

1 The benefit of a complete reference genome for cancer structural 2 variant analysis

3 Luis F Paulin^{1,*}, Jeremy Fan^{2,*}, Kieran O'Neill², Erin Pleasance², Vanessa L. Porter^{2,3,4}, Steven
4 J.M Jones^{2,3,#}, Fritz J. Sedlazeck^{1,5,6,#}

5
6 1: Human Genome Sequencing Center Baylor College of Medicine, Houston, TX, USA

7 2: Canada's Michael Smith Genome Sciences Centre at BC Cancer, Vancouver, BC, Canada

8 3: Department of Medical Genetics, University of British Columbia, Vancouver, BC, Canada

9 4: Michael Smith Laboratories, University of British Columbia, Vancouver, BC, Canada

10 5: Department of Molecular and Human Genetics, Baylor College of Medicine, Houston, TX, USA

11 6: Department of Computer Science, Rice University, Houston, TX, USA

12

13 Equal contributions: *,#

14 Corresponding: sjones@bcgsc.ca, fritz.sedlaeck@bcm.edu

15 Abstract

16 The complexities of cancer genomes are becoming more easily interpreted due to
17 advancements in sequencing technologies and improved bioinformatic analysis. Structural
18 variants (SVs) represent an important subset of somatic events in tumors. While detection of
19 SVs has been markedly improved by the development of long-read sequencing, somatic variant
20 identification and annotation remains challenging.

21

22 We hypothesized that use of a completed human reference genome (CHM13-T2T) would
23 improve somatic SV calling. Our findings in a tumour/normal matched benchmark sample and
24 two patient samples show that the CHM13-T2T improves SV detection and prioritization
25 accuracy compared to GRCh38, with a notable reduction in false positive calls. We also
26 overcame the lack of annotation resources for CHM13-T2T by lifting over CHM13-T2T-aligned
27 reads to the GRCh38 genome, therefore combining both improved alignment and advanced
28 annotations.

29

30 In this process, we assessed the current SV benchmark set for COLO829/COLO829BL across
31 four replicates sequenced at different centers with different long-read technologies. We
32 discovered instability of this cell line across these replicates; 346 SVs (1.13%) were only
33 discoverable in a single replicate. We identify 49 somatic SVs, which appear to be stable as
34 they are consistently present across the four replicates. As such, we propose this consensus set

35 as an updated benchmark for somatic SV calling and include both GRCh38 and CHM13-T2T
36 coordinates in our benchmark. The benchmark is available at: [10.5281/zenodo.10819636](https://zenodo.org/record/10819636)

37
38 Our work demonstrates new approaches to optimize somatic SV prioritization in cancer with
39 potential improvements in other genetic diseases.

40

41 **Introduction**

42 Advancements in long-read sequencing technologies (LRS) have delivered unprecedented
43 insights into cancer genomics (Aganezov et al. 2020; Akagi et al. 2023; Fujimoto et al. 2021;
44 Choo et al. 2023; Thibodeau et al. 2020). Structural variants (SVs), defined as insertions,
45 deletions and rearrangements larger than 50 base pairs (bp), represent an important subset of
46 somatic driver events in tumours (Espejo Valle-Inclan et al. 2022; Aganezov et al. 2020).
47 Examples include the amplification of oncogenes (eg. *ERBB2*) (Nattestad et al. 2018),
48 generation of oncogenic fusion genes (e.g. *BCR-ABL* (Druker et al. 2001)), and silencing of
49 tumour suppressor genes (e.g. *TP53* (Hernández Borrero and El-Deiry 2021)). Accurate
50 detection of somatic SVs and characterization of their functional effects is thus critical for
51 appropriate cancer diagnosis and selection of therapy. Short-read whole genome sequencing is
52 routinely used for this purpose (Pleasant et al. 2020; Tsang et al. 2021), and can detect the
53 majority of SVs (Choo et al. 2023). However, short-read sequencing is often unable to fully
54 characterize the function of oncogenic SVs (Thibodeau et al. 2020; Pleasant et al. 2020).
55 Even worse, short-read sequencing reports a high proportion of falsely identified SVs, which can
56 be misleading (e.g. repeat expansions as translocations) (Thibodeau et al. 2020; Sedlazeck et
57 al. 2018; Mahmoud et al. 2024, 2019). Cell lines (eg. SKBR3) have been analyzed to assess the
58 complexity behind these oncogenic somatic changes. For example, Akagi et al. recently
59 identified virus-mediated progression in head and neck cancer samples through the generation
60 of unstable human papillomavirus (HPV) integrated molecular structures (Akagi et al. 2023).
61 Further research is needed to assess the stability of SVs in cancer genomes, which can help
62 direct priority when analyzing genomics data in heterogeneous cancer samples.

63

64 A key challenge in disease research is variant prioritization to identify potential causative
65 variants of a condition. This is even more amplified in cancer patients since, depending on the
66 type of cancer, the number of mutations are increased compared to normal tissue. Thus,
67 researchers leverage normal controls (often blood) from patients to identify somatic variants that

68 could potentially be driving the cancer (Mandelker and Ceyhan-Birsoy 2020). This process is
69 efficient but also often complicated by multiple factors such as tumor purity, availability of non-
70 cancerous tissue, gene annotation accuracy, and variant comparison/representation (English et
71 al. 2022; Salzberg 2019; Yoshihara et al. 2013). In particular, variant prioritization is impacted
72 by falsely identified variants in either tumor or normal samples. The cause of these falsely
73 identified variants can be attributed to misinterpretation of mapped reads, coverage fluctuations,
74 low quality mapping in tandem repeat regions, and unresolved regions in the reference genome.
75 To improve these, multiple advancements have been proposed. The first complete reference
76 genome CHM13-T2T recently became available, which showed improvements in variant calling
77 (Nurk et al. 2022; Aganezov et al. 2022), and even corrected mis-represented medically
78 relevant genes (Behera et al. 2023). The CHM13-T2T genome added approximately 200 mega
79 base pairs (Mbp) of resolved sequence to the GRCh38 reference build, thus closing reference
80 gaps and completing the centromere and telomere sequences (Nurk et al. 2022). These
81 repetitive regions can play a key role in cancer progression as they include hundreds of protein
82 coding genes that so far are not well understood (English et al. 2023). Perhaps more
83 importantly, they can lead to the onset of genome instability and thus the formation of variants
84 or inclusions of viruses (e.g. HPV) (Porter et al. 2023; Akagi et al. 2023). Annotating these
85 newly assembled regions is challenging, but such annotations are necessary to understand the
86 impact of mutations in these regions. While there are significant improvements in variant calling
87 when SVs are detected against the CHM13-T2T reference (Aganezov et al. 2022), the benefits
88 of using the complete genome reference have yet to be determined for somatic variants in a
89 cancer context. Additionally, since CHM13-T2T is still being annotated, the annotated reference
90 genomes GRCh37 and GRCh38 still hold value for ranking and characterizing SVs (Collins et
91 al. 2020; Tanner et al. 2024). Using both a complete and annotated reference genome would
92 allow researchers to identify novel oncogenic candidate variants across different cancer types.

93
94 Bioinformatics methodologies are continuously undergoing improvements to better utilize these
95 novel reference genomes and enable improved detection of novel alleles, genes, and the
96 variants and mutations that impact them (Smolka et al. 2024; Chen et al. 2024; Majidian et al.
97 2023). One such improvement was a novel method (Leviosam2) (Chen et al. 2024) that lifts
98 over the read alignments between two reference genomes, thus gaining benefits of both and
99 resulting in improved variant detection overall. This contrasts to previous methods that instead
100 lifted over variant calls, which often lead to false positive variant calls, especially with larger
101 rearrangements or duplications (Aganezov et al. 2020). We recently developed the SV caller

102 Sniffles2 to identify SVs using long-read sequencing data (Smolka et al. 2024). While Sniffles2
103 has been applied in germline contexts such as neurological and mendelian disorders, its utility
104 for cancer is still unproven, and somatic mutation detection often requires additional steps to
105 reduce the false discovery rates. Furthermore, widely available benchmark samples are
106 important instruments for benchmarking approaches (Zook et al. 2020; Olson et al. 2023;
107 Majidian et al. 2023; Wagner et al. 2022). Over the past years several groups have proposed
108 benchmarks for normal and cancer samples such as SKBR3 and COLO829 (Craig et al. 2016;
109 Espejo Valle-Inclan et al. 2022). Nevertheless, these require constant vetting as reference
110 genome changes and novel technologies change the downstream results that constitute the
111 established benchmark.

112
113 In this work, we investigated novel advancements in SV calling and reference genome mapping
114 to improve variant prioritization of somatic structural rearrangements. To accomplish this, we
115 analyzed four replicates of the tumour/normal benchmark samples COLO829/COLO829BL at
116 different laboratories to investigate the genome stability of these samples and also investigate
117 how mapping to a complete reference genome (CHM13-T2T) influences the established
118 benchmark. This highlighted a reduction of falsely identified SVs across all
119 COLO829/COLO829BL samples when using the CHM13-T2T reference. This observation was
120 also replicated in POG patient cancer samples where we observed a slight decrease in the
121 number of somatic SV. Nevertheless, the lack of annotations on CHM13-T2T complicates
122 variant prioritization. To overcome this, we propose a liftover approach of aligned reads that
123 combines the benefits from both reference genome versions to deliver less falsely identified
124 variants on GRCh38, further benefiting from the vastly available annotation resources. In
125 addition to this work, we further investigated the genome instability of COLO829/COLO829BL
126 and thus the risk of utilizing previous postulated benchmarks. We highlight multiple falsely
127 reported and missed somatic SVs for this cell line on GRCh38. Furthermore, we propose a
128 corrected, stable benchmark based on the four replicates of COLO829/COLO829BL sequenced
129 across GRCh38 and CHM13-T2T. This should improve future cancer studies that use the
130 CHM13-T2T reference genome for their analyses, which we found improves variant detection
131 and prioritization overall.

132

133 Results

134 *An updated COLO829 structural variation benchmark*

135 We investigated the previously established COLO829/COLO829BL benchmark (Espejo Valle-
136 Inclan et al. 2022), as recent reports highlighted discrepancies between the current benchmark
137 and other sequencing data (Smolka et al. 2024; Shiraishi et al. 2023). We utilized four
138 independent COLO829/COLO829BL samples from four sequencing centres profiled using the
139 ONT and PacBio long read platforms, including the reads from the most recently established
140 benchmark (Espejo Valle-Inclan et al. 2022) (**Supplementary Table 1**). **Figure 1A** shows the
141 steps followed to obtain somatic SV calls that were used to investigate the aforementioned
142 benchmark dataset (see Methods). Briefly, we aligned the four tumor/normal pairs to the
143 GRCh38 reference genome, then identified SVs using Sniffles2's population call/merge strategy
144 (Smolka et al. 2024). We identified 19,684 SVs shared among all samples (**Figure 1B**), of which
145 45 were unique to the COLO829 cancer cell lines. Notably, the COLO829/COLO829BL samples
146 by Valle-Inclan (labeled VAI) had the largest numbers of unique SVs ($n_{\text{tumor}}=200$, $n_{\text{normal}}=193$)
147 while the tumor-normal pair sequenced with the PacBio Revio platform (labeled PBR) has the
148 lowest number ($n_{\text{tumor}}=38$, $n_{\text{normal}}=65$), which could be attributed to either differences in
149 sequencing technology (Mahmoud et al. 2024) (older version of ONT MinION vs. PacBio Revio)
150 or cell-line divergence. Interestingly, we observed that a small number of SVs were shared only
151 by the tumor/normal pairs at individual sequencing centers (i.e detected in both COLO829 and
152 COLO829BL but only in one of the four datasets). The lowest number of
153 COLO829/COLO829BL shared SVs by sequencing center was 15 from Canada's Michael Smith
154 Genome Sciences Centre (labeled as GSC) and the highest was 449 from the Valle-Inclan
155 (ONT 271, PBR 115). Finally, there were 8,659 SVs that were shared between two to seven
156 samples (tumor or normal) from different platforms and sequencing centers, which were omitted
157 from subsequent analysis (**Supplementary Table 2**). We filtered for SVs that were present in
158 the four cancer samples with a variant allele frequency (VAF) $\geq 10\%$ and not present in any of
159 the normal samples. Subsequently, the remaining somatic 44 SVs were manually reviewed in
160 IGV. We denoted this dataset as COLO829-GRCh38 (**Supplementary Table 3**).

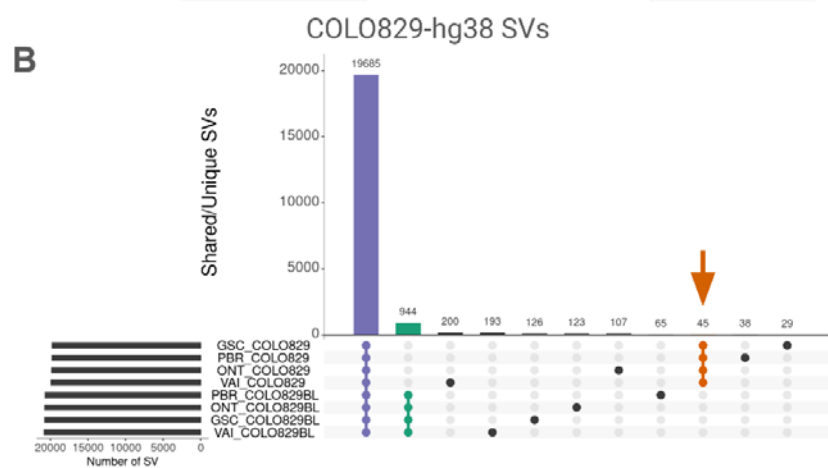
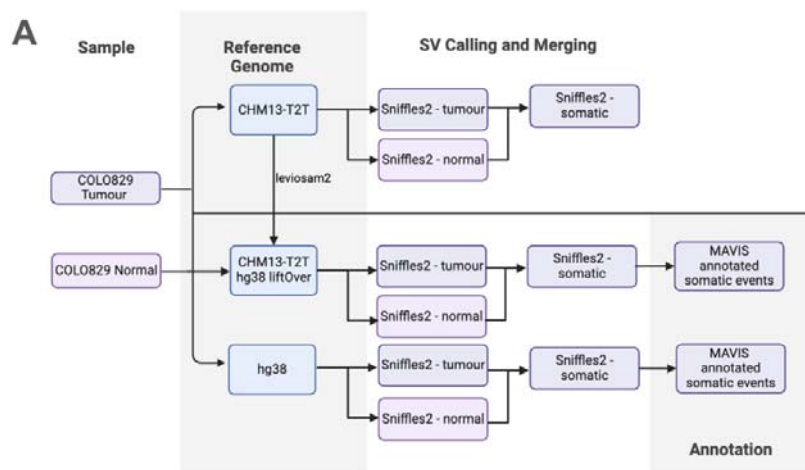
161
162 We compared these 44 SVs called in data from all four centers to the previously established SV
163 benchmark for COLO829 (62 SVs, see Methods) by Valle-Inclan et al. (Espejo Valle-Inclan et al.
164 2022). **Table 1** shows the benchmark summary. In comparison to established benchmark SVs,

165 we categorized 38/44 SVs as true positives (TP), six as false positives (FP), and 24 as false
166 negatives (FN) (**Supplementary Table 4**). We proceeded with an in-depth analysis of the FP
167 and FN calls to understand how the previous benchmark compared to the full set of data from
168 four centres. A manual inspection of FP events revealed that 5/6 were falsely classified as FPs
169 as they showed clear evidence in samples from multiple centres (**Supplementary Table 5**,
170 **Supplementary Figure 1**). For example, the variant INS.C2M6, a heterozygous 69 bases
171 insertion, which involved *PMS2*, was clearly detected in all four cancer samples
172 (**Supplementary Figure 1B**) but was not reported in the previous benchmark (Espejo Valle-
173 Inclan et al. 2022). For the 24 SVs classified as FN, we performed a genotyping experiment
174 (Chander et al. 2019) (also known as force-calling, see Methods, **Supplementary Table 6**) in
175 order to assess if there was any signal in the data to support such SVs. Here, five SVs that were
176 initially missed by our analysis had strong read evidence in all four COLO829 replicates (**Figure**
177 **1D**, **Supplementary Figure 2**), one FN SV only was detected in a control sample with one read
178 thus was discarded by our filter that removed SVs with any evidence in the normal samples
179 (**Supplementary Figure 3**). In addition, ten FN SVs had variant allele frequencies (VAF) below
180 our threshold of 10% and therefore were not included (VAF 1.4-6.2%, **Supplementary Figure**
181 **4**), additionally, four FN SVs could only be detected in 1-3 COLO829 replicates, and thus were
182 filtered out (**Supplementary Figure 5**), and finally, four SVs had no read support for the SVs
183 across all four cancer replicates (**Figure 1C**, **Supplementary Figure 5C**). Interestingly, we
184 missed a well defined deletion present in all four COLO829 replicates (benchmark SV
185 truthset_56_1). Further investigations revealed that a filter applied during SV calling prevented
186 the detection of that particular SV (**Supplementary Figure 2C**). This filter is applied to remove
187 potential erroneous SVs based on the minimum coverage (COV_MIN) needed to detect an SV.
188 In total, 43 somatic SVs detected in SV calling in COLO829-GRCh38, excluding one false
189 positive, and 5 somatic SVs called from force-called genotyping and manual review had strong
190 read evidence in datasets from all four centres.
191

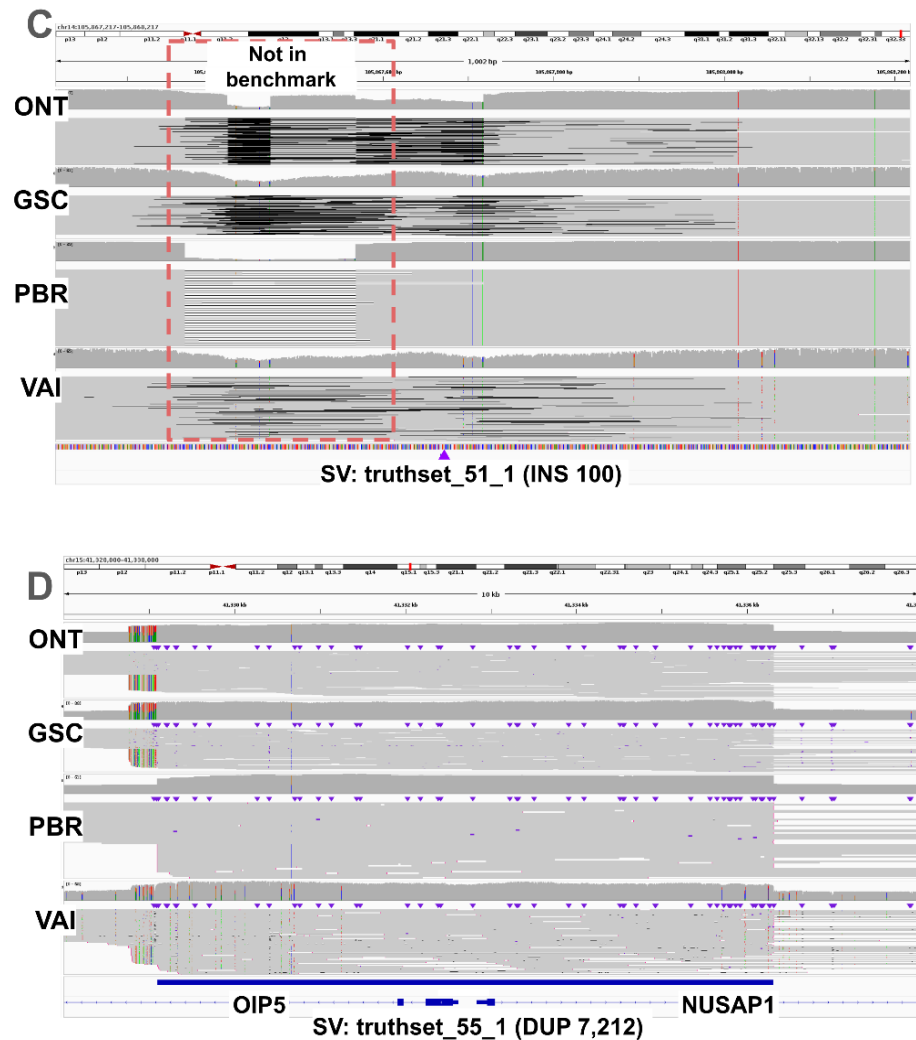
Datset	Version	P	TP	FP	FN	Precision	Recall
COLO829-38	Initial	62	38	6	24	86.36%	61.29%
	Genotyped/IGV	49	44	1	5	97.78%	89.80%
COLO829-T2T	Initial	62	38	5	n/a	88.37%	n/a
	Genotyped/IGV	n/a	43	0	n/a	100.00%	n/a

COLO829-lifted	Initial	62	36	5	26	87.80%	58.06%
	Genotyped/IGV	49	45	0	4	100.00%	91.84%

192 **Table 1.** Benchmark results for the three COLO829 datasets based on the reference genome. For each dataset we
 193 present the initial evaluation (labeled Initial) and the in-depth analysis/evaluation that includes genotyping, and
 194 manual inspection in IGV.
 195



196



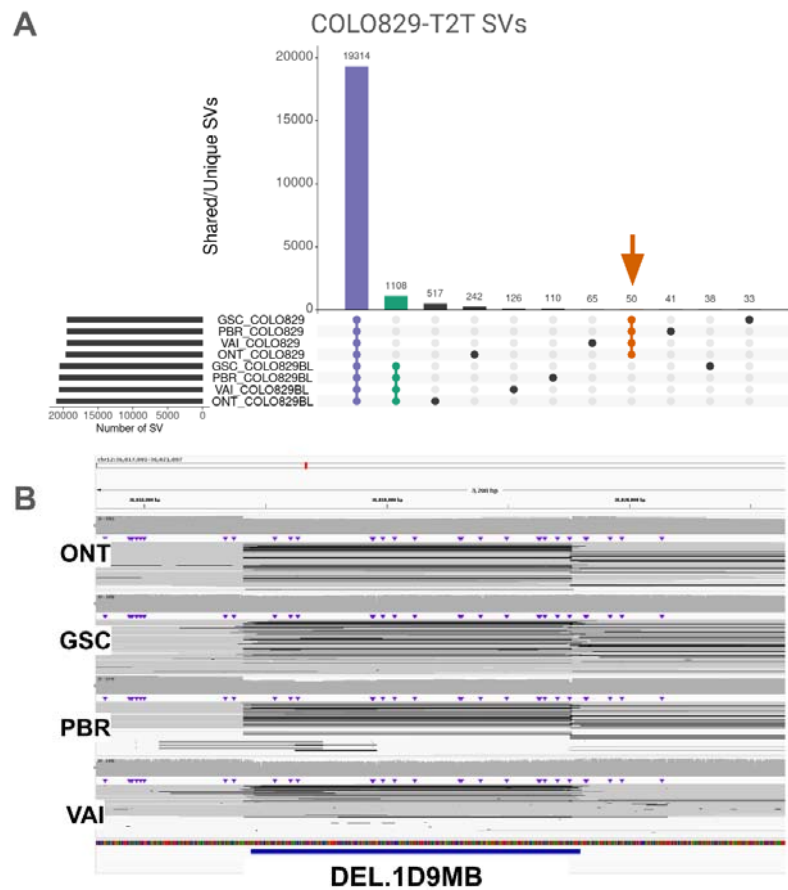
197
 198 **Figure 1. A)** Schematic overview of our variant calling and comparison methods using CHM13-T2T vs
 199 GRCh38. **B)** Upset plot of the shared and unique SVs of the COLO829-GRCh38 dataset across four
 200 COLO829 tumor/normal pairs. Violet events shared by all samples; orange, somatic events shared by
 201 tumour samples only; green, germline events shared by blood normal samples only. **C)** IGV screenshot of
 202 an SV cataloged as missing (truthset_51_1, large purple triangle) according to the benchmark by Valle-
 203 Inclan et al. We did not detect any reads supporting it. Moreover, we detected a DEL in all four samples
 204 which is missing from the benchmark. **D)** IGV screenshot of a DUP that was assigned two distinct SV
 205 types (DUP and INS), and thus cataloged as missing according to the benchmark by Valle-Inclan et al.
 206 Manual inspection in IGV showed the SV.

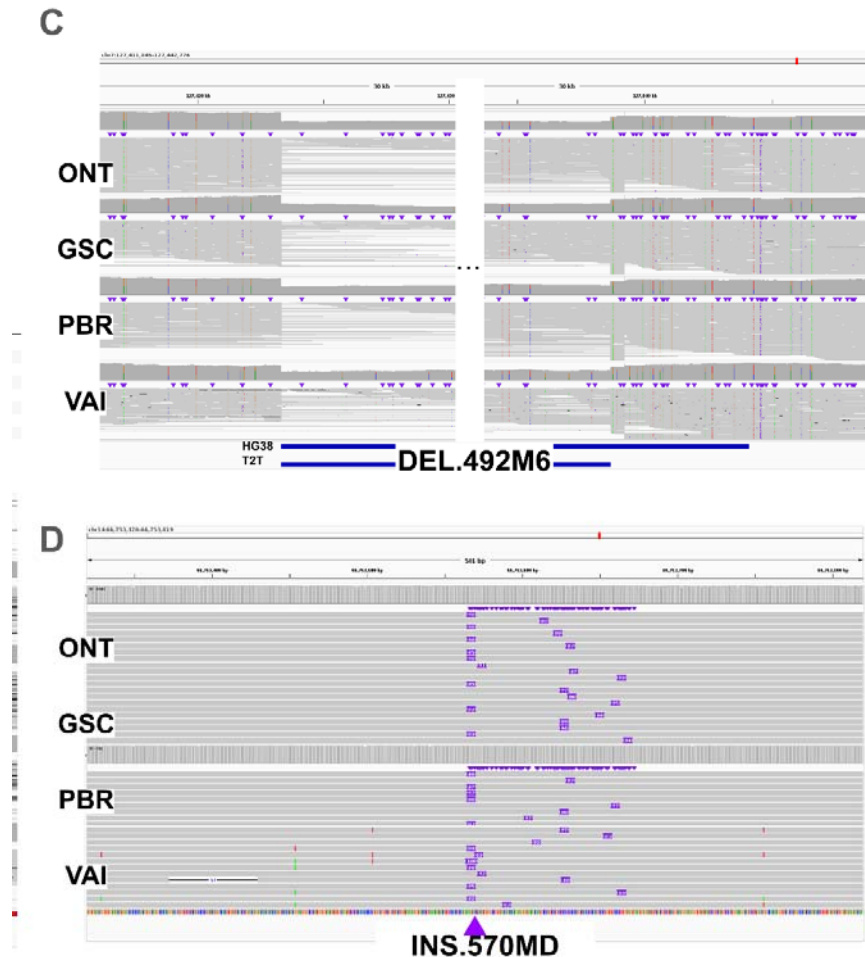
207 ***Leveraging a complete genome reference for cancer structural variation*** 208 ***benchmarking***

209 The recent introduction of CHM13-T2T showed improved variant calling (Aganezov et al. 2022;
210 Nurk et al. 2022) and corrections of a few medically important genes (Behera et al. 2023). Thus,
211 we examined the impact of using CHM13-T2T for structural variant calling in cancer research.
212 Using the CHM13-T2T, we observed a high number of SVs in the telomeric and centromeric
213 regions of the chromosomes, the detection of which was aided by a complete reference
214 (**Supplementary Figure 6**). **Supplementary Table 8** shows the comparison of SVs detected in
215 centromeric and telomeric regions between using the CHM13-T2T reference (labeled as
216 COLO829-T2T) and COLO829-GRCh38. From COLO829-T2T, 18 chromosome have higher
217 number of SVs in centromeric and telomeric regions (average 16.55% increase), while six
218 chromosomes have higher number of SVs in centromeric and telomeric regions in COLO829-
219 GRCh38 (average 1.76% increase). Next, when somatic SVs are identified in the COLO829-
220 T2T dataset (**Supplementary Table 7**), using the same methodology as described for the
221 benchmark SV set, we do not observe an enrichment in such regions (red dots in
222 **Supplementary Figure 6**), suggesting these events are mostly germline variation. We
223 observed a deviation in the ratio of insertions (INS) and deletions (DEL) with an INS:DEL ratio of
224 1:1.4 in CHM13-T2T and 1:0.87 in GRCh38. The switched INS:DEL ratio has been described
225 before and is thought to be caused by too small tandem repeats reported in GRCh38 (Aganezov
226 et al. 2022). Furthermore, we observed a substantial decrease in the number of
227 interchromosomal events (BND) from 225 in GRCh38 to 83 in CHM13-T2T for COLO829. As
228 interchromosomal events are not expected to be found in normal samples, this reduction largely
229 reflects a decrease in noise in variant calling.

230
231 To enable a comparison between the GRCh38 and CHM13-T2T based SV calls, we linked SV
232 calls of the same SV type and chromosome using the read names supporting each SV. Since
233 there is no COLO829 SV benchmark with CHM13-T2T coordinates, we used the read names to
234 determine which SVs correspond to the benchmark based on the previous genotyping. **Figure**
235 **2A** shows the number of SVs shared across the sample replicates. The total number of
236 germline SVs was slightly smaller in CHM13 compared to GRCh38 (19,314 and 19,685
237 respectively), while the number of events considered somatic and shared between tumour
238 samples only was slightly larger (50 and 45 SV, respectively). Unlike the GRCh38 calls, here
239 the samples from the Oxford Nanopore Open Data project (labeled as ONT) had the largest

240 numbers of unique SVs (517, 242 for COLO829/COLO829BL) while the COLO829/COLO829BL
241 sequenced by the GSC has the lowest number (33, 38 for COLO829/COLO829BL). Finally,
242 there were 17,395 SVs that were shared between two to seven samples (tumor or normal) from
243 different platforms and sequencing centers (**Supplementary Table 9**). These SVs were not
244 included in the analysis. We manually reviewed the 50 SVs detected in tumour samples only
245 from all four centres and identified 43 somatic SVs. Out of which 38 were initially classified as
246 TP and five as FP relative to the comparison of the COLO829-GRCh38 somatic SV benchmark
247 by Valle-Inclan et al (see Methods). We identified two interesting cases where the somatic SV is
248 present in both GRCh38 and CHM13, but is called as being a different size due to differences in
249 the references. SV DEL.492M6 differed in size by ~6kb (**Figure 2C**) to the corresponding SV in
250 GRCh38 SV truthset_28_1 from the benchmark by Valle-Inclan et al), overlaps with the
251 glutamate metabotropic receptor *GRM8* and SV DEL.508M6 (corresponding to truthset_30_1 in
252 the benchmark) differed in size by 126kb which in the case of CHM13-T2T SV is not hitting any
253 genetic element but the GRCh38 does overlap with an Olfactory Receptor *OR2A1* and its
254 antisense RNA *OR2A1-AS1*. This difference is a consequence of the size difference between
255 both SVs and not given by differences in the annotation. Next, from the five somatic SVs initially
256 classified as FP according to Valle-Inclan et al , SV DEL.1D9MB is uniquely detectable using
257 the CHM13-T2T reference genome as it is positioned in the centromere of chromosome 12
258 (**Figure 2B**). Manual inspection of the remaining four FP showed that they actually represent
259 TPs, as we detected supporting reads for all four sequenced COLO829 samples
260 (**Supplementary Table 10, Figure 2D, Supplementary Figure 7**). For example, **Figure 2D**
261 shows INS.570MD, a 97 bp insertion that can be clearly observed in all four cancer samples.
262 Next, we used the UCSC genome browser (Kent et al. 2002) to annotate the somatic SVs for
263 COLO829-T2T. Here, we observed similar results to the annotation of somatic SVs in GRCh38
264 coordinates (20/31, 64.5% shared genes), with the addition of nine characterized genetic
265 elements and the centromere (**Supplementary Table 7B**).

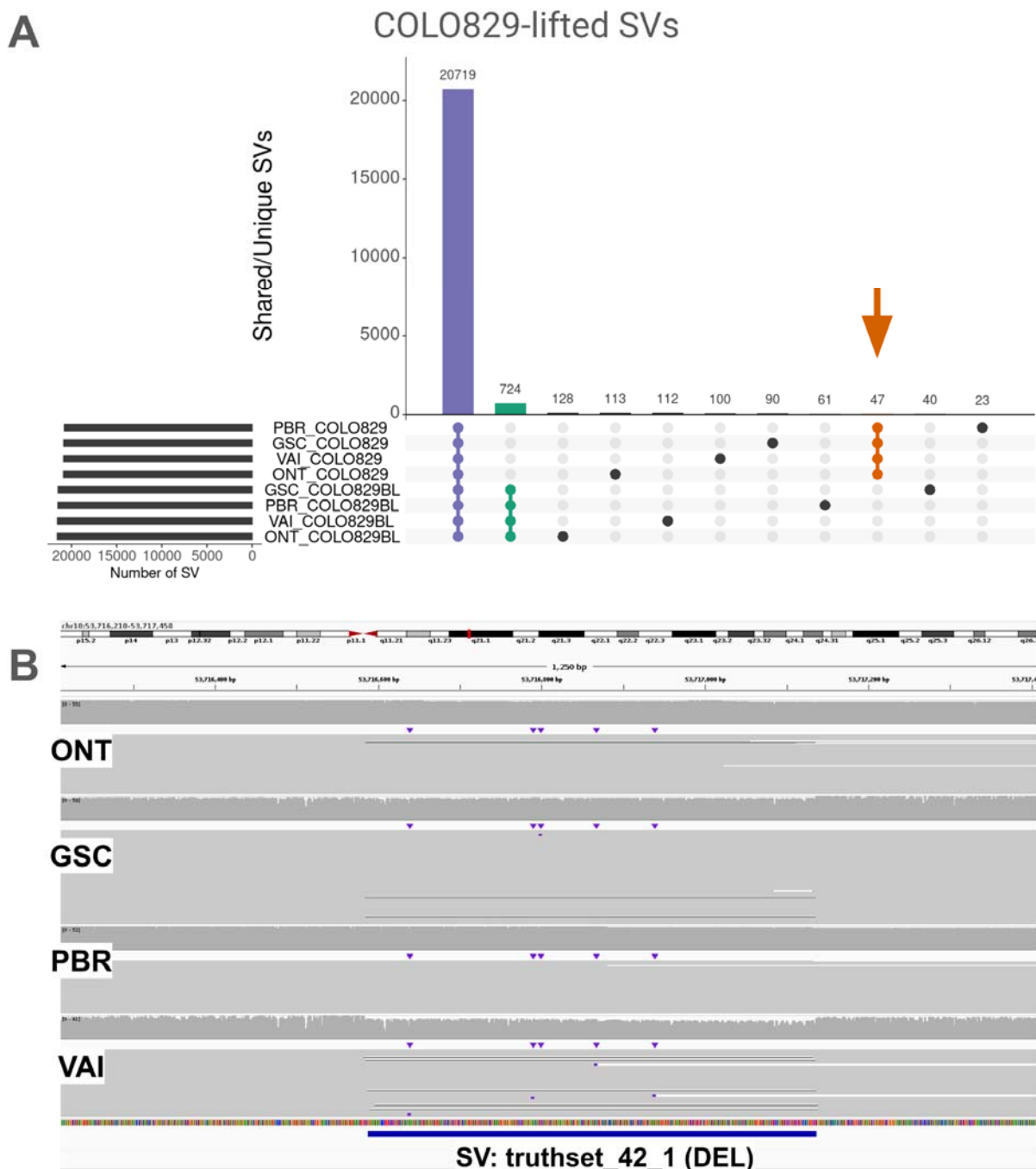




267
268 **Figure 2. A)** Upset plot of the shared and unique SVs of the four COLO829/COLO829BL tumor/normal
269 pairs aligned to the CHM13-T2T reference genome. **B)** IGV screenshot of an SV that was reported with
270 different sizes in CHM13-T2T and GRCh38. **B)** IGV screenshot of an SV detected in the centromere of
271 chromosome 12 on CHM13-T2T. **D)** IGV screenshot of an SV cataloged as FP according to the
272 benchmark by Valle-Inclan et al, which we are able to detect in all four cancer samples.

273
274 Although SV detection is improved (i.e. no FP SVs detected), using CHM13-T2T limits variant
275 prioritization and interpretation as population frequencies are not available on this reference,
276 and GRCh38 has more informative annotation databases (Collins et al. 2020; Nicholas et al.
277 2022; Chowdhury et al. 2022). Thus, we used recent advances in liftover of alignments to take
278 advantage of both the improved mapping using a CHM13-T2T reference (Chen et al. 2024) and
279 the years of annotation and curation of the GRCh38 reference genome. Briefly, we took the
280 CHM13-T2T read alignments and used LevioSAM2 (Chen et al. 2024) (v0.4.1) to liftover the
281 alignments from CHM13-T2T to GRCh38 coordinates. We denoted this dataset as “COLO829-
282 lifted” (**Supplementary Table 11**). **Figure 3A** shows an upset plot with the number of SVs

283 shared among all samples, those unique to the cancer cell line COLO829 (marked with an
284 arrow), the blood control COLO829BL and the unique SV per sample. We compared the
285 COLO829-lifted calls to the benchmark by Valle-Inclan et al, alongside the dataset COLO829-
286 GRCh38 that was produced earlier. **Table 1** shows how leveraging the CHM13-T2T genome
287 improves SV calling and decreases both the number of SV classified as FP and FN when
288 compared to the native GRCh38 alignment, after genotyping and manual review
289 (**Supplementary Table 12**). For the case of the SVs labeled as FP in comparison to the
290 benchmark by Valle-Inclan et al, we were able to identify all the SVs as TP with high confidence
291 in all the cancer samples (**Supplementary Table 13, Supplementary figure 8**). One interesting
292 example is DEL.6A4M6, which is a deletion that was also called in COLO829-GRCh38 but with
293 a greater length (46 bases in COLO829-lifted vs. 138 bases in COLO829-GRCh38); however,
294 manual inspection in IGV shows the 46 base DEL only demonstrating the more accurate result
295 obtained when incorporating T2T alignment. Next, when analysing the FN we identified two SVs
296 that were also missed in the COLO829-GRCh38 (benchmark SVs: truthset_12_1 and
297 truthset_56_1, **Supplementary Figure 9**). In both cases we applied filters during SV calling that
298 removed them from the final call set (COV_CHANGE for the case of SV:truthset_12_1 and
299 COV_MIN for SV:truthset_56_1). Then, two SVs labeled as FN were detected as INS in our
300 callset and as DUP in the benchmark. These two types are often hard to distinguish as a
301 duplication is an insertion of the same sequence next to itself (Mahmoud et al. 2019). The rest
302 of the FN are very similar the COLO829-GRCh38 results (**Supplementary Figure 9-12**) with
303 the difference of truthset_23_1 being missed in the GRCh38 dataset but not in the liftover and
304 truthset_38_1 being classified as low-allele frequency (VAF < 10%) in the GRCh38 and present
305 in the liftover.
306



307
 308 **Figure 3. A)** Upset plot of the shared and unique SVs of the four COLO829 tumor/normal pairs aligned to
 309 the CHM13-T2T reference genome and lifted over to GRCh38 (**Supplementary Table 14**). **B)** IGV
 310 screenshot of an SV that was labeled as FN according to the benchmark by Valle-Inclan et al. Our
 311 analysis calls removed these SVs due to lack of evidence in all samples (**Supplementary Table 13**)
 312
 313 Overall, the approach of incorporating alignment to CHM13-T2T followed by liftover analysis
 314 eliminated all FP somatic SV calls for COLO829. Such a reduction in FP has important impacts

315 particularly when analysing patient samples, where a reduction in FP can have implications for
316 interpretation and necessary follow-up analysis. The SV liftover approach allows for combining
317 the improved accuracy gained by incorporating CHM13-T2T alignments (precision in COLO829-
318 GRCh38 97.78% vs COLO829-T2T/COLO829-lifted 100.00%) with the biological and clinical
319 annotations available on GRCh38. Additionally, we observed a drastic reduction in BND (225
320 and 130 respectively for COLO829-GRCh38 and COLO829-lifted, 83 in COLO829-T2T), both in
321 tumor and normal samples, indicating improvements in germline SV calling which has
322 implications for variant prioritization in disease research. We used MAVIS (Reisle et al. 2019) to
323 annotate our proposed somatic SVs for the COLO829 cell-line in GRCh38 space. From the 49
324 somatic SVs detected in COLO829-lifted (45 TP and 4 FN), 35 SVs overlap with 26 known
325 genes, of which 25 are related to cancer. Some examples are *FHIT* (near a fragile site), tumor
326 suppressors *ITIH5*, *MAGI2*, *PTEN*, *WWOX*, and cell-proliferation control gene *TMX3* (Bellon et
327 al. 2021; Cao et al. 2020; Yehia et al. 2023; Husanie et al. 2022; Zhang et al. 2019).

328

329 Since these 49 somatic SVs could be detected across all 4 replicates of COLO829, we suggest
330 using this call set as a refined benchmark for this important tumor/normal control cell line
331 (**Supplementary Table 15**). The use of data from multiple centres and independent cell line
332 passages is particularly important as we demonstrate differences between samples that are
333 likely to in part represent instability and evolution of this cell line at the SV level in addition to the
334 previously described instability at the SNV level (Craig et al. 2016). Thus, accurate benchmarks
335 that incorporate multiple replicates are necessary to reduce possible sources of error introduced
336 by a single sample.

337 ***Utilizing a complete genome reference for cancer SV detection***

338 Our results show the benefit of using the CHM13-T2T reference for SV analysis in cancer.
339 Furthermore, we introduced and assessed a liftover approach that leverages the benefits of
340 both reference genomes to improve the detection and annotation of somatic SVs in
341 COLO829/COLO829BL. Additionally, variant prioritization can be challenging as most SVs are
342 not cancer drivers, but instead passenger mutations picked up in the process of cancer
343 evolution. Cancer samples have diversity in tumour content (fraction of sequenced DNA derived
344 from tumour cells compared to normal cells) which can result in somatic SVs with low VAF.
345 Because of the overall complexity of SV analysis in cancer, a reduction in false positive calls is
346 greatly beneficial. Given these results, we next investigated the use of the CHM13-T2T
347 reference genome for the analysis of two different cancer patient samples. In this pilot

348 experiment, we used two samples with high whole genome coverage (52x and 54x) and
349 matching normal controls (blood and skin) (**Supplementary Table 1**) (O'Neill et al. 2024). When
350 using CHM13-T2T alignment followed by GRCh38-liftover analysis, we observed fewer somatic
351 SVs when compared to GRCh38 (**Supplementary Table 17A-D**).

352

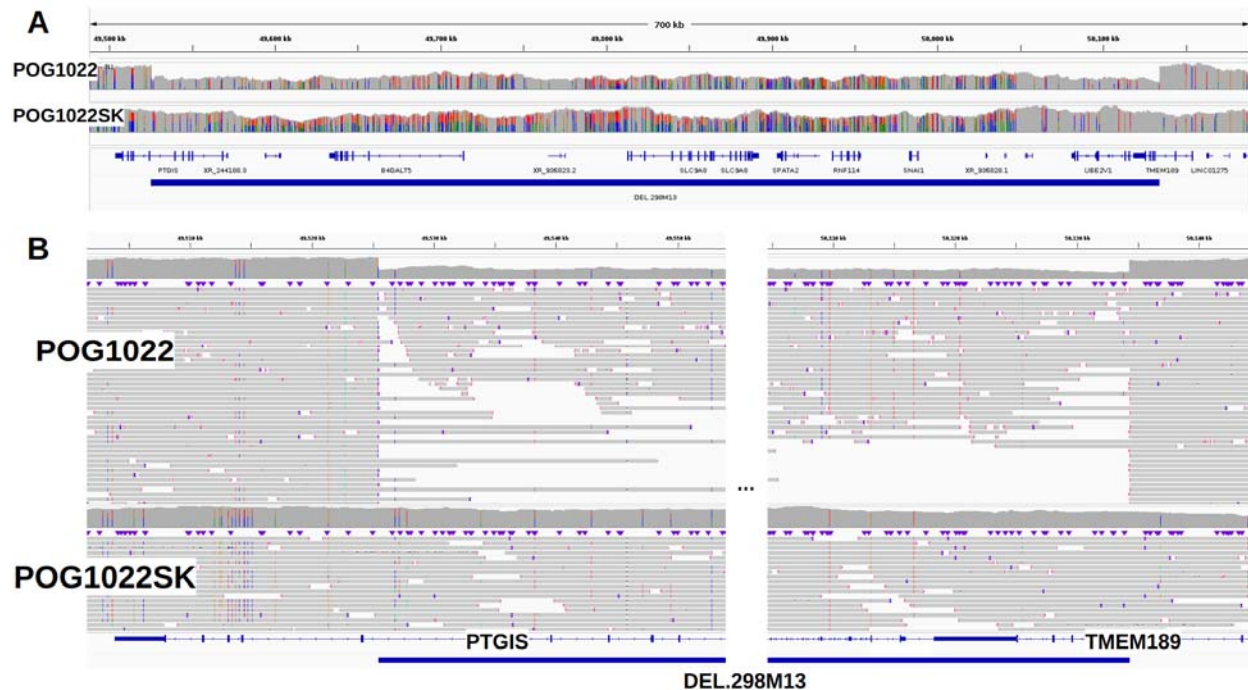
353 For POG044, we observed 193 SVs in the liftover analysis compared to 201 using the GRCh38
354 reference (**Supplementary Table 17A and 17B**), with a similar INS to DEL ratio (1.22:1 in the
355 liftover, 1.23:1 in GRCh38). We observed a reduction in the number of INS, DEL and BND (96,
356 79, 13, to 100, 81, 16, respectively) and no difference in DUP or INV. For sample POG1022, we
357 detected an increase in the number of DEL in both the liftover and the GRCh38 analysis
358 (**Supplementary Table 17C and 17D**). Also, the INS to DEL ratio was skewed towards DEL
359 (0.81:1 and 0.77:1 for liftover and the GRCh38 respectively). In a more detailed inspection we
360 detected a high number of DEL in pericentromeric regions of chromosomes 4, 10 and 13 (eight,
361 36 and 14, respectively). Also, we detected no change in the number of BND, and a decrease in
362 the rest of the SV types.

363

364 We annotated SVs from the liftover analysis (in GRCh38) using MAVIS, but overall did not
365 observe any SV that would be predicted to be causative in tumour formation. MAVIS predicted a
366 non-synonymous SV on chr1 in the sample POG044 (**Supplementary Table 18**). This SV
367 affects the neuroblastoma breakpoint family gene *NBPF20* which has been associated with
368 several types of cancer. This region was affected by several deletion events, moreover most of
369 them show low mapping quality and thus were not used during SV calling, specially in the
370 control sample (**Supplementary Figure 13**).

371 For sample POG1022, three non-synonymous coding events were categorized by MAVIS
372 (**Supplementary Table 19**). The first one, a 24.2Kb inversion in chromosome 2 which affects
373 the immunoglobulin kappa variable gene *IGKV3-11* (**Supplementary Figure 14**) which
374 expression has been observed in myelomas (<https://www.proteinatlas.org/ENSG00000241351-IGKV3-11/pathology>). The second, a 8.9Kb deletion also in chromosome 2 which affects the
375 ankyrin repeat domain gn *ANKRD36* (**Supplementary Figure 15**) which has been used as a
376 biomarker of disease progression in Leukemia (Iqbal et al. 2021). Upon further inspection, we
377 detected a larger deletion in the same region (12kb) with some reads supporting the SV in the
378 control. The third non-synonymous coding event was a 609kb deletion in chromosome 20 which
379 was interestingly only detectable in the lift over analysis and not in the GRCh38 alignment.
380 Visual inspection of the region in both the liftover analysis and GRCh38 alignment shows that
381

382 the latter has a weaker signal, although it is present. This SV affects the prostaglandin synthase
383 *PTGIS* (**Figure 4**) which has been linked with various cancers like prostate cancer (Qiao et al.
384 2023), colorectal cancer (Ding et al. 2023) and cancer-free trichothiodystrophy (Lombardi et al.
385 2021).



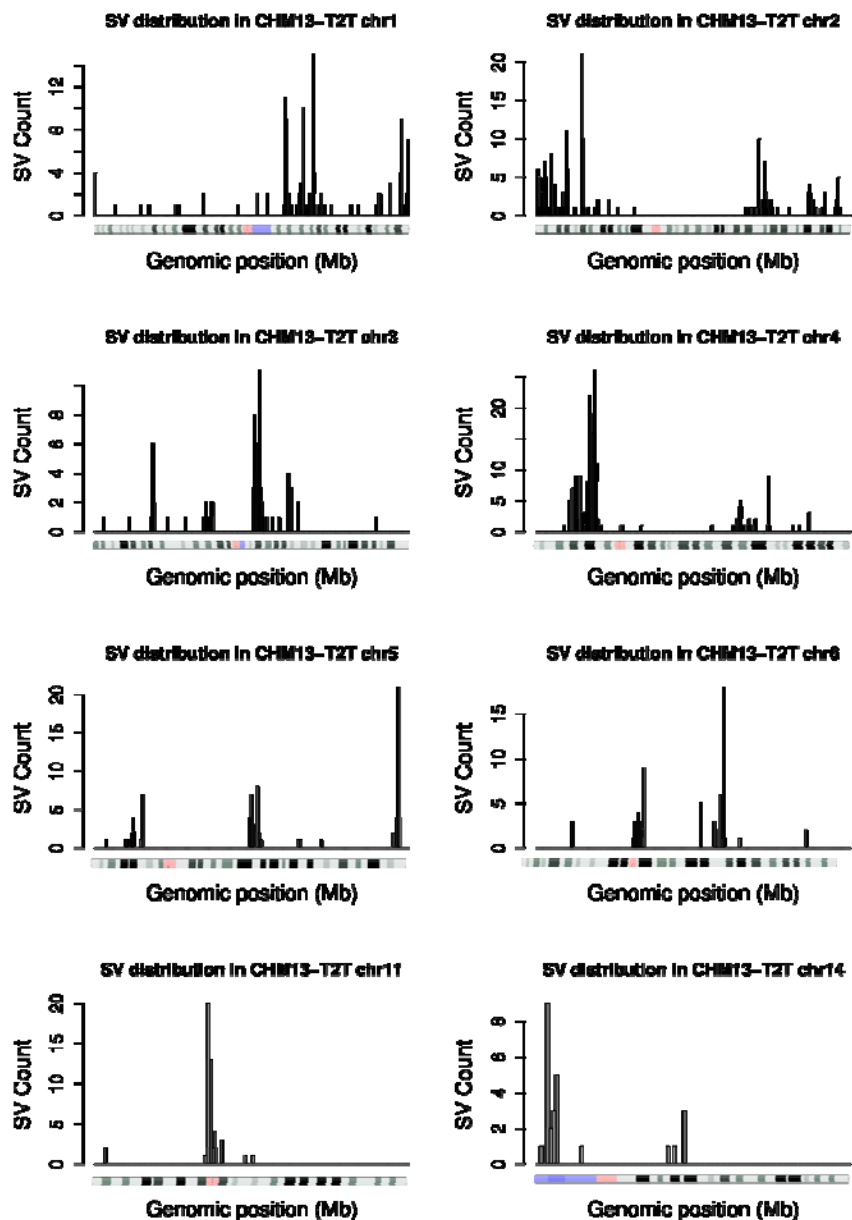
386
387 **Figure 4. A)** Coverage pattern for a 609kb somatic deletion (DEL.298M13) detected in POG1022 that
388 affects *PTGIS*. (**Supplementary Table 18**). **B)** IGV screenshot that zooms into the breakpoints of the
389 somatic deletion DEL.298M13.

390
391 Closer inspection of the 140 somatic SVs from POG044 that were annotated by MAVIS (44
392 were filtered or collapsed) showed that 96 (68.6%) were not somatic, meaning we could detect
393 reads supporting the SV in the control samples. From those 30 were caused by SV of the same
394 type that were collapsed in the tumor sample but not in the control causing a difference in size
395 and thus two different SV (**Supplementary Figure 16A**), 23 were in low complexity regions
396 where we could observe multiple SVs of the same type (**Supplementary Figure 16B**) and 18
397 have read support in the control with mapping quality lower than out threshold ($MQ \geq 20$) and
398 thus were not used during SV calling (**Supplementary Figure 16C**). For sample POG1022,
399 MAVIS annotated 184 somatic SVs (100 were filtered or collapsed) from which 129 (70.1%) had
400 reads supporting the SV in the control sample. From those, 64 occurred in low complexity

401 regions (**Supplementary Figure 17A**) and 46 were caused by collapsed SV of the same type
402 (difference in size **Supplementary Figure 17B**).

403 Next, we investigated the SVs detected when using the CHM13-T2T. While we found CHM13-
404 T2T somatic events in both POG044 (n = 407) and POG1022 (n = 433) (**Supplementary Table**
405 **20**), upon manual review, we could match 55 to GRCh38 in POG044 and 83 in POG1022 based
406 on the read names. Moreover, even when we observed an increase in the number of INS and
407 DEL using CHM13-T2T, in both cases we observed a reduction in the number of BNDs, while
408 the rest of the SV types were comparable. Similar to our previous analysis with the COLO829
409 cell-line, we observed the INS:DEL ratio skewed towards deletions.

410
411 A closer inspection of the somatic events in these two samples from the POG cohort, we
412 observed that most were not somatic, but lay in low complexity regions (i.e centromere/telomere
413 but not exclusively). We detected reads with support of the SVs in the control samples
414 (POG044BL and POG1022SK respectively) but with mapping quality below or established
415 threshold, such that the combined with the mis-alignments, they caused INS, and mostly DEL
416 SVs that were labeled as somatic. Thus, new advances in algorithms to overcome these
417 previously inaccessible regions are needed. Furthermore, when removing the SV that were not
418 in the aforementioned repetitive regions (centromere, telomere), we detected a lower number of
419 SV than using the GRCh38 (**Figure 5, Supplement figure 18-19**)



420

421

422 **Figure 5.** Somatic SV in cancer sample POG1022 along the CHM13-T2T genome. Shown are
423 chromosomes 1-6, 11 and 14 where we detected a high number of somatic SVs, much of which were
424 located in repetitive, low complexity regions.

425

426 Discussion

427 In this work we assessed the benefits of using a complete reference genome for SV calling and
428 prioritization in tumors. For this we examined four replicates of COLO829/COLO829BL together

429 with two tumor samples aligned to both CHM13-T2T and GRCh38. We saw lower false-positive
430 somatic SV calls when using CHM13-T2T as the reference. To enable taking advantage of
431 CHM13-T2T mapping, while also benefiting from the richer clinical annotations available for
432 GRCh38, we demonstrated the efficacy of lifting over alignments using LevioSAM (Chen et al.
433 2024). Furthermore, we developed a new approach to trace variants across reference genome
434 versions using read names instead of coordinates, which is less sensitive to different allelic
435 representations. Lastly, we identified inconsistencies among replicates of the existing
436 COLO829/COLO829BL benchmark set, and propose new consensus somatic SV benchmarks
437 for both GRCh38 and CHM13-T2T that addresses these inconsistencies. Our suggested
438 approach should lead to more accurate and less laborious cancer SV analysis, while the
439 improved benchmark provides a more accurate testbed for assessment of cancer SV callers.

440
441 While the advent of long reads holds promise to improve the detection of complex alleles
442 together with resolution in tandem repeats, it also highlights the problem around variant
443 prioritizations or annotation. In contrast to short reads, we often see many novel variants (eg.
444 SV) that aren't part of public databases such as gnomadSV (Collins et al. 2021), making
445 prioritization and annotation more difficult. To overcome this, analysis often relies on tumor-
446 normal comparison, which streamlines the detection of causative variants (Mandelker and
447 Ceyhan-Birsoy 2020). Short reads tumour-normal analysis often indicates large (>10kbp)
448 causative SV (Choo et al. 2023) but also has higher false positive and false negative rates
449 (Mahmoud et al. 2019; Aganezov et al. 2022), which can hinder rapid SV prioritization during
450 clinical workup of tumors. An example is the reporting of repeat expansions as
451 translocation/BND events (Sedlazeck et al. 2018). We find that using long reads on different
452 reference genomes appears to reduce the false positive calls and can thus improve variant
453 prioritization. In this paper, we observed an important decrease in the number of BND when
454 utilizing CHM13-T2T compared to the GRCh38 genome. Sniffles2 falsely called 130 SVs for
455 GRCh38, 84 for CHM13-T2T and 83 for CHM13-T2T lifted over to GRCh38. False positive SV
456 calls are often caused by misassembled and incomplete reference regions (Mahmoud et al.
457 2019)]; the completeness of CHM13-T2T reduces these false positives. However, clinical
458 annotations of SVs remain incomplete for CHM13-T2T. We therefore suggest a liftover
459 mapping approach that combines the increased accuracy of CHM13-T2T together with the utility
460 of annotations across GRCh38 itself. Furthermore, this liftover approach does not double the
461 analytical time nor the costs for analysis as it utilizes a chain file to rapidly liftover the read
462 alignments without additional pairwise alignments needed (Chen et al. 2024). This also permits

463 downstream utilization of annotation databases and approaches such as gnomadSV (Collins et
464 al. 2021) and STIX (Chowdhury et al. 2022) to further filter and improve variant prioritization.

465

466 COLO829 and its matched control COLO829BL are well established cancer cell lines, which
467 provide advantages for testing and benchmarking of sequencing and analysis approaches
468 (Craig et al. 2016; Espejo Valle-Inclan et al. 2022; Pleasance et al. 2010). However, our study
469 clearly demonstrates that caution should be used when interpreting and comparing results from
470 a single analysis, sequencing run, or even biological sample. When we compared our
471 benchmark to the most up to date benchmark from Valle-Inclan et al (Espejo Valle-Inclan et al.
472 2022), we identified multiple SVs that did not appear in any other replicates of COLO829
473 sequenced at other centers, and some of which were not identified in re-analysis of the original
474 Nanopore data. For example, four SV present in the benchmark from Valle-Inclan could not be
475 identified in any of the cancer replicates, which includes a replicate from the aforementioned
476 benchmark. Moreover, ten SVs could not be identified in all cancer replicates or had low VAF
477 (<10%) which complicates the interpretation of the results. This suggests caution should be
478 used when interpreting results in comparison to the previous COLO829/COLO829BL
479 benchmark (Espejo Valle-Inclan et al. 2022). Furthermore, across the different replicates, we
480 could demonstrate and manually validate differences between cancer replicate samples (i.e SVs
481 uniquely identified per replicate), similar to changes documented for somatic SNVs observed in
482 different COLO829 passages (Craig et al. 2016). These changes clearly highlight the instability
483 of COLO829/COLO829BL and need to be taken into consideration across benchmarks.
484 Otherwise, there is potential for incorrect interpretation of the accuracy of sequencing and
485 analysis techniques evaluated against the benchmark. To aid this, we identified a core set of
486 SVs that appear to be maintained stably across independent sequencing replicates, and we
487 propose this set to be used as the updated version of the COLO829/COLO829BL benchmark
488 for the detection of somatic cancer SVs. It is interesting to note that some differences in the
489 replicates of COLO829/COLO829BL can be attributed to technical artifacts, as the data set
490 includes Nanopore data generated with prior versions of the approach that might suffer from
491 base calling biases (ie. deletions) that have since been improved (Kolmogorov et al. 2023).
492 Nevertheless, biological artifacts are clearly present given the evolution of the cell line.

493

494 Our work demonstrates new approaches to optimize somatic SV prioritization in cancer with
495 potential improvements in other genetic diseases. We demonstrate this over patient samples
496 but further over COLO829 where we introduce a new benchmark due to its variability. Given all

497 these artifacts it's still important to note the importance of the widely available
498 COLO829/COLO829BL cell line for technology and analytical development.
499

500 **Methods**

501 ***Samples***

502 We used the COLO829 cancer cell-line (melanoma) and its germline control COLO829BL
503 (Blood, B lymphoblast) to assess somatic SV calling with whole genome long-read sequencing.
504 We sequenced the tumor/normal samples with different technologies (ONT PromethION, ONT
505 MinION and PacBio Revio) at different sites: Canada's Michael Smith Genome Sciences
506 Centre, BC, Canada (labeled GSC) sequenced with ONT PromethION; Oxford Nanopore
507 Technologies (labeled ONT) sequenced with ONT PromethION; Pacific Biosciences of
508 California, Inc (labeled REV) sequenced with ONT PacBio Revio; Center for Molecular Medicine
509 and Oncode Institute, UMC Utrecht, Utrecht, the Netherlands (labeled VAI) sequenced with ONT
510 MinION.

511
512 Cancer samples were derived from the Personal Oncogenomics (POG) program (Pleasance et
513 al. 2020), clinical trial NCT02155621, approved by and conducted under the University of British
514 Columbia – BC Cancer Research Ethics Board (H12-00137, H14-00681). POG1022 is a tumour
515 sample (diploid) from a metastatic diffuse large B-cell lymphoma, with the control sample taken
516 from a skin biopsy (POG1022SK). POG044 is a tumor sample (diploid) taken from a recurrent
517 anaplastic oligodendroglioma, with a blood control sample (POG044BL). **Supplementary Table**
518 **1** summarizes the samples used in this study and includes links for accessing the data. Two
519 human reference genomes were utilized: GRCh38 (hg38) and CHM13-T2T (hs1).

520

521 ***Alignment***

522 Minimap2 (Li 2021) (version 2.24-r1122) was used to align the long reads to the human
523 genome. We utilized two different versions of the human genome: GRCh38 and CHM13 version
524 2. We used default parameters, with the output to be in the SAM format (-a) and preset for
525 Oxford Nanopore reads (-x map-ont). Additionally, we converted the alignment to BAM format,
526 sorted and indexed using SAMTools (Danecek et al. 2021) (version 1.16.1).

527

528 ***LiftOver***

529 We used LevioSAM2 (Chen et al. 2024) (version 0.4.1) to liftover alignments from CHM13-T2T
530 into GRCh38. We ran LevioSAM2 using the provided chain files for CHM13-T2Tv2. We
531 differentiated between sequencing technologies (ONT and PacBio) in the configuration file and
532 mapping presets for minimap2 (map-ont and ont_all.yaml for ONT and map-hifi and
533 pacbio_all.yaml for PacBio). Additionally, we used specific allowed gaps (-g) and edit distance (-
534 H) values for each technology: -g 1500 -H 6000 for ONT and -g 1000 -H 100 for PacBio.

535

536 ***SV calling***

537 Sniffles2 (Smolka et al. 2024) (version 2.2) was used to call structural variants (SV). Each
538 sample had three reference backgrounds: GRCh38, CHM13-T2T, and the liftover alignment
539 (CHM13-T2T to GRCh38). For both GRCh38 alignments (native and liftover), we added a
540 tandem repeat annotation file during the run (--tandem-repeats file.bed). This file is provided
541 alongside Sniffles (<https://zenodo.org/records/8121996>). In all cases the reference genome was
542 provided (--reference) and the SNF file was produced (--snf file). The rest of the parameters
543 were left as default. Next, we performed population SV calling with Sniffles2 using the SNF files
544 produced in the previous step. Each population call was done by reference genome, such that
545 we produced three population files, one for GRCh38, one for CHM13 and one for the liftover. In
546 this step Sniffles2 was used with default parameters.

547

548 ***Upset plot***

549 We utilized the UpSetR packager (Conway et al. 2017) to produce upset plots that represent the
550 number of shared SVs among samples. We used a custom script to extract the support vector
551 provided in the VCF file (SUPP_VEC), which encodes the presence/absence of a variant in a
552 given sample. **Supplementary tables 2, 9 and 14** contain the necessary information to produce
553 the upset plots.

554

555 ***Somatic SV detection COLO289***

556 We used the fully-genotyped population VCF file generated in the previous step to assess the
557 SVs that were unique to the cancer samples. We used the support vector provided in the VCF
558 file (SUPP_VEC) to assess the presence/absence of each variant. We used bcftools (version
559 1.16) (Li 2011) to extract SVs that were tumor-only (using the --include option, example:

560 bcftools view --include "SUPP_VEC = '11110000'"). We denoted a somatic-cancer variant if the
561 variant was only present in the cancer samples, had a VAF \geq 10%, and had a minimum of 10
562 supporting reads. Initially, we excluded any SV that had any read support in any of the control
563 replicates. Manual inspection overrode cases where a single read was detected in a control
564 sample.

565

566 ***SV annotation***

567 Post-processing of SVs of the GRCh38 and CHM13-T2T to GRCh38 call sets was conducted
568 with MAVIS (Reisle et al. 2019) (version 3.1.0). Briefly, we collapsed duplicate SVs, and merged
569 SVs by breakpoint proximity (100 bp) and type. Next, we used the RefSeq curated gene
570 annotation track in the UCSC genome browser (Kent et al. 2002) to annotate the impacted
571 genes from the somatic SVs for COLO829-T2T. SV subtype-specific analyses also incorporated
572 RepeatMasker (version 4.1) (Tarailo-Graovac and Chen 2009) to annotate the events. Events
573 flagged as non-synonymous coding variants by MAVIS were manually reviewed.

574

575 ***SV benchmark***

576 We compared the two COLO829 somatic SV datasets that are in GRCh38 coordinates
577 (COLO829-GRCh38 and COLO829-lifted) to a published COLO829 SV benchmark by Valle-
578 Inclan et al. (Espejo Valle-Inclan et al. 2022), COLO829-VAI hereafter. The COLO829-VAI
579 benchmark consists of 68 SVs, with representation from all five SV types (INS, DEL, DUP, INV,
580 BND). We removed six SVs whose length was smaller than 50 bp (default reported by
581 Sniffles2), leaving 62 SVs in the COLO829-VAI benchmark. We used BEDtools (version 2.31)
582 (Quinlan and Hall 2010) to compare the SV coordinates of the COLO829-VAI benchmark to our
583 COLO829-GRCh38 and COLO829-lifted somatic SV datasets.

584

585 For the case of the COLO829-T2T dataset, we included the --output-rnames parameter in
586 Sniffles to output the read names supporting each SV. We then used these read names along
587 with the chromosome and the SV type as a proxy for matching the SV from CHM13-T2T to
588 GRCh38 coordinates to perform a partial benchmark (FP only). Only for one case we used the
589 same procedure to investigate a FN call in the COLO829-T2T dataset (DEL in chr16), because
590 Sniffles made the call and reported the read names. However, the SV was removed from the
591 final call set by the COV_MIN filter, thus it was still considered a FN.

592 ***SV genotyping (force-calling)***

593 We used the Sniffles2 --genotype-vcf option to look for all the FN calls in the COLO829-
594 GRCh38 and COLO829-lifted datasets. This option takes a VCF as input (including the FN calls
595 for our case) and uniquely searches for the SV present in the VCF input and updates the
596 genotype according to what Sniffles detects in the alignment file. For the case of BNDs, we
597 additionally took all the reads overlapping with the coordinates provided by the COLO829-VAI
598 benchmark and searched for reads that had supplementary alignments (SA; flag 2048) and
599 compared the coordinates from the SA tag in the alignment to the "CHR2" value from the INFO
600 field of the COLO829-VAI benchmark. We matched chromosomes and allowed for a 10kb
601 distance between the positions.

602

603 ***Somatic SV detection POG samples***

604 We used the fully-genotyped population VCF file generated by merging tumor and normal
605 samples using Sniffles2 SNF files (sniffles --input tumor.snf normal.snf --vcf merge.vcf.gz).
606 Once merged, we used the support vector provided in the VCF file (SUPP_VEC) to assess the
607 presence/absence of each SV. We used bcftools (version 1.16) to extract SVs that were tumor-
608 only (using the --include option, example: bcftools view --include "SUPP_VEC = '10'"). We
609 denoted a somatic-cancer variant if the SV was only present in the cancer samples, a VAF >=
610 10% and a minimum of 10 read support.

611 **Data access**

612 **Supplementary Table 1** summarizes the samples used in this study and includes links for
613 accessing the data. The new COLO829/COLO829BL proposed benchmark and merge VCF
614 files that include the eight tumor/normal samples can be found at [10.5281/zenodo.10819636](https://zenodo.org/record/10819636).
615 Nutty, a Sniffles2 companion app for parsing the VCF was used <https://github.com/lfpaulin/nutty>
616 It contains commands to reproduce the COLO829 and POG analysis.

617 **Acknowledgements**

618 This study was in part supported by funding from the Canada Research Chairs Program, Terry
619 Fox Research Institute Marathon of Hope and the British Columbia Cancer Foundation. FJS,
620 LFP was supported by NIH (UM1DA058229, 1UG3NS132105-01, 1U01HG011758-01). This
621 study was conducted with the financial support of The Terry Fox Research Institute and the

622 Terry Fox Foundation. The views expressed in the publication are the views of the authors and
623 do not necessarily reflect those of the Terry Fox Research Institute or the Terry Fox Foundation.

624 **Conflict of interest**

625 The following authors disclose relevant potential competing interests: Kieran O'Neill, Vanessa
626 Porter, Luis F Paulin and Steven J.M. Jones received travel funding from Oxford Nanopore
627 Technologies to present at conferences in 2022 and/or 2023. Fritz J Sedlazeck receives
628 research support from ONT, Pacbio, Illumina and Genentech. Luis F Paulin received research
629 support from Genentech from 2021 to 2023.

630

631 **References**

- 632 Aganezov S, Goodwin S, Sherman RM, Sedlazeck FJ, Arun G, Bhatia S, Lee I, Kirsche M,
633 Wappel R, Kramer M, et al. 2020. Comprehensive analysis of structural variants in breast
634 cancer genomes using single-molecule sequencing. *Genome Res* **30**: 1258–1273.
- 635 Aganezov S, Yan SM, Soto DC, Kirsche M, Zarate S, Avdeyev P, Taylor DJ, Shafin K, Shumate
636 A, Xiao C, et al. 2022. A complete reference genome improves analysis of human genetic
637 variation. *Science* **376**: eabl3533.
- 638 Akagi K, Symer DE, Mahmoud M, Jiang B, Goodwin S, Wangsa D, Li Z, Xiao W, Dunn JD, Ried
639 T, et al. 2023. Intratumoral Heterogeneity and Clonal Evolution Induced by HPV
640 Integration. *Cancer Discov* **13**: 910–927.
- 641 Behera S, LeFaive J, Orchard P, Mahmoud M, Paulin LF, Farek J, Soto DC, Parker SCJ, Smith
642 AV, Dennis MY, et al. 2023. FixItFelix: improving genomic analysis by fixing reference
643 errors. *Genome Biol* **24**: 31.
- 644 Bellon M, Bialuk I, Galli V, Bai X-T, Farre L, Bittencourt A, Marçais A, Petrus MN, Ratner L,
645 Waldmann TA, et al. 2021. Germinal epimutation of Fragile Histidine Triad (FHIT) gene is
646 associated with progression to acute and chronic adult T-cell leukemia diseases. *Mol*
647 *Cancer* **20**: 86.
- 648 Cao Z, Ji J, Wang F-B, Kong C, Xu H, Xu Y-L, Chen X, Yu Y-W, Sun Y-H. 2020. MAGI-2
649 downregulation: a potential predictor of tumor progression and early recurrence in Han
650 Chinese patients with prostate cancer. *Asian J Androl* **22**: 616–622.
- 651 Chander V, Gibbs RA, Sedlazeck FJ. 2019. Evaluation of computational genotyping of structural
652 variation for clinical diagnoses. *Gigascience* **8**.
653 <http://dx.doi.org/10.1093/gigascience/giz110>.
- 654 Chen N-C, Paulin LF, Sedlazeck FJ, Koren S, Phillippy AM, Langmead B. 2024. Improved
655 sequence mapping using a complete reference genome and lift-over. *Nat Methods* **21**: 41–
656 49.

- 657 Choo Z-N, Behr JM, Deshpande A, Hadi K, Yao X, Tian H, Takai K, Zakusilo G, Rosiene J, Da
658 Cruz Paula A, et al. 2023. Most large structural variants in cancer genomes can be
659 detected without long reads. *Nat Genet* **55**: 2139–2148.
- 660 Chowdhury M, Pedersen BS, Sedlazeck FJ, Quinlan AR, Layer RM. 2022. Searching thousands
661 of genomes to classify somatic and novel structural variants using STIX. *Nat Methods* **19**:
662 445–448.
- 663 Collins RL, Brand H, Karczewski KJ, Zhao X, Alföldi J, Francioli LC, Khera AV, Lowther C,
664 Gauthier LD, Wang H, et al. 2020. A structural variation reference for medical and
665 population genetics. *Nature* **581**: 444–451.
- 666 Collins RL, Brand H, Karczewski KJ, Zhao X, Alföldi J, Francioli LC, Khera AV, Lowther C,
667 Gauthier LD, Wang H, et al. 2021. Author Correction: A structural variation reference for
668 medical and population genetics. *Nature* **590**: E55.
- 669 Conway JR, Lex A, Gehlenborg N. 2017. UpSetR: an R package for the visualization of
670 intersecting sets and their properties. *Bioinformatics* **33**: 2938–2940.
- 671 Craig DW, Nasser S, Corbett R, Chan SK, Murray L, Legendre C, Tembe W, Adkins J, Kim N,
672 Wong S, et al. 2016. A somatic reference standard for cancer genome sequencing. *Sci*
673 *Rep* **6**: 24607.
- 674 Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, Whitwham A, Keane T,
675 McCarthy SA, Davies RM, et al. 2021. Twelve years of SAMtools and BCFtools.
676 *Gigascience* **10**. <http://dx.doi.org/10.1093/gigascience/giab008>.
- 677 Ding H, Wang K-Y, Chen S-Y, Guo K-W, Qiu W-H. 2023. Validating the role of PTGIS gene in
678 colorectal cancer by bioinformatics analysis and in vitro experiments. *Sci Rep* **13**: 16496.
- 679 Druker BJ, Sawyers CL, Kantarjian H, Resta DJ, Reese SF, Ford JM, Capdeville R, Talpaz M.
680 2001. Activity of a specific inhibitor of the BCR-ABL tyrosine kinase in the blast crisis of
681 chronic myeloid leukemia and acute lymphoblastic leukemia with the Philadelphia
682 chromosome. *N Engl J Med* **344**: 1038–1042.
- 683 English AC, Menon VK, Gibbs RA, Metcalf GA, Sedlazeck FJ. 2022. Truvari: refined structural
684 variant comparison preserves allelic diversity. *Genome Biol* **23**: 271.
- 685 English A, Dolzhenko E, Jam HZ, Mckenzie S, Olson ND, De Coster W, Park J, Gu B, Wagner
686 J, Eberle MA, et al. 2023. Benchmarking of small and large variants across tandem
687 repeats. *bioRxiv*. <http://dx.doi.org/10.1101/2023.10.29.564632>.
- 688 Espejo Valle-Inclan J, Besselink NJM, de Bruijn E, Cameron DL, Ebler J, Kutzera J, van
689 Lieshout S, Marschall T, Nelen M, Priestley P, et al. 2022. A multi-platform reference for
690 somatic structural variation detection. *Cell Genom* **2**: 100139.
- 691 Fujimoto A, Wong JH, Yoshii Y, Akiyama S, Tanaka A, Yagi H, Shigemizu D, Nakagawa H,
692 Mizokami M, Shimada M. 2021. Whole-genome sequencing with long reads reveals
693 complex structure and origin of structural variation in human genetic variations and somatic
694 mutations in cancer. *Genome Med* **13**. <https://pubmed.ncbi.nlm.nih.gov/33910608/>
695 (Accessed March 1, 2024).

- 696 Hernández Borrero LJ, El-Deiry WS. 2021. Tumor suppressor p53: Biology, signaling pathways,
697 and therapeutic targeting. *Biochim Biophys Acta Rev Cancer* **1876**: 188556.
- 698 Husanie H, Abu-Remaileh M, Maroun K, Abu-Tair L, Safadi H, Atlan K, Golan T, Aqeilan RI.
699 2022. Loss of tumor suppressor WWOX accelerates pancreatic cancer development
700 through promotion of TGF β /BMP2 signaling. *Cell Death Dis* **13**: 1074.
- 701 Iqbal Z, Absar M, Akhtar T, Aleem A, Jameel A, Basit S, Ullah A, Afzal S, Ramzan K, Rasool M,
702 et al. 2021. Integrated Genomic Analysis Identifies ANKRD36 Gene as a Novel and
703 Common Biomarker of Disease Progression in Chronic Myeloid Leukemia. *Biology* **10**.
704 <http://dx.doi.org/10.3390/biology10111182>.
- 705 Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, Zahler AM, Haussler D. 2002. The
706 human genome browser at UCSC. *Genome Res* **12**: 996–1006.
- 707 Kolmogorov M, Billingsley KJ, Mastoras M, Meredith M, Monlong J, Lorig-Roach R, Asri M,
708 Alvarez Jerez P, Malik L, Dewan R, et al. 2023. Scalable Nanopore sequencing of human
709 genomes provides a comprehensive view of haplotype-resolved variation and methylation.
710 *Nat Methods* **20**: 1483–1492.
- 711 Li H. 2011. A statistical framework for SNP calling, mutation discovery, association mapping and
712 population genetical parameter estimation from sequencing data. *Bioinformatics* **27**: 2987–
713 2993.
- 714 Li H. 2021. New strategies to improve minimap2 alignment accuracy. *Bioinformatics* **37**: 4572–
715 4574.
- 716 Lombardi A, Arseni L, Carriero R, Compe E, Botta E, Ferri D, Uggè M, Biamonti G, Peverali FA,
717 Bione S, et al. 2021. Reduced levels of prostaglandin I synthase: a distinctive feature of the
718 cancer-free trichothiodystrophy. *Proc Natl Acad Sci U S A* **118**.
719 <http://dx.doi.org/10.1073/pnas.2024502118>.
- 720 Mahmoud M, Gobet N, Cruz-Dávalos DI, Mounier N, Dessimoz C, Sedlazeck FJ. 2019.
721 Structural variant calling: the long and the short of it. *Genome Biol* **20**: 246.
- 722 Mahmoud M, Huang Y, Garimella K, Audano PA, Wan W, Prasad N, Handsaker RE, Hall S,
723 Pionzio A, Schatz MC, et al. 2024. Utility of long-read sequencing for All of Us. *Nat*
724 *Commun* **15**: 837.
- 725 Majidian S, Agostinho DP, Chin C-S, Sedlazeck FJ, Mahmoud M. 2023. Genomic variant
726 benchmark: if you cannot measure it, you cannot improve it. *Genome Biol* **24**: 221.
- 727 Mandelker D, Ceyhan-Birsoy O. 2020. Evolving Significance of Tumor-Normal Sequencing in
728 Cancer Care. *Trends Cancer Res* **6**: 31–39.
- 729 Nattestad M, Goodwin S, Ng K, Baslan T, Sedlazeck FJ, Rescheneder P, Garvin T, Fang H,
730 Gurtowski J, Hutton E, et al. 2018. Complex rearrangements and oncogene amplifications
731 revealed by long-read DNA and RNA sequencing of a breast cancer cell line. *Genome Res*
732 **28**: 1126–1135.
- 733 Nicholas TJ, Cormier MJ, Quinlan AR. 2022. Annotation of structural variants with reported

- 734 allele frequencies and related metrics from multiple datasets using SVAFotate. *BMC*
735 *Bioinformatics* **23**: 490.
- 736 Nurk S, Koren S, Rhie A, Rautiainen M, Bizkadze AV, Mikheenko A, Vollger MR, Altemose N,
737 Uralsky L, Gershman A, et al. 2022. The complete sequence of a human genome. *Science*
738 **376**: 44–53.
- 739 Olson ND, Wagner J, Dwarshuis N, Miga KH, Sedlazeck FJ, Salit M, Zook JM. 2023. Variant
740 calling and benchmarking in an era of complete human genome sequences. *Nat Rev*
741 *Genet* **24**: 464–483.
- 742 O'Neill K, Pleasance E, Fan J, Akbari V, Chang G, Dixon K, Csizmok V, MacLennan S, Porter
743 V, Galbraith A, et al. 2024. Long-read sequencing of an advanced cancer cohort resolves
744 rearrangements, unravels haplotypes, and reveals methylation landscapes. *medRxiv*
745 2024.02.20.24302959. <https://www.medrxiv.org/content/10.1101/2024.02.20.24302959v1>
746 (Accessed March 7, 2024).
- 747 Pleasance ED, Cheetham RK, Stephens PJ, McBride DJ, Humphray SJ, Greenman CD, Varela
748 I, Lin M-L, Ordóñez GR, Bignell GR, et al. 2010. A comprehensive catalogue of somatic
749 mutations from a human cancer genome. *Nature* **463**: 191–196.
- 750 Pleasance E, Titmuss E, Williamson L, Kwan H, Culibrk L, Zhao EY, Dixon K, Fan K, Bowlby R,
751 Jones MR, et al. 2020. Pan-cancer analysis of advanced patient tumors reveals
752 interactions between therapy and genomic landscapes. *Nat Cancer* **1**: 452–468.
- 753 Porter VL, O'Neill K, MacLennan S, Corbett RD, Ng M, Culibrk L, Hamadeh Z, Iden M, Schmidt
754 R, Tsaih S-W, et al. 2023. Genomic structures and regulation patterns at HPV integration
755 sites in cervical cancer. *bioRxiv*. <http://dx.doi.org/10.1101/2023.11.04.564800>.
- 756 Qiao D, Liu Y, Lei Y, Zhang C, Bu Y, Tang Y, Zhang Y. 2023. rRNA-Derived Small RNA rsRNA-
757 28S Regulates the Chemoresistance of Prostate Cancer Cells by Targeting PTGIS. *Front*
758 *Biosci* **28**: 102.
- 759 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic
760 features. *Bioinformatics* **26**: 841–842.
- 761 Reisle C, Mungall KL, Choo C, Paulino D, Bleile DW, Muhammadzadeh A, Mungall AJ, Moore
762 RA, Shlafman I, Coope R, et al. 2019. MAVIS: merging, annotation, validation, and
763 illustration of structural variants. *Bioinformatics* **35**: 515–517.
- 764 Salzberg SL. 2019. Next-generation genome annotation: we still struggle to get it right. *Genome*
765 *Biol* **20**: 92.
- 766 Sedlazeck FJ, Rescheneder P, Smolka M, Fang H, Nattestad M, von Haeseler A, Schatz MC.
767 2018. Accurate detection of complex structural variations using single-molecule
768 sequencing. *Nat Methods* **15**: 461–468.
- 769 Shiraishi Y, Koya J, Chiba K, Okada A, Arai Y, Saito Y, Shibata T, Kataoka K. 2023. Precise
770 characterization of somatic complex structural variations from tumor/control paired long-
771 read sequencing data with nanomonsv. *Nucleic Acids Res* **51**: e74.

- 772 Smolka M, Paulin LF, Grochowski CM, Horner DW, Mahmoud M, Behera S, Kalef-Ezra E,
773 Gandhi M, Hong K, Pehlivan D, et al. 2024. Detection of mosaic and population-level
774 structural variants with Sniffles2. *Nat Biotechnol*. [http://dx.doi.org/10.1038/s41587-023-](http://dx.doi.org/10.1038/s41587-023-02024-y)
775 02024-y.
- 776 Tanner A, Sagoo MS, Mahroo OA, Pulido JS. 2024. Genetic analysis of ocular tumour-
777 associated genes using large genomic datasets: insights into selection constraints and
778 variant representation in the population. *BMJ Open Ophthalmol* **9**.
779 <http://dx.doi.org/10.1136/bmjophth-2023-001565>.
- 780 Tarailo-Graovac M, Chen N. 2009. Using RepeatMasker to identify repetitive elements in
781 genomic sequences. *Curr Protoc Bioinformatics* **Chapter 4**: 4.10.1–4.10.14.
- 782 Thibodeau ML, O'Neill K, Dixon K, Reisle C, Mungall KL, Krzywinski M, Shen Y, Lim HJ, Cheng
783 D, Tse K, et al. 2020. Improved structural variant interpretation for hereditary cancer
784 susceptibility using long-read sequencing. *Genet Med* **22**.
785 <https://pubmed.ncbi.nlm.nih.gov/32624572/> (Accessed March 1, 2024).
- 786 Tsang ES, Grisdale CJ, Pleasance E, Topham JT, Mungall K, Reisle C, Choo C, Carreira M,
787 Bowlby R, Karasinska JM, et al. 2021. Uncovering Clinically Relevant Gene Fusions with
788 Integrated Genomic and Transcriptomic Profiling of Metastatic Cancers. *Clin Cancer Res*
789 **27**: 522–531.
- 790 Wagner J, Olson ND, Harris L, McDaniel J, Cheng H, Fungtammasan A, Hwang Y-C, Gupta R,
791 Wenger AM, Rowell WJ, et al. 2022. Curated variation benchmarks for challenging
792 medically relevant autosomal genes. *Nat Biotechnol* **40**: 672–680.
- 793 Yehia L, Plitt G, Tushar AM, Joo J, Burke CA, Campbell SC, Heiden K, Jin J, Macaron C,
794 Michener CM, et al. 2023. Longitudinal Analysis of Cancer Risk in Children and Adults With
795 Germline PTEN Variants. *JAMA Netw Open* **6**: e239705.
- 796 Yoshihara K, Shahmoradgoli M, Martínez E, Vegesna R, Kim H, Torres-Garcia W, Treviño V,
797 Shen H, Laird PW, Levine DA, et al. 2013. Inferring tumour purity and stromal and immune
798 cell admixture from expression data. *Nat Commun* **4**: 2612.
- 799 Zhang X, Gibhardt CS, Will T, Stanisz H, Körbel C, Mitkovski M, Stejerean I, Cappello S,
800 Pacheu-Grau D, Dudek J, et al. 2019. Redox signals at the ER-mitochondria interface
801 control melanoma progression. *EMBO J* **38**: e100871.
- 802 Zook JM, Hansen NF, Olson ND, Chapman L, Mullikin JC, Xiao C, Sherry S, Koren S, Phillippy
803 AM, Boutros PC, et al. 2020. A robust benchmark for detection of germline large deletions
804 and insertions. *Nat Biotechnol* **38**: 1347–1355.