

RESEARCH

Somrit: The Somatic Retrotransposon Insertion Toolkit

Alister V. D'Costa^{1,2} and Jared T. Simpson^{1,2,3*}

*Correspondence:

jared.simpson@oicr.on.ca

¹Ontario Institute for Cancer Research, Toronto, Canada

Full list of author information is available at the end of the article

Abstract

Mobile elements, such as retrotransposons, have the ability to express and re-insert themselves into the genome, with over half the human genome being made up of mobile element sequence. Somatic mobile element insertions (MEIs) have been shown to cause disease, including some cancers. Accurate identification of where novel retrotransposon insertion events occur in the genome is crucial to understand the functional consequence of an insertion event. In this paper we describe somrit, a modular toolkit for detecting somatic MEIs from long reads aligned to a reference genome. We identify the initial read-to-reference mapping step as a potential source of error when the insertion is similar to a nearby repeat in the reference genome and develop a consensus-realignment procedure to resolve this. We show how somrit improves the sensitivity of detection for rare somatic retrotransposon insertion events compared to existing tools, and how the local realignment procedure can reduce false positive translocation calls caused by mis-mapped reads bearing MEIs.

Somrit is openly available at: <https://github.com/adcosta17/somrit>

Keywords: retrotransposon; somatic; structural variation; nanopore

1

2

3 Background

4 Mobile elements are DNA sequences that can change genomic position and re-insert
5 themselves into the genome [1]. A large fraction of the human genome is composed

6 of mobile element sequence with thousands of identified copies [1, 2, 3]. While many
7 of these copies are partial fragments reflecting ancient insertion events that can no
8 longer actively move [2, 3, 4], some more recent copies retain the ability to be
9 expressed and re-insert themselves [5, 6]. Retrotransposons are a class of mobile
10 elements that includes LINE-1 (L1), Alu and SVA elements [7, 8, 9]. Full length
11 human LINE-1 elements are ~6kbp in length and encode proteins for retrotrans-
12 position, allowing for their re-insertion into the genome via an RNA intermediate
13 and reverse transcription [10, 11]. Smaller SVA (~2000bp) and Alu (~300bp) ele-
14 ments rely on the LINE-1 retrotransposition mechanism for re-integration into the
15 genome [12, 13, 14]. LINE-1 insertions usually occur at LINE-1 endonuclease recog-
16 nition motifs [15, 16], often include a target-site duplication (TSD) [15, 17] and a
17 poly-A tail [18] and may contain genomic flanking sequence from the LINE-1 ele-
18 ment of origin [19, 20]. Due to their mobile nature the exact number and location
19 of retrotransposons in the genome varies from person to person, with any individ-
20 ual having some inherited copies not found in the human reference, and possibly
21 somatic copies present in a subset of cells [21].

22 While often only a handful of retrotransposon copies in any given individual retain
23 the ability to be expressed and re-inserted back into the genome, their expression
24 has been linked to disease progression. Prior research has shown somatic insertion
25 of LINE-1 elements activates oncogenes and directly drives cancer progression in
26 some colorectal cancers [22]. Somatic insertion of LINE-1 elements may alter gene
27 expression, including a slowing of DNA translation possibly affecting the expres-
28 sion of tumor suppressor genes [10]. Due to the large amount of mobile element
29 sequence in the genome, retrotransposon insertions have the potential to gener-
30 ate chromosomal rearrangements including deletions, duplications, inversions and
31 translocations, as they may mislead homologous recombination repair pathways to
32 cause non-allelic homologous recombination events [23, 24]. These larger changes in

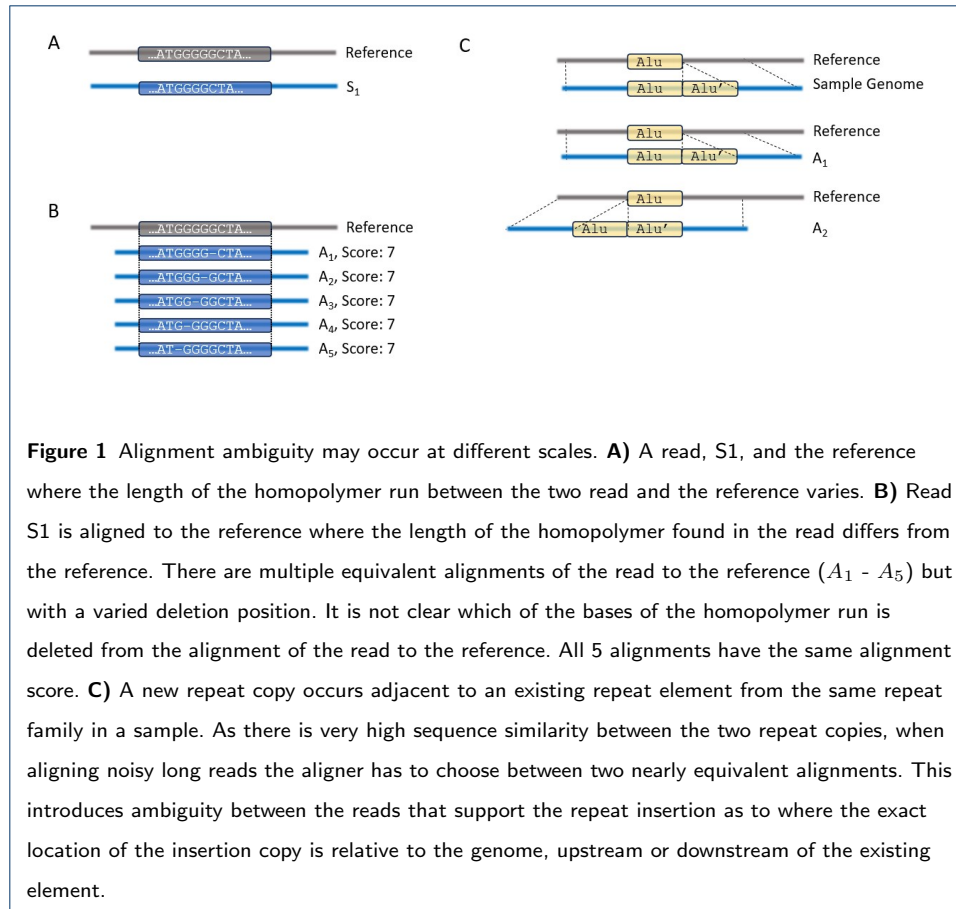
33 the genome can contribute to the loss of tumor suppressors, activation of oncogenes
34 or the generation of fusion proteins that may drive cancer progression [25].

35 Tools for detecting germline and somatic mobile element insertions (MEIs), in-
36 cluding retrotransposons, have been developed for short reads including MELT[26],
37 TraFiC-mem [27], RetroSeq [28] and xTea [29]. While effective, short reads have
38 limited repeat resolution for large insertions, insertions containing varying numbers
39 of repeat copies and insertions into existing repetitive sequence being particularly
40 problematic and hard to detect [30, 31, 32]. Long read technologies like the Ox-
41 ford Nanopore (ONT) and Pacific Biosciences (PacBio) instruments can generate
42 sequencing reads exceeding 10kbp. These reads can therefore fully span a retro-
43 transposon insertion with flanking sequence allowing the genomic location of the
44 insertion to be identified [33] (e.g. a full length ~6kbp LINE-1 element can be fully
45 contained within a 10kbp read with 4kb of flanking sequence available to inform
46 the location of the repeat). This has prompted the development of tools to detect
47 mobile element insertions from long reads such as tldr [34] and xTea-Long [29].
48 These tools have mainly been designed to detect polymorphic repeats that present
49 as heterozygous and homozygous variants within an individual genome and hence
50 they often require multiple reads to support an insertion call. When looking at
51 somatic variation, such as in a tumor, insertions may occur at very low frequen-
52 cies and hence be supported by only a single (or very few) reads depending on the
53 variant allele frequency within the cellular population. Methods designed to detect
54 polymorphic variation may miss these somatic insertion events due to their very
55 low read support. Additionally, many large insertion events in long reads may not
56 be correctly identified by current state of the art long read aligners.

57 While *de novo* assembly of diploid genomes is becoming the gold-standard method
58 for detecting structural variants, including retrotransposon insertions, most meth-
59 ods currently rely on mapping reads to a reference genome [35]. Hence, having high

60 quality read-reference alignments is crucial to detecting MEIs. Aligners such as
61 minimap2 [36] are designed to tolerate large gaps in the alignment by using affine
62 gap-scoring penalties [36, 37, 38, 39]. Despite these scoring schemes, the alignments
63 may be truncated just prior to the insertion event, lowering the read support for the
64 insertion. Further, we have found that the location of the insertion event introduced
65 by the aligner may differ read-to-read when the inserted sequence is similar to a
66 repeat copy already existing in the reference, an effect we term *repeat alignment*
67 *ambiguity*, and also observed by [40]. This problem is analogous to the classical case
68 of aligning two sequences with different lengths of homopolymer runs. In that case,
69 the exact base that has been inserted/deleted is not known, so the placement of
70 the gap is ambiguous with multiple alignments having the same alignment score
71 (illustrated in **Figure 1B**). In our case, the repeat could be placed either before or
72 after the existing element and the aligner's choice may depend solely on the pat-
73 tern of matches/mismatches caused by sequencing errors (**Figure 1C**). Later in the
74 Results section, we quantify how often this artifact occurs as a function sequence
75 divergence between the repetitive elements.

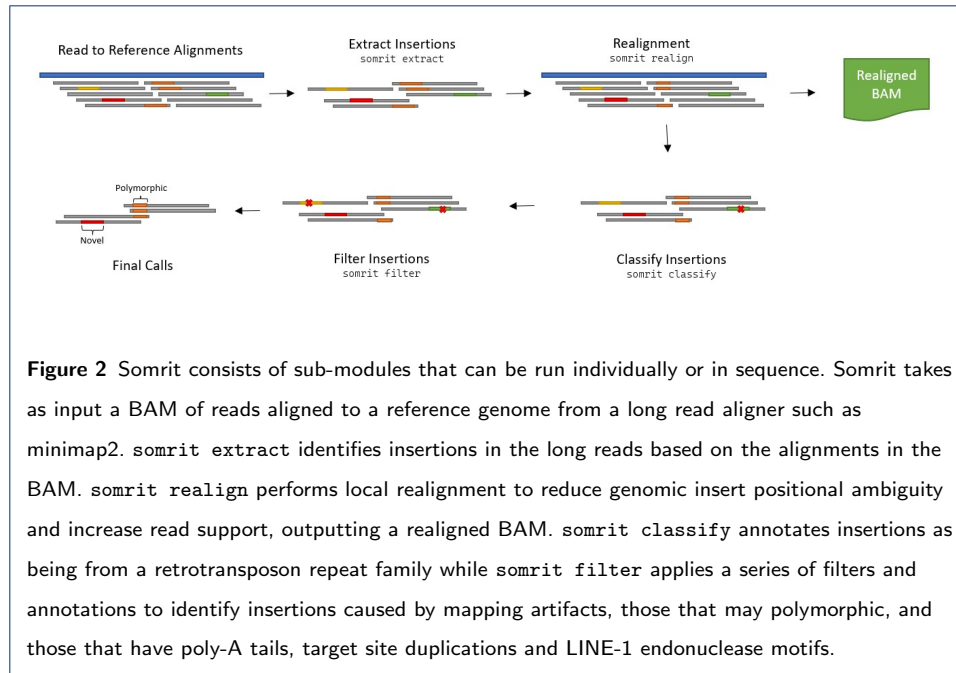
76 To address both repeat alignment ambiguity and alignments truncated due to
77 mobile element insertions we developed somrit, the somatic retrotransposon inser-
78 tion toolkit, to detect novel somatic retrotransposon insertion events and MEIs
79 from long reads mapped to a reference genome. Somrit is a modular toolkit consist-
80 ing of subprograms with standard input/output files. Importantly, it has steps not
81 found in traditional SV detection workflows aimed to recover insertions that may
82 be missed due to alignment truncation, and to resolve repeat alignment ambiguity.
83 In this work we first describe somrit, providing an overview of each sub-module and
84 then show how somrit can be used to detect novel somatic MEIs and help avoid
85 false positive translocations from general purpose SV callers.



86 Methods

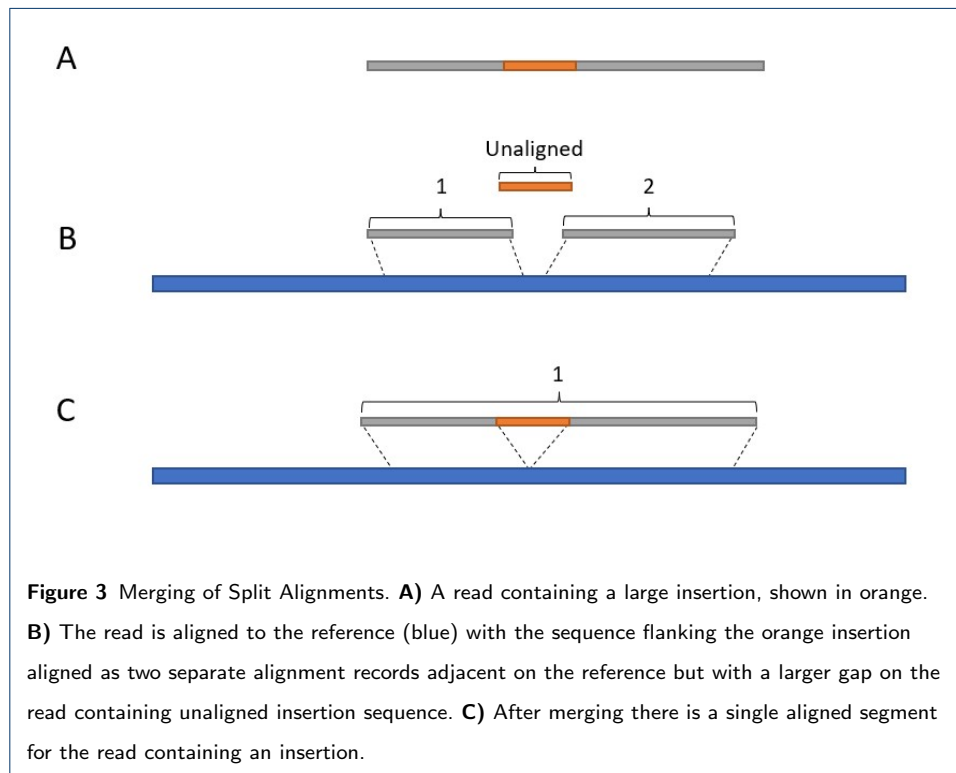
87 Somrit contains individual sub-modules designed to be run as standalone tools or
88 as part of a larger workflow. **Figure 2** shows the somrit modules in the order they
89 would normally be run to call somatic retrotransposon insertions.

90 `somrit extract`. Somrit's first step is to extract candidate retrotransposon in-
91 sertions from the reads aligned to the reference genome. We consider two cases. In
92 the simple case, reads containing long insertions (by default, 50bp) with a minimum
93 flanking anchor sequence (500bp) are exported to a tsv file. Second, we attempt to
94 recover alignments that were erroneously split due to the presence of a large in-
95 sersion within the read, shown in **Figure 3B**. Let $q.d, t.d$ be the distance between
96 the pair of alignments on the query (read) and target (reference), respectively. We
97 merge the pair of alignments when $q.d \geq 100$ and $t.d \leq 100$ by writing the first BAM



98 record with a new CIGAR string and deleting the second BAM record (**Figure 3C**).

99 The coordinates of these insertions are also output to the TSV file.



100 `somrit realign`. Next, we perform local realignment around candidate insertions
101 to reduce alignment ambiguity and increase read support. The explicit goal of `somrit`
102 is to detect novel repetitive insertions that have high sequence similarity to existing
103 repeat copies within the reference genome. This can make it difficult for the read
104 mapper to identify the correct insertion location when the insert happens to occur
105 in a region already containing a copy of the repeat. In this case the output of
106 `somrit extract` may have different representations of the same insertion event
107 across multiple reads. As the level of read support is a key parameter for structural
108 variant calling this can cause false negatives, or worse, the caller might identify
109 multiple separate insertions. `somrit realign` aims to reconcile the alignments of
110 all reads carrying an insertion and recover supporting reads entirely missed by the
111 mapping and extract steps. This process is inspired by the predominant approach for
112 small variant calling, which generates candidate haplotypes containing combinations
113 of variants [41], [42], [43]. Here, we apply the same idea to large insertions found
114 from long reads. `somrit realign` focuses on insertions at least n (default $n = 50bp$)
115 employing a process similar to that of Iris [44], a tool for refining the position of
116 structural variants in long reads, and SVJedi-Graph [45].

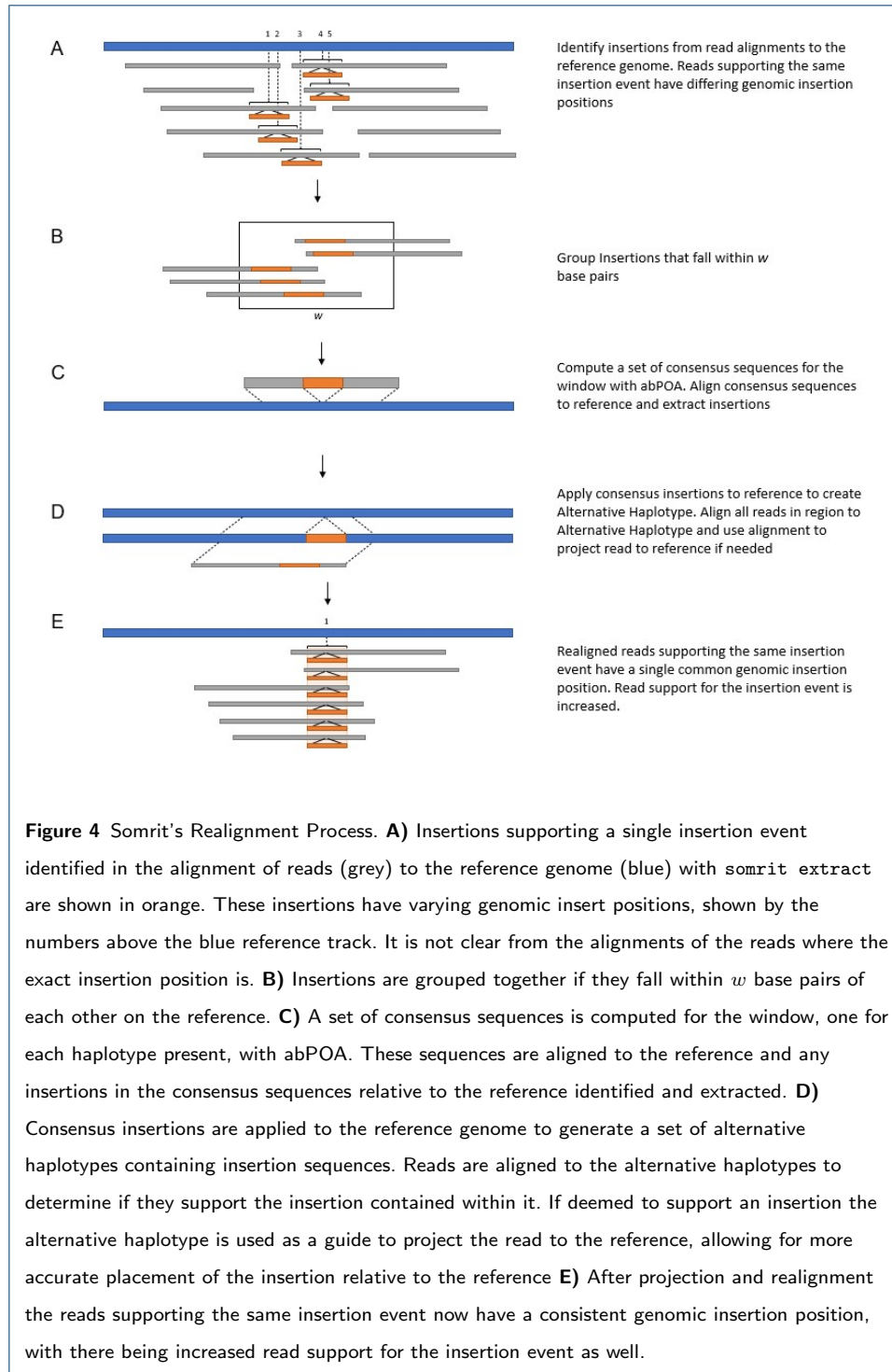
117 The `realign` module contains two steps: realignment and alignment projection.
118 In the realignment step insertions identified by `somrit extract` (**Figure 4A**) are
119 grouped based on genomic position into 1000bp windows (**Figure 4B**). Adjacent
120 windows that contain insertions are merged together up to a max window size
121 (default 25000bp). A set of consensus sequences (default=3, one for each germline
122 haplotype for assumed diploid samples, and one for a germline haplotype containing
123 the putative somatic insertion) is generated for each window from the insertion-
124 supporting reads using abPOA [46]. Each consensus sequence is aligned back to
125 the reference sequence for this window to identify a refined insertion position and
126 sequence (**Figure 4C**). For each refined insert identified (minimum size 50bp), the

127 insert is spliced into the reference to generate an alternative haplotype sequence
128 that only contains a single insertion. Each read in the window is then aligned to the
129 original reference as well as each alternative haplotype for the window. If the read's
130 alignment score against the alternative haplotype is greater than the alignment
131 score to the reference we note the read as supporting the insertion and flag it for
132 projection (**Figure 4D**).

133 Once all reads have been tagged with the insertions they support, we calculate
134 a new read-to-reference alignment in a step we call alignment projection. For each
135 read a new haplotype is constructed by splicing in all insertions supported by that
136 read. The read is then aligned to this haplotype, and the haplotype is aligned to the
137 reference genome. We then iterate over the pair of read-to-haplotype and haplotype-
138 to-reference CIGAR strings to determine the read-to-reference alignment. The BAM
139 record for the read is then updated based on this projected alignment (**Figure 4E**).
140 If a read is not selected for projection the original BAM record(s) for the read are
141 retained. In addition to an updated BAM, `somrit realign` outputs an updated
142 tsv with the coordinates and sequences of insertions after realignment.

143 `somrit classify`. The set of refined insertions are then assigned to a retro-
144 transposon repeat family. Each insert's sequence is aligned to a library of known
145 human retrotransposon consensus sequences compiled from Tubio et al [27] (avail-
146 able: <https://gitlab.com/mobilegenomesgroup/TraFiC>) and DFAM [47] using min-
147 imap2's mappy API. Inserts that have no mapping to a retrotransposon consensus
148 sequence with quality higher than 20 are unassigned, otherwise the insert is assigned
149 to the repeat family with the highest alignment score.

150 `somrit filter`. The final step applies annotations and filters to the classified
151 repeats by appending new columns to the TSV record similar to VCF filter columns.
152 These filters include:



153

- `IN_CONTROL_SAMPLE`: If `somrit` is run with multiple samples and one is designated a matched normal control sample (e.g. for tumour/normal pairs),

154

155 this filter is used to identify which insertions are also found in the designated
156 control sample within ± 500 bp.

157 • IN_CENTROMERE and IN_TELOMORE: insertions that fall within a cen-
158 tromeric or telomeric region respectively based on a bed file provided by the
159 user.

160 • LOW_MAPPING_QUALITY: Insertions found in a genomic window (\pm -
161 500bp of the insertion position) where the average mapping quality over all
162 reads aligned in the region is ≤ 20 .

163 • MIN_READS: This filter flags insertions that do not have the user-specified
164 number of supporting reads (by default, 1).

165 • IN_SECONDARY_MAPPING: If reads supporting the insertion have multiple
166 alignments with a mapping quality ≥ 20 that overlap the insertion position,
167 the insertion is flagged.

168 • POLYMORPHIC: This filter flags insertions that appear to be polymorphic
169 germline variation between the individual and the reference rather than so-
170 matic variation, based on the fraction of insertion supporting reads relative to
171 all reads aligned in a genomic window. Let $f(w)$ be the fraction of insertion
172 supporting reads within a window of w bp. An insertion is flagged as polymor-
173 phic if $f(500) > 0.8$ or $f(200) > 0.5$ or $f(100) > 0.3$. The varied window sizes
174 and fraction cutoffs were used as some genomic regions varied in coverage. A
175 larger window may contain a number of reads that align within the window
176 but do not overlap the insertion position, with parts of the window flanking
177 the insertion having higher coverage than the area around the insertion itself.

178 In addition to these filters, each putative insertion is annotated with features
179 expected of real retrotransposon insertions:

- 180 • Annotate Poly-A Tail: the insertion contains a Poly-A/T tail ≥ 10 bp within
181 50bp of either the start or end of the insert sequence. If present the Poly-A/T
182 sequence is listed in the output column.
- 183 • Annotate TSD: retrotransposon insertions have characteristic sequence du-
184 plication generated as part of the re-insertion process. A local dynamic pro-
185 gramming alignment is used to identify duplicated sequence at least 5bp in
186 length between the start or end of the insertion and the genomic region. If a
187 duplication of at least 5bp is found, the TSD sequence is listed in the output
188 column.
- 189 • Annotate Motif: the reference sequence 2bp upstream and 4bp downstream of
190 the identified insertion position is extracted for comparison to the canonical
191 LINE-1 endonuclease recognition motif sequence. The motif sequence is listed
192 in the output column.

193 Implementation and Pipeline

194 The tsv file generated by `somrit filter` is the final output of the program, with
195 the inserts passing all filters considered the final called somatic insertions. Som-
196 rit is implemented python (`extract`, `classify` and `filter` modules) and C++
197 (`realign`). The code for all modules, a documentation of parameters, a tutorial,
198 and a snakemake file to automate the process of running all 4 module sequentially
199 are available at <https://github.com/adcosta17/somrit>.

200 Results

201 In this section we first quantify how often repeat alignment ambiguity occurs. Next
202 we evaluate somrit's ability to detect both polymorphic and somatic insertions
203 from simulated and real nanopore data, comparing somrit to existing tools for both
204 tasks and showing how somrit's use of local realignment to reduce repeat positional
205 ambiguity and increase read support improved its ability to detect MEIs compared
206 to other tools. Finally, we finally show how realignment around MEIs can reduce

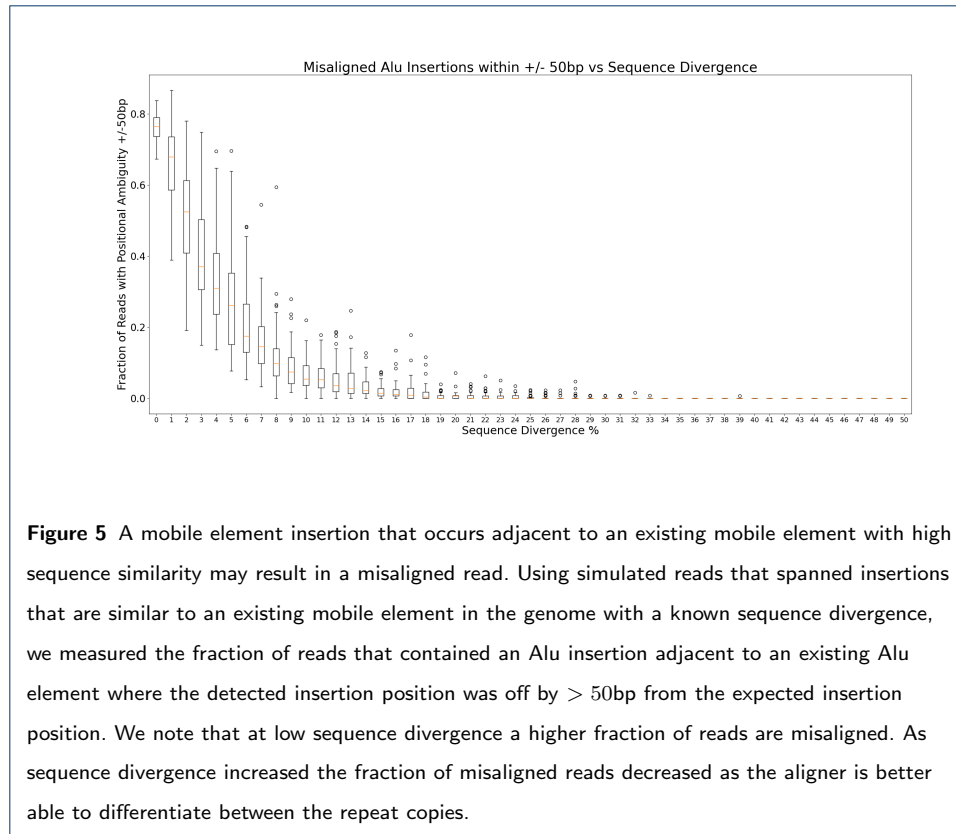
207 the number of false positive translocation calls from general purpose structural
208 variation detection tools.

209 0.1 Repeat Alignment Ambiguity

210 We first quantified how often repeat alignment ambiguity occurs, specifically when
211 a novel mobile element insertion occurs adjacent to an existing copy of the same
212 mobile element in the genome. Using the RepeatMasker [48] annotation of GRCh38
213 we identified existing Alu elements at least 250bp in length. We then generated an
214 insertion by randomly selecting an Alu element present in GRCh38, modifying it to a
215 set level of sequence divergence and then inserting it back into the genome beside the
216 original copy. This process simulates the insertion of an identical or near identical
217 mobile element directly next to an existing mobile element. We then simulated long
218 reads using pbsim2 [49], mean read length of 30kb and per-base accuracy of 95%,
219 that supported the insertion and aligned the reads back to GRCh38 using minimap2
220 [36]. We parsed these read alignments to identify where the aligner had placed the
221 insertion and compared it to the expected position where we had made the insertion.
222 We generated 100 insertions for sequence divergence from 0-50%, in steps of 1%.
223 If a read was mapped > 50 bp away from the expected insertion position it was
224 considered to be misaligned. We see from the results shown in **Figure 5** that at
225 low levels of sequence divergence there is a high fraction of reads that misaligned,
226 exhibiting repeat alignment ambiguity, as the aligner cannot differentiate between
227 the existing copy and the new insertion copy. As sequence divergence increases the
228 proportion of reads whose insertion is incorrectly placed decreases as the aligner is
229 better able to differentiate between the repeat copies.

230 Detection of Polymorphic Insertions

231 Polymorphic insertions, defined here as variants between an individual's inherited
232 genome and the reference, make up the vast majority of large insertions (≥ 50 bp)
233 found by SV callers. To evaluate the performance of somrit in detecting polymorphic



234 retrotransposon insertions compared to existing tools, xTea-Long and tldr, we used
235 publicly available data downloaded from the Human Pan-genome Reference Con-
236 sortium (HPRC)[35]. For samples HG00438, HG00621, HG00673, HG00735 and
237 HG00741 we downloaded raw Oxford Nanopore (ONT) reads, the accompanying
238 diploid hifiasm assembly [50] and matching RepeatMasker [48] annotation of the
239 assembled contigs. The ONT reads for each of the five samples were downsampled
240 to set coverage levels, with three replicates drawn for each coverage level. These
241 read sets were then aligned to GRCh38 and the resulting BAM files passed as input
242 to somrit, xTea-Long and tldr.

243 *Generating ground truth calls*

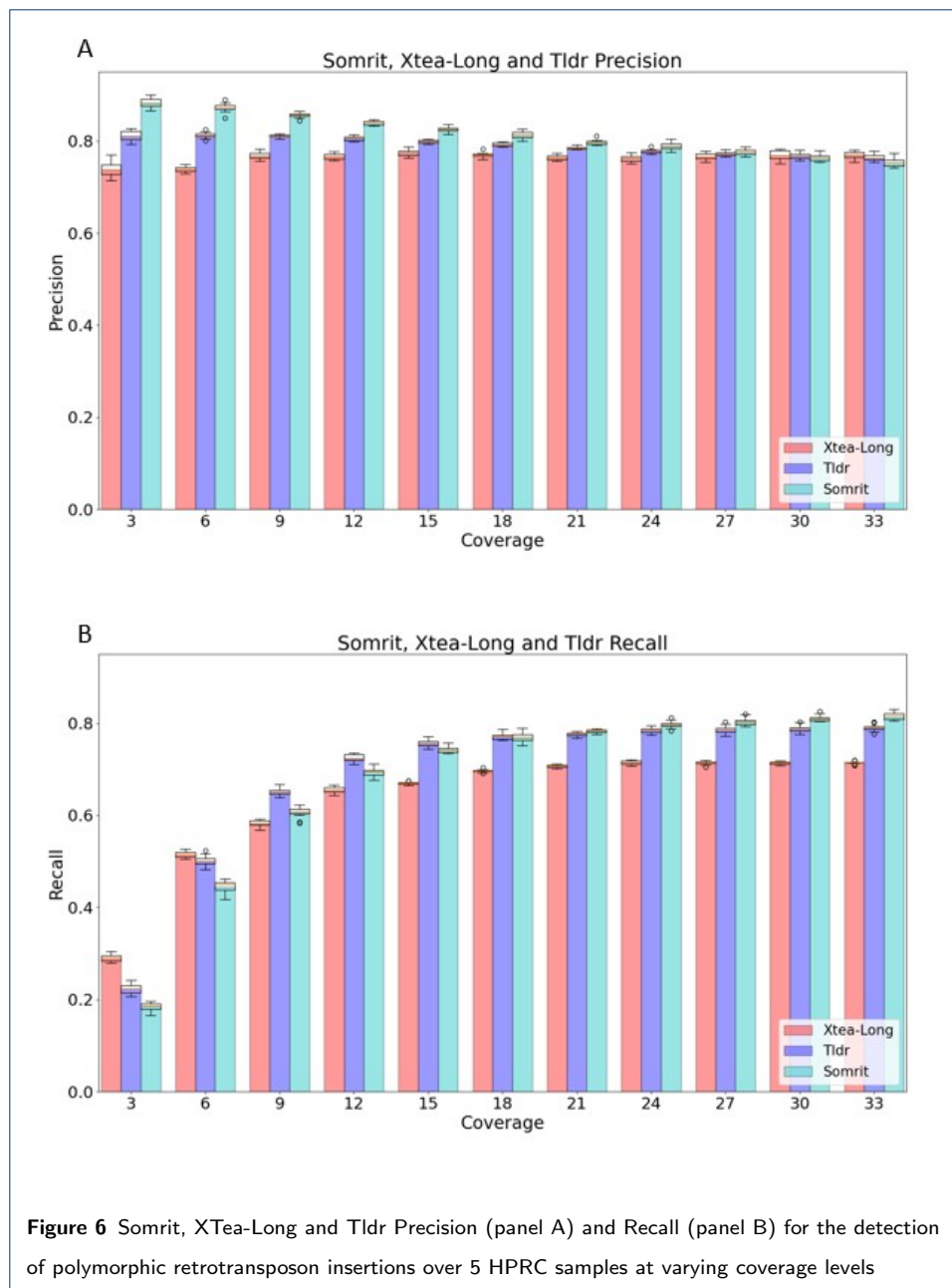
244 We used the high quality diploid assemblies for each HPRC sample to derive a set
245 of ground truth insertions. For each sample the maternal and paternal contigs were
246 aligned to GRCh38 with minimap2 (v2.22-r1101, preset `asm5`). Insertions at least

247 100bp in length on contig alignments with mapping quality ≥ 20 were identified
248 and extracted from the alignment CIGAR strings. These insertions were then an-
249 notated with repeat families using the RepeatMasker [48] annotation of the contigs.
250 Insertions where at least half the insertion sequence on the contig was annotated
251 to a single retrotransposon repeat family by RepeatMasker were assigned to that
252 repeat family. These insertions, their repeat family annotation (if any), and their
253 reference coordinate are our truth set for the subsequent evaluation.

254 *Comparing somrit, tldr and xTea-Long*

255 We ran somrit with default settings except for increasing the minimum insertion
256 supporting read threshold to 3 and requiring at least 1000bp of flanking sequence
257 on each side of an insertion on the read. Also, we did not use somrit's polymor-
258 phic filtering step. We ran tldr (v1.2.2) and xTea-Long (v0.19) using their default
259 parameters, except for also requiring at least 1000bp of flanking sequence for tldr.
260 We compared the retrotransposon insertion calls made by each tool to the ground
261 truth described in the previous section. A called insertion was considered a true
262 positive if it is within 500bp of an insertion call from the same retrotransposon
263 repeat family in the truth set. All other called insertions were considered false posi-
264 tives. Any insertions in the truth set where we did not find a called insertion within
265 500bp annotated to the same retrotransposon repeat family were considered false
266 negatives.

267 **Figure 6** shows the precision and recall for all three tools at different coverage
268 levels, with 3 replicates per sample per coverage level. Somrit has higher precision
269 at lower coverage and slightly lower precision at higher coverage compared to xTea-
270 Long and tldr. At higher levels of coverage false insertions from mapping artifacts
271 may have their read support increased beyond the threshold of 3 supporting reads,
272 resulting in false positive calls. xTea-Long had higher recall than somrit and tldr at
273 the lowest coverage levels, while tldr and somrit had higher recall than xTea-Long



274 at higher coverage levels, with somrit's recall slightly higher than tldr at the highest
275 coverage levels. While tldr considers both reads that fully span an insertion event
276 and reads whose alignments are clipped at the insertion event as supporting the
277 insertion, somrit only looks at reads with a spanning insertion. This may contribute
278 to the observed recall for the somrit being lower than tldr at lower coverage levels.

279 Detection of simulated somatic insertions

280 While somrit can be used to detect polymorphic insertions it is designed primarily
281 to detect somatic insertions. This is more challenging than detecting polymorphic
282 insertions as the read support may be much lower, even down to a single read. While
283 tools like tldr have been previously used to detect rare somatic retrotransposon
284 insertions[51], they are not designed for detecting insertions from a single read,
285 which is one design goal of somrit. To quantify somrit's ability to detect novel
286 somatic retrotransposon insertion events we ran somrit on a set of simulated novel
287 somatic retrotransposon insertions. For comparison we ran tldr and xTea-Long with
288 the minimum read support lowered to 1. We also ran the general purpose SV caller
289 Sniffles2 (v2.0) [52] in its somatic detection mode.

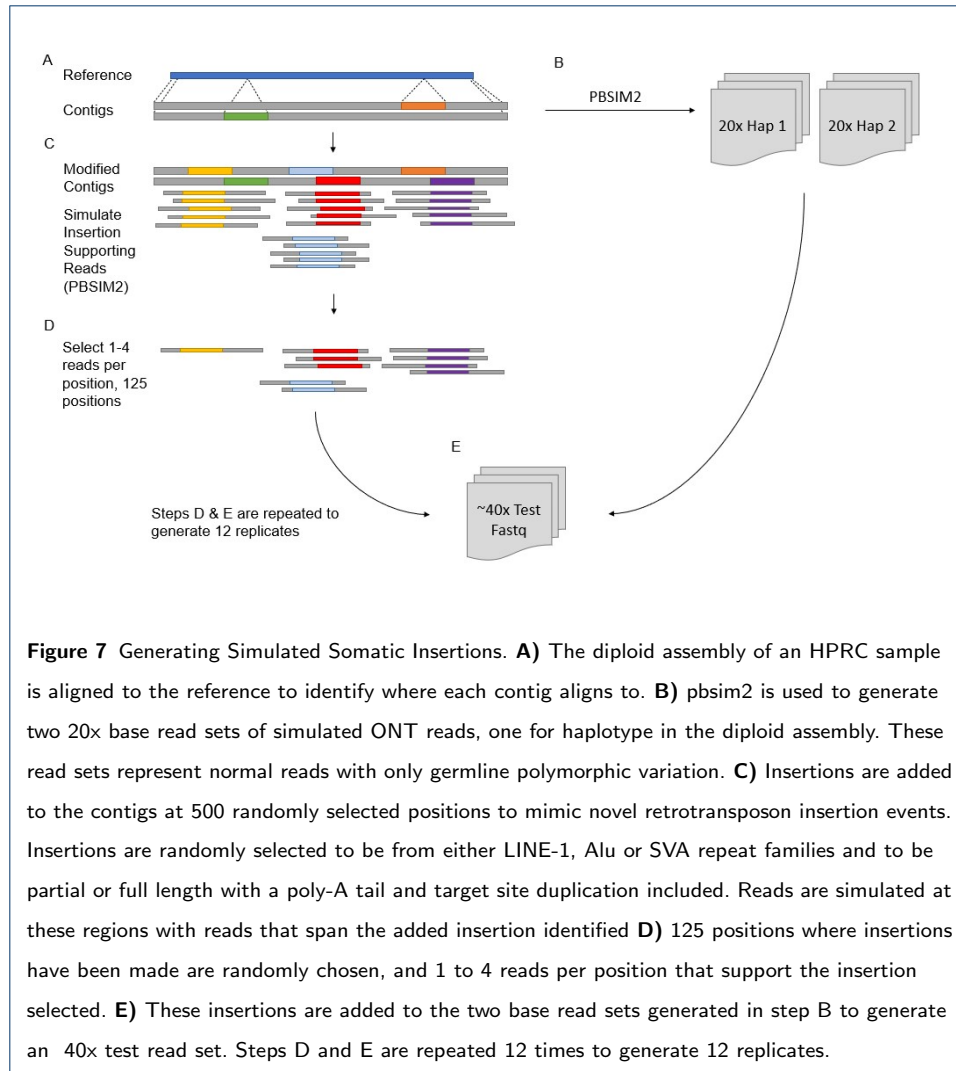
290 *Generating Simulation Data*

291 We simulated long reads with pbsim2[49] from the diploid assembly for 4 of the
292 5 HPRC samples: HG00438, HG00621, HG00735 and HG00741. We simulated 20x
293 coverage from both the maternal and paternal contigs (40x total), to act as a base-
294 line read set free from somatic variation (**Figure 7A** and **B**). Next, we randomly
295 selected 500 positions on the assembly contigs for the location of somatic inser-
296 tions. For each selected position we randomly choose a retrotransposon repeat fam-
297 ily (LINE-1, Alu or SVA) and insert length. For half of the selected positions a full
298 length insertion is selected with the length of the remainder drawn uniformly be-
299 tween 100 bp and the full repeat length. Using a consensus sequence for the repeat
300 family selected [27][47], we truncated the sequence if needed removing bases from
301 the 5' end, generated a poly-A tail at the 3' end between 10 and 40 bp and added
302 the modified sequence to the contig at the selected position. We also generated a
303 target site duplication (TSD) to mimic the real genomic insertion process. We then
304 simulated long reads from each contig in a 50kbp flanking region around the inser-
305 tion position. We recorded which reads fully spanned the insertion with non-insert

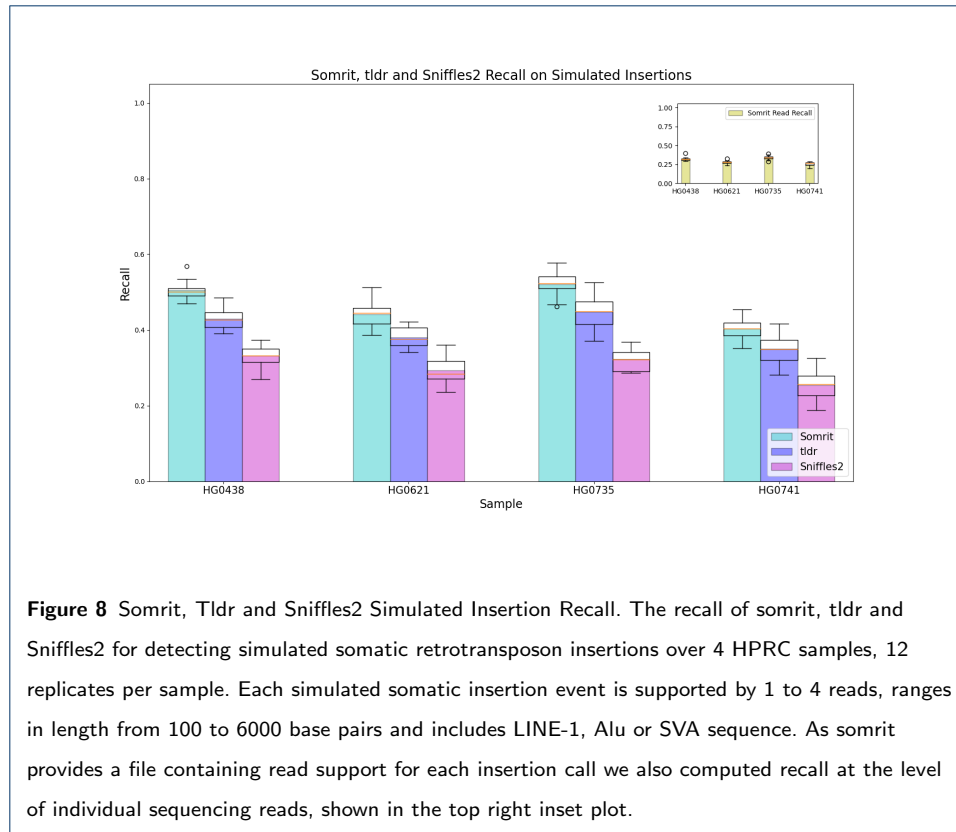
306 flanking sequence. For each position on the contig where we added an insertion we
307 noted the insert sequence, repeat family, poly-A tail length, TSD, and the expected
308 location of the insertion on GRCh38, and the list of read names deemed to sup-
309 port the insertion event (**Figure 7C**). In order to generate each test read set that
310 contained simulated somatic retrotransposon insertion events we started with the
311 simulated baseline read set for the sample and randomly selected 125 of the 500
312 positions where we had simulated a novel somatic insertion event. For each of these
313 125 positions we randomly selected between 1 to 4 supporting reads and added them
314 to the base read set to generate a test read set (**Figure 7D and E**). This process
315 is repeated to generate 12 replicates per sample, each with a randomly generated
316 set of 125 positions.

317 *Detecting simulated insertion events*

318 For each replicate we mapped the reads to GRCh38 with minimap2 and then ran
319 tldr, Sniffles2 and somrit as described above to identify somatic insertions in the
320 same 4 HPRC samples. We ran tldr and somrit in their default settings but with
321 the minimum read flank size set to 1000bp and minimum read support set to 1. We
322 compared the calls made by each tool to the truth data for each replicate, noting the
323 insertion events that were detected as passing retrotransposon insertions by somrit
324 and tldr as well any insertions detected by Sniffles2 (as Sniffles2 does not annotate
325 insertions as being from a retrotransposon repeat family). Passing insertion calls
326 made within 500bp of an expected simulated insertion position with the same repeat
327 family annotation were considered true positives. If no passing insertion with the
328 same repeat annotation was found within 500bp of a simulated insertion position,
329 the simulated insertion position was considered a false negative. We additionally ran
330 xTea-Long in its default settings but as xTea-Long is not designed for the detection
331 of somatic insertions it was unable to detect almost all the simulated insertions.
332 Thus we do not report the results of xTea-Long in this analysis.



333 We computed recall for each tool and over all replicates for the 4 HPRC samples
334 (**Figure 8**). Somrit had the highest recall, followed by Sniffles2. As somrit generates
335 calls for individual reads so we also calculate a read-level recall (the proportion of
336 reads, rather than positions, that have an insertion that were called by somrit;
337 **Figure 8 inset**). Even though somrit outperformed all other tools, its best recall
338 for any sample did not exceeded 60%, indicating the difficulty of detecting insertions
339 with minimal read support. Due to the repetitive nature of the human genome a
340 repeat insertion event in a long read makes it harder for a long read aligner to
341 correctly align the read. Thus reads with insertions may either have split alignments
342 at the expected insertion position where realignment is unable to increase the read



343 support, or may not align to the expected region of the reference, or the reference
344 at all. All tools evaluated rely on alignment BAM files as input and thus are limited
345 by the shortcomings of current long read aligners when aligning reads with repeat
346 insertions.

347 Identifying novel L1-mCherry retrotransposon insertions in Nanopore reads

348 In a recent analysis by Gerdes et al [51] HeLa cells were treated with a plasmid vector
349 containing a modified mouse LINE-1 fused with an mCherry reporter. Successful
350 integration of this modified vector into the HeLa cells results in a novel insertion
351 of the L1-mCherry construct sequence in the cells. This is an ideal experiment to
352 test the performance of somatic retrotransposon insertion detection tools as the
353 mCherry sequence allows novel insertions to be definitively identified. Using the
354 ONT data generated by Gerdes et al for these samples we evaluated somrit's ability
355 to call these previously identified insertions.

356 We downloaded reads from the HeLa/L1-mCherry experiment that Gerdes et al
357 deposited in the ENA. The construct generated insertions of up to 10.2kbp once
358 integrated into the HeLa genome[51]. Each of the five read sets are WGS nanopore
359 sequencing run of a HeLa cell line expanded over 3-5 passages from single L1-
360 mCherry insertion harboring colonies, barcoded and pooled in equal amounts be-
361 fore being sequenced with a single PromethION flow cell [51]. It is expected that
362 individual L1-mCherry insertions will appear as somatic insertion events occurring
363 in a small fraction of cells in the sample, with possibly just a single read supporting
364 the insertion event.

365 We ran somrit and tldr on the 5 samples, with minimum read support set to 1
366 and a minimum read flank size of 1000 bp. We ran both tools using the L1-mCherry
367 consensus as the only possible repeat family. Both tools were run in multi-sample
368 mode, where multiple samples are analysed jointly. This allowed for insertion calls to
369 have supporting reads from multiple biological replicate cell line colonies, making it
370 easier to identify novel somatic events that only occur in a single sample rather than
371 those that represent any possible polymorphic variation seen across all samples. As
372 the L1-mCherry sequence contained both a full length LINE-1 sequence and the
373 mCherry protein coding gene we explicitly filtered the final detected insertion calls
374 to ensure they aligned with at least one base pair to the non LINE-1 sequence of
375 the L1-mCherry construct.

Table 1 L1-mCherry insertions detected by Somrit and Tldr. The total number of L1-mCherry insertions, the number of insertions unique to each tool, the number of insertions shared between the tools and the number of insertions shared with Gerdes et al's call set.

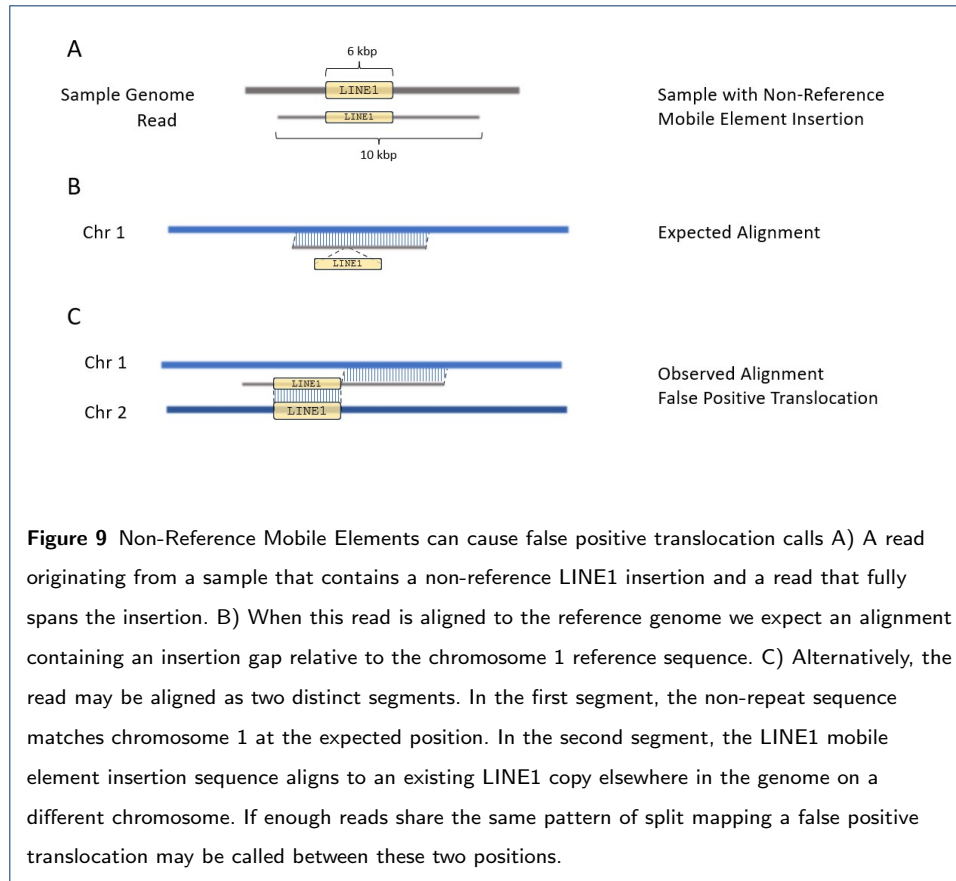
	Total Inserts	Unique To Tool	Shared Between Tools	Shared with Gerdes et al
somrit	67	10	57	35/41
tldr	62	5	57	40/41

376 The results in **Table 1** show that somrit is able to identify additional insertions
377 containing L1-mCherry sequence beyond what was initially detected by Gerdes et
378 al. Each of the insertion calls made by somrit used ≤ 5 reads, with 37 insertions
379 being detected with just a single supporting read. Both tools detected insertions

380 that were unique to the tool. The 5 insertions unique to tldr represent the 5 of
381 the 6 insertions of the Gerdes et al set that somrit did not identify. Of these 5
382 somrit was able to identify 3 as being annotated to the L1-mCherry sequence, but
383 these insertions were flagged for not having enough flanking sequence in the read
384 alignment. Manual inspection of the 10 insertions unique to somrit showed that 8
385 of the 10 were mapped to the 3' end of the L1-mCherry construct sequence, with
386 target site duplications and poly-A tails or mapped to non-LINE1 sequence in the
387 L1-mCherry construct, indicating they are likely to be novel insertions of the L1-
388 mCherry construct sequence into the cells. The remaining two insertions called only
389 by somrit mapped mainly to the LINE1 portion of the L1-mCherry construct, with
390 only a small fraction of the insertion sequence aligning to non-LINE-1 sequence in
391 the construct, with no mCherry specific sequence identified. Thus these insertions
392 are likely false positives.

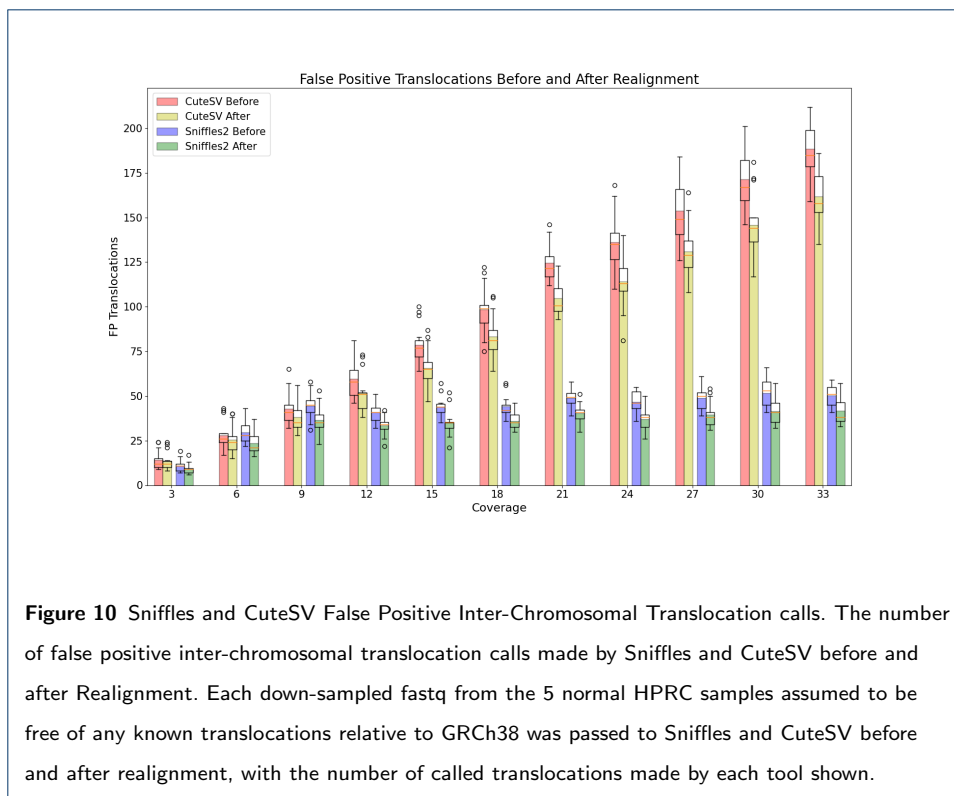
393 Repeat realignment reduces false positive translocation calls

394 While somrit is primarily designed to detect mobile element insertions, local re-
395 alignment with `somrit realign` may be useful in reducing false positive calls from
396 general purpose SV callers such as Sniffles2 and CuteSV. As the human genome
397 contains many copies of mobile repetitive elements such as retrotransposons and
398 the exact location of these elements varies between individuals and the reference
399 genome, some sequencing reads that partially cover a non-reference mobile repet-
400 itive element may appear to have a split mapping when aligned to the reference
401 genome. This split mapping occurs as the aligner may map the non-repetitive se-
402 quence correctly but maps the portion of the read containing the mobile element
403 sequence to an existing repetitive element copy elsewhere in the genome. If a num-
404 ber of reads are misaligned in this way a general purpose SV caller may incorrectly
405 interpret this as a translocation. We propose that `somrit realign` may help reduce
406 the number of false positive translocation calls induced by this effect (**Figure 9**).



407 To evaluate how `somrit realign` could be used to reduce false positive transloca-
408 tion calls we first ran Sniffles2 and CuteSV on the aligned reads for each HPRC sam-
409 ple at various read depths, noting the number of inter-chromosomal translocation
410 calls made. As the HPRC samples are generated from lymphoblastoid cell lines clas-
411 sified as karyotypically normal, thus free of any known inter-chromosomal translo-
412 cation events, we considered any inter-chromosomal translocation calls made by the
413 tools as false positives. We then detected candidate insertions (`somrit extract`)
414 for realignment (`somrit realign`) to generate a new BAM file. We then used the
415 realigned bam as input into Sniffles2 and CuteSV, noting the number of called
416 inter-chromosomal translocations after realignment.

417 **Figure 10** shows that in both tools there is a reduction in the number of false
418 positive inter-chromosomal translocation calls made after realignment. We observe
419 up to 41% and 31% reduction in the number of false positive inter-chromosomal



420 translocation calls made by Sniffles2 and CuteSV, respectively, with the effect most
421 noticeable at higher coverage levels.

422 Discussion

423 In this paper we introduce somrit, a toolkit for the identification of somatic retro-
424 transposon insertion events in long reads. We show that somrit is able to detect
425 existing polymorphic MEIs with comparable precision and recall to state-of-the-art
426 tools. We also show that somrit is able to detect somatic MEIs in both simulated
427 and real nanopore data, outperforming other methods at identifying insertions with
428 single read support. In addition, we show that realignment around MEIs can reduce
429 false positive translocation calls in general purpose SV callers.

430 While these results show somrit's effectiveness, they have limitations. Somrit
431 firstly requires a large amount of time and memory to run, more than existing
432 tools for retrotransposon insertion detection. The majority of the computational
433 burden lies with `somrit realign` and the generation of consensus sequences with

434 abPOA. While abPOA uses adaptive banded alignment to reduce the time and
435 memory usage needed to compute a consensus sequence this process is still time
436 and memory intensive. As the time and memory needed to compute the consensus
437 sequences scales with both the number and length of input sequences, limiting the
438 number of read sequences used to generate the consensus sequences at high input
439 coverage can be considered to reduce the time and memory usage.

440 Somrit is also limited in its ability to realign insertions that may be missed by the
441 initial mapping of reads to the reference. If there is an insertion present in a sample,
442 but there is no alignment made by the aligner that introduces an alignment gap for
443 the insertion, somrit is unable to recover this insertion. This becomes problematic
444 for somatic detection where insertions may have low read support and an aligner
445 may clip the alignment of the single read supporting an insertion event, with somrit
446 unable to detect or recover the insertion.

447 While realignment is able to increase the read support for genuine insertion events,
448 in some cases the realignment process may increase the number of reads that support
449 a mapping or alignment artifact. The decision to realign a read using an alternative
450 haplotype containing an insertion as a guide is based on comparing the alignment
451 score between the alternative and reference haplotypes. A higher scoring alignment
452 to the alternative haplotype indicates that the read may support the insertion. If
453 an false alternative haplotype is generated by a mapping artifact, and the read has
454 a marginally higher alignment score to this haplotype, a mapping artifact could be
455 introduced into the read, with the read now supporting a false insertion. We believe
456 this effect can be seen in the decreased precision somrit has compared to other
457 tools at higher levels of coverage for polymorphic insertion detection, as at higher
458 coverage levels there is a greater chance a mapping artifact supported by one or two
459 reads has its read support increased through realignment to three or more reads.

460 More stringent criteria for selecting alternative haplotypes may help alleviate this
461 issue.

462 **Competing interests**

463 J.T.S. receives research funding from Oxford Nanopore Technologies.

464 **Acknowledgements**

465 We thank the members of the Simpson Lab at the Ontario Institute for Cancer Research for their helpful
466 suggestions, ideas and feedback in debugging. A.D. was the recipient of an Ontario Graduate Scholarship. J.T.S
467 receives funding from the Government of Ontario through the Ontario Institute of Cancer Research.

468 **Code Availability**

469 Somrit is available at <https://github.com/adcosta17/somrit>. Scripts to generate the evaluation and analysis
470 performed as well as generate the plots shown in this paper can be found at
471 <https://github.com/adcosta17/somrit-test>.

472 **Author details**

473 ¹Ontario Institute for Cancer Research, Toronto, Canada. ²Department of Computer Science, University of Toronto,
474 Toronto, Canada. ³Department of Molecular Genetics, University of Toronto, Toronto, Canada.

475 **References**

- 476 1. Bourque, G., Burns, K.H., Gehring, M., Gorbunova, V., Seluanov, A., Hammell, M., Imbeault, M., Izsvák, Z.,
477 Levin, H.L., Macfarlan, T.S., *et al.*: Ten things you should know about transposable elements. *Genome biology*
478 **19**(1), 1–12 (2018)
- 479 2. Ostertag, E.M., Kazazian Jr, H.H.: Biology of mammalian L1 retrotransposons. *Annual review of genetics* **35**, 501
480 (2001)
- 481 3. Szak, S.T., Pickeral, O.K., Makalowski, W., Boguski, M.S., Landsman, D., Boeke, J.D.: Molecular archeology
482 of L1 insertions in the human genome. *Genome biology* **3**(10), 1–18 (2002)
- 483 4. Sudmant, P.H., Rausch, T., Gardner, E.J., Handsaker, R.E., Abyzov, A., Huddleston, J., Zhang, Y., Ye, K.,
484 Jun, G., Hsi-Yang Fritz, M., *et al.*: An integrated map of structural variation in 2,504 human genomes. *Nature*
485 **526**(7571), 75–81 (2015)
- 486 5. Penzkofer, T., Dandekar, T., Zemojtel, T.: L1base: from functional annotation to prediction of active line-1
487 elements. *Nucleic acids research* **33**(suppl.1), 498–500 (2005)
- 488 6. Beck, C.R., Collier, P., Macfarlane, C., Malig, M., Kidd, J.M., Eichler, E.E., Badge, R.M., Moran, J.V.: Line-1
489 retrotransposition activity in human genomes. *Cell* **141**(7), 1159–1170 (2010)
- 490 7. Hancks, D.C., Kazazian, H.H.: Roles for retrotransposon insertions in human disease. *Mobile DNA* **7**(1), 1–28
491 (2016)
- 492 8. Hancks, D.C., Kazazian Jr, H.H.: Active human retrotransposons: variation and disease. *Current opinion in*
493 *genetics & development* **22**(3), 191–203 (2012)
- 494 9. Goodier, J.L.: Restricting retrotransposons: a review. *Mobile DNA* **7**(1), 1–30 (2016)
- 495 10. Xiao-Jie, L., Hui-Ying, X., Qi, X., Jiang, X., Shi-Jie, M.: Line-1 in cancer: multifaceted functions and potential
496 clinical implications. *Genetics in Medicine* **18**(5), 431–439 (2016)
- 497 11. Dombroski, B.A., Mathias, S.L., Nanthakumar, E., Scott, A.F., Kazazian Jr, H.H.: Isolation of an active human
498 transposable element. *Science* **254**(5039), 1805–1808 (1991)

- 499 12. Dewannieux, M., Esnault, C., Heidmann, T.: Line-mediated retrotransposition of marked alu sequences. *Nature*
500 *genetics* **35**(1), 41–48 (2003)
- 501 13. Ade, C., Roy-Engel, A.M., Deininger, P.L.: Alu elements: an intrinsic source of human genome instability.
502 *Current opinion in virology* **3**(6), 639–645 (2013)
- 503 14. Ostertag, E.M., Goodier, J.L., Zhang, Y., Kazazian Jr, H.H.: Sva elements are nonautonomous retrotransposons
504 that cause disease in humans. *The American Journal of Human Genetics* **73**(6), 1444–1451 (2003)
- 505 15. Feng, Q., Moran, J.V., Kazazian Jr, H.H., Boeke, J.D.: Human I1 retrotransposon encodes a conserved
506 endonuclease required for retrotransposition. *Cell* **87**(5), 905–916 (1996)
- 507 16. Symer, D.E., Connelly, C., Szak, S.T., Caputo, E.M., Cost, G.J., Parmigiani, G., Boeke, J.D.: Human I1
508 retrotransposition is associated with genetic instability in vivo. *Cell* **110**(3), 327–338 (2002)
- 509 17. Moran, J.V., Holmes, S.E., Naas, T.P., DeBerardinis, R.J., Boeke, J.D., Kazazian Jr, H.H.: High frequency
510 retrotransposition in cultured mammalian cells. *Cell* **87**(5), 917–927 (1996)
- 511 18. Doucet, A.J., Wilusz, J.E., Miyoshi, T., Liu, Y., Moran, J.V.: A 3 poly (a) tract is required for line-1
512 retrotransposition. *Molecular cell* **60**(5), 728–741 (2015)
- 513 19. Pickeral, O.K., Makatowski, W., Boguski, M.S., Boeke, J.D.: Frequent human genomic dna transduction driven
514 by line-1 retrotransposition. *Genome research* **10**(4), 411–415 (2000)
- 515 20. Goodier, J.L., Ostertag, E.M., Kazazian Jr, H.H.: Transduction of 3-flanking sequences is common in I1
516 retrotransposition. *Human molecular genetics* **9**(4), 653–657 (2000)
- 517 21. Ewing, A.D., Kazazian, H.H.: High-throughput sequencing reveals extensive variation in human-specific I1
518 content in individual human genomes. *Genome research* **20**(9), 1262–1270 (2010)
- 519 22. Cajuso, T., Sulo, P., Tanskanen, T., Katainen, R., Taira, A., Hänninen, U.A., Kondelin, J., Forsström, L.,
520 Välimäki, N., Aavikko, M., *et al.*: Retrotransposon insertions can initiate colorectal cancer and are associated
521 with poor survival. *Nature communications* **10**(1), 1–9 (2019)
- 522 23. Han, K., Lee, J., Meyer, T.J., Remedios, P., Goodwin, L., Batzer, M.A.: L1 recombination-associated deletions
523 generate human genomic variation. *Proceedings of the national academy of sciences* **105**(49), 19366–19371
524 (2008)
- 525 24. Sen, S.K., Han, K., Wang, J., Lee, J., Wang, H., Callinan, P.A., Dyer, M., Cordaux, R., Liang, P., Batzer,
526 M.A.: Human genomic deletions mediated by recombination between alu elements. *The American Journal of*
527 *Human Genetics* **79**(1), 41–53 (2006)
- 528 25. Scott, E.C., Gardner, E.J., Masood, A., Chuang, N.T., Vertino, P.M., Devine, S.E.: A hot I1 retrotransposon
529 evades somatic repression and initiates human colorectal cancer. *Genome research* **26**(6), 745–755 (2016)
- 530 26. Gardner, E.J., Lam, V.K., Harris, D.N., Chuang, N.T., Scott, E.C., Pittard, W.S., Mills, R.E., Devine, S.E.,
531 Consortium, .G.P., *et al.*: The mobile element locator tool (melt): population-scale mobile element discovery
532 and biology. *Genome research* **27**(11), 1916–1929 (2017)
- 533 27. Tubio, J.M., Li, Y., Ju, Y.S., Martincorena, I., Cooke, S.L., Tojo, M., Gundem, G., Pipinikas, C.P., Zamora, J.,
534 Raine, K., *et al.*: Extensive transduction of nonrepetitive dna mediated by I1 retrotransposition in cancer
535 genomes. *Science* **345**(6196), 1251343 (2014)
- 536 28. Keane, T.M., Wong, K., Adams, D.J.: Retroseq: transposable element discovery from next-generation
537 sequencing data. *Bioinformatics* **29**(3), 389–390 (2013)
- 538 29. Chu, C., Borges-Monroy, R., Viswanadham, V.V., Lee, S., Li, H., Lee, E.A., Park, P.J.: Comprehensive
539 identification of transposable element insertions using multiple sequencing technologies. *Nature*
540 *communications* **12**(1), 1–12 (2021)
- 541 30. Marsili, L., Duque, K.R., Bode, R.L., Kauffman, M.A., Espay, A.J.: Uncovering essential tremor genetics: The

- 542 promise of long-read sequencing. *Frontiers in neurology* **13** (2022)
- 543 31. Pollard, M.O., Gurdasani, D., Mentzer, A.J., Porter, T., Sandhu, M.S.: Long reads: their purpose and place.
544 *Human molecular genetics* **27**(R2), 234–241 (2018)
- 545 32. Gong, T., Hayes, V.M., Chan, E.K.: Detection of somatic structural variants from short-read next-generation
546 sequencing data. *Briefings in Bioinformatics* **22**(3), 056 (2021)
- 547 33. Jain, M., Olsen, H.E., Paten, B., Akeson, M.: The oxford nanopore minion: delivery of nanopore sequencing to
548 the genomics community. *Genome biology* **17**(1), 1–11 (2016)
- 549 34. Ewing, A.D., Smits, N., Sanchez-Luque, F.J., Faivre, J., Brennan, P.M., Richardson, S.R., Cheetham, S.W.,
550 Faulkner, G.J.: Nanopore sequencing enables comprehensive transposable element epigenomic profiling.
551 *Molecular Cell* **80**(5), 915–928 (2020)
- 552 35. Wang, T., Antonacci-Fulton, L., Howe, K., Lawson, H.A., Lucas, J.K., Phillippy, A.M., Popejoy, A.B., Asri, M.,
553 Carson, C., Chaisson, M.J., et al.: The human pangenome project: a global resource to map genomic diversity.
554 *Nature* **604**(7906), 437–446 (2022)
- 555 36. Li, H.: Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**(18), 3094–3100 (2018)
- 556 37. Jain, C., Rhie, A., Hansen, N.F., Koren, S., Phillippy, A.M.: Long-read mapping to repetitive reference
557 sequences using winnowmap2. *Nature Methods* **19**(6), 705–710 (2022)
- 558 38. Jain, C., Rhie, A., Zhang, H., Chu, C., Walenz, B.P., Koren, S., Phillippy, A.M.: Weighted minimizer sampling
559 improves long read mapping. *Bioinformatics* **36**(Supplement_1), 111–118 (2020)
- 560 39. Ren, J., Chaisson, M.J.: Ira: A long read aligner for sequences and contigs. *PLOS Computational Biology*
561 **17**(6), 1009078 (2021)
- 562 40. Audano, P.A., Beck, C.R.: Small allelic variants are a source of ancestral bias in structural variant breakpoint
563 placement. *bioRxiv*, 2023–06 (2023)
- 564 41. Albers, C.A., Lunter, G., MacArthur, D.G., McVean, G., Ouwehand, W.H., Durbin, R.: Dindel: accurate indel
565 calls from short-read data. *Genome research* **21**(6), 961–973 (2011)
- 566 42. Poplin, R., Ruano-Rubio, V., DePristo, M.A., Fennell, T.J., Carneiro, M.O., Van der Auwera, G.A., Kling, D.E.,
567 Gauthier, L.D., Levy-Moonshine, A., Roazen, D., et al.: Scaling accurate genetic variant discovery to tens of
568 thousands of samples. *BioRxiv*, 201178 (2018)
- 569 43. Garrison, E., Marth, G.: Haplotype-based variant detection from short-read sequencing. *arXiv preprint*
570 *arXiv:1207.3907* (2012)
- 571 44. Kirsche, M., Prabhu, G., Sherman, R., Ni, B., Aganezov, S., Schatz, M.C.: Jasmine: Population-scale structural
572 variant comparison and analysis. *BioRxiv*, 2021–05 (2021)
- 573 45. Romain, S., Lemaitre, C.: Svjedi-graph: improving the genotyping of close and overlapping structural variants
574 with long reads using a variation graph. *Bioinformatics* **39**(Supplement_1), 270–278 (2023)
- 575 46. Gao, Y., Liu, Y., Ma, Y., Liu, B., Wang, Y., Xing, Y.: abpoa: an simd-based c library for fast partial order
576 alignment using adaptive band. *Bioinformatics* **37**(15), 2209–2211 (2021)
- 577 47. Storer, J., Hubley, R., Rosen, J., Wheeler, T.J., Smit, A.F.: The dfam community resource of transposable
578 element families, sequence models, and genome annotations. *Mobile DNA* **12**(1), 1–14 (2021)
- 579 48. Chen, N.: Using repeat masker to identify repetitive elements in genomic sequences. *Current protocols in*
580 *bioinformatics* **5**(1), 4–10 (2004)
- 581 49. Ono, Y., Asai, K., Hamada, M.: Pbsim2: a simulator for long-read sequencers with a novel generative model of
582 quality scores. *Bioinformatics* **37**(5), 589–595 (2021)
- 583 50. Cheng, H., Concepcion, G.T., Feng, X., Zhang, H., Li, H.: Haplotype-resolved de novo assembly using phased
584 assembly graphs with hifiasm. *Nature methods* **18**(2), 170–175 (2021)

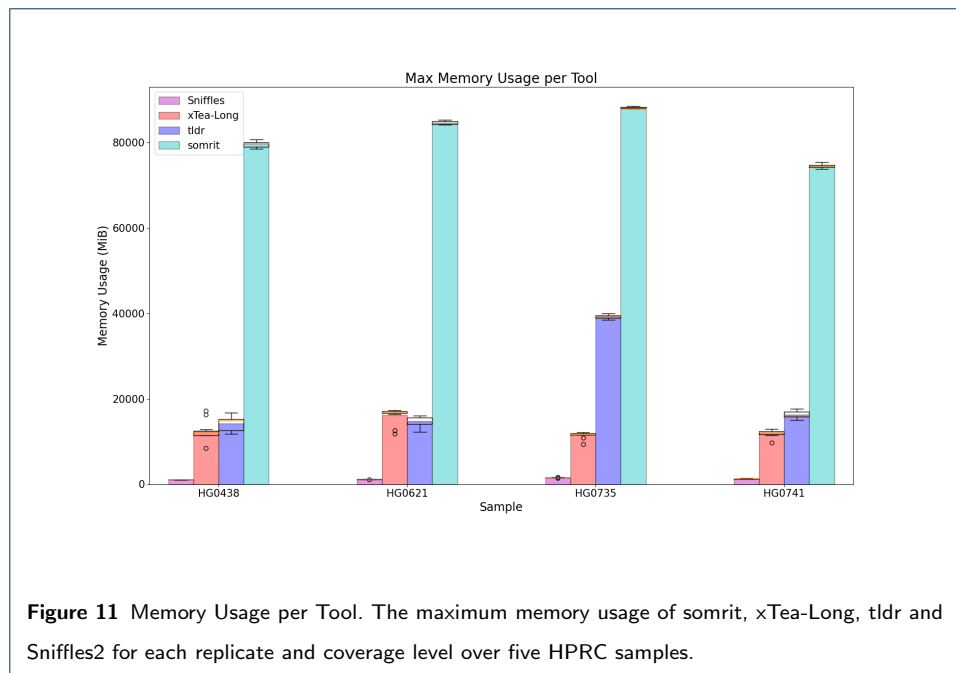
- 585 51. Gerdes, P., Lim, S.M., Ewing, A.D., Larcombe, M.R., Chan, D., Sanchez-Luque, F.J., Walker, L., Carleton,
586 A.L., James, C., Knaupp, A.S., *et al.*: Retrotransposon instability dominates the acquired mutation landscape
587 of mouse induced pluripotent stem cells. *Nature Communications* **13**(1), 1–18 (2022)
- 588 52. Smolka, M., Paulin, L.F., Grochowski, C.M., Mahmoud, M., Behera, S., Gandhi, M., Hong, K., Pehlivan, D.,
589 Scholz, S.W., Carvalho, C.M., *et al.*: Comprehensive structural variant detection: from mosaic to
590 population-level. *Biorxiv*, 2022–04 (2022)

591 Supplementary Information

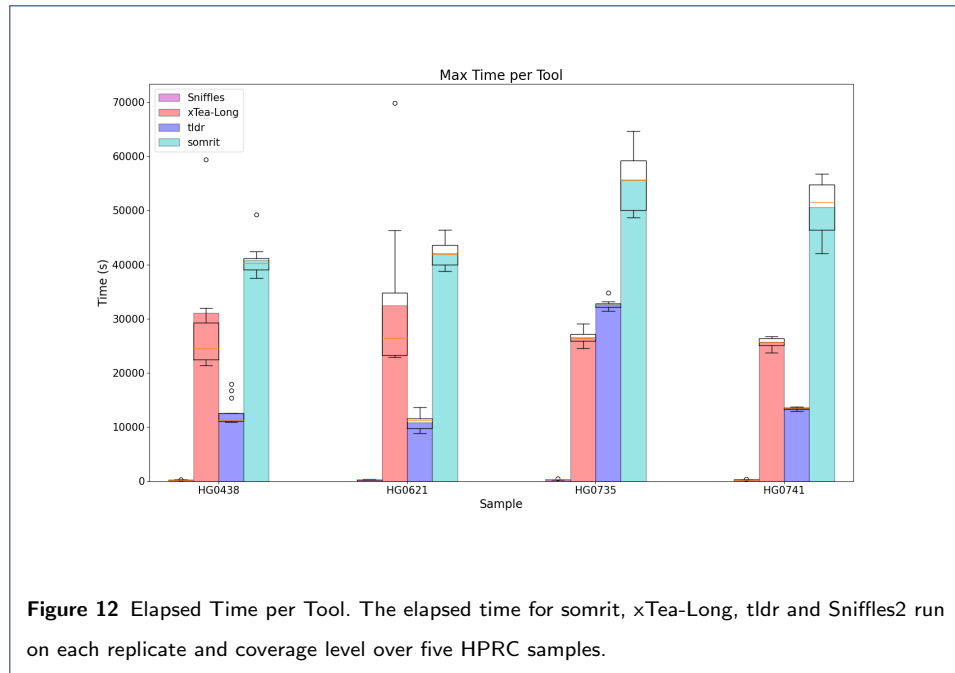
592 Time and Memory Analysis

593 We evaluated how somrit compared to other SV detection methods for computational performance. We noted the
594 total time and memory usage of somrit, xTea-Long, tldr and Sniffles2 runs during the previously mentioned analysis
595 of simulated somatic retrotransposon insertion events on the 4 HPRC samples. The analysis of simulated somatic
596 retrotransposon insertion events used a number of different machines as part of a larger shared computing
597 environment, thus tools were not run on the same machine. While not ideal this approach does allow us to get a
598 range of possible time and memory usages for a tool over different machines with equivalently sized input. For each
599 HPRC sample each tool was run the 12 40x replicates used for the simulated insertion analysis.

600 Figure 11 shows this comparison for memory usage and Figure 12 for time. As somrit consists of multiple individual
601 modules, we noted the time of each step and reported two versions of the time analysis. One version, shown in
602 Figure 12 as somrit total, represents the total time taken if each step is run sequentially with 10 threads. The
603 second version, shown in Figure 12 as somrit ideal, is the total time taken if individual re-alignment jobs for each
604 chromosome are run in parallel with 10 threads each. For somrit memory we took the maximum memory over all
605 modules for a given input fastq.



606 Somrit does have both higher run time overall and higher memory usage than other tools at 40x coverage. If somrit
607 realign is run in parallel per chromosome the total time required for somrit is comparable to that of tldr and



608 xTea-Long. The higher memory usage of somrit is attributed to the consensus generation step of realignment, with
609 abPOA requiring a high amount of memory to generate consensus sequences.