

Clonal assessment of functional mutations in cancer based on a genotype-aware method for clonal reconstruction

Paul Deveau¹⁻³, Leo Colmet Daage², Derek Oldridge⁴⁻⁶, Virginie Bernard⁷, Angela Bellini², Mathieu Chicard², Nathalie Clement², Eve Lapouble⁸, Valérie Combaret⁹, Anne Boland¹⁰, Vincent Meyer¹⁰, Jean-François Deleuze¹⁰, Isabelle Janoueix-Lerosey¹¹, Emmanuel Barillot¹, Olivier Delattre¹¹, John Maris⁴⁻⁶, Gudrun Schleiermacher^{2,12,†,*} and Valentina Boeva^{1,13-16,†,*}

¹Institut Curie, PSL Research University, Mines Paris Tech, INSERM U900, 75005, Paris, France

²Institut Curie, PSL Research University, INSERM U830, Laboratoire RTOP (Recherche Translationnelle en Oncologie Pédiatrique), Département de recherche translationnelle, 75005, Paris, France

³Univ. Paris-Sud, Orsay, France

⁴Division of Oncology, Children's Hospital of Philadelphia, Philadelphia, Pennsylvania, USA

⁵Center for Childhood Cancer Research, Children's Hospital of Philadelphia, Philadelphia, Pennsylvania, USA

⁶Department of Pediatrics, Perelman School of Medicine at the University of Pennsylvania, Philadelphia, Pennsylvania, USA

⁷Institut Curie, PSL Research University, ICGex, 75005, Paris, France

⁸Institut Curie, PSL Research University, Unité de Génétique Somatique, 75005, Paris, France

⁹Centre Léon-Bérard, Laboratoire de Recherche Translationnelle Lyon, France

¹⁰Centre National de Génotypage, Institut de Génomique, CEA, Evry, 91057, France.

¹¹Institut Curie, PSL Research University, INSERM U830, Paris, 75005, France

¹²Institut Curie, PSL Research University, Département de Pédiatrie, Paris, 75005, France

¹³Institut Cochin, Paris, France

¹⁴INSERM U1016, Paris, France

¹⁵CNRS UMR 8104, Paris, France

¹⁶Université Paris Descartes UMR-S1016, 75014, Paris, France

†These authors jointly supervised this work.

*Correspondance should be addressed to valentina.boeva@inserm.fr or gudrun.schleiermacher@curie.fr

Running title: Framework for clonal reconstruction in cancer

Keywords: Clonal inference, Cancer, Algorithms, Neuroblastoma, Whole genome sequencing

Abstract

In cancer, clonal evolution is characterized based on single nucleotide variants and copy number alterations. Nonetheless, previous methods failed to combine information from both sources to accurately reconstruct clonal populations in a given tumor sample or in a set of tumor samples coming from the same patient. Moreover, previous methods accepted as input all variants predicted by variant-callers, regardless of differences in dispersion of variant allele frequencies (VAFs) due to uneven depth of coverage and possible presence of strand bias, prohibiting accurate inference of clonal architecture. We present a general framework for assignment of functional mutations to specific cancer clones, which is based on distinction between passenger variants with expected low dispersion of VAF versus putative functional variants, which may not be used for the reconstruction of cancer clonal architecture but can be assigned to inferred clones at the final stage. The key element of our framework is QuantumClone, a method to cluster variants into clones, which we have thoroughly tested on simulated data. QuantumClone takes into account VAFs and genotypes of corresponding regions together with information about normal cell contamination. We applied our framework to whole genome sequencing data for 19 neuroblastoma trios each including constitutional, diagnosis and relapse samples. We discovered specific pathways recurrently altered by deleterious mutations in different clonal populations. Some such pathways were previously reported (e.g., MAPK and neuritogenesis) while some were novel (e.g., epithelial–mesenchymal transition, cell survival and DNA repair). Most pathways and their modules had more mutations at relapse compared to diagnosis.

1 Introduction

2 The principal cause of cancer is believed to be accumulation of mutations and structural variations (SVs)
3 of the genome. Recently, many efforts have focused on the identification of driver mutations; nonetheless,
4 passenger variants, although they are not directly linked to the disease, may provide additional evidence from
5 which to infer the phylogeny of a tumor and so help uncover the basis for its proliferative activity (Marusyk
6 et al., 2014).

7
8 To understand the role driver mutations play in clonal expansion and cancer progression, it is essential
9 to accurately reconstruct the clonal structure and assign functional variants to it. We define a clone as a
10 cell population that harbors a unique pattern of mutations and SVs. Clones are related to each other and
11 share a common ancestor. A hierarchical phylogenetic tree, which represents the ancestry of clones, can
12 be constructed to reflect the order of appearance of new sets of mutations defining each clone. Each such
13 set of mutations is expected to contain at least one driver mutation or SV giving a selective advantage to
14 the clone compared to its ancestry. A clone can thus have a different behavior from its ancestral clone
15 when facing the same stimuli. With accumulation of driver mutations, clones are likely to gain hallmarks of
16 cancer such as evading growth suppressors, activating invasion and metastasis (Hanahan and Weinberg, 2011).

17
18 High-Throughput Sequencing (HTS) of bulk tumor tissues has allowed uncovering genetic differences at
19 the clonal level in primary and relapse/metastatic tumors. Modern computational methods provide ways to
20 reconstruct the structure of the phylogenetic tree from variant allele frequencies (VAFs) in sequenced reads,
21 where VAF is a proportion of reads supporting each given variant among all reads spanning the position of
22 interest (Fischer et al., 2014; Jiao et al., 2014; Kepler, 2013; Malikić et al., 2015; Miller et al., 2014; Qiao et al.,
23 2014; Schwarz et al., 2014). However, existing methods for clonal reconstruction often neglect information
24 about the genotype of each position, which refers to the paternal or maternal inheritance of a locus and the
25 number of copies of each allele. Accounting for the genotype information is especially crucial in the case of
26 hyper-diploid cancers and cancers with highly rearranged genomes, as the cellular prevalence – measured as
27 the proportion of cancer cells carrying a variant – is linked to VAF through such parameters as copy number
28 of the locus and the number of chromosome bearing the mutation.

29
30 Here we show that by combining the genotype and VAF information it is possible to correctly cluster
31 variants and assign them to specific clones, thus reconstructing the clonal architecture of an individual can-
32 cer. This may be done with our novel method, QuantumClone, designed to reconstruct clones based on both

33 VAF and genotype information. We demonstrate that our algorithm accurately clusters variants on simu-
34 lated data, even when cancer is hyper-diploid or contaminated by normal cells. We also propose a general
35 framework based on QuantumClone to detect driver mutations of clonal evolution. This general approach is
36 applied to 19 neuroblastoma cases; each case includes whole genome sequencing (WGS) data from a sample
37 at diagnosis and relapse. We show that deleterious mutations in neuroblastoma accumulate at relapse in
38 specific pathways such as cell motility (e.g., cell-matrix adhesion and regulation of epithelial–mesenchymal
39 transition, EMT) and cell survival (e.g., PI3K/AKT/mTOR, MAPK or noncanonical Wnt pathways).

40

41 Results

42 The QuantumClone method presented here applies an expectation-maximization (EM) algorithm and allows
43 for accurate inference of clonal structure using VAFs from one or several tumor samples sequenced using
44 WGS. It can analyze variants coming from highly rearranged and hyper-diploid cancer genomes. We exten-
45 sively validated QuantumClone on simulated data, where we compared it with recently published methods
46 (Miller et al., 2014; Roth et al., 2014). We complement QuantumClone with a robust framework for the
47 functional assessment of mutations based on signaling pathway analysis combined with the clonal assignment
48 (Fig. 1).

49

50 The overall framework was applied to WGS neuroblastoma datasets: 19 patients’ primary and relapse
51 samples including 7 new triplets. Novel and previously published samples (Eleveld et al., 2015) have been
52 sequenced at $\sim 100\times$ depth of coverage using Illumina HiSeq 2500 and Complete Genomics sequencing tech-
53 nologies. Application of the QuantumClone-based framework allowed us to discover pathways recurrently
54 altered by mutations in neuroblastoma at diagnosis and relapse.

55

56 Assessment of clonal reconstruction accuracy by QuantumClone

57 For clonal reconstruction using VAFs, we developed an approach that applies an EM algorithm (Fig. 2A,
58 Methods). QuantumClone utilizes genotype information and assigns variants to clones providing the most
59 likely values of cellular prevalence (Fig. 2A, Methods).

60

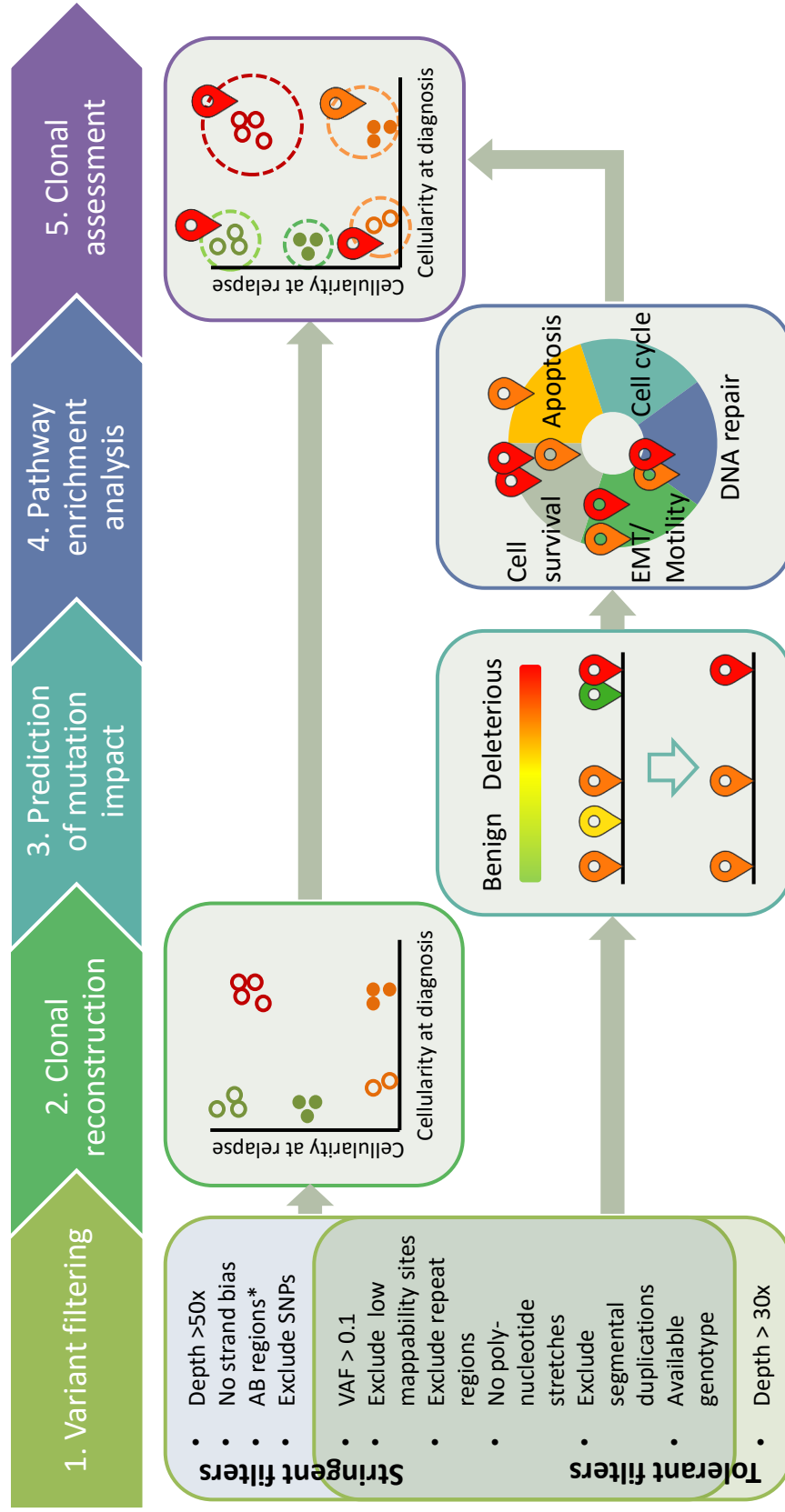


Figure 1: Overview of the general clonal reconstruction workflow: steps 1-5. (1) Variants are filtered to remove false positive calls; stringent filters are used to produce mutations that are further employed for clonal reconstruction (step 2), tolerant filters are used to detect functional mutations (step 3-4). (2) Variants that pass stringent filters and have genotype information assigned to the corresponding genomic loci are used as input to QuantumClone to reconstruct clonal populations. (3) Functional impact of variants passing tolerant filters is assessed. (4) Pathways recurrently altered by deleterious mutations are identified. (5) Finally, possibly damaging mutations belonging to frequently altered pathways are mapped to the reconstructed clones. (*) Stringent filtering keeps mutations located in AB regions only when at least 100 mutations pass this filter.

61 Comparison of QuantumClone with existing methods

62 Using *in silico* data, we compared the performance of QuantumClone, sciClone (Miller et al., 2014), pyClone
63 (Roth et al., 2014) and a generic k -medoids clustering algorithm in inferring clonal structure of a set of
64 tumors derived from the same patient. sciClone is based on variational Bayesian Mixture Models, while py-
65 Clone relies on a hierarchical Bayes statistical model. Partitioning with k -medoids is a more robust version
66 of a widely used k -means clustering algorithm (Kaufman and Rousseeuw, 1987).

67
68 In our simulation experiments, the following parameters were varied within realistic ranges: depth of
69 sequencing ($50 \times$ to $1000 \times$), fraction of contamination by normal cells (from 0 to 50%), number of variants
70 used for the clonal reconstruction (from 50 to 200), number of tumor samples used for each patient (from
71 2 to 5) and number of distinct clones per cancer (from 2 to 9) (Fig. 2B). For each set of parameters, we
72 performed and analyzed 50 independent simulation experiments (Methods). The accuracy of clonal recon-
73 struction was assessed by evaluation of the normalized mutual information (NMI) (Manning et al., 2008). A
74 perfect mutation clustering would result in a NMI value of 1, which corresponds to an identification of the
75 exact number of clones and correct assignment of all the mutations of a clone to the same cluster.

76
77 Our analysis showed that QuantumClone surpasses both published algorithms in experiments with chal-
78 lenging parameter settings: high contamination by normal cells, moderate depth of sequencing or high tumor
79 heterogeneity (Fig. 2B). Indeed, in samples with 50% contamination by normal cells QuantumClone signif-
80 icantly outperformed sciClone and pyClone (p -value = 2.6×10^{-4} and p -value = 3.2×10^{-14} , Welch
81 two sample t-test). While at high values of sequencing depth, all methods provided accurate results, at
82 depth of sequencing of $50 \times$ QuantumClone consistently gave better predictions (p -value = 1.2×10^{-3} and
83 p -value = 7.3×10^{-10} for sciClone and pyClone respectively). In addition, compared to the other meth-
84 ods, QuantumClone took the best advantage of data when multiple samples were provided for the analysis
85 (p -value = 9.4×10^{-3} and p -value = 8.0×10^{-4} for sciClone and pyClone respectively, for simulated tumors
86 with five samples). Also, starting from six clones per tumor, QuantumClone demonstrated significantly better
87 clonal reconstitution accuracy than the other methods (p -value = 4.3×10^{-8} and p -value < 2.2×10^{-16}
88 for sciClone and pyClone respectively).

89

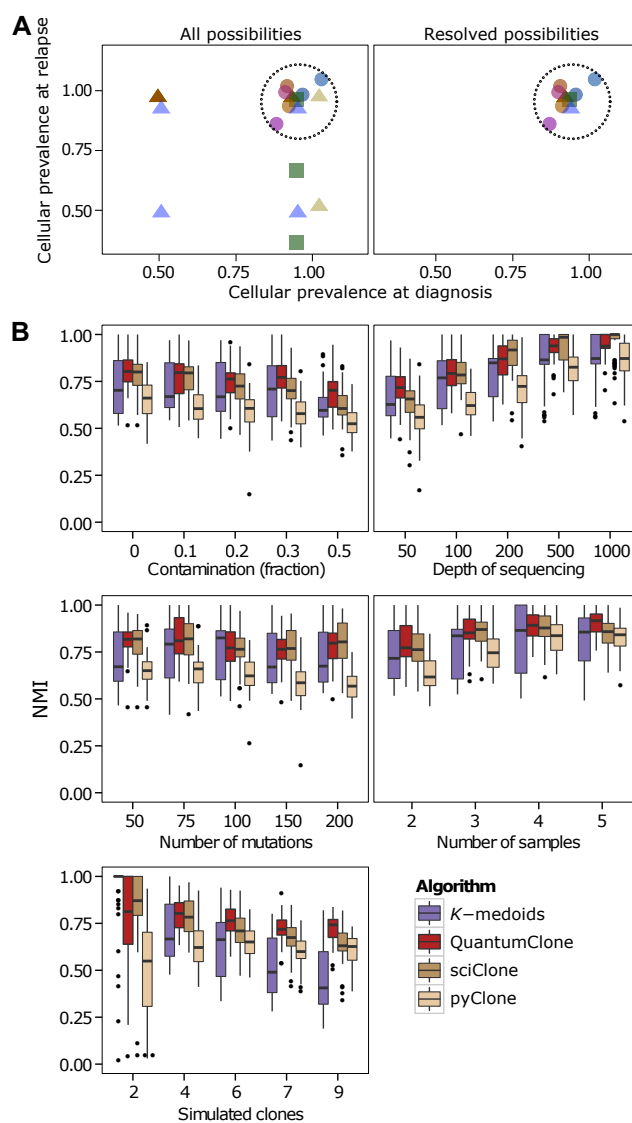


Figure 2: Principle of the QuantumClone algorithm and comparison to published methods. (A) Mutations located in regions of copy number aberration can be present on several chromosomal copies; they can thus be assigned to several cellular prevalence values (left panel). After the Expectation Maximization (EM) step each mutation is attributed the most likely cellular prevalence value (right panel). Each mutation is represented by a specific color. Mutations located in AB regions (circles); mutations located at relapse in regions of gain (squares), mutations located in regions of gain both at diagnosis and relapse (triangles). **(B)** Comparison of QuantumClone to existing methods. Normalized Mutual Information (NMI) is used to assess the quality of clustering on simulated data, with a single parameter varying in each test. QuantumClone (red) shows better performance in difficult settings, i.e., in presence of a high number of clones, low number of input samples, low sequencing depth, and high fraction of contamination by normal cells. Default parameters: two tumor samples without contamination sequenced at 100x; 4 clones; 100 mutations used for clustering.

90 **Assessment of clonal reconstruction accuracy in hyper-diploid cancers or cancers with highly**
91 **rearranged genomes**

92 We expect that in addition to the parameters discussed above, the degree of genome rearrangement and
93 chromosome duplication significantly affects the quality of the mutation clustering and consecutive clonal
94 reconstruction. Indeed, given an observed VAF value, a mutation occurring in a high copy number locus has
95 more possibilities for values of cellular prevalence: a mutation with an observed allele frequency of 25% can
96 only be linked to a cellular prevalence of 50% in a AB locus, while it can arise from cellular prevalence values
97 of 33.3%, 50% or 100% if the genotype is AAAB (Methods).

98

99 In order to validate QuantumClone on diploid and hyper-diploid genomes, we simulated variants in loci
100 of genotype AB, AAB and AABB (Fig. 3). In addition to QuantumClone, we tested the performance of
101 the k -medoids clustering algorithm, and two alternative versions of QuantumClone: QuantumClone-Single
102 and QuantumClone-Alpha. QuantumClone-Single assigned variants to a single copy state, e.g., variants in
103 AAB regions were supposed to only occur on a single chromosome. QuantumClone-Alpha used the same EM
104 algorithm as the default version of QuantumClone but added additional weights to probabilities based on
105 the locus genotype (Methods); e.g., this method suggested that a mutation in an AAB region has 3 times
106 more chances to occur on a single chromosome than on two chromosomes out of three. The default version
107 of QuantumClone assigned equal weights to all possibilities (Methods).

108

109 In all types of regions, QuantumClone and its alternative methods performed better than the baseline
110 k -medoids clustering algorithm (Fig. 3). For AB regions, where VAF of each mutation corresponds to a single
111 possible value of cellular prevalence, the difference in performance between QuantumClone, QuantumClone-
112 Single and QuantumClone-Alpha was not significant; in fact, it was due to random initialization of the EM
113 algorithm.

114

115 In the case of a triploid genome when each mutation was simulated to occur in a single chromosome
116 copy (case AAB*, Fig. 3), QuantumClone provided equally good results as QuantumClone-Single (t-test
117 p -value = 0.63). This validated our EM strategy to automatically select the number of chromosomal copies
118 with a mutation.

119

120 We demonstrated that in a more realistic case when a mutation can have a single or multiple copy status
121 (AAB and AABB regions), QuantumClone performed better than the three other methods. This validated

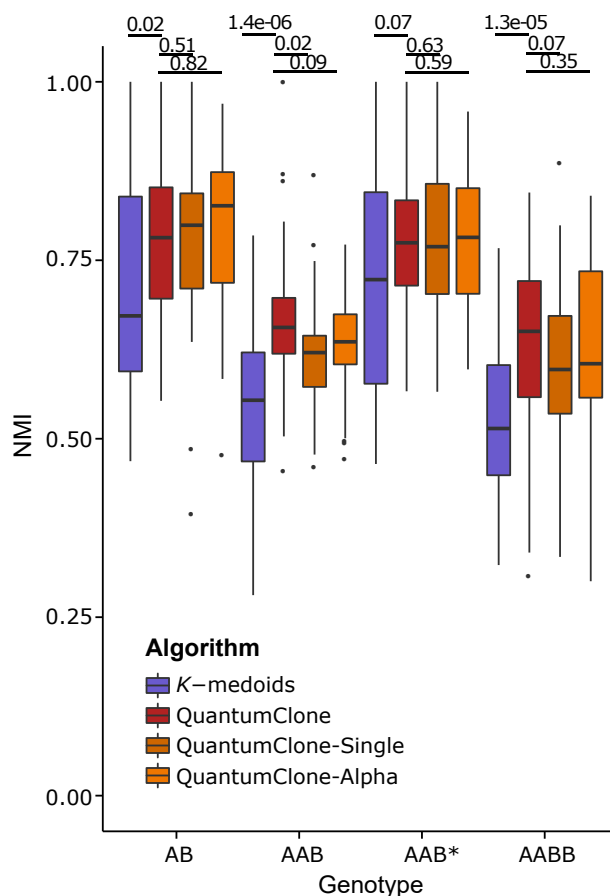


Figure 3: **Quality of clonal reconstruction for mutations located in regions of altered copy number.** QuantumClone-Single assumes that the mutation is always present at a single chromosomal copy; QuantumClone-Alpha uses an alternative algorithm for the selection of the best cellular prevalence given VAF values (Methods). Our comparison shows that the default version of QuantumClone performs as good as QuantumClone-Single in single copy regions but is the best algorithm in the situation when there are multiple possibilities of cellular prevalence for a mutation. P-values are calculated using Welch's t-test. (*) Mutations in AAB* regions are always present at a single copy.

122 our computational strategy for hyper-diploid cancers.

123

124 Overall, validation on simulated data showed that (1) QuantumClone can be applied to cancer samples
125 with hyperploid or rearranged genomes and (2) QuantumClone performs generally better than its peers in
126 difficult settings, e.g., in experiments with low depth of sequencing, when the number of clones is higher than
127 or equal to six, or when the contamination by normal cells is higher than or equal to 0.3.

128

129 **Effect of experimental settings on the clonal reconstruction accuracy**

130 Our analysis allowed us not only to compare QuantumClone to published methods and define the limits of
131 applicability of each method but also to study the effect of experimental settings on the clonal reconstruction
132 accuracy, which can help in the planning of tumor DNA sequencing experiments.

133

134 *Effect of contamination by normal cells.* As expected, the accuracy of clonal reconstruction decreased
135 with the contamination level (Spearman's rank correlation $\rho = -0.29$, $p - value = 8.1 \times 10^{-6}$); here and
136 below, correlation is provided for QuantumClone results only, although the observed trend is usually true for
137 the three other methods.

138

139 *Effect of the number of variants.* Unexpectedly, increasing the number of variants slightly decreased the
140 quality of clustering ($\rho = -0.17$, $p - value = 0.0104$). This effect is due to an artefactual increase in
141 the number of mutation clusters needed to explain larger numbers of observed cellular prevalence values. In
142 other words, independently of the true clonal structure, more variants in the input result in the prediction
143 of higher number of mutation clusters corresponding to clones.

144

145 *Effect of the number of samples.* To reconstruct clones existing in patients' tumors with higher fidelity,
146 recent studies advocate for sequencing of multiple samples obtained from the same patient, either by us-
147 ing different time points or different sites of the tumor (Schwarz et al., 2014). Indeed, we demonstrated
148 that increasing the number of available samples significantly improved the quality of clonal reconstruction
149 ($\rho = 0.40$, $p - value = 6.3 \times 10^{-6}$).

150

151 *Effect of sequencing depth.* Variance in VAF estimation is expected to go down with the increase in
152 sequencing depth, resulting in better mutation clustering. As expected, we observed a significant positive

153 correlation between depth of sequencing and the quality of clonal reconstruction ($\rho = 0.67$, p -value <
154 2.2×10^{-16}).

155

156 *Effect of the number of distinct clones present.* We observe that the quality of clonal reconstruction de-
157 clined with the increase in tumor heterogeneity, i.e., the number of distinct clones present in patients' tumors
158 ($\rho = -0.27$, p -value = 8.2×10^{-5}). As previously mentioned, in heterogeneous samples, QuantumClone
159 showed better performance than the other tested methods.

160

161 **Creating a robust framework for clonal assignment of functional mutations**

162 We proposed a novel concept of reconstruction of the clonal architecture in cancer combined with the attri-
163 bution of functional mutations (potential drivers) to identified clones (Fig. 1). The approach is based on
164 the different usage of '*functional*' variants that potentially affect cell phenotype and '*support*' variants that
165 are used to define clones. *Support* variants can be either drivers or passengers; however, they should have
166 high depth of coverage ($> 50\times$ in our implementation), have no strand bias and should not coincide with
167 annotated single-nucleotide polymorphisms (SNPs). As we showed in the simulation studies (Fig. 2B) only
168 a limited number of support variants are needed for an accurate clonal reconstruction. Therefore, in most
169 cases, we can even afford to limit the set of support variants to those falling in regions of genotype A and
170 AB; in such regions VAF values directly determine values of mutation cellular prevalence (Methods). *Support*
171 variants, because they have a lower variance of observed VAF compared with other variants, are applied to
172 define clones, i.e., *support* variants serve as input to QuantumClone or an alternative method. *Functional*
173 mutations are defined as variants with deleterious properties that affect either genes reported in the Cancer
174 Census List (Futreal et al., 2004) or genes from gene modules/signaling pathways that are recurrently affected
175 by mutations in a given cancer type (Methods). At the last step of our framework, functional mutations are
176 mapped to the clonal structure inferred from support variants based on the likelihood values.

177

178 The QuantumClone R package includes functions for both clonal reconstruction using *support* variants
179 and assignment of *functional* mutations to the defined clones.

180

181 We propose the following options to be used in the analysis framework. Deleterious mutations can be
182 determined using SIFT (Ng and Henikoff, 2003), PolyPhen-2 (Adzhubei et al., 2013) and FunSeq2 (Khurana
183 et al., 2013). Gene module enrichment analysis may be performed using the R package ACSNMiner (Deveau

184 P., Barillot E., Boeva V., Zinovyev A., Bonnet E., *In Press*) using maps and modules of the Atlas of Can-
185 cer Signalling Networks (ACSN) (Kuperstein et al., 2015) completed with the user-defined/cancer-specific
186 modules. For the specific case of neuroblastoma, we created a ‘Neuritogenesis’ map extracted from Molenaar
187 et al. (2012).

188

189 **Characterization of neuroblastoma clonal evolution from diagnosis to relapse:** 190 **application of the QuantumClone-based framework**

191 We applied our framework to investigate the clonal composition of neuroblastoma primary and relapse tumors
192 and study its clonal evolution. We characterized clonal structure of tumors of 22 neuroblastoma patients
193 (clinical data available in Suppl. Table 1). We performed WGS of constitutive DNA, diagnosis and relapse
194 tumor samples of each patient with average depth of coverage $\sim 100\times$. Datasets for 15 patients out of
195 22 came from a previously published study (Eleveld et al., 2015). Sequencing was carried out using both
196 Illumina HiSeq 2500 and Complete Genomics platforms. Reads were mapped to the reference hg19 genome
197 using BWA-aln (Li and Durbin, 2009) (Illumina reads) and the internal Complete Genomics mapping tool
198 (Complete Genomics reads). Variant calling was performed using VarScan2 version 2.3.6 (Koboldt et al.,
199 2013).

200

201 The level of contamination by normal cells varied from 0% to 90%, and only data from 19 patients with
202 a contamination level lower than 70% were kept for further analysis (Suppl. Table 1).

203

204 **Application of filters unifies variant call numbers across different sequencing platforms**

205 In order to remove false positive variant calls, we used a set of stringent filters (Fig. 1, Methods). The initial
206 number of variants in the VarScan2 output was highly dependent on the sequencing technology and platform
207 (Suppl. Fig. 1 and 2). The number of variants called for samples sequenced by the Beijing Genomics Institute
208 (BGI) sequencing platform was an order of magnitude higher than the number of mutations called for samples
209 processed by the Centre National de Génotypage (CNG). Application of a set of filters based on read depth
210 of coverage, read mappability, annotated repetitive regions (listed in Fig. 1 and Methods) allowed us to
211 get comparable numbers of variants for further analysis. In the final list, number of mutations per sample
212 correlated with the age of the patient (Fig. 4, Spearman’s $\rho = 0.54$, p -value = 5.3×10^{-4} ; datasets tested
213 include information from all samples kept after the evaluation of normal contamination). The stability of

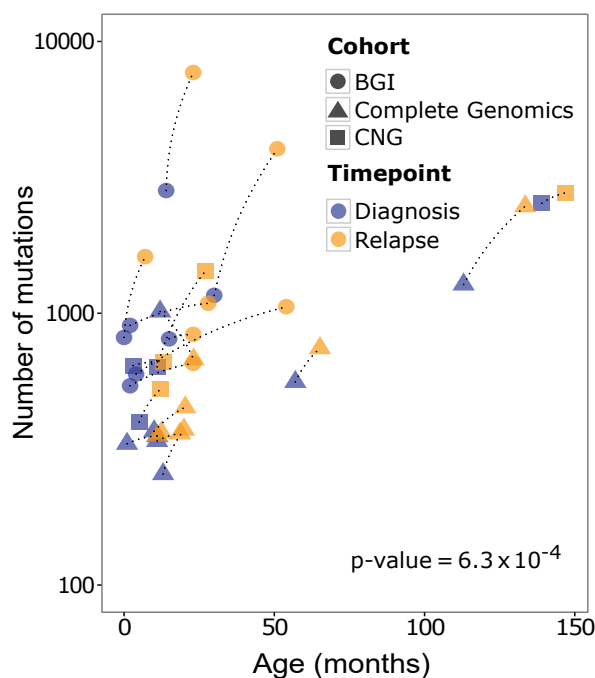


Figure 4: **Statistics on numbers of somatic mutations called using stringent filters for diagnosis and relapse samples from 19 neuroblastoma patients.** The number of somatic mutations is correlated with the age of the patient at the time of the biopsy (Spearman correlation test p -value = 6.3×10^{-4}). Final mutation numbers after the filtering step do not depend on the sequencing center or sequencing technology used. Diagnosis and relapse samples from the same patient are connected by a dotted line.

214 final numbers of predicted variants across sequencing platforms and correlation of these numbers with the
215 age of patients validated our filtering approach. In fact, the presence of such correlation has previously been
216 shown in neuroblastoma (Molenaar et al., 2012). In all but one neuroblastoma patient, there were more
217 variants detected in the relapse sample than in the diagnosis sample with an average two-fold increase.

218 Clonal reconstruction

219 We applied QuantumClone on *support* variants we defined using stringent filters (Fig. 1, Step 2; Fig. 5).
220 Across our cohort, we observed a significant association between the predicted number of clones and the
221 number of mutations per patient (Spearman's $\rho = 0.42$, p -value = 0.011). However, for each given pa-
222 tient, the number of clones at relapse was similar to that at diagnosis, even despite the fact that the relapse
223 samples had about twice as many mutations as the diagnosis samples (number of mutation clusters varied
224 from zero to four with a median of three for both time points).

225

226 We identified mutations coming from the ancestral clone (Figure 6A), i.e. the clone that gave rise to all
227 cells in both diagnosis and relapse samples, in 84% of reconstructed clonal structures (16 out of 19 patients).

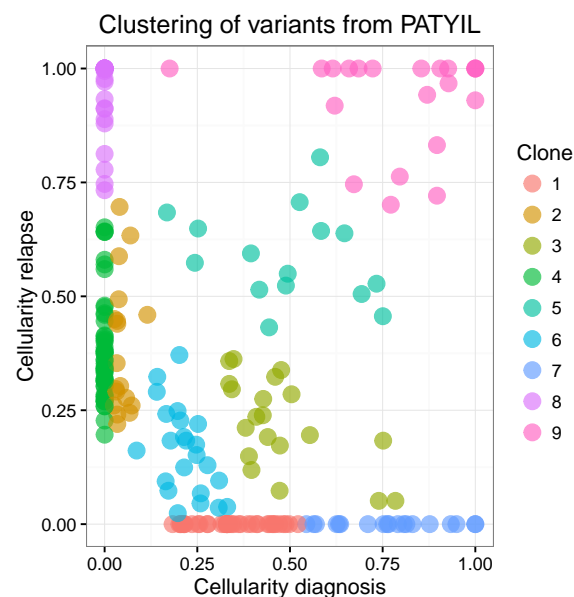


Figure 5: **Visualization of the QuantumClone output on the PATYIL neuroblastoma patient.** Nine mutation clusters corresponding to inferred clones are shown in different colors. The cellular prevalence of each mutation was assessed by QuantumClone by applying the EM algorithm on VAFs using known levels of normal contamination and locus genotype information.

228 The total number of *support* variants attributed to ancestral clones ranged between 3 and 120, with a median
229 of 22.

230 **Annotation of functional mutations in each sample based on the global pathway enrichment** 231 **analysis**

232 In our framework, we assumed that *functional* mutations in a given cancer type, i.e., putative drivers, should
233 target specific signaling pathways or pathway modules (Fig. 1, Step 4). These modules can be identified
234 as those enriched in coding deleterious mutations across the cohort. Thus, we mapped the total of 541
235 deleterious mutations obtained with tolerant filters (Fig. 1, Methods) to the ACSN maps and detected re-
236 currently altered gene modules using the ACSNmineR package. Overall, five general gene maps (apoptosis,
237 DNA repair, EMT/cell motility, cell survival and neuritogenesis) and their 15 gene modules were found to
238 be enriched in mutations (threshold 0.05 on the p-value corrected to account for multiple testing with the
239 Benjamini-Hochberg False Discovery Rate correction, corresponding to the q-value)) (Figure 6B, Supp. Ta-
240 ble 2). Next, deleterious mutations were annotated as *functional* when corresponding genes were included
241 in the enriched pathways, or when such genes belonged to the Cancer Census list. The resulting number of
242 *functional* mutations per sample varied from 1 to 23.

243

244 At this step, among the cell survival modules, the highest enrichment in putative driver mutations was
245 observed for the MAPK pathway ($q - value \leq 7.2 \times 10^{-5}$). In addition, we detected significant enrichment
246 in *functional* mutations of both canonical and non-canonical WNT pathways ($q - value \leq 3.8 \times 10^{-3}$ and \leq
247 1.49×10^{-2} , respectively), and of the PI3K/AKT/mTOR and Hedgehog gene modules ($q - value \leq 1.3 \times 10^{-2}$
248 and $\leq 4.5 \times 10^{-4}$, respectively). As for the modules of other maps, genes coding for the EMT regulators
249 were also significantly affected by the deleterious mutations in our cohort of relapsed neuroblastoma patients
250 ($q - value \leq 1.2 \times 10^{-6}$).

251

252 As a control, we ran ACSNMiner on synonymous variants passing the same filters as the non-synonymous
253 ones (1060 synonymous mutations in 771 genes). Enrichment analysis provided no modules enriched in
254 synonymous variants (*minimal p-value* > 0.28), thus confirming the biological significance of the discovered
255 enrichment in functional mutations for 20 ACSN modules and maps. Strikingly, similar enrichment analysis
256 of intronic variants showed recurrently affected pathways similar to those affected by coding deleterious
257 mutations, such as EMT-motility (*Odds ratio* = 5.3, $p - value < 10^{-20}$) or cell survival (*Odds ratio* = 4.0,
258 $p - value < 10^{-20}$), highlighting the possible role of intronic SNVs in neuroblastoma tumorigenesis.

259 **Assignment of *functional* mutations to the identified clonal structure**

260 Using the results of the mapping of *functional* mutations on the clonal structure detected for each patient
261 by QuantumClone (Fig. 1, Step 5), we annotated mutations as (i) those belonging to expanding clones, (ii)
262 those belonging to shrinking clones, and (iii) those belonging to ancestral clones (Fig. 6A). Overall, 53%,
263 31% and 4,8% of all *functional* mutations fell in these three categories.

264

265 For the majority of samples (16 out of 19 patients), we could not detect *functional* mutations in the
266 ancestral clone. But, interestingly, samples with identified ancestral *functional* mutations contained variants
267 in the *MYC* and *AKT2* genes, which are known to act as drivers in many cancers. Moreover, the ancestral
268 clones detected in our samples, often contained a higher proportion of putative driver mutations affecting
269 the DNA repair, EMT/cell motility and other pathways than any other clones including those expanding at
270 relapse (Fig. 6B and 6C).

271 **Analysis of pathways enriched in *functional* mutations in shrinking and expanding clones**

272 Assignment of mutations to clones shrinking or expanding after the treatment resulted in the identification
273 of 33 and 56 possible driver mutations in these clone types, respectively. Expanding clones had more dele-
274 terious mutations targeting genes from four general maps (DNA repair, EMT/cell motility, cell survival and

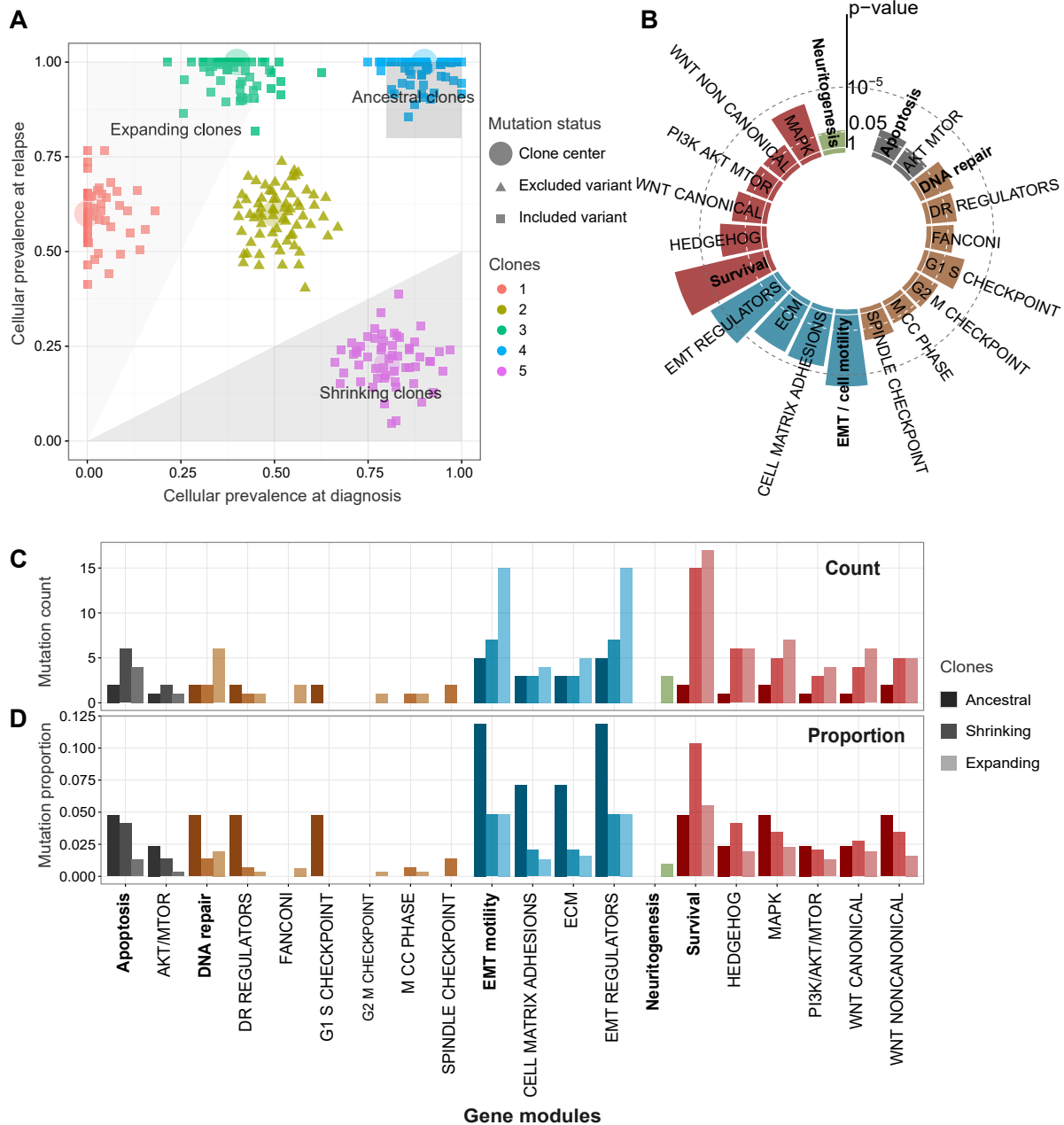


Figure 6: **Annotation of clones in neuroblastoma and pathway enrichment analysis.** (A) Illustration of the rule for assignment of mutations to (i) the ancestral clone (cellular prevalence of the mutation cluster exceeds 80% both at diagnosis and relapse), (ii) clones expanding after the treatment (cellular prevalence of the mutation cluster increases at least two-fold at relapse) and (iii) shrinking clones (cellular prevalence of such mutation clusters decreases at least two-fold). (B) Pathways enriched in deleterious mutations in neuroblastoma. General ACSN maps and the manually added neurogenesis pathway are shown in bold; map sub-modules are shown in the color of the corresponding map. The “Cell Cycle” ACSN map does not show enrichment in mutations and is omitted from the graph. (C) Absolute numbers of deleterious mutations in gene maps and modules in the ancestral clones, and clones expanding and shrinking at relapse. (D) Proportions of *functional* mutations in each module over the total number of detected deleterious mutations in the ancestral, expanding and shrinking clones.

275 neuritogenesis) (Fig. 6C). Similarly, in these clones, most of the corresponding gene modules (e.g., MAPK,
276 WNT canonical or PI3K/AKT/mTOR) were also more frequently targeted. Although the absolute number
277 of *functional* mutations in the expanding clones was about twice as high as in the shrinking ones, the pro-
278 portion of such mutations among all deleterious ones was not significantly different between the shrinking
279 and expanding clones (Fig. 6D).

280

281 In addition to the MAPK pathway mutations previously reported for the relapse neuroblastoma samples
282 (e.g., mutations in the *NRAS* gene (Eleveled et al., 2015)), in the expanding clones we observed coding variants
283 affecting additional MAPK pathway genes such as *NFE2L1*, *PIM1*, *FLT4* and *RSK*. These genes have been
284 shown to be involved in tumorigenesis of prostate (Yu et al., 2015), colorectal (Xiao et al., 2014), breast
285 (Malinen et al., 2013) and other cancers.

286 Discussion

287 Here we propose a pathway-based framework to detect functional mutations in cancer samples and associate
288 the mutations to their corresponding clonal structure. The central part of our framework is represented by
289 the QuantumClone method, which allows reconstruction of clonal populations based on both variant allele
290 frequencies and genotype information. QuantumClone showed stable results on simulated data significantly
291 outperforming other methods in difficult settings such as highly contaminated samples, heterogeneous tumors
292 and relatively low depth of sequencing coverage.

293

294 Our analysis framework is based on two central ideas. First, high-reliability passenger mutations must
295 be used to reconstruct the clonal structure of tumors samples; then, low coverage functional mutations (with
296 high variance in VAFs) should be mapped onto the inferred clonal structure. Second, we suggest to limit the
297 set of functional mutations to those in genes known to be associated with cancer (e.g., Cancer Census genes)
298 or to those in genes from gene modules/pathways frequently disrupted in a given cancer type (Fig. 1).

299

300 We apply the proposed analysis framework to decipher clonal structure in neuroblastoma and assign to
301 clones possible driver mutations. We detect 15 pathways as being altered by mutations in neuroblastoma.
302 We identify genes associated with DNA repair, cell motility, apoptosis and survival to be enriched in func-
303 tional mutations in neuroblastoma. For relapsed neuroblastoma samples, we recover the previously reported
304 enrichment of mutations in the MAPK signaling pathway (Eleveled et al., 2015), while complementing this
305 knowledge with discovery of accumulation of functional mutations at the relapse in such functional gene

306 modules as PI3K/AKT/mTOR, WNT, Hedgehog signaling and modules consisting of genes responsible for
307 cell-matrix adhesion and epithelial–mesenchymal transition (EMT).

308

309 We highlight the lack of enrichment in mutations of the cell cycle pathway in this pediatric cancer. This
310 could be explained by the biological context of neuroblastoma that may already favor proliferation, which
311 accompanies organism development, thus limiting the necessity to disrupt cell cycle mechanisms. However,
312 this observation should be in the future confirmed on a larger cohort.

313

314 For the majority of neuroblastoma patients, we did not identify driver mutations in the ancestral clone.
315 This is in line with the current understanding of neuroblastoma as a type of cancer driven by copy number
316 alterations. In fact, it has been known that the aggressiveness of neuroblastoma is highly associated with
317 changes of the chromosomal copy number profile (MYCN amplification, 1p, 3p, 11q deletions, partial gain
318 of chr17) (Janoueix-Lerosey et al., 2010). And indeed, copy number profiling often detected gain and losses
319 in these recurrently affected chromosomal regions both in the diagnosis and relapse samples from our cohort
320 (Suppl. Fig. 3).

321

322 For most of our samples, we did not succeed in reliably reconstructing the phylogeny of clonal evolution
323 based on cellular prevalence values for identified mutation clusters (clones). Some contradictions between
324 cellular prevalence values between diagnosis and relapse, as well as disappearance at relapse of many potential
325 driver mutations seemingly present in the ancestral clone at diagnosis, may be due to tumor heterogeneity
326 and the fact that biopsies were taken from different tumor sites. This situation has been termed "illusion
327 of clonality" (Bruin et al., 2014). However, the fact that for most of the samples we observed a number of
328 shared mutations between diagnosis and relapse and that copy number breakpoints were consistent between
329 the two time points ensures that there is a common phylogeny between diagnosis and relapse in neuroblas-
330 toma (Bollet et al., 2008).

331

332 In the application of our framework to neuroblastoma sequencing data, we excluded information about
333 SVs and indels. The reason for this was that the analysis of clonal structure is based on the number of se-
334 quencing reads supporting each genetic variant. While we suppose that the number of reads with a mismatch
335 mutation is proportional to the number of DNA molecules harboring this variant, we expect that due to read
336 mapping issues the fraction of reads indicating an indel or a large SV will be generally lower than the actual
337 proportion of DNA molecules with the rearrangement. Eviction of large and small SVs seemingly resulted
338 in a decrease in sensitivity of the detection of genetic driver events. In our neuroblastoma data, a large

339 proportion of observed clones did not contain any predicted driver mutation. In the future, the sensitivity
340 can be improved by using higher depth of coverage data and combining the paired-end datasets with reads
341 produced with the mate-pair protocol or with long PacBio reads.

342

343 The proposed framework can be applied in the future to any type of cancer. The pre-requirements are
344 sufficient number of candidate mutations (at least 50 mutations per sample) and a minimal read depth of
345 coverage of 50×. These requirements are usually met by WGS or whole exome sequencing datasets. Our
346 simulation results show that increasing the number of mutations used for clonal reconstruction above 50 does
347 not improve significantly the clonal reconstruction accuracy provided that mutations specific for every clone
348 are present in the input.

349

350 **Methods**

351 **Datasets**

352 *Patient selection and collection of tumor samples.* The inclusion criteria for this study were histopathologi-
353 cal confirmation of neuroblastoma at original diagnosis and the presence of biopsy material from a subsequent
354 relapse specimen. Patients were included in this study after an informed consent was obtained from par-
355 ents or guardians, with oversight from the ethics committees 'Comité de Protection des Personnes Sud-Est
356 IV', reference L07-95/L12-171, and 'Comité de Protection des Personnes Ile-de-France', reference 0811728 in
357 France, the review board at the Children's Hospital of Philadelphia and review boards at other Children's
358 Oncology Group sites that submitted samples for patients on this study in the United States. In total we
359 obtained material for 22 neuroblastoma patients (tumor tissue at diagnosis, relapse and constitutional DNA,
360 Suppl. Table 1).

361

362 *Whole-genome sequencing of neuroblastoma samples.* In the framework of this study, we carried out
363 Illumina paired-end sequencing for 7 novel neuroblastoma patients (corresponding to 21 samples). Data for
364 15 patients were taken and reanalyzed from the previous study (Eleveld et al., 2015). DNA from 7 patients
365 from the previous study and 7 new ones have been sequenced using Illumina HiSeq 2500 instruments to an av-
366 erage depth of coverage of 80× by Beijing Genomics Institute (BGI) and the Centre National de Génotypage
367 (CNG) respectively. For 8 patients out of 15 previously reported, whole-genome sequencing was performed
368 by Complete Genomics with an average read depth of coverage of 50×. DNA material for each patient

369 (lymphocytes, primary tumors and relapse tumors) was in each case sequenced using the same sequencing
370 platform (see Suppl. Table 1 for more detail).

371

372 *Data processing.* Sequenced reads were mapped to the human genome hg19 using BWA and the internal
373 Complete Genomics tools for Illumina and Complete Genomics datasets respectively. Reads from datasets
374 sequenced using the Illumina platform were realigned around indels with the Genome Analysis ToolKit
375 (GATK) (McKenna et al., 2010), followed by a base recalibration. Due to the inherent structure of Complete
376 Genomics reads, which contain an effective deletion relative to their corresponding genomic library, the indel
377 realignment step was skipped for the Complete Genomics samples.

378

379 **Variant calling and filtering**

380 Mutations were called using VarScan2 (Koboldt et al., 2013). Two sets of variants were created for each
381 patient (see Fig. 1) using tolerant and stringent filtering options. The ‘stringent’ set was further used for
382 clonal reconstruction, while the ‘tolerant’ one was used for inference of recurrently altered pathways.

383

384 Tolerant filters for somatic mutations included those on minimal depth of coverage ($30\times$), minimal per-
385 centage of reads supporting the mutation (10%). In addition, mutations were required to be located in regions
386 of high local mappability (36 bp mappability), outside of repeat and duplicated genomic regions (assessed
387 by the UCSC repeat and segmental duplication region tracks). We further filtered mutations that created a
388 stretch of four or more identical nucleotides. Finally, we only kept mutations located in regions where the
389 genotype evaluated by Control-FREEC was available.

390

391 To obtain a set of high confidence mutations, we required the minimal depth of coverage of $50\times$. We
392 filtered out variants corresponding to polymorphisms present in more than 1% of the population (snp138,
393 1000Genomes, esp6500) except if it was a known cancer related variant (COSMIC database for coding and
394 non-coding mutations). Strand bias was also tested by the Fisher exact test (in addition to the test in
395 VarScan2) to reduce the number of informative mutations to a maximum of 500. We restricted the analysis
396 to AB regions in case when after such filtering we kept more than 100 mutations.

397

398 Copy number analysis

399 Copy number alterations in patients were detected using the Control-FREEC method (Boeva et al., 2012)(ver-
400 sion 7.2) (Suppl. Fig. 3). We selected the main ploidy value so that the predicted copy number and B-allele
401 frequency profiles were consistent. Control-FREEC also provided estimations of the level of contamination
402 by normal cells, which, after manual confirmation, was further used for the clonal reconstruction.

403

404 Three samples with the estimated proportion of contamination by normal cells higher than 70% were
405 excluded from the further analysis (NB0784:diagnosis, NB1434:relapse and NB1471:relapse).

406

407 Comparison of clonal reconstruction between QuantumClone and existing meth- 408 ods

409 *Data simulation.* In silico validation data were generated using the QuantumCat method from package
410 QuantumClone (version 0.15.12.10). QuantumCat simulates genomic mutations, copy number alterations
411 and corresponding VAFs. It relies on the following set of rules:

- 412 1. A binary phylogenetic tree is created to simulate the clonal architecture of the tumor. The mutation
413 cellular prevalence values correspond to the nodes and leaves of the phylogenetic tree.
- 414 2. Cellular prevalence values of mutations from each clone are independent across tumor samples. However,
415 the cellular prevalence of each clone should be always coherent with the phylogenetic tree.
- 416 3. The allelic copy number of all mutation loci was set to AB in the tests carried out to compare Quan-
417 tumClone, sciClone and pyClone (Figure 2B). For QuantumClone validation on triploid and tetraploid
418 genomes (Figure 3), the number of chromosomal copies bearing each mutation was randomly assigned
419 between one and the number of A-alleles for the locus considered. Generation of the genotype, num-
420 ber of chromosomal copies, normal contamination and cellular prevalence of a mutation allows for the
421 computation of the exact VAF, which is the cellular prevalence (taking into account the contamination
422 by normal cells) times the number of copies of the mutations divided by the number of copies of the
423 locus in each cell. On the other hand, the observed VAF is determined by the ratio of the number
424 of reads supporting the mutations divided by the read depth of coverage. Local depth of coverage at
425 each given position was generated by the negative binomial distribution centered on the target depth
426 of sequencing, fitted on experimental data. The number of reads supporting a mutation was simulated
427 from the binomial distribution with the probability of success equal to the exact VAF.

428 *Program versions and parameters.* We used SciClone version 1.0.7 with the following changes to the default
429 parameters: maximal number of clusters was set to 10 and the minimal depth of coverage was set to 0.

430

431 PyClone version 0.12.9 was used with the following parameters : 10 iterations of the Markov chain Monte
432 Carlo, alpha and beta parameters in the Beta base measure for Dirichlet Process set to 1, concentration prior
433 shape set to 1 and the rate parameter in the Gamma prior on the concentration parameter set to 0.001. We
434 used the default Beta binomial distribution with precision parameter set to 1000, prior shape set to 1, rate
435 to 0.0001, and proposal precision set to 0.01.

436

437 We used an implementation of the k -medoids algorithm provided by the R package “fpc”, version 2.1.10,
438 with a range of clusters between 2 and 10.

439

440 For clonal reconstruction, QuantumClone version 0.15.12.10 was used with default parameters except for
441 the the maximal number of clusters, which was set to 10. For clonal reconstruction of neuroblastoma data,
442 mutations used to compute centers of clusters (corresponding to clones) were selected using the stringent set
443 of filters. Copy number information from Control-FREEC (version 7.2) was also passed to the algorithm as
444 well as the predicted value of contamination by normal cells.

445

446 For simulated data, quality of clustering was assessed by using Normalized Mutual Information (NMI),
447 which is given for a group of clones Ω and a group of reconstructed clusters \mathbb{C} :

$$NMI_{(\Omega, \mathbb{C})} = -2 \times \frac{\sum_k \sum_j \frac{|\omega_k \cap c_j|}{N} \log \left(\frac{N \times |\omega_k \cap c_j|}{|\omega_k| |c_j|} \right)}{\sum_k \frac{|\omega_k|}{N} \log \left(\frac{|\omega_k|}{N} \right) + \sum_j \frac{|c_j|}{N} \log \left(\frac{|c_j|}{N} \right)}, \quad (1)$$

448 where N is the number of mutations observed, $|\omega_k|$ the number of mutations in clone k , and $|c_j|$ the
449 number of mutations attributed to cluster j .

450

451 Clonal reconstruction

452 In this section we describe QuantumClone, a method we have developed for the clonal reconstruction of a
453 tumor. QuantumClone performs clustering of cellular prevalence values of mutations defined by:

$$\theta = \frac{VAF \times N_{Ch}}{NC \times P}, \quad (2)$$

454 where θ is the cellular prevalence, N_{Ch} the number of copies of the corresponding locus, NC the number
 455 of chromosomal copies bearing the mutation, and P the tumor purity. For instance, only in case of a purely
 456 diploid tumor without loss of heterozygosity (LOH) regions, with no contamination of the sample by normal
 457 cells, the cellular prevalence is equal to $2 \times VAF$. The latter assumption has been frequently used in cancer
 458 studies (Schramm et al., 2015). As we do not have information about the number of chromosomal copies
 459 bearing a mutation, our approach was to compute each possible value of cellular prevalence associated with
 460 the mutation VAF. For example, a mutation can have a VAF of 1/3 in a locus of genotype AAB when it is
 461 present in 100% of tumor cells on a single chromosome copy and when it is present in 50% of tumor cells
 462 on two chromosomes. Yet the latter case is rather improbable. Each mutation thus corresponds to several
 463 possible values of cellular prevalence; each solution is associated with a value of NC . In order to address
 464 the problem of non-uniqueness of solution, we use an EM algorithm based on the probability to observe
 465 a specific number of reads confirming a mutation given the number of reads overlapping the position, the
 466 contamination and the cellularity of a clone. In more detail, we attribute to each possibility a probability to
 467 observe f reads supporting the variant given that the latter belongs to a clone of cellular prevalence θ , based
 468 on a binomial distribution:

$$P(f|\theta) = \binom{d}{f} \left(\frac{\theta \times NC(1-c)}{N_{Ch}} \right)^f \times \left(\frac{1 - (\theta \times NC(1-c))}{N_{Ch}} \right)^{d-f}, \quad (3)$$

469 where

- 470 • d is the depth of coverage of the variation
- 471 • f is the number of reads supporting the variant
- 472 • c is the sample contamination by normal cells

473 We can then write the log likelihood function to maximize:

$$L = \sum_{i \in \text{mutations}} \sum_{k \in \text{clones}} \sum_{s \in \text{samples}} \sum_{p \in \text{possibilities}(i)} \omega_{(i,p)} t_{(i,k)} \log(P_{i,s,p}(f_{i,s,p}|\theta_{k,s})), \quad (4)$$

where $\omega_{i,p}$ are weights of the possibility computed for a corresponding genotype $xAyB$ (major allele A is present x times and the minor allele B is present y times):

$$\omega_{i,p} = \prod_{s \in \text{samples}} \frac{\binom{x_s}{NC_{i,s,p}} + \binom{y_s}{NC_{i,s,p}}}{2^{N_{Ch_s}}}$$

474 By adding weights that for each variant sum to one, we favour mutations with the lowest number of copies.

475 Each mutation is then attributed to its most likely possibility, which is the possibility with highest probability
476 to belong to a clone. In the situation described above (a variant in a AAB region with the VAF of 1/3), this
477 approach would assign probabilities of 2/3 and 1/2 to the presence of the mutation in 100% and 50% of cells
478 respectively. However, if there is a second mutation present, for example, in a locus of genotype AB with a
479 VAF of 1/2 and thus having unambiguously cellular prevalence of 100%, the first mutation will have a high
480 density of probability for a cellular prevalence of 100% and our approach will assign both mutations to the
481 same cluster corresponding to the same cellular prevalence (100%).

482

483 The number of clones is determined by minimization of the Bayesian Information Criterion (BIC). Priors
484 can be provided by the user, randomly generated or determined by the k -medoids clustering on mutations in
485 A and AB sites when the latter contain enough mutations.

486

QuantumClone-Single and QuantumClone-Alpha variants of QuantumClone. To test accuracy of predic-
tions for the number of copies with a variant and selection of the most likely mutation cellular prevalence, we
designed two alternatives of the QuantumClone method: QuantumClone-Single and QuantumClone-Alpha.
QuantumClone-Single is a modification of QuantumClone that assigns to all variants a single copy state.
QuantumClone-Alpha uses the same EM algorithm as the default version of QuantumClone, except for the
selection of the most likely mutation where it chooses the possibility maximizing the quantity q :

$$q_{f,\theta,p} = \binom{N_{Ch}}{NC} \times P(f|\theta).$$

487 **Analysis of mutation enrichment in signaling pathways and gene modules**

488 ACSNMiner (Deveau P., Barillot E., Boeva V., Zinovyev A., Bonnet E., *In Press*) version 0.16.01.29 was
489 used to compute gene modules and pathways enriched in deleterious mutations. Gene modules included
490 by default in ACSNMiner come from the manually curated Atlas of Cancer Signalling Networks (ACSN)
491 (Kuperstein et al., 2015). In addition to the ACSN modules, we calculated mutation enrichment in a set of
492 neurogenesis genes frequently mutated in neuroblastoma (Molenaar et al. (2012), Suppl. Table 8). We called
493 deleterious mutations as stop-gain mutation or variants that were predicted as possibly damaging or deleteri-
494 ous by SIFT (Ng and Henikoff, 2003), PolyPhen-2 (Adzhubei et al., 2013), or FunSeq2 (Khurana et al., 2013).

495

496 To get a list of genes to use as an input to ACSNMiner, we pooled mutations from all neuroblastoma
497 patients; genes mutated at least once were included in the final list. Modules with a p-value lower than 0.05

498 after Benjamini-Hochberg correction were considered as enriched.

499

500 **Data access**

501 The whole-genome sequencing data have been deposited at the European Genome-phenome Archive (EGA)
502 under accession number EGAS00001001184 for the French cases sequenced at BGI and under accession num-
503 ber EGAS00001001825 for the French cases sequenced at CNG. Sequence data for the US cases are available
504 in the database of Genotypes and Phenotypes (dbGaP) under accession number phs000467.

505

506 QuantumClone is available at <https://github.com/DeveauP/QuantumClone/> and can be downloaded as
507 an R package from the CRAN repository.

508

509 **Acknowledgments**

510 GS and her team were supported by the Annenberg Foundation and the Nelia and Amadeo Barletta Founda-
511 tion. Funding was also obtained from SiRIC/INCa (Grant INCa-DGOS-4654) and from the CEST of Institut
512 Curie. This study was also funded by the Associations Enfants et Santé, Association Hubert Gouin Enfance
513 et Cancer, Les Bagouz à Manon, Les amis de Claire. VB and her team were supported by the ATIP-Avenir
514 Program, the ARC Foundation and the "Who Am I?" Project. EB was supported by the ABS4NGS project
515 of the French Program 'Investissement d'Avenir'. Sequencing of French samples was carried out in a collab-
516 oration of Institut Curie with CEA/IG/CNG financed by France Génomique infrastructure, as part of the
517 program "Investissements d'Avenir" from the Agence Nationale pour la Recherche (contract ANR-10-INBS-
518 09). JM and his team were supported in part by US National Institutes of Health grants RC1MD004418
519 to the TARGET consortium, and CA98543 and CA180899 to the Children's Oncology Group. In addition,
520 this project was funded in part with Federal funds from the National Cancer Institute, National Institutes
521 of Health, under Contract No. HHSN261200800001E. The content of this publication does not necessarily
522 reflect the views of policies of the Department of Health and Human Services, nor does mention of trade
523 names, commercial products, or organizations imply endorsement by the U.S. Government.

524 **DISCLOSURE DECLARATION**

525 We have no conflict of interest to declare.

References

- Adzhubei, I., Jordan, D. M., and Sunyaev, S. R., 2013. Predicting Functional Effect of Human Missense Mutations Using PolyPhen-2. *Current protocols in human genetics / editorial board, Jonathan L. Haines ... [et al.]*, **0** 7:Unit 7.20.
- Boeva, V., Popova, T., Bleakley, K., Chiche, P., Cappo, J., Schleiermacher, G., Janoueix-Lerosey, I., Delattre, O., and Barillot, E., 2012. Control-FREEC: a tool for assessing copy number and allelic content using next-generation sequencing data. *Bioinformatics*, **28**(3):423–425.
- Bollet, M. A., Servant, N., Neuvial, P., Decraene, C., Lebigot, I., Meyniel, J.-P., Rycke, Y. D., Savignoni, A., Rigaille, G., Hupé, P., *et al.*, 2008. High-Resolution Mapping of DNA Breakpoints to Define True Recurrences Among Ipsilateral Breast Cancers. *Journal of the National Cancer Institute*, **100**(1):48–58.
- Bruin, E. C. d., McGranahan, N., Mitter, R., Salm, M., Wedge, D. C., Yates, L., Jamal-Hanjani, M., Shafi, S., Murugaesu, N., Rowan, A. J., *et al.*, 2014. Spatial and temporal diversity in genomic instability processes defines lung cancer evolution. *Science*, **346**(6206):251–256.
- Eleveld, T. F., Oldridge, D. A., Bernard, V., Koster, J., Daage, L. C., Diskin, S. J., Schild, L., Bentahar, N. B., Bellini, A., Chicard, M., *et al.*, 2015. Relapsed neuroblastomas show frequent RAS-MAPK pathway mutations. *Nature Genetics*, **47**(8):864–871.
- Fischer, A., Vázquez-García, I., Illingworth, C. J. R., and Mustonen, V., 2014. High-Definition Reconstruction of Clonal Composition in Cancer. *Cell Reports*, **7**(5):1740–1752.
- Futreal, P. A., Coin, L., Marshall, M., Down, T., Hubbard, T., Wooster, R., Rahman, N., and Stratton, M. R., 2004. A census of human cancer genes. *Nature Reviews Cancer*, **4**(3):177–183.
- Hanahan, D. and Weinberg, R. A., 2011. Hallmarks of Cancer: The Next Generation. *Cell*, **144**(5):646–674.
- Janoueix-Lerosey, I., Schleiermacher, G., and Delattre, O., 2010. Molecular pathogenesis of peripheral neuroblastic tumors. *Oncogene*, **29**(11):1566–1579.
- Jiao, W., Vembu, S., Deshwar, A. G., Stein, L., and Morris, Q., 2014. Inferring clonal evolution of tumors from single nucleotide somatic mutations. *BMC Bioinformatics*, **15**(1):35.
- Kaufman, L. and Rousseeuw, P., 1987. Clustering by means of medoids. *Statistical Data Analysis Based on the L1-Norm and Related Methods*, :405–416.
- Kepler, T. B., 2013. Reconstructing a B-cell clonal lineage. I. Statistical inference of unobserved ancestors. *F1000Research*, .

- Khurana, E., Fu, Y., Colonna, V., Mu, X. J., Kang, H. M., Lappalainen, T., Sboner, A., Lochovsky, L., Chen, J., Harmanci, A., *et al.*, 2013. Integrative Annotation of Variants from 1092 Humans: Application to Cancer Genomics. *Science (New York, N.Y.)*, **342**(6154):1235587.
- Koboldt, D. C., Larson, D. E., and Wilson, R. K., 2013. Using VarScan 2 for Germline Variant Calling and Somatic Mutation Detection. *Current protocols in bioinformatics / editorial board, Andreas D. Baxevanis ... [et al.]*, **44**:15.4.1–15.4.17.
- Kuperstein, I., Bonnet, E., Nguyen, H.-A., Cohen, D., Viara, E., Grieco, L., Fourquet, S., Calzone, L., Russo, C., Kondratova, M., *et al.*, 2015. Atlas of Cancer Signalling Network: a systems biology resource for integrative analysis of cancer data with Google Maps. *Oncogenesis*, **4**(7):e160.
- Li, H. and Durbin, R., 2009. Fast and accurate short read alignment with Burrows–Wheeler transform. *Bioinformatics*, **25**(14):1754–1760.
- Malikic, S., McPherson, A. W., Donmez, N., and Sahinalp, C. S., 2015. Clonality Inference in Multiple Tumor Samples using Phylogeny. *Bioinformatics*, **31**(9):1349–1356.
- Malinen, M., Jääskeläinen, T., Pelkonen, M., Heikkinen, S., Väisänen, S., Kosma, V.-M., Nieminen, K., Mannermaa, A., and Palvimo, J. J., 2013. Proto-oncogene PIM-1 is a novel estrogen receptor target associating with high grade breast tumors. *Molecular and Cellular Endocrinology*, **365**(2):270–276.
- Manning, C. D., Raghavan, P., and Schütze, H., 2008. *Introduction to information retrieval*. Cambridge University Press, New York.
- Marusyk, A., Tabassum, D. P., Altrock, P. M., Almendro, V., Michor, F., and Polyak, K., 2014. Non-cell-autonomous driving of tumour growth supports sub-clonal heterogeneity. *Nature*, **514**(7520):54–58.
- McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernysky, A., Garimella, K., Altshuler, D., Gabriel, S., Daly, M., *et al.*, 2010. The Genome Analysis Toolkit: A MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Research*, **20**(9):1297–1303.
- Miller, C. A., White, B. S., Dees, N. D., Griffith, M., Welch, J. S., Griffith, O. L., Vij, R., Tomasson, M. H., Graubert, T. A., Walter, M. J., *et al.*, 2014. SciClone: Inferring Clonal Architecture and Tracking the Spatial and Temporal Patterns of Tumor Evolution. *PLoS Comput Biol*, **10**(8):e1003665.
- Molenaar, J. J., Koster, J., Zwijnenburg, D. A., van Sluis, P., Valentijn, L. J., van der Ploeg, I., Hamdi, M., van Nes, J., Westerman, B. A., van Arkel, J., *et al.*, 2012. Sequencing of neuroblastoma identifies chromothripsis and defects in neuritogenesis genes. *Nature*, **483**(7391):589–593.

- Ng, P. C. and Henikoff, S., 2003. SIFT: predicting amino acid changes that affect protein function. *Nucleic Acids Research*, **31**(13):3812–3814.
- Qiao, Y., Quinlan, A. R., Jazaeri, A. A., Verhaak, R., Wheeler, D. A., and Marth, G. T., 2014. Subclone-Seeker: a computational framework for reconstructing tumor clone structure for cancer variant interpretation and prioritization. *Genome Biology*, **15**(8):443.
- Roth, A., Khattra, J., Yap, D., Wan, A., Laks, E., Biele, J., Ha, G., Aparicio, S., Bouchard-Côté, A., and Shah, S. P., *et al.*, 2014. PyClone: statistical inference of clonal population structure in cancer. *Nature Methods*, **11**(4):396–398.
- Schramm, A., Köster, J., Assenov, Y., Althoff, K., Peifer, M., Mahlow, E., Odersky, A., Beisser, D., Ernst, C., Henssen, A. G., *et al.*, 2015. Mutational dynamics between primary and relapse neuroblastomas. *Nature Genetics*, **47**(8):872–877.
- Schwarz, R. F., Trinh, A., Sipos, B., Brenton, J. D., Goldman, N., and Markowitz, F., 2014. Phylogenetic Quantification of Intra-tumour Heterogeneity. *PLoS Computational Biology*, **10**(4):e1003535.
- Xiao, X., Liu, Z., Wang, R., Wang, J., Zhang, S., Cai, X., Wu, K., Bergan, R. C., Xu, L., and Fan, D., *et al.*, 2014. Genistein suppresses FLT4 and inhibits human colorectal cancer metastasis. *Oncotarget*, **6**(5):3225–3239.
- Yu, G., Lee, Y.-C., Cheng, C.-J., Wu, C.-F., Song, J. H., Gallick, G. E., Yu-Lee, L.-Y., Kuang, J., and Lin, S.-H., 2015. RSK Promotes Prostate Cancer Progression in Bone through ING3, CKAP2 and PTK6-mediated Cell Survival. *Molecular cancer research : MCR*, **13**(2):348–357.