consider the magnitude of the
449 entity in the bins compared, which does not matter in this case since they are binary, hence the branch
450 lengths in the phylogenetic trees will be in units of number of aberrations. No bias in the estimation of
451 branch lengths were seen for ML in this study.

452 Using the function *stem*, a column is added in which the tree will be rooted. The stem can be chosen as
453 a cell containing all mutations shared between the subclones or a cell containing no alterations thus
454 representing a normal cell. The event matrix is then transformed into phyDat format using the function
455 *phydatevent*. This is the data class needed for phylogenetic analysis using the R package phangorn³⁴.

456 In the next step the maximum likelihood and maximum parsimony algorithms were used in order to
457 reconstruct phylogenetic trees based on the event matrices. The maximum likelihood trees were

458 reconstructed using the pml algorithm in the package phangorn³⁴. First a Hamming distance matrix was
459 calculated from the event matrix which was used to obtain an initial tree given by the neighbor joining
460 method. The initial tree as well as the initial event matrix was used as input variables in the pml
461 algorithm. This function returns an object containing the tree parameters, the data as well as the
462 likelihood for that phylogenetic tree. In order to optimize the tree parameters further, the function
463 optim.pml was used in combination with the Jukes Cantor model. The tree was subsequently rooted in
464 a constructed cell having all the events shared between all subclones. Since the model used is time
465 reversible the choice of the root does not influence the computed likelihood³⁵. The tree was visualized
466 using the ggtree package in R³⁶.

467 The maximum parsimony trees were constructed using the parsimony ratchet algorithm (pratchet) in the
468 R package phangorn with the Fitch algorithm³⁷. Using the acctran algorithm the branch lengths and the
469 ancestral character probability distributions were obtained. The trees were rooted in a cell containing no
470 alterations.

471

472 **Performance testing using simulated bulk sampling data**

473 In order to further assess the reliability of the algorithm, it was evaluated using simulated bulk
474 genotyping data using a basic 3D-lattice based, stochastic model of tumor growth. The simulation is
475 initiated imagining a cell having one single genetic aberration, representing a stem event, which will be
476 conveyed to all cells comprising the virtual tumor. In each time unit one cell can proliferate. When a
477 cell has been chosen for proliferation a certain inherent mutation frequency determines whether the cell
478 will mutate or not. If not, two cells identical to the mother cell are obtained, otherwise a stochastic
479 genetic copy number aberration algorithm is used to randomly select a genetic alteration (copy number
480 alteration), conjointly considering the chromosome boundaries and sizes. Making use of a random
481 number generator, a chromosome is randomly selected, then a start and end position and finally the type
482 of event. The result is one cell identical to the mother and one cell harboring one additional genetic
483 aberration. The spatial orientation of the cells is also considered where each cell is assumed to have 26

484 neighbors and two cells cannot occupy the same position. The position of the second daughter cell is
 485 randomly selected among the available neighbor positions. The simulation was conducted generating
 486 40,000 cells giving a simple 3D lattice structure illustrating the spatial intratumoral heterogeneity for
 487 three different mutation frequencies (**Figure 4a**). Note that only the fractions of cells harboring a certain
 488 genetic alteration is of importance in this model and not any absolute numbers. The goal is not to
 489 correctly simulate tumor growth but to obtain a random mixture of entities to be demixed using
 490 DEVOLUTION. The mutation frequency chosen, 10/40,000, 50/40,000 and 100/40,000, is not
 491 biologically accurate but chosen such that a certain number of mutations will appear during the
 492 simulation, resulting in segment files resembling the MCF-distributions seen in the patient cases.

493 Virtual biopsies were drawn from the set of simulated entities while varying the number of biopsies
 494 from 1 to 10. The biopsies were drawn randomly from different parts of the simulated entities. Values
 495 for x, y and the z coordinates were randomly chosen while fulfilling

$$496 \quad x \in [x_{min}, x_{max}], y \in [y_{min}, y_{max}], z \in [z_{min}, z_{max}]$$

$$497 \quad \sqrt{x^2 + y^2 + z^2} \in \left[\frac{r_{mean}}{2}, r_{mean} \right] \in \mathbb{R}$$

498 where

$$499 \quad r_{mean} = \frac{x_{max} + x_{min} + y_{max} + y_{min} + z_{max} + z_{min}}{6}$$

500 Based on the position chosen, the entities within a radius of 2 units were extracted and the fraction of
 501 cells harboring each of the alterations found was calculated. The data from each cell can be used to
 502 create artificial bulk sampling with MCFs as well as single cell data. Hence, the segment files can be
 503 analyzed while having the true subclones at hand for comparison. DEVOLUTION was evaluated using
 504 1-10 biopsies for the three tumors.

505 **Statistics and reproducibility**

506 For phylogenetic tree characteristics, significance was tested using the Mann-Whitney U-test (two-
507 tailed). The stem lengths, total branch lengths and number of branches between individual patients are
508 assumed to be independent of each other and not normal distributed.

509 All simulation points were repeated at least three times and the standard deviation are illustrated with
510 bars in Figure 4 c and d.

511 **Data availability**

512 All data generated or analyzed during this study are included in this article as supplementary files 1-4.

513 **Code availability**

514 The code is freely available and relies on R 4.0.2 or later. Setup instructions and dependencies can be
515 found on github.

516 <https://github.com/NatalieKAndersson/DEVOLUTION>

517

518 **Author contributions**

519 N.A. and D.G. conceived and designed the project, N.A., D.G. and S.C. developed the methodology,
520 A.V., J.K. and D.G. did tumor biopsy data acquisition, N.A. and D.G. analyzed and interpreted the data,
521 N.A., D.G., S.C., A.V. and J.K. contributed towards the manuscript.

522 **Competing interests**

523 The authors declare no competing interests.

524 **Acknowledgements**

525 We would like to thank the Swegene Centre for Integrative Biology at Lund University (SCIBLU) for
526 assistance.

527 **References**

- 528 1 Nowell, P. C. The clonal evolution of tumor cell populations. *Science* **194**, 23-28,
529 doi:10.1126/science.959840 (1976).
- 530 2 Quail, D. F. & Joyce, J. A. Microenvironmental regulation of tumor progression and metastasis.
531 *Nat Med* **19**, 1423-1437, doi:10.1038/nm.3394 (2013).
- 532 3 Komarova, N. L., Burger, J. A. & Wodarz, D. Evolution of ibrutinib resistance in chronic
533 lymphocytic leukemia (CLL). *Proc Natl Acad Sci U S A* **111**, 13906-13911, doi:10.1073/pnas.1409362111 (2014).
- 534 4 Misale, S. *et al.* Emergence of KRAS mutations and acquired resistance to anti-EGFR therapy in
535 colorectal cancer. *Nature* **486**, 532-536, doi:10.1038/nature11156 (2012).
- 536 5 Leder, K. *et al.* Fitness conferred by BCR-ABL kinase domain mutations determines the risk of
537 pre-existing resistance in chronic myeloid leukemia. *PloS one* **6**, e27682, doi:10.1371/journal.pone.0027682
538 (2011).
- 539 6 Cresswell, G. D. *et al.* Intra-Tumor Genetic Heterogeneity in Wilms Tumor: Clonal Evolution and
540 Clinical Implications. *EBioMedicine* **9**, 120-129, doi:10.1016/j.ebiom.2016.05.029 (2016).
- 541 7 Andor, N. *et al.* Pan-cancer analysis of the extent and consequences of intratumor heterogeneity.
542 *Nat Med* **22**, 105-113, doi:10.1038/nm.3984 (2016).
- 543 8 Mengelbier, L. H. *et al.* Intratumoral genome diversity parallels progression and predicts outcome
544 in pediatric cancer. *Nat Commun* **6**, 6125, doi:10.1038/ncomms7125 (2015).
- 545 9 Martelotto, L. G., Ng, C. K. Y., Piscuoglio, S., Weigelt, B. & Reis-Filho, J. S. Breast cancer intra-
546 tumor heterogeneity. *Breast Cancer Research : BCR* **16**, 210, doi:10.1186/bcr3658 (2014).
- 547 10 Villamon, E. *et al.* Genetic instability and intratumoral heterogeneity in neuroblastoma with
548 MYCN amplification plus 11q deletion. *PloS one* **8**, e53740, doi:10.1371/journal.pone.0053740 (2013).
- 549 11 Williams, M. J., Werner, B., Barnes, C. P., Graham, T. A. & Sottoriva, A. Identification of neutral
550 tumor evolution across cancer types. *Nat Genet* **48**, 238-244, doi:10.1038/ng.3489 (2016).
- 551 12 Gerlinger, M. *et al.* Intratumor heterogeneity and branched evolution revealed by multiregion
552 sequencing. *N Engl J Med* **366**, 883-892, doi:10.1056/NEJMoa1113205 (2012).
- 553 13 Sottoriva, A. *et al.* A Big Bang model of human colorectal tumor growth. *Nat Genet* **47**, 209-216,
554 doi:10.1038/ng.3214 (2015).
- 555 14 Andersson, N. *et al.* Extensive Clonal Branching Shapes the Evolutionary History of High-Risk
556 Pediatric Cancers. *Cancer research* **80**, 1512-1523, doi:10.1158/0008-5472.CAN-19-3468 (2020).
- 557 15 Alves, J. M., Prieto, T. & Posada, D. Multiregional Tumor Trees Are Not Phylogenies. *Trends*
558 *Cancer* **3**, 546-550, doi:10.1016/j.trecan.2017.06.004 (2017).
- 559 16 Mallory, X. F., Edrisi, M., Navin, N. & Nakhleh, L. Methods for copy number aberration detection
560 from single-cell DNA-sequencing data. *Genome biology* **21**, 208, doi:10.1186/s13059-020-02119-8 (2020).
- 561 17 Karlsson, J. *et al.* Four evolutionary trajectories underlie genetic intratumoral variation in
562 childhood cancer. *Nat Genet* **50**, 944-950, doi:10.1038/s41588-018-0131-y (2018).
- 563 18 Deshwar, A. G. *et al.* PhyloWGS: reconstructing subclonal composition and evolution from
564 whole-genome sequencing of tumors. *Genome biology* **16**, 35, doi:10.1186/s13059-015-0602-8 (2015).
- 565 19 Roth, A. *et al.* PyClone: statistical inference of clonal population structure in cancer. *Nat Methods*
566 **11**, 396-398, doi:10.1038/nmeth.2883 (2014).
- 567 20 Miller, C. A. *et al.* SciClone: inferring clonal architecture and tracking the spatial and temporal
568 patterns of tumor evolution. *PLoS Comput Biol* **10**, e1003665, doi:10.1371/journal.pcbi.1003665 (2014).

569 21 Strino, F., Parisi, F., Micsinai, M. & Kluger, Y. TrAp: a tree approach for fingerprinting subclonal
570 tumor composition. *Nucleic Acids Res* **41**, e165, doi:10.1093/nar/gkt641 (2013).

571 22 Ben-David, U. & Amon, A. Context is everything: aneuploidy in cancer. *Nat Rev Genet* **21**, 44-
572 62, doi:10.1038/s41576-019-0171-x (2020).

573 23 Vogelstein, B. *et al.* Cancer genome landscapes. *Science* **339**, 1546-1558,
574 doi:10.1126/science.1235122 (2013).

575 24 Heim, S. & Mitelman, F. *Cancer cytogenetics: chromosomal and molecular genetic aberrations*
576 *of tumor cells*. (John Wiley & Sons, 2015).

577 25 Hamming, R. W. Error detecting and error correcting codes. *Bell Labs Technical Journal* **29**, 147-
578 160 (1950).

579 26 Ha, G. *et al.* TITAN: inference of copy number architectures in clonal cell populations from tumor
580 whole-genome sequence data. *Genome Res* **24**, 1881-1893, doi:10.1101/gr.180281.114 (2014).

581 27 Oesper, L., Mahmoody, A. & Raphael, B. J. THetA: inferring intra-tumor heterogeneity from
582 high-throughput DNA sequencing data. *Genome biology* **14**, R80, doi:10.1186/gb-2013-14-7-r80 (2013).

583 28 Popic, V. *et al.* Fast and scalable inference of multi-sample cancer lineages. *Genome biology* **16**,
584 91, doi:10.1186/s13059-015-0647-8 (2015).

585 29 Niknafs, N., Beleva-Guthrie, V., Naiman, D. Q. & Karchin, R. SubClonal Hierarchy Inference
586 from Somatic Mutations: Automatic Reconstruction of Cancer Evolutionary Trees from Multi-region Next
587 Generation Sequencing. *PLoS Comput Biol* **11**, e1004416, doi:10.1371/journal.pcbi.1004416 (2015).

588 30 Zare, H. *et al.* Inferring clonal composition from multiple sections of a breast cancer. *PLoS*
589 *Comput Biol* **10**, e1003703, doi:10.1371/journal.pcbi.1003703 (2014).

590 31 Maley, C. C. *et al.* Classifying the evolutionary and ecological features of neoplasms. *Nat Rev*
591 *Cancer* **17**, 605-619, doi:10.1038/nrc.2017.69 (2017).

592 32 Turajlic, S., Sottoriva, A., Graham, T. & Swanton, C. Resolving genetic heterogeneity in cancer.
593 *Nat Rev Genet* **20**, 404-416, doi:10.1038/s41576-019-0114-6 (2019).

594 33 Ester, M., Kriegel, H.-P., Sander, J. & Xu, X. in *Proceedings of the Second International*
595 *Conference on Knowledge Discovery and Data Mining* 226–231 (AAAI Press, Portland, Oregon, 1996).

596 34 Schliep, K. P. phangorn: phylogenetic analysis in R. *Bioinformatics* **27**, 592-593,
597 doi:10.1093/bioinformatics/btq706 (2011).

598 35 Felsenstein, J. Evolutionary trees from DNA sequences: a maximum likelihood approach. *J Mol*
599 *Evol* **17**, 368-376, doi:10.1007/BF01734359 (1981).

600 36 Yu, G., Smith, D. K., Zhu, H., Guan, Y. & Lam, T. T. Y. ggtree: an R package for visualization
601 and annotation of phylogenetic trees with their covariates and other associated data. *Methods in Ecology and*
602 *Evolution* **8**, 28-36 (2017).

603 37 Nixon, K. C. The parsimony ratchet, a new method for rapid parsimony analysis. *Cladistics* **15**,
604 407-414 (1999).

605

606

607

608 **Supplementary datasets**

609 **S. Data 1 Childhood cancer data set.** All segment files for NB, WT and RMS along with the clustering
610 obtained by the software. Segments are annotated similar to Supplementary Figure 1.

611 **S. Data 2 The event matrices.** The files named “before” illustrates the allocation within each biopsy
612 while the files denoted “after” are the final event matrices.

613 **S. Data 3 Simulation results.** The segment file for each simulation, the biopsy positions randomly
614 chosen, which unique cells are present in these biopsies, the event matrix and evaluation of which
615 alterations are correctly allocated using the software.

616 **S. Data 4 Summary of all simulations.** Tables illustrating the number of alterations in total, alterations
617 with a clone size > 10 % and the number of alterations correctly allocated.

Figures

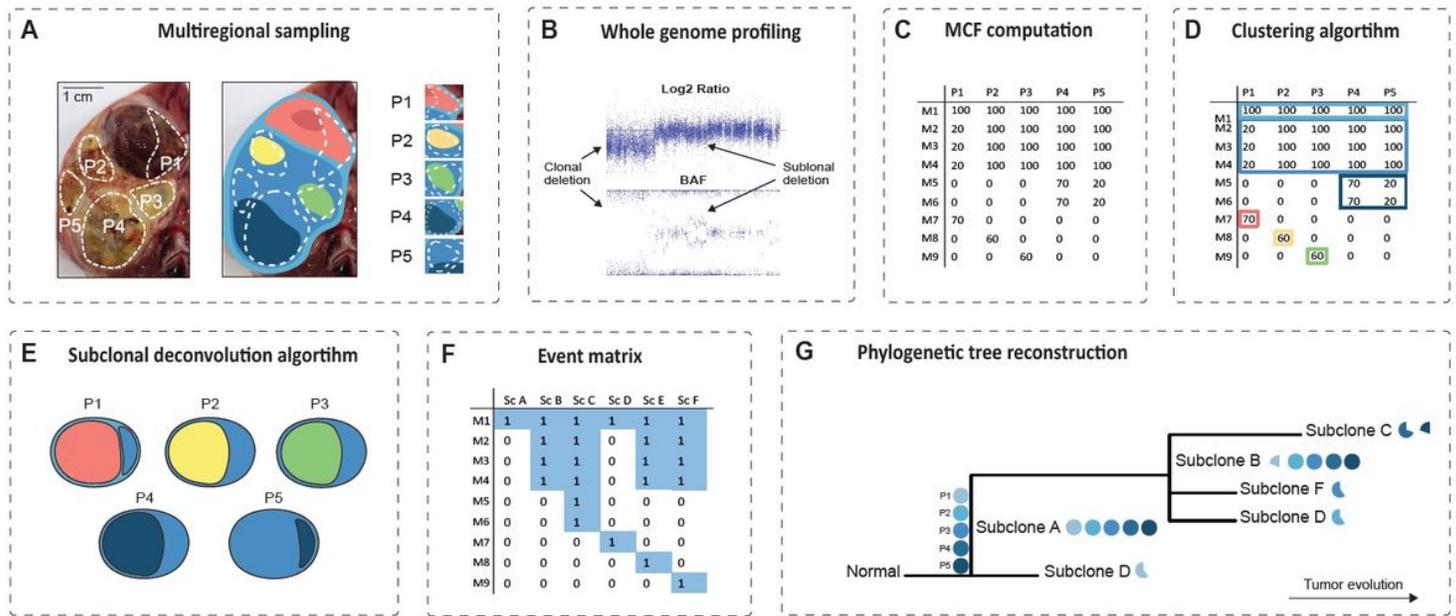


Figure 1

Overview of the methodological outline a) An example of multiregional sampling to obtain biopsies P1-P5 from a Wilms tumor. The tumor is composed of several subclones with distinct genomic profiles, exemplified by the schematic genomic landscape in the rightmost panel where each color unifies cells with identical sets of genetic alterations. The photograph is adapted from a previous publication¹⁴ and the colors do not represent true sets of genetic alterations. b) Acquisition of genomic tumor data for each sample, exemplified by copy number analysis by SNP-array. c) The whole genome profiling data can be used to compute the mutated clone fraction (MCF) illustrating the proportion of cells in each biopsy harboring a certain genetic alteration (M1-9 in the left column). d) A clustering algorithm is employed to identify genetic alterations that seem to follow each other in size across samples. e) A subclonal deconvolution algorithm determines the temporal order of these clusters by considering the information obtained throughout all samples while minimizing the occurrence of parallel evolution and back mutations. f) The proposed solution for the temporal order of subclones (Sc) is integrated into an event matrix. g) This event matrix can be used to generate phylogenetic trees with either the maximum likelihood or the parsimony method. At the stem, all available biopsies for the patient are visualized as filled circles with the biopsy name (P1-P5). At the branches of the tree the subclones can be seen along with pie charts illustrating in which biopsies, and in what fraction of the tumor cells they appear.

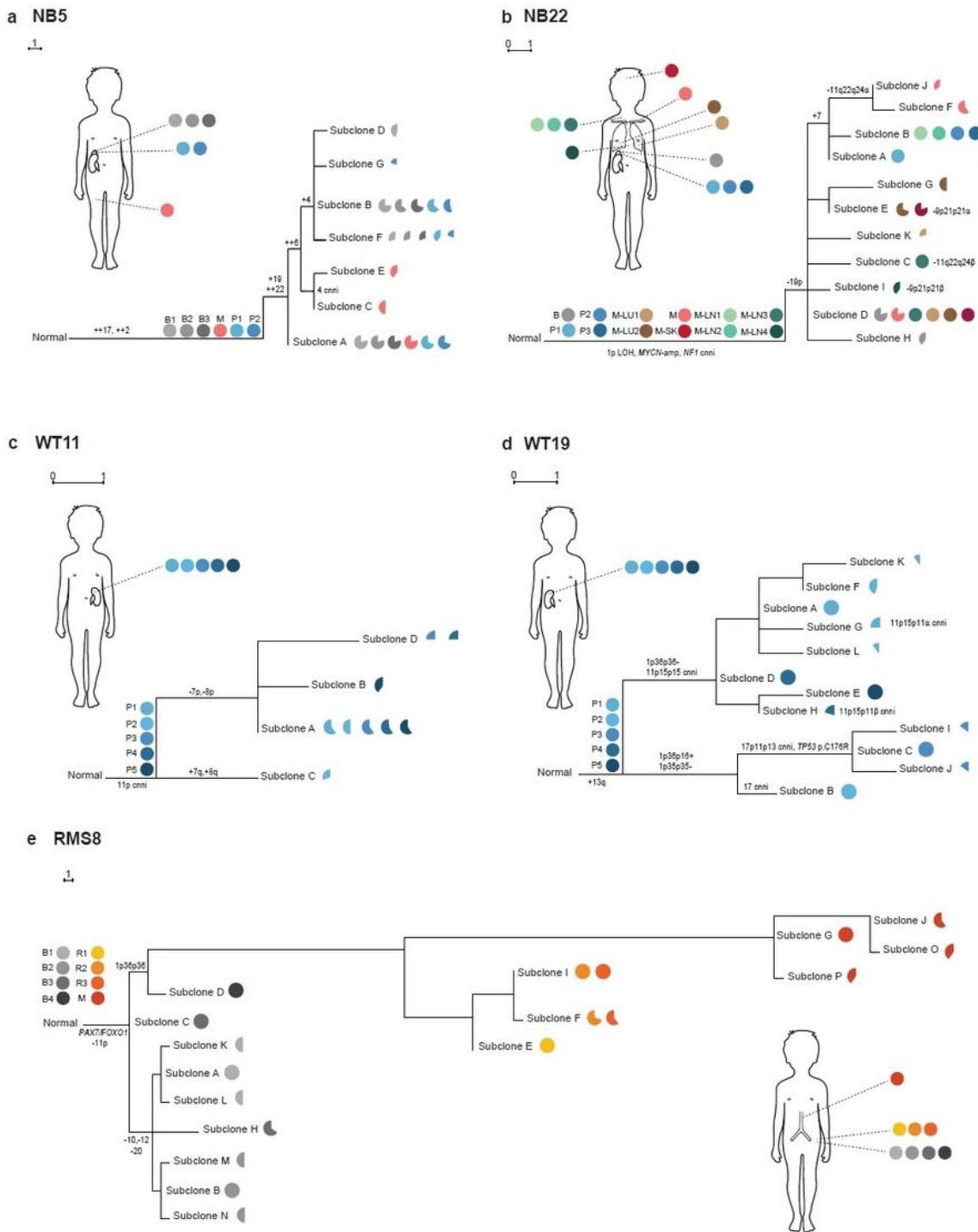


Figure 2

Phylogenetic trees of childhood cancers. At the stem of the maximum parsimony trees illustrated here, the biopsies available from each patient are denoted. The genetic alterations belonging to the stem are present in all cells in all samples, indicated by filled pies. The endpoints represent cell populations harboring distinct genomic profiles (subclones), whose fractions per sample are visualized by pie charts. The scale bar indicates the distance corresponding to one genetic aberration. Gains and losses of

chromosomes or segments of chromosomes that are characteristic of each tumor subtype are indicated by + and – signs. a) In NB5 samples are available from the primary tumor before treatment (B1-B3), a synchronous metastasis (M), and the primary tumor post treatment (P1-P2). The metastasis must have originated from a subclone harboring the stem events only and another subclone with the copy number profile seen in subclone A, indicating polyclonal seeding. The metastasis also has a private copy number neutral imbalance (cnni) of chromosome 4. b) NB22 also shows evidence of polyclonal seeding. Samples are from the primary tumor before treatment (B), the primary tumor post treatment (P1-P3), metastases to the lung (M-LU1-2), to the lymph nodes (M-LN1-4), the skull (M-SK) and from the area around the clavicle (M). The stem harbors a 1p cnni, MYCN-amplification and a NF1 deletion. Greek letters denote different structural alterations targeting partly overlapping regions. c) WT11, shows subclones present across multiple primary tumor areas (P1-P5). d) WT19 display a similar distribution of subclones as WT11 across the post treatment

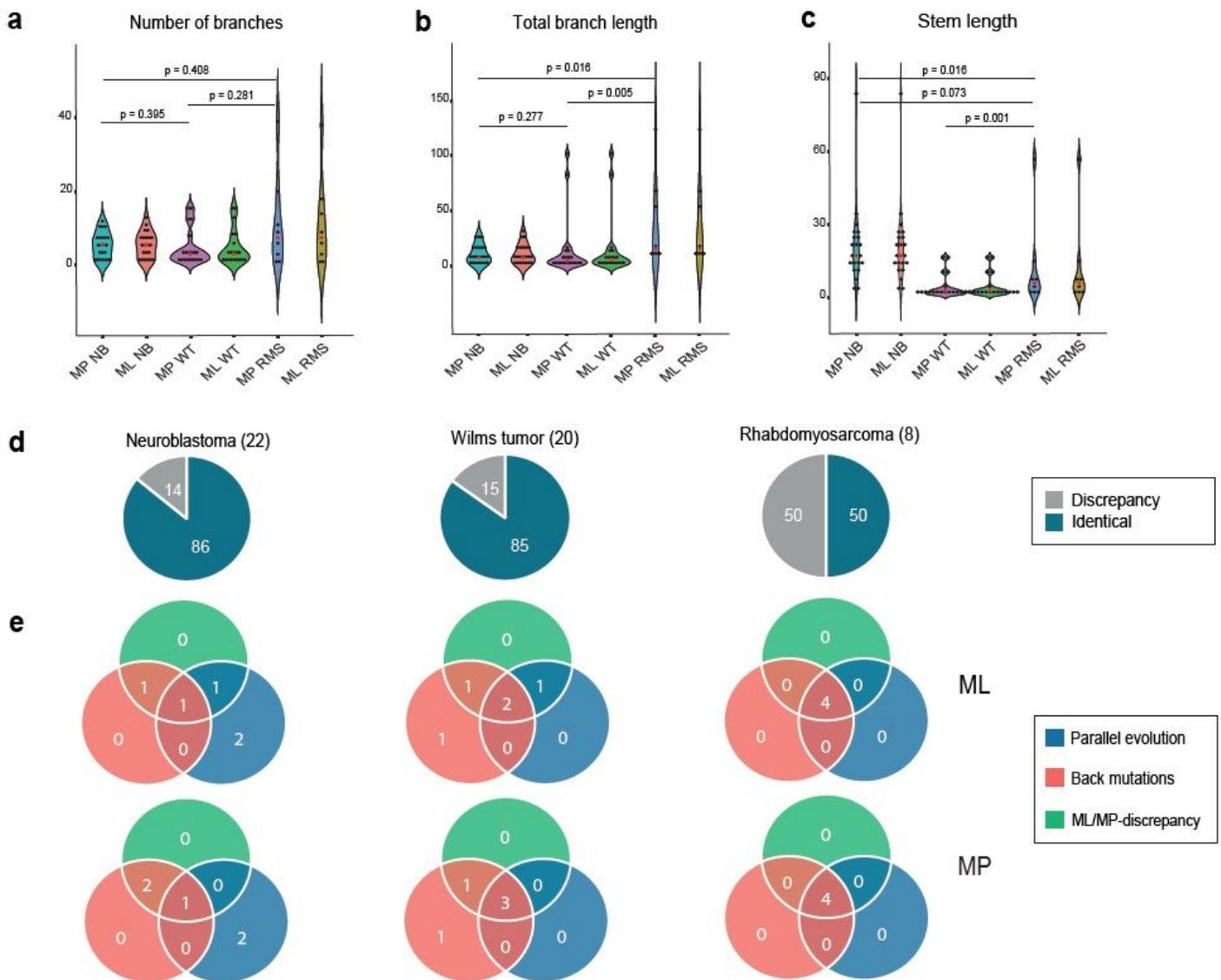


Figure 3

Structural properties of the generated phylogenetic trees. Violin plots of a) the number of branches, b) total branch lengths and c) total stem lengths for neuroblastoma (NB), Wilms tumor (WT) and rhabdomyosarcoma (RMS) using either the maximum likelihood (ML) or parsimony (MP) method for phylogenetic reconstruction. Significance represented by P-values were calculated using the two-sided Mann-Whitney U-test. d) In 14% of neuroblastomas, 15% of Wilms tumors and 50% of the 8 rhabdomyosarcomas analyzed, the phylogenetic trees obtained using the ML and MP methods differed from one another. NB1, NB24, WT1, WT2, WT3 and WT5 are not included in this calculation since they did not display any private genetic alterations. Hence only event matrices could be generated but not phylogenetic trees. If including these cases, the proportions would be 12.8% for NB and 12.5% for WT. e) Venn diagrams of how often discrepancies between ML/MP trees, back mutations, and parallel evolution occurred together in the same case/tree. Numbers indicate the number of tumors where a particular contradiction or combination of contradictions occurred.

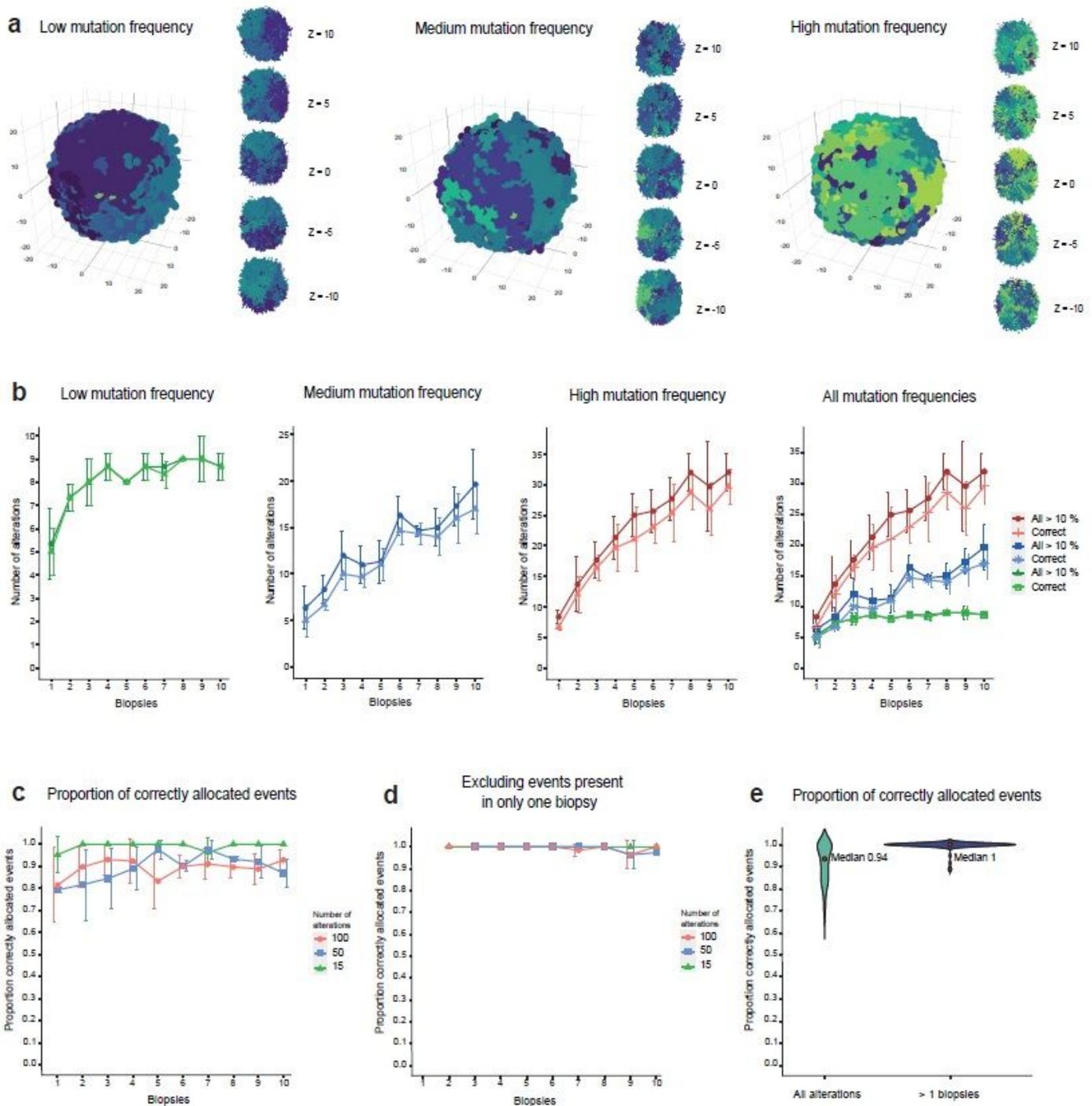


Figure 4

Properties of the simulated bulk genotyping data set and performance measure of DEVOLUTION. a) Visualization of three simulated tumors with increasing mutation frequency. To the right of each tumor are cross-sections at five positions (z). Each color represents a subclone harboring a unique genetic profile. b) When increasing the number of biopsies, more genetic alterations and hence subclones, are identified. In addition, the algorithm is also able to identify more subclones correctly. Each point in the graph is the mean of three consecutive measurements and the error bars consequently the standard deviation. c) The proportion of genetic alterations correctly allocated \pm SD when increasing the number of

biopsies from 1-10 for three different mutation frequencies resulting in 10 (green), 50 (blue) and 100 (red) genetic alterations present in the virtual tumor. Each point is represents the mean of three consecutive measurements. d) The proportion of correctly allocated genetic alterations when excluding genetic alterations that were only found in a single biopsy. e) Violin plot showing the spread of the proportion of correctly allocated genetic alterations when including all alterations and when excluding the genetic alterations only found in one biopsy.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [S.Figure12.pdf](#)
- [S.Figure3.pdf](#)
- [S.Figure3text.pdf](#)
- [S.Figure4.pdf](#)
- [S.Figure4text.pdf](#)
- [S.Figure5.pdf](#)
- [S.Figure6.pdf](#)
- [S.Figure7.pdf](#)