

Detecting Somatic Mutations Without Matched Normal Samples Using Long Reads

Jared T. Simpson ^{*1,2,3}

¹Ontario Institute for Cancer Research, Toronto, Canada

²Department of Molecular Genetics, University of Toronto, Toronto, Canada

³Department of Computer Science, University of Toronto, Toronto, Canada

Abstract

DNA sequencing of tumours to identify somatic mutations has become a critical tool to guide the type of treatment given to cancer patients. The gold standard for mutation calling is comparing sequencing data from the tumour to a matched normal sample to avoid mis-classifying inherited SNPs as mutations. This procedure works extremely well, but in certain situations only a tumour sample is available. While approaches have been developed to find mutations without a matched normal, they have limited accuracy or require specific types of input data (e.g. ultra-deep sequencing). Here we explore the application of single molecule long read sequencing to calling somatic mutations without matched normal samples. We develop a simple theoretical framework to show how haplotype phasing is an important source of information for determining whether a variant is a somatic mutation. We then use simulations to assess the range of experimental parameters (tumour purity, sequencing depth) where this approach is effective. These ideas are developed into a prototype somatic mutation caller, *smrest*, and its use is demonstrated on two highly mutated cancer cell lines. Finally, we argue that this approach has potential to measure clinically important biomarkers that are based on the genome-wide distribution of mutations: tumour mutation burden and mutation signatures.

1 Introduction

Large scale efforts to sequence thousands of cancer genomes has been a major undertaking since the invention of high-throughput DNA sequencing instruments. These projects have catalogued driver mutations (Weinstein et al. 2013; Bailey et al. 2018; Rheinbay et al. 2020; “Pan-cancer analysis of whole genomes” 2020), defined mutational processes and the signatures they leave on the genome (Alexandrov, Nik-Zainal, et al. 2013; Alexandrov, Kim, et al. 2020; Nik-Zainal et al. 2012), tracked the evolutionary trajectories of tumours (Shah et al. 2009; Sottoriva, Spiteri, et al. 2013; Sottoriva, Kang, et al. 2015; Gerstung et al. 2020), uncovered the originating cell and tissue types of tumours (Jiao et al. 2020; Hendrikse et al. 2022) and discovered mutation-based biomarkers that can guide treatment choice (Chan et al. 2019; Zhao et al. 2019; André et al. 2020). Underpinning these studies was the development of both the high throughput sequencing instruments and analysis methods that can accurately detect mutated bases from the sequenced reads. Unlike in typical human genome sequencing projects, where one typically wants to discover variation within a sequenced genome compared to a reference genome, cancer genome projects aim to discover *somatic* mutations that occurred during the development and progression of a tumour. This requires finding genetic

*jsimpson@oicr.on.ca

38 variation within an individual, where a subset of cells contain a mutated copy of a chromosome
39 with respect to the chromosome inherited from the individual’s parents (see **Figure1a**).

40 The gold standard method of detecting somatic mutations is to compare sequencing data from
41 the tumour to a matched “normal” sample that is assumed to be representative of the individual’s
42 inherited genome. This method, referred to as tumour-normal calling, is highly accurate and widely
43 used (Koboldt et al. 2009; Larson et al. 2012; Saunders et al. 2012; Cibulskis et al. 2013; Fang et al.
44 2021). Tumour-only calling methods have also been developed for situations where a normal sample
45 is not available (e.g. for biobanked tissue samples where a normal was not collected, or to simplify
46 clinical workflows). These methods typically rely on analysis of the fraction of reads supporting a
47 candidate mutation, as this may differ from the fraction supporting inherited heterozygous SNPs,
48 usually coupled with extensive filtering of the candidate calls against databases of variants known
49 to occur in the human population (Smith et al. 2016; Kalatskaya et al. 2017; Sun et al. 2018). While
50 this strategy can be effective in certain situations, like intermediate purity tumours that are very
51 deeply sequenced (Sun et al. 2018), it is inherently less accurate than tumour-normal pair calling,
52 and filtering against population database raises issues of bias for underrepresented populations
53 (Nassar et al. 2022).

54 Thus far, cancer genome sequencing projects have primarily relied on highly accurate short read
55 sequencing technologies. Long read sequencing technologies from Oxford Nanopore Technologies
56 (ONT) and Pacific Biosciences (PacBio) are increasingly accurate (Sereika et al. 2022; Kolesnikov
57 et al. 2023), and improvements to instrument throughput have expanded the range of possible
58 applications. Long read sequencing is now the gold standard for genome assembly (Rhie et al.
59 2021; Nurk et al. 2022). Both ONT and PacBio sequencers can measure the genome and epigenome
60 simultaneously by directly detecting base modifications (Flusberg et al. 2010; Laszlo et al. 2013;
61 Schreiber et al. 2013; Simpson et al. 2017). Recently, a tumour-normal calling approach has been
62 developed for long reads that has comparable accuracy to short read sequencing (Zheng, Su, et al.
63 2023).

64 In this paper, we explore the potential for using long read sequencing to perform tumour-only
65 mutation calling. The advantage of long reads for this problem is that reads can often be assigned
66 to individual haplotypes, transforming the problem of detecting somatic mutations from potentially
67 small shifts in the variant allele fraction, into detecting the presence of two or more bases within
68 a single haplotype, as proposed in the mosaic variant detection method by Darby et al. (2019)
69 for 10X Genomics linked reads. In this work we formalize the problem and use simulations to
70 assess the applicability of this approach as a function of key experimental parameters (sequencing
71 depth, tumour purity, sequencing error rate). Then, we develop a prototype mutation caller, *smrest*
72 (somatic mutation rate estimator), for real long read data and demonstrate its use on cancer cell
73 lines. Finally we present the intended application of this tool, which is to discover genome-wide
74 mutation patterns that can be used to guide therapy choice, such as tumour mutation burden or
75 mutation signatures.

76 **2 Results**

77 **2.1 Overview of Method and Feasibility**

78 Somatic mutations are by definition mosaic; they occur in a subset of cells within the human body.
79 When a particular cell becomes cancerous the complement of somatic mutations contained within
80 that originating cell, and subsequent mutations that occur during the tumour’s expansion, rise in

81 frequency (**Figure 1a**). These mutations can be detected by comparing DNA sequences from a
82 tumour sample to a blood or non-cancerous tissue from the same individual (tumour-normal calling,
83 reviewed in Xu 2018). The subject of this paper however is calling mutations in absence of the
84 matched normal sample, referred to as “tumour-only” calling. This problem has been studied for
85 short read sequencing and, most relevant to this work, Darby et al. 2019 introduced the idea of
86 using haplotype phasing patterns for 10X Genomics Linked Reads.

87 Importantly for these prior approaches, and essential to this work, is that real tumour samples are
88 typically mixtures of both cancerous and healthy cells (**Figure 1b**). The fraction of cancerous
89 cells is referred to as tumour purity or tumour cellularity which we will denote α . The presence
90 of normal cells provides evidence of the allele that the individual inherited from their parent. In
91 a sequencing experiment this may shift the variant allele fraction (VAF; the proportion of reads
92 supporting the putative mutation) away from the ratio expected of heterozygous SNPs (0.5 for copy
93 number balanced autosomes, as each parental haplotype is equally likely to be sampled by a read).
94 The ability to detect mutations from the shift in VAF alone is strongly dependent on sequencing
95 depth, sequencing error rate and tumour purity (**Figure 1c**). For example, when tumours are
96 nearly pure ($\alpha \approx 1$) somatic mutations are nearly indistinguishable from heterozygous SNPs and
97 at the other extreme, where few cells in the sample derive from the tumour ($\alpha \approx 0$), somatic
98 mutations will look like sequencing errors. If the parental haplotype (maternal or paternal) of each
99 read is known, then the problem of identifying somatic mutations simplifies to identifying whether
100 there is sufficient evidence of a mixture of alleles within a haplotype (one inherited allele from the
101 normal cells, and the mutant allele from the cancerous cells). **Figure 1d** illustrates how phased
102 reads bearing evidence of a variant can help identify whether it is a heterozygous SNP, somatic
103 mutation or sequencing error.

104 To explore the feasibility of somatic mutation calling from long read sequencing we developed sta-
105 tistical models for phased and unphased data to calculate the posterior probability that a given
106 position of the genome harbors a somatic mutation based on the number of reads supporting
107 the reference and alternate allele, along with the tumour’s purity and sequencing error rate (see
108 **Methods**). We then generated simulated data according to this model to explore the relation-
109 ship between sequencing depth, error rate and purity (**Figure 1e**). In all parameter sets tested
110 haplotype-phased data provides equal or superior variant calling accuracy. For a fixed depth and
111 sequencing error rate this enables accurate somatic mutation calling across a wider range of tumour
112 purity. For example, at 80x sequencing depth with 1% error rate the somatic mutation calling F1
113 exceeded 0.9 for simulated samples with purity in the range 0.23 – 0.84 when the data was phased
114 but only 0.42 – 0.52 when the data is not phased.

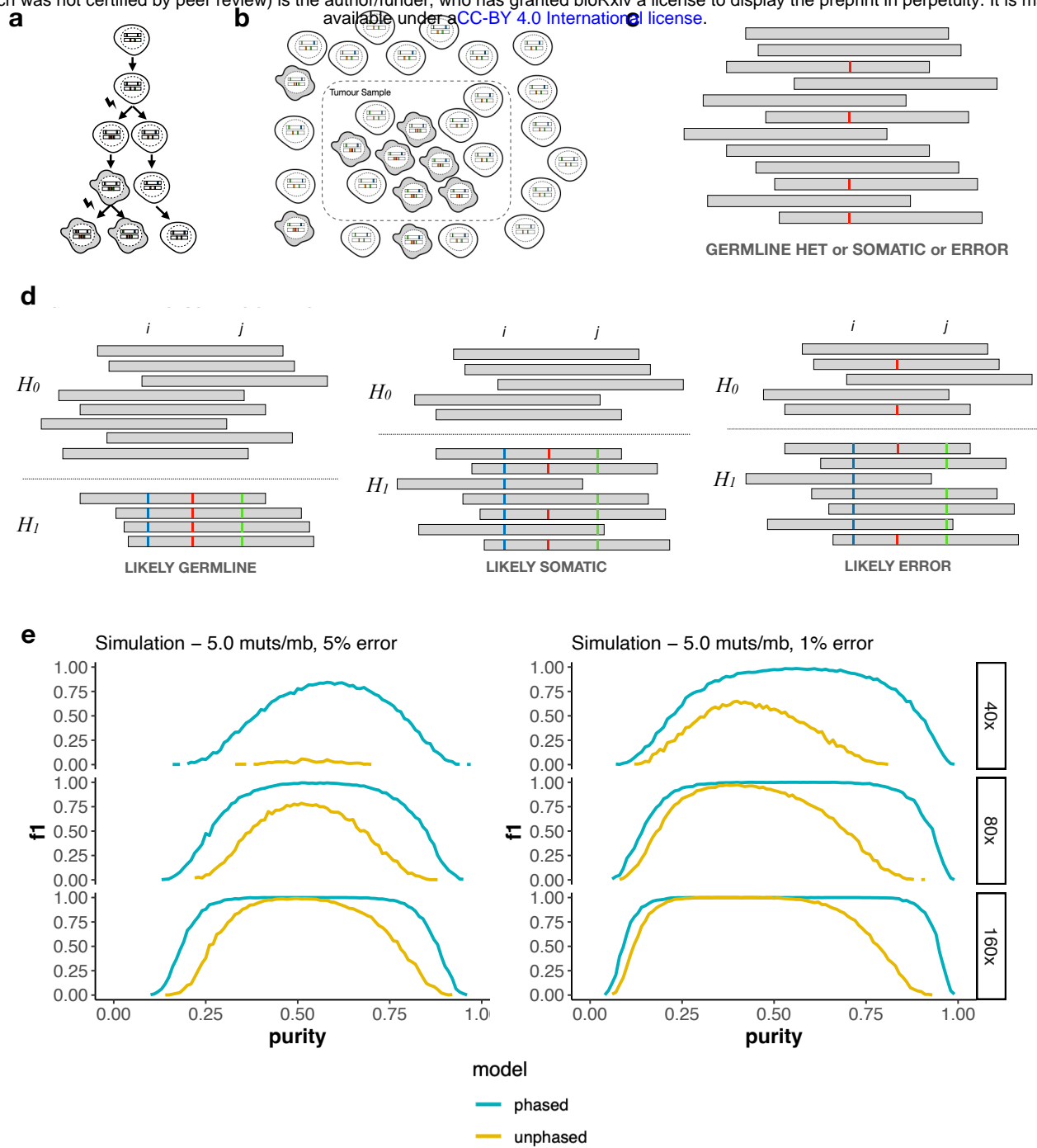


Figure 1: **a**. Cell lineages accumulate somatic mutations over time. When a healthy cell (white) becomes malignant (grey) the mutations contained within that lineage rise in frequency. Additional mutations within the tumour population may generate subclones, where the mutations are contained in a subset of the tumour cell population. **b**. Tumour samples are typically a mixture of cancerous cells and normal cells. The proportion of tumour cells is called the tumour *purity* and denoted α . When extracting DNA from the tumour sample the somatic mutations will present as mosaics with only some of the sequenced reads supporting the mutation. **c**. Variant callers that do not use phasing information need to make a decision based on the number of reads (grey bars) that support the alternate allele (red) and reference allele. Depending on the number of reads supporting the alternate allele it may be ambiguous whether the position is a heterozygous SNP, sequencing error or somatic mutation. **d**. If the haplotype of each read (H_0 , H_1) can be determined using nearby heterozygous SNPs (positions i and j with blue and green variants, respectively) the arrangement of reads supporting the alternate allele (red) can help determine whether it is a somatic mutation. We expect all phased reads to support the alternate allele when the position is a heterozygous SNP (left) or to support a mixture of reference and alternate alleles when it is a somatic mutation (middle). When the position contains a sequencing error we do not expect the evidence of the error to segregate by haplotype (right). **e**. Simulation results demonstrating that using haplotype phasing information can classify somatic mutations more accurately across a range of tumour purity, read depth and error rates.

115 2.2 Mutation Calling on Single Molecule Long Reads

116 The simulations presented above demonstrate that haplotype phasing can improve somatic mu-
117 tation calling accuracy. These simulations assume ideal data however, where every read can be
118 correctly assigned to a haplotype and the reads are perfectly aligned to the reference genome. As a
119 further proof of principle, we developed and benchmarked a mutation caller, called *smrest* (somatic
120 mutation rate estimator) for real data. Briefly, this program first detects heterozygous SNPs, phases
121 them, partitions reads into haplotypes (haplotagging), then calls mutations using a procedure that
122 extends the simulation model with a probabilistic alignment framework that is widely used in other
123 variant callers (Albers et al. 2011; H. Li 2011; Garrison and Marth 2012; Poplin et al. 2018; Cooke,
124 Wedge, and Lunter 2021). Finally, putative somatic variants are filtered to remove likely artifacts
125 (e.g. from mapping errors, or systematic sequencing errors) and the resulting mutations that pass
126 quality control are output as a VCF file. The mutation caller also generates a BED file describ-
127 ing the regions of the genome that were deemed callable (both haplotypes detected with sufficient
128 sequencing depth). A high-level presentation of the mutation calling procedure is provided here,
129 with complete details in the **Methods**.

130 **Genotyping and Phasing.** The input to standard haplotype phasing algorithms is a set of
131 heterozygous SNPs (Patterson et al. 2015). While variant calling and genotyping long read data
132 from human genomes can now be done with high accuracy (Zheng, S. Li, et al. 2022) cancer genomes
133 may have copy number imbalances that shift the read support of heterozygous SNPs away from
134 the expectation of equal support for the paternal and maternal alleles. In the worst case where one
135 parental haplotype is entirely lost, known as loss-of-heterozygosity (LOH), the evidence of the lost
136 allele will only come from the normal cells within a sample and hence be a function of the tumour’s
137 purity. To account for these factors we developed a genotyper for known variant sites that relaxes
138 the assumption of copy number balance. This method is designed to be conservative and prefer
139 false negatives over false positives, as the latter is more likely to introduce erroneous haplotype
140 assignments. The set of heterozygous SNPs found by this procedure is used as input to *whatshap*
141 (Patterson et al. 2015) for phasing.

142 **Mutation Calling and Filtering.** The phased VCF file, and the original BAM file containing
143 the reads mapped to the human reference genome, are input into the mutation calling program
144 *smrest call*. This program haplotags each read, discovers candidate variants, then calculates class
145 probabilities for each candidate using a read-haplotype likelihood model (see **Methods**). Finally,
146 summary statistics are gathered for each output call to facilitate mutation filtering, e.g. when
147 variants show evidence of strand bias.

148 2.3 Mutation Calling in COLO829

149 To characterize the performance of *smrest* for long read tumour-only somatic mutation calling we
150 first analyzed COLO829, a melanoma cell line with a matched normal (COLO829BL). In all exper-
151 iments we mixed reads from the tumour and normal together, without knowing the origin of each
152 read, to simulate tumour-only sequencing with a controlled purity and sequencing depth. COLO829
153 has a high mutation rate of ≈ 14 mutations per megabase (Titmuss et al. 2022) due to ultravio-
154 let light damage (Plesance et al. 2010), and is frequently used to benchmark the performance of
155 sequencing technologies and analysis methods (Arora et al. 2019; Espejo Valle-Inclan et al. 2022).
156 We downloaded Oxford Nanopore R10.4.1 reads for both COLO829 and COLO829BL from the
157 Oxford Nanopore Open Datasets Collection. The tumour and normal sample were sequenced to
158 95x and 57x sequencing depth, respectively. The reads were prepared with a mixture of standard

159 and ultra-long library preparations with 7.5x and 7.1x coverage of reads exceeding 100kbp in length,
160 simplifying haplotype phasing. We also downloaded Illumina short reads for this pair of samples.
161 Access information for all datasets used in this work are provided in **Table 1**.

162 To provide context to the performance of smrest we also ran ClairS (Zheng, Su, et al. 2023), a
163 recently developed tumour-normal pair (abbreviated TN hereafter) caller based on neural networks
164 on the long read data. For short reads we ran Mutect2 (Benjamin et al. 2019), which is designed for
165 TN calling but also supports tumour-only (abbreviated TO) calling via filtering against population
166 databases and a panel-of-normals to remove sequencing artifacts. It is included here to compare
167 the performance of smrest against a short read method that is commonly used when a matched
168 normal sample is not available.

169 To prepare the datasets both the tumour and normal cell line samples were downsampled to coverage
170 of 20x or 40x. For the TO methods the read sets were merged together to give a single sample with
171 equal coverage of the tumour and normal (50% purity). The tumour purity parameter to smrest
172 was set to this known value ($\alpha = 0.5$). For the TN callers the tumour and normal read sets were
173 provided as separate inputs.

174 The mutation calls for all programs were compared to an externally curated call set (**Methods**)
175 that was considered the ground truth. We discarded all mutations (in the truth data, or any call set
176 from any program) with predicted VAF $< 10\%$ as low VAF mutations cannot be reliably called at
177 the sequencing depths used for this analysis. In addition, we only considered substitution variants.
178 As cell lines accumulate mutations (Petljak et al. 2019) that will appear as somatic mutations in
179 tumour-only analysis we identified putative mutations in COLO829BL using the short read datasets
180 and Mutect2 in TN mode by swapping the tumour and normal inputs, and filtering the results
181 to avoid calling inherited SNPs as mutations in loss-of-heterozygosity regions (**Methods**). Any
182 called mutations found in this COLO829BL mutation list were ignored for subsequent analysis. We
183 characterized the accuracy of each program on each dataset in the regions of the genome designated
184 as “high confidence” (denoted HC, **Figure 2a**), and the subset of the HC regions that could be
185 successfully called and phased by smrest (**Figure 2b**) as precision-recall curves stratified by a
186 variant’s QUAL score (for smrest and ClairS) or TLOD (for mutect2-TN and mutect2-TO).

187 As expected the variant callers that use a tumour-normal strategy perform extremely well, even
188 for the lowest depth dataset (20x/20x). For tumour-only calling the mutect2-based approach can
189 achieve high recall, but with limited precision due to the difficulty of identifying somatic muta-
190 tions from short reads even when provided with a database of population variants. In contrast,
191 smrest achieves high precision without a matched normal and without filtering against population
192 databases. This method is not a replacement for tumour-normal calling however. When the entire
193 set of high-confidence regions is considered (spanning 2,080Mbp of human reference GRCh38) recall
194 ranges from 0.61 (20x/20x) to 0.86 (40x/40x). Sensitivity is primarily determined by the ability to
195 phase the genome, as recall improves to 0.87-0.96 when only successfully phased high-confidence
196 regions are considered. This highlights the critical need for genome-wide phasing if this approach
197 is to be used for comprehensive somatic mutation detection.

198 Our method was primarily developed, tested, and debugged on COLO829 data sequenced with
199 Oxford Nanopore long reads. To assess whether our approach generalizes, we additionally analyzed
200 COLO829 PacBio HiFi data, and a highly mutated breast cancer cell line (HCC1395/HCC1395BL)
201 sequenced with both ONT R10.4.1 and PacBio HiFi reads. The accuracy of each mutation call
202 set was calculated as above, however here we computed the intersection of the ONT and PacBio-
203 HiFi phased regions to ensure all call sets were analyzed over the same regions of the genome. For

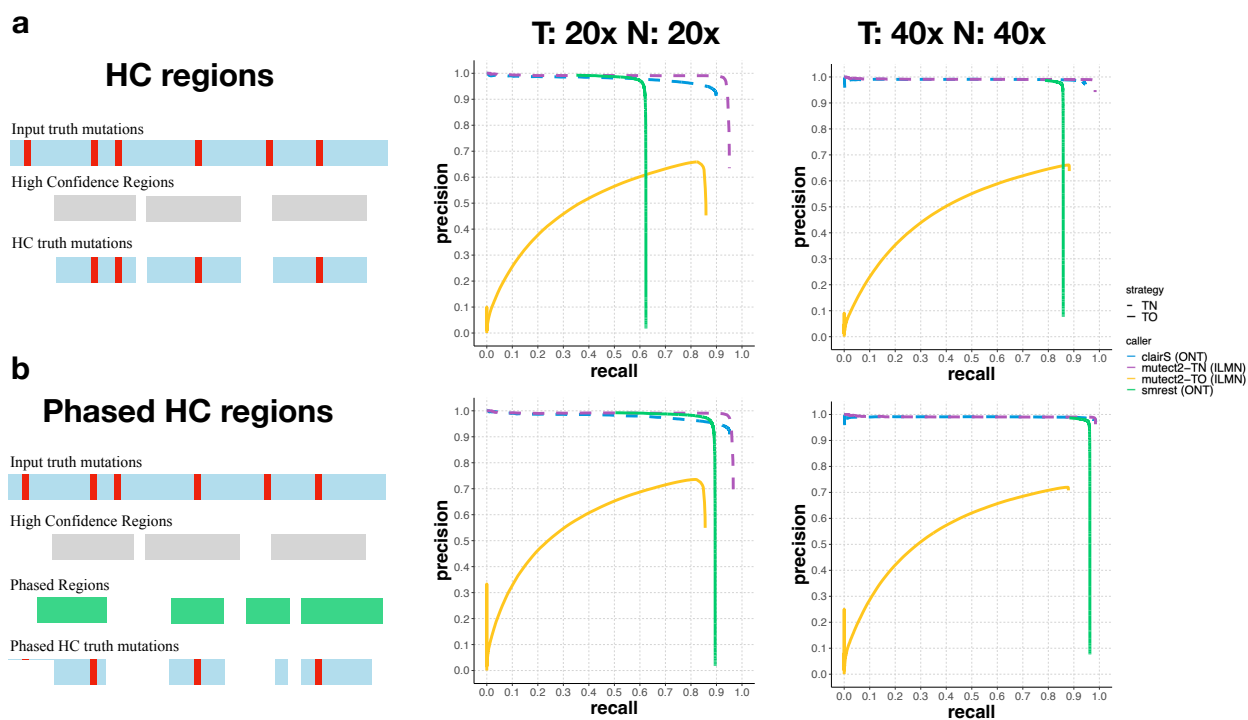


Figure 2: Somatic mutation calling performance on COLO829. Short and long read datasets were prepared with 20x tumour and 20x normal coverage and 40x tumour and 40x normal coverage. Each data set was provided as input to smrest (green solid line), clairS (blue dashed line), mutect2 in TN mode (purple dashed line) or mutect2 in TO mode (yellow solid line). Precision/recall curves were calculated using an external call set as ground truth. Two subsets of the mutation calls were analyzed, one consisting of calls that lie in regions deemed High Confidence (**Panel a**, top row) and one consisting of calls that lie in the HC regions that were successfully phased by smrest at the given coverage (**Panel b**, bottom row).

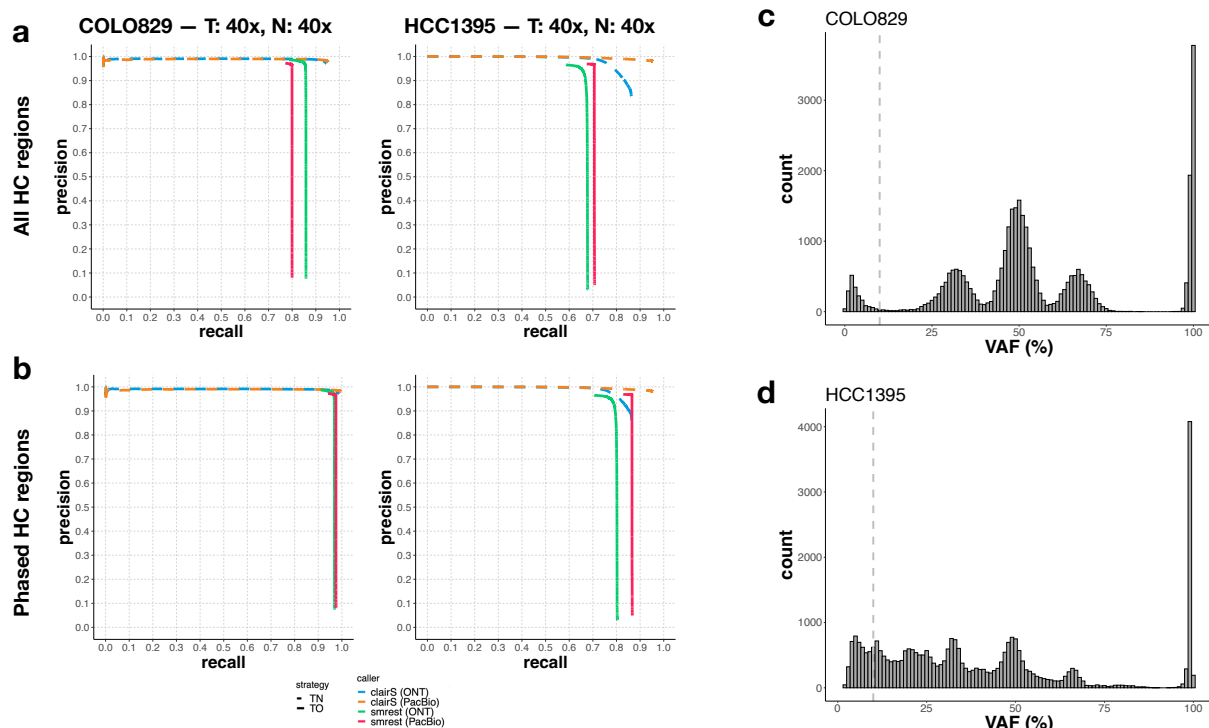


Figure 3: Somatic mutation calling performance on 40x tumour and normal coverage of COLO829 (left column of panel **a.** and **b.** and HCC1395 (right column) with ONT and PacBio data. As in the previous figure, callset accuracy was calculated over the High Confidence regions (**a.**), or the phased subset of the HC regions (**b.**). Each data set was provided as input to smrest (ONT: green solid line, PacBio: red solid line) or clairS (ONT: blue dashed line, PacBio: orange dashed line). The distribution of variant allele frequencies for the truth call set for each sample is shown in (COLO829: **c.**, HCC1395: **d.**), with the lower VAF threshold for inclusion within the analysis annotated as a vertical dashed line.

204 COLO829, the Oxford Nanopore call set from smrest had slightly higher recall than the PacBio call
 205 set (**Figure 3a**), likely due to the presence of ultra-long reads that simplify the construction of long
 206 range haplotypes. When computing accuracy over the phased regions of the genome the ONT and
 207 PacBio data were both highly accurate. ClairS performed very well with data from each technology.
 208 The HCC1395/HCC1395BL cell line is more challenging than COLO829 due to the presence of a
 209 long tail of low VAF mutations (**Figure 3c,d**) possibly indicating multiple subclones. While smrest
 210 was able to precisely call mutations for both the ONT and PacBio datasets, the PacBio dataset had
 211 higher recall, likely due to the higher accuracy allowing lower frequency mutations to be identified.
 212 This trend is also seen in ClairS where the recall on PacBio data was somewhat higher than on
 213 ONT reads.

214 2.4 Assessing the Effect of Tumour Purity

215 In our simulations (**Figure 1e**) we observed that variant calling performance is a function of tumour
 216 purity. We therefore performed a series of experiments to assess this effect on the real datasets. Here
 217 we selected tumour purity in the range 0-100%, then computed the number of COLO829 (tumour)
 218 and COLO829BL (normal) reads needed to reach a total of 40x or 80x coverage at the selected
 219 purity. Each mutation calling program was run and analyzed as described above. For this analysis
 220 smrest was not provided the known value of tumour purity, it was left at its default value of $\alpha = 0.5$.

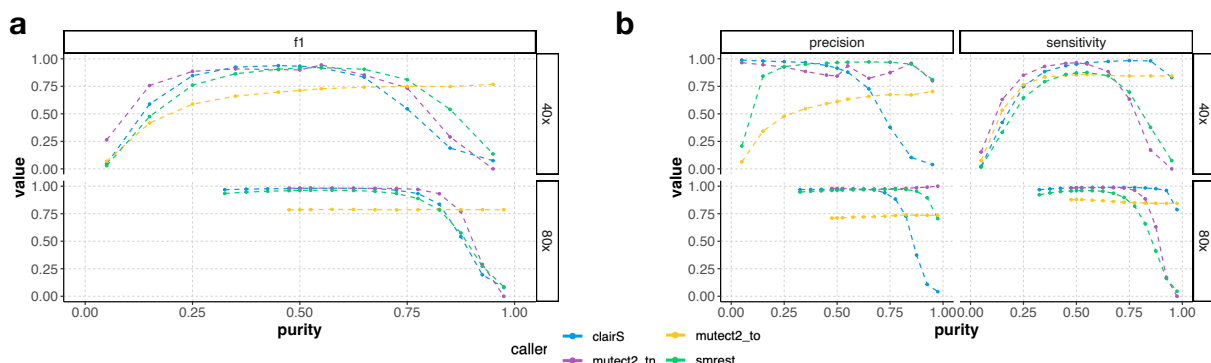


Figure 4: Somatic mutation calling accuracy (**Panel a:** F1, **Panel b:** precision and sensitivity) metrics as a function of tumour purity for 40x (top facet within each panel) and 80x (bottom) coverage sequencing of COLO829/COLO829BL with smrest (green), clairS (blue), mutect2-TN (purple) or mutect2-TO (yellow)

221 In addition, only the phased HC regions were assessed. The results are shown in **Figure 4**. As
222 expected and seen in the simulations, the performance of all mutation calling approaches is limited
223 at extreme values of tumour purity as at low purity there is insufficient tumour depth to confidently
224 identify mutations and at high purity there is insufficient normal depth to confidently say what
225 the inherited allele is. Our program achieves comparable performance to the TN callers and in
226 particular maintains high precision across a wide range of tumour purities, at both 40x and 80x
227 coverage. In contrast, short read TO calling fails to achieve high precision. These results highlight
228 the importance of sequencing depth as 80x coverage allows a much wider range of (simulated)
229 sample purities to be reliably analyzed, as predicted by our simulations.

230 2.5 Estimating Tumour Mutation Burden and the Mutation Spectrum

231 Cancer genome sequencing is increasingly used to help guide treatment choices. While there are
232 now many known point mutations, indels and structural variants that may indicate response or
233 resistance to certain therapies (Krysiak et al. 2023), an emerging class of biomarker is based on the
234 overall pattern of mutation across the genome. Perhaps most prominently, the tumour's mutation
235 burden (TMB; the number of observed coding mutations per megabase of coding sequence) is used
236 for predicting response to immunotherapies, for example pembrolizumab (Marabelle et al. 2020)
237 for tumours classified as TMB-high (≥ 10 mutations/MB, Marcus et al. 2021). PARP inhibitors
238 are used to treat patients with mutations in the homologous repair pathways (Patel, Sarkaria,
239 and Kaufmann 2011; Miller et al. 2020), which can be detected with high accuracy using mutation
240 signatures (Davies et al. 2017). As both TMB and mutation signatures are genome-wide phenomena
241 they can be determined by sequencing only a subset of the genome (Milbury et al. 2022).

242 While we have demonstrated that in certain situations our approach can be highly sensitive, the
243 requirement of phasing every base of the genome that might harbour a mutation of interest may
244 limit the application of this approach for general purpose somatic mutation detection. Therefore
245 our primary intended use case is to detect biomarkers that can be found from accurate mutation
246 calling in defined regions of the genome. In our case, we propose to calculate TMB from the
247 phased subset of the high confidence regions, by dividing the number of called mutations by the
248 total size of the genome phased. To assess the feasibility of this strategy we performed additional
249 mixture experiments where reads from COLO829 and COLO829BL were again merged together

250 into a single sample with a defined purity (30%-70%, in steps in 10%) and depth (10x to 82x, in
 251 steps of 8x). Here we considered all mutations that passed our filtering criteria with a minimum
 252 quality score of 20 to be called. In this analysis we included a healthy blood sample openly released
 253 by Oxford Nanopore to assess the false positive rate in a sample that is assumed to be free from
 254 somatic mutations. In **Figure 5a**, we observe that the estimated mutation burden converges to the
 255 expected value of 14 mutations per megabase of analyzed sequenced (Titmuss et al. 2022), with the
 256 exception of the 30% purity sample which underestimates TMB at 13.1 muts/MB at the maximum
 257 analyzed depth of 74x. In the healthy blood sample few mutations were called (maximum of 0.2
 258 muts/MB at 66x).

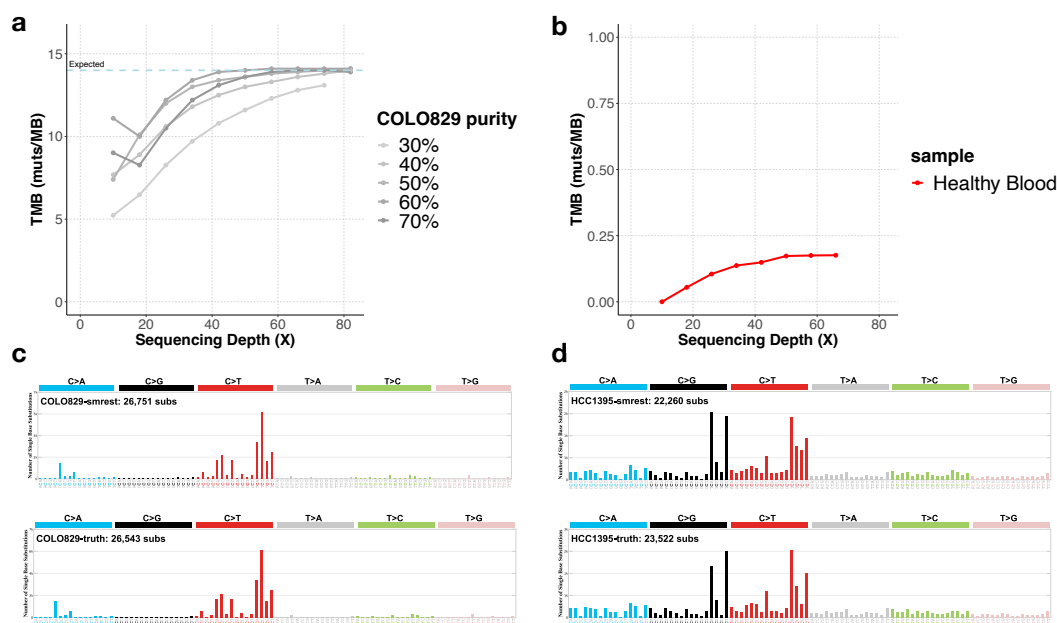


Figure 5: **Panel a** presents the results of estimating tumour mutation burden (TMB) on COLO829 with varying tumour purity (each line series, with darker lines having higher purity) and sequencing depth (x-axis). **Panel b** shows a healthy blood sample (red) as a negative control where few mutations are expected (note different y-axis range than **panel a**). **Panels c** (COLO829) and **d** (HCC1395) show the sequence contexts of called mutations (known as the mutation spectrum) for smrest (top facet) and the ground truth call set (bottom facet) from the 40x/40x experiments.

259 Mutation signature profiling uses a vector of 96 unique sequence contexts (six different mutation
 260 types, each with a single preceding and following base) to determine the relative contribution of
 261 mutagenic processes, like UV damage or DNA repair deficiencies (Alexandrov, Nik-Zainal, et al.
 262 2013; Davies et al. 2017). As this analysis requires accurate determination of mutation counts for
 263 each sequence context, we assessed the spectrum of mutations found by our program compared to
 264 the spectrum of the ground truth data for the 40x/40x datasets for COLO829 (**Figure 5b**) and
 265 HCC1395 (**Figure 5c**). The mutation spectrum from smrest is highly similar to that of the ground
 266 truth data (cosine similarity > 0.99) and predominantly consists of C>T mutations in the context
 267 TC>TT as expected of a sample with UV damage. Similarly, the mutation spectrum for HCC1395
 268 is consistent with the ground truth (cosine similarity > 0.98).

269 3 Discussion

270 In this work we analyzed the potential for calling somatic mutations using single molecule long read
271 sequencing of tumour samples without matched normal samples. Our results suggest that when
272 certain conditions are met - most importantly when tumour purity is within a band determined
273 by sequencing depth - that accurate mutation calling is possible. However, there are limitations
274 to this study that will need to be further explored. The samples sequenced here are all cell lines,
275 where sufficient DNA quantity for long read (and ultra long read) sequencing is easily achieved.
276 Sequencing real tumour samples, particularly solid tumours, will be more challenging and require
277 extensive protocol optimization, or accepting a shorter read length, which will impact the amount
278 of the genome that can be phased and called using this approach. In addition, the cell lines we
279 used are both highly mutated and hence favourable for calculating accuracy statistics given the very
280 large number of true positive mutations found. This approach must also be assessed on tumour
281 samples with a lower mutation rate, although the analysis of the healthy blood sample presented
282 in Figure 5 suggests a low false positive rate.

283 The algorithms used here are a proof-of-concept that can be improved in a number of ways. Most
284 notably, we do not incorporate allele specific copy number in our mutation classification model
285 unlike in other methods (e.g. Sun et al. 2018). Also, we used a default purity parameter rather than
286 jointly estimating purity and copy number (Carter et al. 2012; Cameron, Baber, Shale, Papenfuss,
287 et al. 2019) as is common in many other approaches. Similarly we do not attempt to infer a cancer
288 cell fraction distribution for subclonal mutations. We use `whatshap` to phase the heterozygous
289 SNPs but this is not designed to account for allele specific copy number changes in cancer, which is
290 an additional source of information. Future work aims to incorporate these improvements, as well
291 as support other mutation types like short indels, and the inference of microsatellite instability,
292 which can be predictive of response to checkpoint inhibitors (K. Li et al. 2020).

293 4 Methods

294 4.1 Simulations

295 4.1.1 Definitions and notation

- 296 • μ := somatic mutation rate (probability a given base of the genome is mutated in the tumour)
- 297 • π := heterozygous SNP rate (probability a given base is a SNP, here fixed at 1/1000)
- 298 • α := tumour purity
- 299 • ϵ := sequencing error rate
- 300 • a_i := number of reads containing the alternate allele at position i of the genome
- 301 • r_i := number of reads containing the reference allele at position i of the genome
- 302 • $H_j[i]$:= the i -th base on haplotype j ($S_j[i]$ defined similarly)

303 4.1.2 Simulated data generation for clonal tumours

304 For each simulation the average sequencing depth (λ) and tumour purity (α) are input parameters.
305 To generate the simulated data for a parameter setting (λ , α) the following procedure is used.
306 First, the genome size G is set to a constant value (here, $G=100,000,000$) then two haplotype

307 strings (H_0, H_1) are initialized from the first G bases of human chromosome 2. Each haplotype
 308 base was randomly mutated with probability π to simulate germline SNPs. A somatic copy of each
 309 haplotype (S_0, S_1) was made and mutated at rate μ .

310 For every position i , we draw the total sequencing depth d_i from a Poisson(λ) distribution and
 311 partition the depth across the two haplotypes by drawing $d_{i,0} \sim \text{Binom}(d_i, 0.5)$ and setting $d_{i,1} =$
 312 $d_i - d_{i,0}$. Then we assign a base to each read. If a position contains a somatic mutation ($S_j[i] \neq H_j[i]$)
 313 then the base is set to $S_j[i]$ with probability α (the chance of sampling a cancer cell from the mixture)
 314 otherwise it is set to $H_j[i]$. Finally, the drawn base is randomly changed to one of the three other
 315 bases with probability ϵ to simulate random sequencing errors. The simulated data is aggregated
 316 into a vector containing the number of observed A, C, G and T bases. These vectors were passed
 317 into the classifiers described below.

318 4.1.3 Simulated data classifiers

319 **Unphased data.** The classifier for unphased data considers three possibilities for each position i ,
 320 $c_i \in \{\text{somatic}, \text{het}, \text{reference}\}$, using the number of reads supporting the reference base, denoted
 321 r_i , and the number of reads supporting the non-reference base with highest read count, a_i . It is
 322 possible for a_i to be 0 if all reads support the reference base. The classifier for phased data is
 323 similar but the classifications and observed data are defined per haplotype j ($c_{i,j}, r_{i,j}, a_{i,j}$).

324 We calculate the posterior probability a site contains a somatic mutation after observing a_i alternate
 325 bases and r_i reference bases as:

$$P(c_i = \text{somatic} | a_i, r_i) = \frac{P(a_i, r_i | c_i = \text{somatic})P(c_i = \text{somatic})}{\sum_c P(a_i, r_i | c_i = c)P(c_i = c)} \quad (1)$$

326 The likelihood term accounts for sequencing errors using a binomial observation model. To observe
 327 a read with an alternative base the read must either sample a mutated haplotype (with probability
 328 $\alpha/2$) with a correct basecall ($1 - \epsilon$), or sample a non-mutated haplotype (with probability $1 - \alpha/2$)
 329 with a base that erroneously supports the alternative allele (ϵ). By summing these cases we have
 330 the chance of observing an alternative base given the position contains a somatic mutation:

$$P(a_i, r_i | c_i = \text{somatic}) = \text{Bin}(a_i, a_i + r_i, (1 - \epsilon)\frac{\alpha}{2} + \epsilon(1 - \frac{\alpha}{2})) \quad (2)$$

331 The likelihoods for the heterozygous and reference classifications are similar:

$$P(a_i, r_i | c_i = \text{het}) = \text{Bin}(a_i, a_i + r_i, \frac{(1 - \epsilon)}{2} + \frac{\epsilon}{2}) = \text{Bin}(a_i, a_i + r_i, \frac{1}{2}) \quad (3)$$

$$P(a_i, r_i | c_i = \text{reference}) = \text{Bin}(a_i, a_i + r_i, \epsilon) \quad (4)$$

332 To complete the calculation we use the priors specified by the parameters described above:

$$P(c_i = \text{somatic}) = \mu \quad (5)$$

$$P(c_i = \text{het}) = \pi \quad (6)$$

$$P(c_i = \text{reference}) = 1 - \mu - \pi \quad (7)$$

333 4.1.4 Phased data

334 When the sequencing data can be phased by assigning each read to a haplotype the problem becomes
335 simpler. To illustrate, consider a position where the reference base is T and we have observed 17
336 reads with T and 12 reads with C. Under the previous model it is plausible that the position is
337 heterozygous C/T and the haplotype bearing the reference allele was sampled more often simply
338 by chance. If the haplotype for each read is known however, we can directly calculate whether
339 the haplotype with the non-reference allele contains sufficient evidence of the reference allele from
340 contaminating non-cancerous cells to classify the position as somatic. Using the example above
341 suppose we have partitioned the reads into two haplotypes where H_1 contains 11 observations of C
342 and 5 observations of T and H_2 contains 1 observation of C and 12 of T. It now appears plausible
343 that the position is a T>C somatic mutation with the 5 observations of T due to contaminating
344 normal cells and the individual's (inherited) genotype at this position is T/T. Clearly this only
345 works when the sample is not purely tumour ($\alpha < 1$).

346 The classifier from the previous section can be modified to support phased data:

$$P(c_{i,j} = \text{somatic} | a_{i,j}, r_{i,j}) = \frac{P(a_{i,j}, r_{i,j} | c_{i,j} = \text{somatic})P(c_{i,j} = \text{somatic})}{\sum_c P(a_{i,j}, r_{i,j} | c_{i,j} = c)P(c_{i,j} = c)} \quad (8)$$

347 The likelihoods for **somatic** and **reference** classes are similar to above with the factor of 2 dropped
348 in the **somatic** case (as now the haplotype for each read is known):

$$P(a_{i,j}, r_{i,j} | c_{i,j} = \text{somatic}) = \text{Bin}(a_{i,j}, a_{i,j} + r_{i,j}, (1 - \epsilon)\alpha + \epsilon(1 - \alpha)) \quad (9)$$

$$P(a_{i,j}, r_{i,j} | c_{i,j} = \text{reference}) = \text{Bin}(a_{i,j}, a_{i,j} + r_{i,j}, \epsilon) \quad (10)$$

349 In the case for **het** however, all reads are expected to contain the alternate base, with any reference
350 reads being sequencing errors:

$$P(a_{i,j}, r_{i,j} | c_{i,j} = \text{het}) = \text{Bin}(a_{i,j}, a_{i,j} + r_{i,j}, 1 - \epsilon) \quad (11)$$

351 4.1.5 Implementation

352 The complete simulation procedure is implemented in the function `sim_pileup` in `simulation.rs`
353 in the `smrest` software package.

354 4.2 Single Molecule Long Read Somatic Mutation Calling

355 The mutation calling procedure for real data is derived from the models presented above. A major
356 difference however is that instead of counting the number of reads supporting the reference and
357 alternate allele, which can be inaccurate if the read-to-reference alignment provided in the BAM file
358 is unreliable, we use a likelihood-based calculation that is common with many other variant callers
359 (Albers et al. 2011; Garrison and Marth 2012; Poplin et al. 2018; Cooke, Wedge, and Lunter 2021).
360 The procedure we use is derived from the Longshot variant caller (Edge and Bansal 2019) with the
361 modification that the core read-haplotype likelihood calculation is changed from Longshot’s Hidden
362 Markov Model to `kprobaln` from HTSlib (H. Li 2011), which incorporates quality scores. In early
363 experiments we found (data not shown) that this change significantly improved somatic mutation
364 calling accuracy.

365 **Calculating Read-Haplotype Likelihoods.** In the following methods we rely on calculating
366 the probability of observing a certain sequencing read, R_k , given a known or assumed haplotype
367 sequence, H , denoted $P(R_k|H)$ and referred to as the read-haplotype likelihood. The calculation
368 of $P(R_k|H)$ typically uses the forward algorithm on a hidden Markov model parameterized with
369 gap open, gap extension and substitution probabilities that model the properties of the sequencing
370 platform that generated the read. Longshot also calculates a read-allele likelihood by summing
371 the read-haplotype likelihoods over all haplotypes that contain a particular allele at a particular
372 variant site. We denote this as $P(R_k|V_i^R)$ and $P(R_k|V_i^A)$ for the reference and alternate alleles of a
373 candidate variant at position i . We will denote the set of reads that are informative at site i (they
374 cross the reference position and pass all QC filters) as \mathcal{R}_i .

375 It is inefficient to perform the forward algorithm over the entire length of each read, so Longshot
376 constrains the calculation to short windows surrounding a position of interest. This procedure
377 involves assembling a set of haplotypes containing combinations of candidate variant alleles, then
378 evaluating $P(R_k|H)$ for each one. We directly use Longshot’s code for performing these calculations
379 and refer to the methods in Edge and Bansal 2019 for further details.

380 **Genotyping.** The genotyping procedure takes as input a BAM file and a VCF file containing
381 known SNPs in the human population. In the experiments presented in this manuscript we used
382 gnomAD v3 (Karczewski et al. 2020) biallelic SNPs that have allele frequency $> 0.1\%$. First, the
383 VCF and BAM file are provided to Longshot’s `extract_fragments` function to calculate read-allele
384 likelihoods. Next we calculate genotype likelihoods for each site in the VCF file. Unlike in typical
385 genotyping applications, where reads at heterozygous sites have equal chance of being drawn from
386 each allele, we account for unbalanced copy number. In the following, let f be the probability of
387 drawing a read from the haplotype with higher copy number and G_i be the genotype at position i .
388 The genotype likelihoods for the homozygous REF and ALT cases are straightforward:

$$P(\mathcal{R}_i|G_i = 0/0) = \prod_{R_k \in \mathcal{R}_i} P(R_k|V_i^R) \quad (12)$$

$$P(\mathcal{R}_i|G_i = 1/1) = \prod_{R_k \in \mathcal{R}_i} P(R_k|V_i^A) \quad (13)$$

389 To calculate the heterozygous genotype likelihood we need to account for uncertainty whether the
390 read came from the haplotype bearing the reference or alternate allele, and which one has higher
391 copy number:

$$\begin{aligned}
 P(\mathcal{R}_i | G_i = 0/1) &= \frac{1}{2}P(\mathcal{R}_i | (V_i^R, V_i^A)) + \frac{1}{2}P(\mathcal{R}_i | (V_i^A, V_i^R)) \\
 &= \frac{1}{2} \prod_{R_k \in \mathcal{R}_i} P(R_k | V_i^R) f + P(R_k | V_i^A) (1 - f) \\
 &\quad + \frac{1}{2} \prod_{R_k \in \mathcal{R}_i} P(R_k | V_i^A) f + P(R_k | V_i^R) (1 - f)
 \end{aligned}$$

392 We fix $f = 0.25$ for the experiments in this paper. Future work could fit f for each copy number
 393 segment of the genome in an iterative procedure that calls an initial set of heterozygous SNPs, esti-
 394 mates f , then repeats. Genotype probabilities are calculated using priors that expect approximately
 395 1 in 1,000 sites to be a variant, with 2/3 of variant sites being heterozygous.

396 The genotype with the highest probability is assigned for each site. Any sites that have evidence
 397 of sequencing strand bias (Guo et al. 2012) are left uncalled.

398 **Phasing.** The VCF file output by the genotyping procedure and the BAM file is provided to
 399 `whatshap` (Patterson et al. 2015) for phasing. Default parameters are used, except for the addition
 400 of `--ignore-read-groups`. The output is a phased VCF file.

401 **Haplotagging.** The somatic mutation calling algorithm requires partitioning the input reads by
 402 haplotype. Most phasing software, including `whatshap`, can produce a BAM file where reads are an-
 403 notated with a prediction of which haplotype they originate from, a process known as haplotagging.
 404 In `smrest` we adopt the same procedure, but produce the read-to-haplotype assignments as needed
 405 to avoid the time and space required to produce a large BAM file. The output of `whatshap` is set
 406 of phased heterozygous SNPs. For simplicity we treat the phased haplotypes as a pair of strings B_j
 407 ($j \in \{0, 1\}$) where $B_j[i] \in \{R, A\}$ is the allele at position i of haplotype j . This formulation neglects
 408 the segmentation of the genome into multiple phased blocks, where the phase of adjacent blocks is
 409 unknown. This is handled however by our quality control procedure (see below) to discard reads
 410 that cross phase block boundaries.

411 The haplotype assignments can be viewed as a vector A where $A[k] \in \{0, 1, -\}$ is the assignment
 412 of read k to one of the two haplotypes, or unassigned ($-$). The posterior probability that read k
 413 originates from haplotype j is:

$$P(A[k] = j | R_k) = \frac{P(R_k | A[k] = j) P(A[k] = j)}{P(R_k | A[k] = 0) P(A[k] = 0) + P(R_k | A[k] = 1) P(A[k] = j)}$$

414 Letting \mathcal{I}_k be the set of heterozygous sites that are found in read k , the likelihood is computed as
 415 the product of the read-variant likelihoods calculated by Longshot:

$$P(R_k | A[k] = j) = \prod_{i \in \mathcal{I}_k} P(R_k | V_i^{B_j[i]})$$

416 The read is assigned to the most probable haplotype and a quality score (log-scaled probability the
 417 haplotype assignment is incorrect) is calculated. If this quality score is less than 20, or if the alleles

418 supported by the read mismatch more than 10% of the alleles in its assigned haplotype, the read is
419 left unassigned and not used for somatic mutation calling. This conservative approach avoids the
420 difficulty of assigning a haplotype to reads that may span multiple phased blocks. Better treatment
421 of these reads is an avenue for future work.

422 **Somatic mutation calling.** The somatic mutation calling procedure uses the phased VCF file
423 and the original input BAM. For efficiency and parallelization mutations are calculated over 10Mbp
424 windows of the genome. First, reads within the calling window are assigned to a haplotype (if
425 possible) using the procedure described above. Next, a set of candidate somatic variants is found.
426 A position i is considered *callable* if both haplotypes have at least 10x sequencing depth and the
427 sum of haplotype depth is not greater than 400x. The callable positions are recorded and output
428 as a BED file. Next, the most frequently observed non-reference base on one of the haplotypes is
429 found. If this base is seen in more than 10% of the reads on the haplotype, and at least 3 times on
430 the haplotype, the position and base are recorded on the list of candidate variants. This procedure
431 is necessarily very permissive and generates a large list of candidate variants, very few of which are
432 expected to be actual somatic mutations. This list of candidate variants is input into Longshot's
433 `extract_fragments` algorithm as described above to calculate read-allele likelihoods.

434 Next, every candidate is classified as a somatic mutation, an inherited SNP or reference allele. The
435 general procedure follows the classification described in the simulation section where a reference
436 classification expects all reads to match the reference allele (V_i^R), the heterozygous SNP classification
437 expects all reads to match the alternate allele (V_i^A) and the somatic mutation expects a mixture of
438 reference and alternative alleles. Here the read-variant likelihoods are used and the model is also
439 extended to account for subclonal mutations. The calculations are performed for each haplotype
440 separately. Letting $\mathcal{R}_{i,j}$ be the set of reads on haplotype j that are informative about position
441 i :

$$P(\mathcal{R}_{i,j} | c_{i,j} = \text{reference}) = \prod_{R_k \in \mathcal{R}_{i,j}} P(R_k | V_i^R) \quad (14)$$

$$P(\mathcal{R}_{i,j} | c_{i,j} = \text{het}) = \prod_{R_k \in \mathcal{R}_{i,j}} P(R_k | V_i^A) \quad (15)$$

442 For the somatic class, assuming for the moment all mutations are clonal:

$$P(\mathcal{R}_{i,j} | c_{i,j} = \text{somatic}) = \prod_{R_k \in \mathcal{R}_{i,j}} P(R_k | V_i^R)(1 - \alpha) + P(R_k | V_i^A)\alpha \quad (16)$$

443 **Modelling subclonal mutations.** The methods described thus far assume that every mutation is
444 clonal and contained in every cancer cell however real tumours have subclonal mutations. We define
445 the *cancer cell fraction*, denoted by ϕ_i , as the proportion of tumour cells that carry a mutation at
446 position i (for convenience $\phi_i = 0$ when position i does not have a somatic mutation). For clonal
447 variants $\phi_i = 1$, a variant is *subclonal* when $0 < \phi_i < 1$. Typically in real cancers a subset of
448 mutations will be clonal, those that arose in the founding lineage of the tumour, so we model the
449 cancer cell fraction distribution as a mixture where a proportion of variants are clonal, denoted ρ ,
450 and the frequencies for the remaining $1 - \rho$ variants are drawn from a Beta distribution:

$$P(\phi_i = x; \rho, a, b) = \begin{cases} \rho & x = 1 \\ (1 - \rho)\text{Beta}(x; a, b) & 0 < x < 1 \end{cases}$$

451 If ϕ_i is known, it would be straightforward to modify the somatic mutation classifier as the chance
452 of sampling a mutated tumour cell is $\alpha\phi_i$. The likelihood then becomes:

$$P(R_k | c_{i,j} = \text{somatic}, \phi_i) = P(R_k | V_i^R)(1 - \alpha\phi_i) + P(R_k | V_i^A)\alpha\phi_i \quad (17)$$

453 ϕ_i is not known however, so we integrate it out:

$$P(R_k | c_{i,j} = \text{somatic}) = \int_0^1 P(R_k | c_{i,j} = \text{somatic}, \phi_i) P(\phi_i; \rho, a, b) d\phi_i \quad (18)$$

454 In our experiments we set $\rho = 0.95, a = 2, b = 2$ for the CCF distribution. The integration is
455 numerically approximated in discrete bins of ϕ_i over the range $[0, 1.0]$.

456 **Filtering variants.** After identifying somatic mutations we apply filters to remove mutations that
457 either break assumptions of our model or have features that are indicative of problematic regions
458 of the genome. The filters currently used are:

- 459 • **MaxOtherHaplotypeObservations:** We expect somatic mutations to appear on only one of
460 the two haplotypes, so use the non-called haplotype as an internal control. If the variant
461 appears on this haplotype more than 2 times, or in more than 20% of the reads, this filter is
462 applied.
- 463 • **MinObsPerStrand:** This filter is applied when the variant does not appear in reads from both
464 the forward and reverse sequencing strands, as this is commonly indicative of sequencing
465 artifacts.
- 466 • **PossibleAlignmentArtifact:** Our model assumes that the reads used for variant calling
467 (those in $\mathcal{R}_{i,j}$) are reliably mapped and aligned. While Longshot applies QC filters to the reads
468 it uses, primarily mapping quality thresholds, some read alignments may still be erroneous.
469 In particular, we found that reads spanning structural variant breakpoints can have regions
470 with very high mismatch rates, which can be called as somatic mutations, so this filter is
471 applied to remove such calls.
- 472 • **LowQual:** the PHRED-scaled mutation quality score falls below the calling threshold.
- 473 • **MinHaplotypeDepth:** the depth on either haplotype falls below the calling threshold of 10x
474 coverage.
- 475 • **StrandBias:** As in most variant callers we calculate a strand bias p-value to identify possible
476 sequencing artifacts (Guo et al. 2012).

477 **Implementation.** smrest is implemented in Rust and available under the MIT license on github:
478 <https://github.com/jts/smrest>. The repository contains a Snakemake (Köster and Rahmann
479 2018) pipeline that automates the entire process, starting from a BAM file containing reads mapped
480 to the human reference genome.

481 4.3 Experiments

482 **Data Access and Preparation.** FASTQ, BAM or CRAM files were downloaded from public
483 repositories:

Sample	Technology	Location/Accessions	Reference
COLO829	ONT	colo829_2023.04/COLO829/	-
COLO829BL	ONT	colo829_2023.04/COLO829BL/	-
COLO829	PacBio-HiFi	revio/2023Q2/COLO829	-
COLO829BL	PacBio-HiFi	revio/2023Q2/COLO829-BL	-
COLO829	Illumina	ERR2752450	Cameron, Baber, Shale, Valle-Inclan, et al. 2021
COLO829BL	Illumina	ERR2752449	Cameron, Baber, Shale, Valle-Inclan, et al. 2021
HCC1395	ONT	SRR25005626	Zheng, Su, et al. 2023
HCC1395BL	ONT	SRR25005625	Zheng, Su, et al. 2023
HCC1395	PacBio-HiFi	revio/2023Q2/HCC1395	-
HCC1395BL	PacBio-HiFi	revio/2023Q2/HCC1395-BL	-
HCC1395	Illumina	WGS_IL_T_1.bwa.dedup.bam	Xiao et al. 2021
HCC1395BL	Illumina	WGS_IL_N_1.bwa.dedup.bam	Xiao et al. 2021
CliveOME	ONT	cliveome_kit14_2022.05/	-

Table 1: Summary of datasets used in this manuscript

484 BAM or CRAM files that were already mapped to GRCh38 were used directly, otherwise the reads
485 were mapped with minimap2 or bwa mem for long and short reads, respectively. The raw signal
486 data for the CliveOME data set was downloaded and basecalled with `wf-basecalling` using model
487 `dna_r10.4.1_e8.2_400bps v4.1.0` in sup mode.

488 **Software Versions.** The following software tools were used in this work:

Software	Version	Reference
minimap2	2.24-r1122	H. Li 2018
bwa	0.7.17-r1188	H. Li 2013
samtools	1.16	H. Li et al. 2009
bcftools	1.16	Danecek et al. 2021
clairS	0.1.16	Zheng, Su, et al. 2023
mutect2	4.4.0.0	Benjamin et al. 2019
whatshap	1.7	Patterson et al. 2015
bedtools	2.30.0	Quinlan and Hall 2010
snakemake	7.19.1	Köster and Rahmann 2018

Table 2: Summary of software used in this manuscript

489 **Downsampling.** To prepare BAM files with a specified coverage level we first computed the total
490 number of bases contained in the full depth BAM with `samtools stats`, calculated the proportion
491 of reads needed to reach the specified coverage, then generated a new BAM by passing this value
492 to the `-s` argument of `samtools view`.

493 **ClairS Mutation Calling.** ClairS was run using singularity as described in the README and
494 provided with a tumour and normal BAM file using the `--tumour-bam-fn` and `--normal-bam-fn`
495 arguments. The `--platform` argument was set to `ont_r10_dorado_4khz` for ONT data or `pacbio-hifi`

496 for PacBio-HiFi data.

497 **Mutect2 Tumour-Normal Mutation Calling.** To generate a raw VCF file the GATK `mutect2`
498 command was run with the `panel-of-normals` argument set to `1000g_pon.hg38.vcf.gz` and germline
499 resource set to `af-only-gnomad.hg38.vcf.gz`. To force multi-nucleotide variants to be called
500 as individual SNV the `--max-mnp-distance 0` argument was provided. The raw mutation calls
501 were filtered with the `FilterMutectCalls` using the output of the `CalculateContamination` com-
502 mand.

503 **Mutect2 Tumour-only Mutation Calling.** To call mutations in tumour-only mode, `Mutect2`
504 was run as above but provided with a single BAM file containing downsampled reads from both
505 the tumour and normal, and the `--normal-sample-name` argument was omitted.

506 **smrest Mutation Calling.** `smrest` was run using the snakemake pipeline implementing the
507 procedure described in the previous section. A single BAM file, containing reads mixed from
508 the tumour and normal cell lines, was provided as input. In all experiments the tumour purity
509 parameter was set to 0.5.

510 **Truth Data.** The COLO829 truth mutation set was derived from the NovaSeq VCF file provided
511 by the New York Genome Centre (Arora et al. 2019). This VCF file was processed to split MNVs
512 into individual SNVs using `bcftools norm -a`. The HCC1395 truth mutation set was downloaded
513 from the SEQC2 FTP site.

514 **Genome stratification.** To evaluate mutation calling performance for all tools, we restricted
515 the analysis to high confidence (HC) regions of the genome. For COLO829 we defined the high
516 confidence regions by first using `bedtools` to compute the union between GIAB’s `alldifficultregions`
517 and `HG001_v4.2.1_complexandSVs` BED files (Zook et al. 2014), and ENCODE’s `hg38-blacklist.v2`
518 BED (Amemiya, Kundaje, and Boyle 2019). We then took the complement of this union BED file
519 to define the HC regions. For HCC1395 we intersected this BED file with the BED file provided
520 with the truth data from the SEQC FTP.

521 **Identifying cell line artifacts.** As tumour-only mutation calling may identify mutations acquired
522 in cell culture, which aren’t present in the truth data we used, we removed these calls from our
523 analysis. First, we ran `Mutect2` in TN mode but provided the normal sample name in place of the
524 tumour sample name to generate a list of putative mutations in the normal cell line. As loss-of-
525 heterozygosity regions in the tumour would be called as mutations using this procedure, we filtered
526 out any variant call that was within 5,000bp of an annotated population variant (using the `POPAF`
527 field in `Mutect2`’s VCF) to generate the final list of normal cell line artifacts. Any mutation call
528 matching this list was not included in accuracy calculations.

529 **Analyzing Called Mutations.** A mutation call set is annotated using the truth data as follows.
530 Each mutation in either the call set or truth set is represented by its chromosome, position, reference
531 allele and alternate allele. The union between the call set and truth set is taken and each record
532 output in a TSV file. Each record is annotated with whether it is contained in the truth set, the
533 call set (excluding hard filtered calls, with the exception of `LowQual` calls as these are retained
534 for calculating precision-recall curves) and within the regions specified within the specified BED
535 file (either the HC regions described above, or the phased subset of these regions). Each record is
536 also annotated with the mutation caller’s confidence (the `QUAL` field for `smrest` and `clairS`, `TLOD`
537 for `Mutect2`), the `VAF` for the truth mutation or called mutation (if applicable) and any filters
538 applied in the VCF file. Accuracy statistics (F1, precision, sensitivity) and precision-recall curves
539 are calculated from this file after removing records where `VAF` is below the analysis threshold

540 (10%), outside of the regions specified in the BED file or present in the list of normal cell artifacts.
541 When calculating F1, precision or sensitivity a minimum QUAL of 20 was used for smrest calls.
542 Default values were used (PASS calls) for clairS and mutect2.

543 **Estimating Tumour Mutation Burden.** Tumour mutation burden was estimated by counting
544 the number of QC-PASS mutation calls with a minimum QUAL score of 20, and dividing this num-
545 ber by the total number of bases where mutation calling was performed (from the BED file output
546 by smrest). The CliveOME sample listed in **Table 1** was the healthy blood sample control.

547 **Mutation Signatures.** The mutation spectrum was extracted and plotted using SigProfilerMa-
548 trixGenerator (Bergstrom et al. 2019).

549 **Implementation.** The source code for the mutation caller is provided at [https://github.com/](https://github.com/jts/smrest)
550 [jts/smrest](https://github.com/jts/smrest). The code used to generate all results in this manuscript is provided as a Snakemake
551 pipeline and associated python scripts at <https://github.com/jts/smrest-analysis-pipeline>.

552 Acknowledgements

553 The author thanks Joanna Pineda for prior work on somatic mutation calling using phased 10X
554 Genomics Linked Reads (<https://github.com/jopineda/10xtrim>), and Felix Beaudry and Tom
555 Ouellette for comments on a draft version of this manuscript. The author also thanks Philip
556 Zuzarte, Jim Shaw, Chris Wright, Alvin Ng, Matthew Loose and Winston Timp for discussions
557 related to this manuscript.

558 The author is supported by the Ontario Institute for Cancer Research through funds provided
559 by the Government of Ontario, the Government of Canada through Genome Canada and Ontario
560 Genomics (OGI-136 and OGI-201) and the National Human Genome Research Institute (NHGRI
561 project 5R01HG009190).

562 Conflict of Interest

563 J.T.S. receives research funding from Oxford Nanopore Technologies (ONT) and has received travel
564 support to attend and speak at meetings organized by ONT, and is on the Scientific Advisory Board
565 of Day Zero Diagnostics.

566 References

- 567 Weinstein, John N. et al. (Oct. 2013). “The Cancer Genome Atlas Pan-Cancer analysis project”.
568 en. In: *Nature Genetics* 45.10. Number: 10 Publisher: Nature Publishing Group, pp. 1113–1120.
569 ISSN: 1546-1718. DOI: 10.1038/ng.2764. URL: <https://www.nature.com/articles/ng.2764>
570 (visited on 02/05/2024).
- 571 Bailey, Matthew H. et al. (2018). “Comprehensive characterization of cancer driver genes and
572 mutations”. In: *Cell* 173.2. ISBN: 0092-8674 Publisher: Elsevier, 371–385. e18.
- 573 Rheinbay, Esther et al. (2020). “Analyses of non-coding somatic drivers in 2,658 cancer whole
574 genomes”. In: *Nature* 578.7793. ISBN: 0028-0836 Publisher: Nature Publishing Group UK Lon-
575 don, pp. 102–111.
- 576 *Pan-cancer analysis of whole genomes* (2020). In: *Nature* 578.7793. ISBN: 0028-0836 Publisher:
577 Nature Publishing Group UK London, pp. 82–93.

- 578 Alexandrov, Ludmil B., Serena Nik-Zainal, et al. (Jan. 2013). “Deciphering signatures of mutational
579 processes operative in human cancer”. eng. In: *Cell Reports* 3.1, pp. 246–259. ISSN: 2211-1247.
580 DOI: 10.1016/j.celrep.2012.12.008.
- 581 Alexandrov, Ludmil B., Jaegil Kim, et al. (2020). “The repertoire of mutational signatures in human
582 cancer”. In: *Nature* 578.7793. ISBN: 0028-0836 Publisher: Nature Publishing Group UK London,
583 pp. 94–101.
- 584 Nik-Zainal, Serena et al. (May 2012). “Mutational processes molding the genomes of 21 breast
585 cancers”. eng. In: *Cell* 149.5, pp. 979–993. ISSN: 1097-4172. DOI: 10.1016/j.cell.2012.04.024.
- 586 Shah, Sohrab P. et al. (2009). “Mutational evolution in a lobular breast tumour profiled at single
587 nucleotide resolution”. In: *Nature* 461.7265. ISBN: 0028-0836 Publisher: Nature Publishing Group
588 UK London, pp. 809–813.
- 589 Sottoriva, Andrea, Inmaculada Spiteri, et al. (2013). “Intratumor heterogeneity in human glioblas-
590 toma reflects cancer evolutionary dynamics”. In: *Proceedings of the National Academy of Sciences*
591 110.10. ISBN: 0027-8424 Publisher: National Acad Sciences, pp. 4009–4014.
- 592 Sottoriva, Andrea, Haeyoun Kang, et al. (2015). “A Big Bang model of human colorectal tumor
593 growth”. In: *Nature genetics* 47.3. ISBN: 1546-1718 Publisher: Nature Publishing Group, pp. 209–
594 216.
- 595 Gerstung, Moritz et al. (2020). “The evolutionary history of 2,658 cancers”. In: *Nature* 578.7793.
596 ISBN: 0028-0836 Publisher: Nature Publishing Group UK London, pp. 122–128.
- 597 Jiao, Wei et al. (2020). “A deep learning system accurately classifies primary and metastatic cancers
598 using passenger mutation patterns”. In: *Nature communications* 11.1. ISBN: 2041-1723 Publisher:
599 Nature Publishing Group UK London, p. 728.
- 600 Hendrikse, Liam D. et al. (2022). “Failure of human rhombic lip differentiation underlies medul-
601 loblastoma formation”. In: *Nature* 609.7929. ISBN: 0028-0836 Publisher: Nature Publishing
602 Group UK London, pp. 1021–1028.
- 603 Chan, Timothy A. et al. (2019). “Development of tumor mutation burden as an immunotherapy
604 biomarker: utility for the oncology clinic”. In: *Annals of Oncology* 30.1. ISBN: 0923-7534 Pub-
605 lisher: Elsevier, pp. 44–56.
- 606 Zhao, Pengfei et al. (2019). “Mismatch repair deficiency/microsatellite instability-high as a predic-
607 tor for anti-PD-1/PD-L1 immunotherapy efficacy”. In: *Journal of hematology & oncology* 12.1.
608 ISBN: 1756-8722 Publisher: BioMed Central, pp. 1–14.
- 609 André, Thierry et al. (2020). “Pembrolizumab in microsatellite-instability–high advanced colorectal
610 cancer”. In: *New England Journal of Medicine* 383.23. ISBN: 0028-4793 Publisher: Mass Medical
611 Soc, pp. 2207–2218.
- 612 Koboldt, Daniel C. et al. (Sept. 2009). “VarScan: variant detection in massively parallel sequencing
613 of individual and pooled samples”. eng. In: *Bioinformatics (Oxford, England)* 25.17, pp. 2283–
614 2285. ISSN: 1367-4811. DOI: 10.1093/bioinformatics/btp373.
- 615 Larson, David E. et al. (Feb. 2012). “SomaticSniper: identification of somatic point mutations in
616 whole genome sequencing data”. eng. In: *Bioinformatics (Oxford, England)* 28.3, pp. 311–317.
617 ISSN: 1367-4811. DOI: 10.1093/bioinformatics/btr665.
- 618 Saunders, Christopher T. et al. (July 2012). “Strelka: accurate somatic small-variant calling from se-
619 quenced tumor–normal sample pairs”. In: *Bioinformatics* 28.14, pp. 1811–1817. ISSN: 1367-4803.
620 DOI: 10.1093/bioinformatics/bts271. URL: [https://doi.org/10.1093/bioinformatics/
621 bts271](https://doi.org/10.1093/bioinformatics/bts271) (visited on 02/05/2024).
- 622 Cibulskis, Kristian et al. (2013). “Sensitive detection of somatic point mutations in impure and
623 heterogeneous cancer samples”. In: *Nature biotechnology* 31.3. ISBN: 1087-0156 Publisher: Nature
624 Publishing Group US New York, pp. 213–219.

- 625 Fang, Li Tai et al. (Sept. 2021). “Establishing community reference samples, data and call sets
626 for benchmarking cancer mutation detection using whole-genome sequencing”. eng. In: *Nature*
627 *Biotechnology* 39.9, pp. 1151–1160. ISSN: 1546-1696. DOI: 10.1038/s41587-021-00993-6.
- 628 Smith, Kyle S. et al. (Mar. 2016). “SomVarIUS: somatic variant identification from unpaired tissue
629 samples”. In: *Bioinformatics* 32.6, pp. 808–813. ISSN: 1367-4803. DOI: 10.1093/bioinformatics/
630 btv685. URL: <https://doi.org/10.1093/bioinformatics/btv685> (visited on 02/05/2024).
- 631 Kalatskaya, Irina et al. (June 2017). “ISOWN: accurate somatic mutation identification in the
632 absence of normal tissue controls”. en. In: *Genome Medicine* 9.1, p. 59. ISSN: 1756-994X. DOI:
633 10.1186/s13073-017-0446-9. URL: <https://doi.org/10.1186/s13073-017-0446-9> (visited
634 on 02/05/2024).
- 635 Sun, James X. et al. (Feb. 2018). “A computational approach to distinguish somatic vs. germline ori-
636 gin of genomic alterations from deep sequencing of cancer specimens without a matched normal”.
637 eng. In: *PLoS computational biology* 14.2, e1005965. ISSN: 1553-7358. DOI: 10.1371/journal.
638 pcbi.1005965.
- 639 Nassar, Amin H. et al. (Oct. 2022). “Ancestry-driven recalibration of tumor mutational burden and
640 disparate clinical outcomes in response to immune checkpoint inhibitors”. eng. In: *Cancer Cell*
641 40.10, 1161–1172.e5. ISSN: 1878-3686. DOI: 10.1016/j.ccell.2022.08.022.
- 642 Sereika, Mantas et al. (July 2022). “Oxford Nanopore R10.4 long-read sequencing enables the
643 generation of near-finished bacterial genomes from pure cultures and metagenomes without short-
644 read or reference polishing”. en. In: *Nature Methods* 19.7. Number: 7 Publisher: Nature Publishing
645 Group, pp. 823–826. ISSN: 1548-7105. DOI: 10.1038/s41592-022-01539-7. URL: [https://www.
646 nature.com/articles/s41592-022-01539-7](https://www.nature.com/articles/s41592-022-01539-7) (visited on 02/05/2024).
- 647 Kolesnikov, Alexey et al. (Sept. 2023). “Local read haplotagging enables accurate long-read small
648 variant calling”. In: *bioRxiv*, p. 2023.09.07.556731. DOI: 10.1101/2023.09.07.556731. URL:
649 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC10515762/> (visited on 02/05/2024).
- 650 Rhie, Arang et al. (Apr. 2021). “Towards complete and error-free genome assemblies of all vertebrate
651 species”. en. In: *Nature* 592.7856. Number: 7856 Publisher: Nature Publishing Group, pp. 737–
652 746. ISSN: 1476-4687. DOI: 10.1038/s41586-021-03451-0. URL: [https://www.nature.com/
653 articles/s41586-021-03451-0](https://www.nature.com/articles/s41586-021-03451-0) (visited on 02/05/2024).
- 654 Nurk, Sergey et al. (Apr. 2022). “The complete sequence of a human genome”. In: *Science (New*
655 *York, N. Y.)* 376.6588, pp. 44–53. ISSN: 0036-8075. DOI: 10.1126/science.abj6987. URL: [https:
656 //www.ncbi.nlm.nih.gov/pmc/articles/PMC9186530/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9186530/) (visited on 02/05/2024).
- 657 Flusberg, Benjamin A. et al. (June 2010). “Direct detection of DNA methylation during single-
658 molecule, real-time sequencing”. en. In: *Nature Methods* 7.6. Number: 6 Publisher: Nature Pub-
659 lishing Group, pp. 461–465. ISSN: 1548-7105. DOI: 10.1038/nmeth.1459. URL: [https://www.
660 nature.com/articles/nmeth.1459](https://www.nature.com/articles/nmeth.1459) (visited on 02/05/2024).
- 661 Laszlo, Andrew H. et al. (Nov. 2013). “Detection and mapping of 5-methylcytosine and 5-hydroxymethylcytosine
662 with nanopore MspA”. In: *Proceedings of the National Academy of Sciences* 110.47. Publisher:
663 Proceedings of the National Academy of Sciences, pp. 18904–18909. DOI: 10.1073/pnas.
664 1310240110. URL: [https://www.pnas.org/doi/full/10.1073/pnas.
665 1310240110](https://www.pnas.org/doi/full/10.1073/pnas.1310240110) (visited on 02/05/2024).
- 666 Schreiber, Jacob et al. (Nov. 2013). “Error rates for nanopore discrimination among cytosine,
667 methylcytosine, and hydroxymethylcytosine along individual DNA strands”. eng. In: *Proceedings*
668 *of the National Academy of Sciences of the United States of America* 110.47, pp. 18910–18915.
669 ISSN: 1091-6490. DOI: 10.1073/pnas.1310615110.
- 670 Simpson, Jared T. et al. (Apr. 2017). “Detecting DNA cytosine methylation using nanopore se-
671 quencing”. eng. In: *Nature Methods* 14.4, pp. 407–410. ISSN: 1548-7105. DOI: 10.1038/nmeth.
672 4184.

- 673 Zheng, Zhenxian, Junhao Su, et al. (Aug. 2023). *ClairS: a deep-learning method for long-read*
674 *somatic small variant calling*. en. Pages: 2023.08.17.553778 Section: New Results. DOI: 10.1101/
675 2023.08.17.553778. URL: [https://www.biorxiv.org/content/10.1101/2023.08.17.](https://www.biorxiv.org/content/10.1101/2023.08.17.553778v1)
676 553778v1 (visited on 02/05/2024).
- 677 Darby, Charlotte A. et al. (Aug. 2019). “Samovar: Single-Sample Mosaic Single-Nucleotide Variant
678 Calling with Linked Reads”. eng. In: *iScience* 18, pp. 1–10. ISSN: 2589-0042. DOI: 10.1016/j.
679 isci.2019.05.037.
- 680 Xu, Chang (Feb. 2018). “A review of somatic single nucleotide variant calling algorithms for next-
681 generation sequencing data”. In: *Computational and Structural Biotechnology Journal* 16, pp. 15–
682 24. ISSN: 2001-0370. DOI: 10.1016/j.csbj.2018.01.003. URL: [https://www.ncbi.nlm.nih.](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5852328/)
683 [gov/pmc/articles/PMC5852328/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC5852328/) (visited on 02/05/2024).
- 684 Albers, Cornelis A. et al. (June 2011). “Dindel: accurate indel calls from short-read data”. eng. In:
685 *Genome Research* 21.6, pp. 961–973. ISSN: 1549-5469. DOI: 10.1101/gr.112326.110.
- 686 Li, Heng (Apr. 2011). “Improving SNP discovery by base alignment quality”. eng. In: *Bioinformatics*
687 (*Oxford, England*) 27.8, pp. 1157–1158. ISSN: 1367-4811. DOI: 10.1093/bioinformatics/btr076.
- 688 Garrison, Erik and Gabor Marth (July 2012). *Haplotype-based variant detection from short-read*
689 *sequencing*. arXiv:1207.3907 [q-bio]. DOI: 10.48550/arXiv.1207.3907. URL: [http://arxiv.](http://arxiv.org/abs/1207.3907)
690 [org/abs/1207.3907](http://arxiv.org/abs/1207.3907) (visited on 05/25/2023).
- 691 Poplin, Ryan et al. (July 2018). *Scaling accurate genetic variant discovery to tens of thousands of*
692 *samples*. en. Pages: 201178 Section: New Results. DOI: 10.1101/201178. URL: [https://www.](https://www.biorxiv.org/content/10.1101/201178v3)
693 [biorxiv.org/content/10.1101/201178v3](https://www.biorxiv.org/content/10.1101/201178v3) (visited on 02/05/2024).
- 694 Cooke, Daniel P., David C. Wedge, and Gerton Lunter (July 2021). “A unified haplotype-based
695 method for accurate and comprehensive variant calling”. eng. In: *Nature Biotechnology* 39.7,
696 pp. 885–892. ISSN: 1546-1696. DOI: 10.1038/s41587-021-00861-3.
- 697 Patterson, Murray et al. (June 2015). “WhatsHap: Weighted Haplotype Assembly for Future-
698 Generation Sequencing Reads”. eng. In: *Journal of Computational Biology: A Journal of Com-*
699 *putational Molecular Cell Biology* 22.6, pp. 498–509. ISSN: 1557-8666. DOI: 10.1089/cmb.2014.
700 0157.
- 701 Zheng, Zhenxian, Shumin Li, et al. (Dec. 2022). “Symphonizing pileup and full-alignment for deep
702 learning-based long-read variant calling”. en. In: *Nature Computational Science* 2.12. Number:
703 12 Publisher: Nature Publishing Group, pp. 797–803. ISSN: 2662-8457. DOI: 10.1038/s43588-
704 022-00387-x. URL: [https://www.nature.com/articles/s43588-](https://www.nature.com/articles/s43588-022-00387-x)
705 [022-00387-x](https://www.nature.com/articles/s43588-022-00387-x) (visited on
05/25/2023).
- 706 Titmuss, Emma et al. (Sept. 2022). “TMBur: a distributable tumor mutation burden approach for
707 whole genome sequencing”. eng. In: *BMC medical genomics* 15.1, p. 190. ISSN: 1755-8794. DOI:
708 10.1186/s12920-022-01348-z.
- 709 Pleasance, Erin D. et al. (Jan. 2010). “A comprehensive catalogue of somatic mutations from a
710 human cancer genome”. en. In: *Nature* 463.7278. Number: 7278 Publisher: Nature Publishing
711 Group, pp. 191–196. ISSN: 1476-4687. DOI: 10.1038/nature08658. URL: [https://www.nature.](https://www.nature.com/articles/nature08658)
712 [com/articles/nature08658](https://www.nature.com/articles/nature08658) (visited on 02/05/2024).
- 713 Arora, Kanika et al. (Dec. 2019). “Deep whole-genome sequencing of 3 cancer cell lines on 2 sequenc-
714 ing platforms”. eng. In: *Scientific Reports* 9.1, p. 19123. ISSN: 2045-2322. DOI: 10.1038/s41598-
715 019-55636-3.
- 716 Espejo Valle-Inclan, Jose et al. (June 2022). “A multi-platform reference for somatic structural
717 variation detection”. eng. In: *Cell Genomics* 2.6, p. 100139. ISSN: 2666-979X. DOI: 10.1016/j.
718 xgen.2022.100139.

- 719 Benjamin, David et al. (Dec. 2019). *Calling Somatic SNVs and Indels with Mutect2*. en. Pages:
720 861054 Section: New Results. DOI: 10.1101/861054. URL: [https://www.biorxiv.org/content/
721 10.1101/861054v1](https://www.biorxiv.org/content/10.1101/861054v1) (visited on 05/25/2023).
- 722 Petljak, Mia et al. (Mar. 2019). “Characterizing Mutational Signatures in Human Cancer Cell Lines
723 Reveals Episodic APOBEC Mutagenesis”. en. In: *Cell* 176.6, 1282–1294.e20. ISSN: 00928674. DOI:
724 10.1016/j.cell.2019.02.012. URL: [https://linkinghub.elsevier.com/retrieve/pii/
725 S0092867419301618](https://linkinghub.elsevier.com/retrieve/pii/S0092867419301618) (visited on 02/05/2024).
- 726 Krysiak, Kilannin et al. (Jan. 2023). “CIViCdb 2022: evolution of an open-access cancer variant
727 interpretation knowledgebase”. In: *Nucleic Acids Research* 51.D1, pp. D1230–D1241. ISSN: 0305-
728 1048. DOI: 10.1093/nar/gkac979. URL: <https://doi.org/10.1093/nar/gkac979> (visited on
729 02/05/2024).
- 730 Marabelle, Aurélien et al. (Oct. 2020). “Association of tumour mutational burden with outcomes
731 in patients with advanced solid tumours treated with pembrolizumab: prospective biomarker
732 analysis of the multicohort, open-label, phase 2 KEYNOTE-158 study”. English. In: *The Lancet
733 Oncology* 21.10. Publisher: Elsevier, pp. 1353–1365. ISSN: 1470-2045, 1474-5488. DOI: 10.1016/
734 S1470-2045(20)30445-9. URL: [https://www.thelancet.com/journals/lanonc/article/
735 PIIS1470-2045\(20\)30445-9/fulltext](https://www.thelancet.com/journals/lanonc/article/PIIS1470-2045(20)30445-9/fulltext) (visited on 06/06/2023).
- 736 Marcus, Leigh et al. (Sept. 2021). “FDA Approval Summary: Pembrolizumab for the Treatment of
737 Tumor Mutational Burden-High Solid Tumors”. eng. In: *Clinical Cancer Research: An Official
738 Journal of the American Association for Cancer Research* 27.17, pp. 4685–4689. ISSN: 1557-3265.
739 DOI: 10.1158/1078-0432.CCR-21-0327.
- 740 Patel, Anand G., Jann N. Sarkaria, and Scott H. Kaufmann (Feb. 2011). “Nonhomologous end join-
741 ing drives poly(ADP-ribose) polymerase (PARP) inhibitor lethality in homologous recombination-
742 deficient cells”. eng. In: *Proceedings of the National Academy of Sciences of the United States of
743 America* 108.8, pp. 3406–3411. ISSN: 1091-6490. DOI: 10.1073/pnas.1013715108.
- 744 Miller, R. E. et al. (Dec. 2020). “ESMO recommendations on predictive biomarker testing for homol-
745 ogous recombination deficiency and PARP inhibitor benefit in ovarian cancer”. eng. In: *Annals
746 of Oncology: Official Journal of the European Society for Medical Oncology* 31.12, pp. 1606–1622.
747 ISSN: 1569-8041. DOI: 10.1016/j.annonc.2020.08.2102.
- 748 Davies, Helen et al. (2017). “HRDetect is a predictor of BRCA1 and BRCA2 deficiency based on
749 mutational signatures”. In: *Nature medicine* 23.4. ISBN: 1078-8956 Publisher: Nature Publishing
750 Group US New York, pp. 517–525.
- 751 Milbury, Coren A. et al. (Mar. 2022). “Clinical and analytical validation of FoundationOne®CDx,
752 a comprehensive genomic profiling assay for solid tumors”. In: *PLoS ONE* 17.3, e0264138. ISSN:
753 1932-6203. DOI: 10.1371/journal.pone.0264138. URL: [https://www.ncbi.nlm.nih.gov/pmc/
754 articles/PMC8926248/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8926248/) (visited on 02/05/2024).
- 755 Carter, Scott L. et al. (May 2012). “Absolute quantification of somatic DNA alterations in human
756 cancer”. eng. In: *Nature Biotechnology* 30.5, pp. 413–421. ISSN: 1546-1696. DOI: 10.1038/nbt.
757 2203.
- 758 Cameron, Daniel L., Jonathan Baber, Charles Shale, Anthony T. Papenfuss, et al. (Sept. 2019).
759 *GRIDSS, PURPLE, LINX: Unscrambling the tumor genome via integrated analysis of structural
760 variation and copy number*. en. Pages: 781013 Section: New Results. DOI: 10.1101/781013. URL:
761 <https://www.biorxiv.org/content/10.1101/781013v1> (visited on 02/07/2024).
- 762 Li, Kai et al. (Jan. 2020). “Microsatellite instability: a review of what the oncologist should know”.
763 In: *Cancer Cell International* 20, p. 16. ISSN: 1475-2867. DOI: 10.1186/s12935-019-1091-8.
764 URL: <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC6958913/> (visited on 05/25/2023).

- 765 Edge, Peter and Vikas Bansal (Oct. 2019). “Longshot enables accurate variant calling in diploid
766 genomes from single-molecule long read sequencing”. eng. In: *Nature Communications* 10.1,
767 p. 4660. ISSN: 2041-1723. DOI: 10.1038/s41467-019-12493-y.
- 768 Karczewski, Konrad J. et al. (May 2020). “The mutational constraint spectrum quantified from
769 variation in 141,456 humans”. eng. In: *Nature* 581.7809, pp. 434–443. ISSN: 1476-4687. DOI:
770 10.1038/s41586-020-2308-7.
- 771 Guo, Yan et al. (Nov. 2012). “The effect of strand bias in Illumina short-read sequencing data”.
772 eng. In: *BMC genomics* 13, p. 666. ISSN: 1471-2164. DOI: 10.1186/1471-2164-13-666.
- 773 Köster, Johannes and Sven Rahmann (Oct. 2018). “Snakemake-a scalable bioinformatics workflow
774 engine”. eng. In: *Bioinformatics (Oxford, England)* 34.20, p. 3600. ISSN: 1367-4811. DOI: 10.
775 1093/bioinformatics/bty350.
- 776 Cameron, Daniel L., Jonathan Baber, Charles Shale, Jose Espejo Valle-Inclan, et al. (July 2021).
777 “GRIDSS2: comprehensive characterisation of somatic structural variation using single breakend
778 variants and structural variant phasing”. In: *Genome Biology* 22.1, p. 202. ISSN: 1474-760X. DOI:
779 10.1186/s13059-021-02423-x. URL: <https://doi.org/10.1186/s13059-021-02423-x>
780 (visited on 02/06/2024).
- 781 Xiao, Wenming et al. (Sept. 2021). “Toward best practice in cancer mutation detection with whole-
782 genome and whole-exome sequencing”. In: *Nature biotechnology* 39.9, pp. 1141–1150. ISSN: 1087-
783 0156. DOI: 10.1038/s41587-021-00994-5. URL: [https://www.ncbi.nlm.nih.gov/pmc/
784 articles/PMC8506910/](https://www.ncbi.nlm.nih.gov/pmc/articles/PMC8506910/) (visited on 02/06/2024).
- 785 Li, Heng (Sept. 2018). “Minimap2: pairwise alignment for nucleotide sequences”. In: *Bioinformatics*
786 34.18, pp. 3094–3100. ISSN: 1367-4803. DOI: 10.1093/bioinformatics/bty191. URL: <https://doi.org/10.1093/bioinformatics/bty191> (visited on 02/06/2024).
- 787 – (May 2013). *Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM*.
788 arXiv:1303.3997 [q-bio]. DOI: 10.48550/arXiv.1303.3997. URL: [http://arxiv.org/abs/1303.
789 3997](http://arxiv.org/abs/1303.3997) (visited on 02/07/2024).
- 790 Li, Heng et al. (Aug. 2009). “The Sequence Alignment/Map format and SAMtools”. In: *Bioin-
791 formatics* 25.16, pp. 2078–2079. ISSN: 1367-4803. DOI: 10.1093/bioinformatics/btp352. URL:
792 <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC2723002/> (visited on 02/06/2024).
- 793 Danecsek, Petr et al. (Feb. 2021). “Twelve years of SAMtools and BCFtools”. In: *GigaScience* 10.2,
794 giab008. ISSN: 2047-217X. DOI: 10.1093/gigascience/giab008. URL: [https://doi.org/10.
795 1093/gigascience/giab008](https://doi.org/10.1093/gigascience/giab008) (visited on 02/06/2024).
- 796 Quinlan, Aaron R. and Ira M. Hall (Mar. 2010). “BEDTools: a flexible suite of utilities for com-
797 paring genomic features”. In: *Bioinformatics* 26.6, pp. 841–842. ISSN: 1367-4803. DOI: 10.1093/
798 bioinformatics/btq033. URL: <https://doi.org/10.1093/bioinformatics/btq033> (visited
799 on 02/07/2024).
- 800 Zook, Justin M. et al. (Mar. 2014). “Integrating human sequence data sets provides a resource
801 of benchmark SNP and indel genotype calls”. en. In: *Nature Biotechnology* 32.3. Number: 3
802 Publisher: Nature Publishing Group, pp. 246–251. ISSN: 1546-1696. DOI: 10.1038/nbt.2835.
803 URL: <https://www.nature.com/articles/nbt.2835> (visited on 02/07/2024).
- 804 Amemiya, Haley M., Anshul Kundaje, and Alan P. Boyle (June 2019). “The ENCODE Blacklist:
805 Identification of Problematic Regions of the Genome”. en. In: *Scientific Reports* 9.1. Number:
806 1 Publisher: Nature Publishing Group, p. 9354. ISSN: 2045-2322. DOI: 10.1038/s41598-019-
807 45839-z. URL: <https://www.nature.com/articles/s41598-019-45839-z> (visited on
808 02/07/2024).
- 809 Bergstrom, Erik N. et al. (Aug. 2019). “SigProfilerMatrixGenerator: a tool for visualizing and
810 exploring patterns of small mutational events”. eng. In: *BMC genomics* 20.1, p. 685. ISSN: 1471-
811 2164. DOI: 10.1186/s12864-019-6041-2.
- 812