

1 Ultra-fast Prediction of Somatic Structural Variations by Reduced Read 2 Mapping via Pan-Genome *k*-mer Sets

3

4 Min-Hak Choi^{1,#}, Jang-il Sohn^{1,2,#}, Dohun Yi^{1,#}, A Vipin Menon¹, Yeon Jeong Kim⁴, Sungkyu Kyung⁵,
5 Seung-Ho Shin⁵, Byunggook Na⁶, Je-Gun Joung⁴, Sungro Yoon⁶, Youngil Koh⁷, Daehyun Baek⁸,
6 Tae-Min Kim⁹, and Jin-Wu Nam^{1,2,3*}

7

8 ¹Department of Life Science, Hanyang University, Seoul 04763, Republic of Korea,

9 ²Research Institute for Convergence of Basic Sciences, Hanyang University, Seoul 04763, Republic
10 of Korea,

11 ³Research Institute for Natural Sciences, Hanyang University, Seoul 04763, Republic of Korea,

12 ⁴Samsung Genome Institute, Samsung Medical Center, 81 Irwon-ro, Gangnam-gu, Seoul, 06351,
13 Republic of Korea,

14 ⁵GENINUS Inc, Seoul 05836, Republic of Korea,

15 ⁶Department of Electrical and Computer Engineering, Seoul National University, 08826, Republic of
16 Korea,

17 ⁷College of Medicine, Seoul National University, Seoul 03080, Republic of Korea,

18 ⁸School of Biological Sciences, Seoul National University, Seoul 08826, Republic of Korea.

19 ⁹Department of Medical Informatics and Cancer Research Institute, College of Medicine, The
20 Catholic University of Korea, Seoul 06591, Republic of Korea.

21

22 [#]These authors contributed equally

23 ^{*}Correspondence: Jin-Wu Nam

24 Tel: +82-2-2220-2428, Fax: +82-2-2298-0319, Email: jwnam@hanyang.ac.kr

25

26 Contributions

27 MHC, JIS, DHY, and AVM performed analyses; MHC, JIS, DHY, and AVM contributed to writing
28 the codes; MHC, JIS, BKN, and SRY contributed to parallel computing; YK and Genius provided
29 validation datasets; YJK and JGJ performed experimental validations; JIS and JWN contributed to
30 writing the manuscript; DHB, TMK, and JWN supervised the project; and JWN conceived the idea.

31

32 **Keywords:** Somatic structural variation, whole genome sequencing, cancer panel sequencing, fusion
33 gene, pan-genome *k*-mer set, reduced mapping

35 **ABSTRACT**

36 Genome rearrangements often result in copy number alterations of cancer-related genes and cause the
37 formation of cancer-related fusion genes. Current structural variation (SV) callers, however, still
38 produce massive numbers of false positives (FPs) and require high computational costs. Here, we
39 introduce an ultra-fast and high-performing somatic SV detector, called ETCHING, that significantly
40 reduces the mapping cost by filtering reads matched to pan-genome and normal *k*-mer sets. To reduce
41 the number of FPs, ETCHING takes advantage of a Random Forest classifier that utilizes six
42 breakend-related features. We systematically benchmarked ETCHING with other SV callers on
43 reference SV materials, validated SV biomarkers, tumor and matched-normal whole genomes, and
44 tumor-only targeted sequencing datasets. For all datasets, our SV caller was much faster ($\geq 15X$) than
45 other tools without compromising performance or memory use. Our approach would provide not only
46 the fastest method for largescale genome projects but also an accurate clinically practical means for
47 real-time precision medicine.

48

49 **Introduction**

50 Chromosomal rearrangements in coding regions and regulatory non-coding elements often cause
51 malignancy of somatic cells. Although structural variations (SVs) occur much less frequently than
52 single nucleotide variations (SNVs), the SVs often have a greater impact on cellular functions and
53 gene expression¹. In particular, large SVs (>1Kbp), which include large insertions (INSs), deletions
54 (DELs), inversions (INVs), duplications (DUPs), and translocations (TRAs), are more often
55 associated with gain- and/or loss-of-function of cancer-related genes and druggable target genes for
56 cancer treatments than are SNVs²⁻⁷. For instance, *ERBB2* amplification in breast cancers (BRCA) ^{8,9},
57 *EML4-ALK* fusion in lung cancer ¹⁰, and *BCR-ABL* fusion in chronic myeloid leukemia ¹¹ are well-
58 known SV-driven cancer drivers and actionable targets for cancer treatments. Hence, the rapid
59 detection of cancer-related SVs is indispensable for companion diagnostics and targeted cancer
60 therapy.

61 So far, a handful of SV callers have been introduced to find germline and somatic SVs in
62 normal and/or tumor samples by using a read-based approach — read-depth¹², discordant read-pairs
63 ¹³, soft-clipped reads¹⁴⁻¹⁶, and their combinations¹⁷⁻²¹ — or by using a *k*-mer-based approach²². Some
64 of them utilize local assembly of reads^{13, 20-22} to precisely detect breakpoints (BPs) and SV types.
65 Regardless of the approach, all current SV callers require genome mapping of all input reads.
66 Although the mapping process is an indispensable step for the confident identification of SVs, it
67 consumes most of the computing time in processing massive whole genome sequencing (WGS) data.
68 For instance, the genome mapping of 30X WGS data from a cancer patient takes ~300 hours with a
69 single thread on a high-performing computer, resulting in delayed diagnosis. Furthermore, SV studies
70 for largescale WGS projects, such as those undertaken by the Pan-Cancer Analysis of Whole
71 Genomes (PCAWG)²³ and the Genome Aggregation Database (gnomAD)²⁴ consortiums, would be
72 only doable by institutes with access to a giant computing facility or expensive cloud computing
73 services.

74 The majority of sequenced reads are reference reads (perfectly matched to the reference
75 genome), which could be dispensable for SV calling. Mapping the reference reads consumes

76 expensive computing time. It may also increase background noise resulting from imprecise and
77 ambiguous alignments of the reads (mainly due to repeats or low-complexity regions) or from
78 unresolved misassemblies of the reference ^{25, 26}. Thus, only mapping informative (non-reference)
79 reads to detect SVs would both reduce computing time and increase accuracy.

80 In general, somatic SV callers use a case-control design that compares tumor (case) SVs with
81 those of matched normals (control) to detect somatic (case-specific) SVs. The absence of matched
82 normal samples may lead to either a failure of SV-calling or a high FP rate spawned by germline SVs.
83 In particular, cancer panel sequencing is frequently carried out using only tumor samples. Using the
84 pan-genome sequences containing all non-medical variations instead of a matched normal sample
85 would help to enhance the accuracy of SV calling in this situation.

86 In this study, we developed ETCHING, an ultra-fast SV detection method. Our approach
87 significantly reduces the number of reads to be mapped by excluding those from the reference and/or
88 pan-genome *k*-mer (PGK) set. This new strategy drastically reduces running time (it is at least ~15
89 times faster than other methods) without compromising performance by taking advantage of machine-
90 learning-based classification to remove FP SVs further. ETCHING displays either comparable to or
91 better accuracy than other state-of-the-art SV detection tools on benchmarking whole genome and
92 panel sequencing datasets as well as reference materials.

93

94 **Results**

95 **Fast prediction of somatic SVs**

96 We report the development of ETCHING (Efficient deTection of CHromosomal rearrangements and
97 fusIoN Genes) – a fast computational SV caller that comprises four stepwise modules: Filter, Caller,
98 Sorter, and Fusion-identifier (Fig. 1a; Supplementary Fig. 1; see Methods for more details). The Filter
99 module uses one of three different filters: a Pan-Genome k -mer (PGK) filter that excludes tumor reads
100 in which all k -mers are present in PGK, a Normal filter that removes those reads in which all k -mers
101 come from normal reads (not using reference genomes), or a combined (PGKN) filter (Fig. 1b). PGK
102 is a unique set of 31-mers from 10 human genome assemblies and nonpathogenic single nucleotide
103 polymorphisms (SNPs) from dbSNP ($\sim 3.9 \times 10^9$ k -mers; Supplementary Fig. 2; Supplementary Table
104 1). This module allows us to collect tumor-specific (TS) reads by filtering reference reads, those with
105 germline variations, and those matched to normal reads. We used The Cancer Genome Atlas (TCGA)
106 BRCA WGS data used in a previous SV study²⁷ for checking the Filter module. Of the BRCA
107 samples, 31 and 9 were selected for training and hold-out test, respectively, by random selection
108 (Methods; Supplementary Table 2). For the hold-out test dataset, the Filter module excluded about
109 96.2% of the reads by PGK, 99.2% by Normal, and 99.4% by PGKN (Fig. 1c). The remaining TS
110 reads clearly present BPs with a sharp decay of read-depth in somatic DEL, DUP, INV, and TRA
111 examples, reminiscent of the chemical etching process (Fig. 1d). This filtration method significantly
112 shortened the mapping process. The mapping time for TS reads from the nine hold-out BRCA WGS
113 datasets with varied coverages (33–68X and 27–56X in tumor and normal samples, respectively) was
114 approximately 300 times faster than that for all reads (Unfiltered) using BWA-MEM²⁸ (Fig. 1e).

115 After mapping TS reads to the reference genome (hg19), the Caller module collects simple-
116 clipped reads to find initial BPs (Supplementary Fig. 3a) and then defines breakends (BNDs) for BP
117 pairs by considering the clipped direction (Supplementary Fig. 3b). The identified BNDs were then
118 assigned to an SV type, such as DEL, DUP, INV, and TRA, according to their position and the
119 clipped direction (Supplementary Fig. 3c; Methods). Next, the Sorter module predicts a confidence
120 score for each SV call using machine learning models pre-trained over the 31 training datasets

121 (Methods). Because there is no ground truth for the TCGA dataset, we instead used a silver standard
122 set of SVs, simultaneously detected by multiple SV callers, during training and evaluation (Methods).
123 Random Forest ²⁹ (RF)-based sorter was chosen as our default SV sorter module (Methods;
124 Supplementary Fig. 4). In the last step, with the predicted SVs, the Fusion-identifier module predicts
125 fusion-gene (FG) candidates (Methods).

126 We compared the running time of ETCHING with those of other SV callers over the hold-out
127 test dataset. The CPU time (running time converted in a single thread) for the entire SV prediction
128 process of ETCHING was at least 15 times less than those of the other SV callers (Fig. 1f). In real
129 (wall-clock) time, ETCHING took 2.2 hours on average, meaning that it was at least 6.6 times faster
130 than the second-fastest caller (Manta), on 30 threads (Supplementary Fig. 5). The ETCHING process
131 not only reduced the running time but also increased the precision of the SV prediction (Fig. 1g). The
132 PGK, Normal, and PGKN filters gradually reduced the number of FP reads with little compromise of
133 the true positive (TP) rates, resulting in better performances (F1-score) with the filters on BRCA
134 WGS and HCC1395 cell line WGS datasets (Fig. 1g; Supplementary Fig. 6). Taken together, these
135 results suggest that ETCHING provides high-performance SV prediction at a faster rate than other SV
136 callers.

137

138 **ETCHING displays robust performance**

139 We next sought to systematically benchmark the performance of ETCHING against the performances
140 of the read-based callers DELLY, LUMPY, Manta, and SvABA, as well as that of a *k*-mer-based
141 caller, novoBreak, over WGS data from the HCC1395 cancer cell line (50X) and its matched normal
142 cell line, HCC1395 BL (30X). Because the HCC1395 dataset also lacks ground-truth SVs, we again
143 used the approach of employing silver standard SVs identified by multiple callers, mentioned above.
144 The precision-recall (PR) curves over varying parameters showed that ETCHING performed more
145 robustly than the other callers, particularly for precision (Fig. 2a). We obtained optimal cutoffs, which
146 were used in the following benchmarking analyses for fair comparisons (Fig. 2a, red indicator;
147 Methods).

148 Because the performances of SV callers tend to be affected by the read-depth³⁰, we then
149 examined the robustness of the SV callers over varying read-depths. For this comparison, we
150 randomly subsampled 40% (20X), 60% (30X), and 80% (40X) of the reads from the HCC1395 cancer
151 line (50X) while keeping the depth of the normal reads fixed, and then performed benchmarking
152 analyses with the optimal cutoffs (Fig. 2b). ETCHING displayed a robust performance, regardless of
153 the read-depth, and showed a slightly increased precision as the read-depth became higher. In contrast,
154 Manta and SvABA presented lower recall rates at low read-depths but performances that were
155 comparable to that of ETCHING at 50X.

156 To compare the performance of ETCHING on primary tumor samples with varying read-
157 depth and tumor purity with those of the other tools, the nine hold-out BRCA samples were again
158 used as the benchmarking dataset (Supplementary Table 2). In this analysis, ETCHING showed
159 results that were superior or comparable to those of other tools, regardless of the SV type (Fig. 2c).
160 Notably, ETCHING robustly predicted all SV types while displaying high F1-scores across samples,
161 compared to other tools. We also benchmarked the SV callers over 33 true SVs from four thyroid
162 cancer (THCA) samples of TCGA as an independent evaluation dataset. The performance of
163 ETCHING was comparable to those of SvABA and novoBreak in terms of the F1-scores (Fig. 2d;
164 Supplementary Fig. 7; Supplementary Table 3).

165 Because the silver standard set of SVs could still include FPs, we selected high-quality (HQ)
166 SVs with depth-difference and connect-pair scores for DEL/DUP and INV/TRA, respectively
167 (Supplementary Fig. 8; see Supplementary Note for details). With HQ SVs, ETCHING still displayed
168 an accuracy that was comparable or superior to that of the other tools (Supplementary Fig. 9).

169

170 **SV prediction of experimentally validated targets**

171 For experimental validation of the SV callers, we newly sequenced the whole genomes of 26 multiple
172 myeloma (MM) samples with matched normal samples (Supplementary Table 4). We first
173 benchmarked the SV callers using the MM samples, and found that ETCHING outperformed the
174 others over a silver standard set of all SV types (Fig. 2e; Supplementary Fig. 9 and 10). Notably, its

175 performance exceeded that of another *k*-mer-based caller, novoBreak, which showed a lower
176 precision, particularly for INV and TRA types.

177 We then evaluated all of the SV callers using known clinical SV biomarkers of MM, such as
178 DELs (in 1q25, *p16*, *RBI*, and *TP53*) and IGH rearrangements (including DELs and TRAs)(Fig. 3a)³¹.
179 Fluorescence *in situ* hybridization (FISH) and karyotype were first examined on the SV biomarkers
180 (Supplementary Table 5). However, because FISH probe sets (Supplementary Table 6) of the SV
181 biomarkers cannot discern focal deletion/duplication from (partial) aneuploidy, the true set of
182 biomarker SVs were selected through a manual curation by considering read-depth changes and
183 unbalanced minor allele frequency (Supplementary Fig. 11a,b) as well as discordant paired-reads in
184 tumor and normal samples for each patient (Supplementary Fig. 11c–e; Methods). The SV set
185 supported by FISH and/or karyotype (excluding aneuploidy) were well overlapped with manually
186 curated SV biomarkers (Fig. 3b). We accordingly benchmarked ETCHING and other SV callers with
187 the manually curated SV biomarkers as a true set. The receiver operating characteristics (ROC)
188 showed that ETCHING displays comparable or slightly better performances than other callers in the
189 SV biomarker level (Fig. 3c). Of 23 curated biomarkers, ETCHING detected 19. When breaking
190 down the *IGH* rearrangements into SV level, known MM target genes, *FGFR3*, *IL6ST*, *CCND3*,
191 *CCND1*, and *IPLL5* were detected as translocation partners by manual curation (Fig. 3b middle;
192 Supplementary Table 7). Of the 38 SVs, ETCHING detected 17 SVs but missed seven including a
193 *p16* DEL (MM17), three *IGH* DELs (MM10, 12, and 18), and three *IGH* TRAs (MM1, 11, and 14)
194 (Fig. 3b).

195 We further searched for SVs related to actionable (cancer-druggable and clinically verified)
196 targets from the OncoKB database³². ETCHING detected five actionable SV targets – *BRCA2* DEL
197 (MM22), *ALK* DUP (SNUH19_MM04), *PIK3CA* DUP (MM15), *AKT1* DUP (MM3), and *NTRK1*
198 DUP (SNUH19_MM01) (Fig. 3b,d). Of the five predicted targets, three targets (excluding those from
199 the MM3 and MM15 patients, which lack tumor DNA quantities) were verified by targeted PCR (Fig.
200 3e). The PCR products expected after amplification of *ALK* DUP (SNUH19_MM04), *BRCA2* DEL

201 (MM22), and *NTRK1* DUPs (SNUH19_MM01) were observed in tumor but not in normal samples,
202 indicating that the SV targets are true cases.

203

204 **SV and FG prediction in cancer panel sequencing**

205 Targeted gene panel sequencing is more relevant than WGS for clinical applications, and clinical
206 laboratories daily produce panel sequencing data with the aim of finding actionable target variations,
207 SNVs, SVs, and FGs. Targeted gene panel sequencing is often applied to detect low-frequency
208 alterations such as somatic SNVs or FGs in cell-free DNA from cancer patients. To test the
209 effectiveness of the SV callers in such clinical situations, we analyzed 56 targeted gene panel
210 sequencing data derived from three types of cell-free DNA (cfDNA) reference material (Methods):
211 Complete Reference (CR), Complete Mutation Mix (CMM), and Mutation Mix v2 (MMv2). Each
212 type contains two or three synthetic FGs with low mutant allele ratios (0.5–5%) and wild-type (WT)
213 alleles from a cell line, GM24385.

214 Because cancer panel sequencing approaches generally lack matched normal data, ETCHING
215 was first set to use a PGK filter to extract TS reads for SV prediction. Other benchmarking tools, with
216 the exception of novoBreak, also predict SVs in the absence of normal data. novoBreak, given its
217 requirement for normal data, used simulated data from the hg19 reference genome (Methods). Note
218 that we ran all tools with default parameters that display a better recall rate for panel sequencing data.

219 This analysis showed that, along with LUMPY and DELLY, ETCHING is one of the top
220 callers in terms of recall over such low mutant allele frequencies (Fig. 4a,b), while showing a
221 moderate level of additional calls in targeted regions (Fig. 4c). Additional calls could be either FP
222 calls or germline SVs from the WT sample. Compared to other tools, ETCHING barely predicted
223 additional calls in non-target regions, indicating a relatively low frequency of FP calls (Fig. 4c, gray).
224 Because the reference materials include WT data that lack mutant alleles, the benchmarking analyses
225 of SV prediction were also performed using the WT data as the normal sample. ETCHING was then
226 set to use a PGKN filter. Unlike the other tools, ETCHING and LUMPY maintained high recall rates
227 (Supplementary Fig. 12) compared to the results obtained without WT data (Fig. 4a,b). This result

228 indicates that ETCHING and LUMPY can effectively remove FPs without compromising the recall
229 rate for targeted gene panel sequencing data, regardless of the presence of matched normal data. Since
230 BreaKmer³³ is specialized for targeted sequencing data, we also tested it on the same dataset.
231 However, BreaKmer failed to report any result (Methods).

232 Then, on the cancer panel sequencing data from formalin-fixed paraffin-embedded (FFPE)
233 and frozen tissues from a previous study of BreaKmer³³, we evaluated the performances of
234 ETCHING and other tools including BreaKmer. The data consists of 105 replicates from 37 samples
235 of different types of cancers (Supplementary Table 8). Because the data included tumor samples
236 without matched normal samples, ETCHING utilized the PGK, rather than the PGKN, filter for this
237 prediction as above. All settings for the other tools were the same as were used for the reference
238 materials. Since the data contains small variants in *FLT3* and *KIT*, we included small variations in this
239 analysis (Methods). We first ran BreaKmer and compared its results to those of the previous study³³.
240 BreaKmer was still very specific, giving only 479 additional calls across all 105 cases. However, it
241 showed a lower recall rate (78 out of 105) than that they reported. It is possibly due to the lack of bait
242 information or using a different version of BreaKmer (Methods). Other tools showed comparable
243 recall rates (94 to 103 out of 105). Although DELLY showed the most sensitive performance, it was
244 at the cost of massive additional calls (about one million). ETCHING found 98 true variants, and its
245 number of additional calls was the lowest level except BreaKmer.

246 ETCHING was one of three tools that were able to detect all eight *FLT3* indels, which
247 appeared in diverse forms including seven cases of DUPs (32-73bp) and one case of small indel (30bp)
248 (Supplementary Fig. 13). Taken together, these results indicate that ETCHING shows high
249 performance for detecting SVs and FGs in both WGS and targeted sequencing data, indicating its
250 general usability.

251

252 **Benchmarking computational efficiency**

253 ETCHING significantly reduced the running time through implementation of the Filter module,
254 resulting in computational speeds that were at least 15 times faster than those of the other tools (Fig.

255 1f). Such fast predictions result from significantly reduced genome mapping of reads. Although
256 novoBreak also takes advantage of the *k*-mer approach to assembly TS contigs with BPs, it requires
257 prior genome mapping of all reads to find read clusters, which is a time-consuming step. To confirm
258 this conclusion, we determined the running times of ETCHING and novoBreak for each step (read
259 filtration, mapping, and SV calling) on the HCC1395 dataset (Fig. 5a). As shown in Fig. 1f, based on
260 its CPU time, ETCHING was approximately 15 times faster than novoBreak, mostly due to a
261 reduction in the mapping time. In fact, most of novoBreak's running time was spent in the mapping
262 step (87%, 283.5 CPU-hours), whereas ETCHING used about 13% of its running time for this step
263 (2.6 CPU-hours; Fig. 5a). Unlike other tools, ETCHING significantly reduces computational costs
264 through its filtration-and-mapping strategy (Fig. 5b). Using multiple processes (30 threads) for
265 parallel computing, ETCHING completed the entire procedure for nine hold-out datasets in 2.2h on
266 average and for HCC1395 in 1.5h (Supplementary Fig. 5 and 14). Application of different numbers of
267 threads showed that the efficiency of ETCHING approached saturation (1.5h) over 25 threads
268 (Supplementary Fig. 15).

269 Computational efficiency reflects both speed and memory usage, which have a trade-off
270 relationship. However, benchmarking the memory usages of SV callers on 20X and 61X tumor
271 samples showed no such relationship (Fig. 5c), which is probably because the memory usage is more
272 dependent on the number of *k*-mers than the sequencing depth. In fact, ETCHING consistently used
273 ~12G RAM, regardless of the size of the input dataset, which is comparable or more efficient than
274 other tools in terms of memory usage. This fixed memory usage is mostly attributable to the size of
275 the PGKN set, which is the least variable. Taken together, these results show that ETCHING is
276 computationally very efficient, yet does not exhibit compromised performance.

277

278 **Discussion**

279 Here, we introduced a high performing and very efficient SV caller, ETCHING, which takes
280 advantage of a scalable PGK set ($>3.9 \times 10^9$ 31-mers). Matched normal samples can extend
281 ETCHING to the PGKN *k*-mer set to enrich reads with somatic variations. *k*-mer counting, and

282 searching for an exact k -mer in the large k -mer set, impose critical challenges on k -mer-based SV
283 callers. ETCHING utilized K-mer Counter (KMC)³⁴ for efficient k -mer counting and employed a
284 parallel roll-encoding method for searching for TS k -mers, allowing a highly efficient k -mer
285 processing method.

286 ETCHING has excellent potential for the prediction of somatic SVs, even without matched
287 normal data. The PGK filter module can remove reads present in pan-genome or containing common
288 variations from tumor sequencing data (Fig. 1b). Although ETCHING may produce FPs, it is still
289 useful in the absence of matched normal data (Fig. 1g; Supplementary Fig. 6). This flexibility will be
290 quite helpful, particularly for clinical sequencing, which often lacks such matched normal data (Fig.
291 4).

292 ETCHING found five additional druggable SV targets (in *ALK*, *NTRK1*, *BRCA2*, *PIK3CA*,
293 and *AKT1*), three of which (*ALK*, *NTRK1*, and *BRCA2*) were validated by PCR analysis, in MM
294 patients who did not carry SV biomarkers. *ALK* amplification is a potential molecular target in several
295 cancers and *ALK* inhibitors could be beneficial to patients carrying such an *ALK* amplification³⁵.
296 Because multiple SV events of DELs, DUPs, and INVs were detected around the *NTRK1* gene in
297 SNUH19_MM01, the two most likely paths for their creation were confirmed by PCR (Figure 3e).
298 Although the *NTRK1-LMNA* fusion is known to be a druggable target, the amplification of 1q23.1,
299 where the *NTRK1* locus resides, has also been proposed as a candidate hotspot in the progression of
300 MM³⁶. Because *BRCA2* loss of function is a known cancer driver, we examined biallelic inactivation
301 of the *BRCA2* gene by searching for somatic or germline SNVs or indels at that locus but confirmed
302 no clinically relevant variations in the other allele.

303 ETCHING can also predict other types of variations, such as germline and *de novo* mutations.
304 With a k -mer set from a reference genome (such as hg19), it can predict germline SVs. If we use k -
305 mers of parental genome sequences, ETCHING can find *de novo* mutations in offspring genomes. The
306 current version of ETCHING predicts FG candidates from DNA sequencing data, but the detection of
307 high-confidence FGs requires transcriptome data, such as RNA-seq. Such detection will be possible,
308 without a need for other FG callers, by using a k -mer set of reference transcriptomes or RNA-seq data

309 from normal samples. Hence, by the selection of an appropriate k -mer set, ETCHING can be a multi-
310 purpose predictor for diverse types of genomic variations and FGs.

311 Although both ETCHING and novoBreak take advantage of TS reads to predict somatic SVs,
312 the main strategy of ETCHING is distinct from that of novoBreak, which collects TS reads by
313 comparing tumor and normal reads after mapping (the mapping-and-filtration approach). Instead,
314 ETCHING uses a filtration-and-mapping approach, which makes ETCHING much faster than
315 novoBreak, by as much as an order of magnitude (Fig. 5; Supplementary Fig. 15). In addition,
316 novoBreak performs a local *de novo* assembly using the resulting TS reads to assemble TS contigs,
317 which is another source of the heavy computational burden. The resulting contigs are aligned to a
318 reference genome to predict SVs and BPs based on the mapping patterns of the contigs. Thus, the risk
319 of misassembly also cannot be neglected. In contrast, ETCHING predicts all possible SVs using split-
320 reads of TS reads and filters FPs by a RF module, achieving a low FP rate.

321 In summary, ETCHING is the fastest method for SV and FG prediction, and this speed has
322 been achieved without compromising its performance or memory usages. We believe that our new
323 approach will not only provide an efficient strategy for predicting various variations in mega-genome
324 projects but will also contribute to real-time clinical applications.

325

326 **Data availability**

327 WGS data from 26 MM samples can be downloaded from the Clinical & Omics Data Archive
328 (CODA; registration number: R002594) of the Korean National Institute of Health. Targeted gene
329 panel sequencing data from reference materials are available at our website
330 (<http://big.hanyang.ac.kr/ETCHING>).

331

332 **Code availability**

333 ETCHING was designed for 64-bit Linux systems. At least 16 GB of RAM is required. We
334 recommend at least 64 GB. All source and binary codes used in the study are available at
335 <http://big.hanyang.ac.kr/ETCHING> and GitHub (<https://github.com/ETCHING-team/ETCHING>).

336

337 **Acknowledgments**

338 We thank all of the BIGLab members, and Professor Sun Kim of Seoul National University for
339 critical reading and comments. This work was supported by the National Research Foundation (NRF)
340 funded by the Ministry of Science & ICT (2014M3C9A3063541 to JWN) and by the Korean Health
341 Technology R&D Project, Ministry of Health and Welfare, Republic of Korea (HI15C3224 to JWN).

342

343 **Ethics declarations**

344 All MM samples used in this study were prepared under the Human Biospecimen Ethics Guidelines
345 and were approved by the Internal Review Board (IRB) of SNUH.

346

347 **Competing interests**

348 None declared

349

350

351 **Figure legends**

352 **Fig. 1.** A schematic overview of ETCHING. **a.** A schematic showing the flow through the ETCHING
353 process, which comprises four stepwise modules (Filter, Caller, Sorter, and Fusion-identifier). **b.** The
354 Filter module collects TS reads containing at least one TS k -mer not present in the k -mer sets (PGK,
355 Normal, or PGKN). **c.** The percentage of TS reads that pass through the PGK, Normal, and PGKN
356 filters. **d.** The mapping patterns of the total tumor reads (unfiltered, gray) and TS reads (filtered, blue)
357 are shown for representative DEL, DUP, INV, and TRA loci via Integrative Genome Viewer. **e.** The
358 mapping times (CPU times) required for the total tumor reads (Unfiltered) and TS reads filtered by
359 PGK, Normal, and PGKN using BWA-MEM. **f.** The total running time (CPU time) of the SV callers.
360 **g.** The precision, recall, and F1-scores of ETCHING with total tumor (Unfiltered) and TS reads
361 collected by PGK, Normal, and PGKN. **(c,e-g)** The analyses were done with nine BRCA WGS
362 datasets. The error bars indicate the first to third quartile range, and the height of the boxes indicate
363 median values.

364

365 **Fig. 2.** Performances of ETCHING and benchmarking SV callers. **a.** PR curves of ETCHING and
366 benchmarking tools on the HCC1395 dataset. The red symbols indicate the points corresponding to
367 optimal parameters. **b.** Precision, recall, and F1-scores of ETCHING and benchmarking tools over
368 sub-sampled data with different sequencing depths from the HCC1395 tumor sample. **c.** Precision,
369 recall, and F1-scores of ETCHING and benchmarking tools for all types of SVs over nine hold-out
370 test datasets of TCGA BRCA samples. Each dot denotes the performance of each tool on a sample.
371 The height of the bar plots indicates the median performance of each tool on nine samples, and the red
372 error bars are the first and third quartiles. **d.** The performances of ETCHING and benchmarking tools
373 on four TCGA THCA samples. Because there were only a few true SVs from each of the samples, we
374 combined them as one value. **e.** The performances of ETCHING and benchmarking tools on MM
375 samples. **d,e.** Otherwise, as in (c).

376

377 **Fig. 3.** Prediction of SVs and FGs by SV callers using MM samples containing known clinical
378 biomarkers and actionable SV targets. **a.** Known clinical SV and FG biomarkers (also known as
379 clinical targets) of MM. The type of SV of known clinical biomarkers are indicated on the appropriate
380 chromosomes. **b.** Summary of manually curated, experiment-supported, and ETCHING-detected SV
381 biomarkers, known MM targets, and actionable targets from OncoKB (tier1, 2, and 3). **c.** ROC curves
382 of ETCHING and benchmarking tools are shown along with accuracies (acc) and F1 scores as an inset.
383 The accuracies and F1 scores were calculated on optimal parameters. **d.** The read-depth landscapes
384 for chromosomes in which clinical biomarkers and targets were found. **e.** Experimental validation of
385 three predicted actionable SV targets by PCR. The blue arrows indicate the expected sizes of the PCR
386 amplicons in the gel images. ‘N’ indicates the normal sample and ‘T’ indicates the tumor sample.
387 (bottom) The dotted lines indicate the junctions formed from tandem DUPs and DELs. The red arrows
388 are the forward and reverse PCR primers.

389

390 **Fig. 4.** SV and FG predictions on targeted gene panel sequencing data. **a.** The TP calls (labeled as
391 ‘Found’ in orange) and false negatives (labeled as ‘Missed’ in gray) of SV callers for cfDNA
392 reference materials – CR, CMM, and MMv2 – with different mutant allele ratios (0.5 to 5.0%; gray to
393 black). CR and CMM include *NCOA4-RET*, *EML4-ALK*, and *CD74-ROSI* FGs, and MMv2 includes
394 *NCOA4-RET* and *TPR-ALK* FGs. **b.** The recall rates of benchmarking SV callers on the reference
395 materials across different mutant allele ratios. **c.** The additional calls in target regions (colors) and
396 non-target regions (gray). **d.** The heatmap summarizes the TPs (labeled as ‘Found’ in orange), false
397 negatives (labeled as ‘Missed’ in gray), and additional calls for 105 cancer panel sequencing datasets.
398 The panels on the right show the total number of TPs and additional calls. The white-to-black gradient
399 indicates the number of additional calls on each SV caller. The color-coded charts (top) indicate
400 cancer types, known alterations, and detection methods. Abbreviations: Diffuse large B-cell
401 lymphoma (DLBCL), desmoplastic small round cell tumor (DSRC), gastrointestinal stromal tumor
402 (GIST), acute lymphoblastic leukemia (ALL), primitive neuroectodermal tumor (PNET), follicular B-

403 cell lymphoma (FL), acute myeloid leukemia (AML), chronic myelogenous leukemia (CML), and
404 lung adenocarcinoma (LA).

405

406 **Fig. 5.** Computational costs of ETCHING and the benchmarking tools. **a.** Stepwise comparison of the
407 CPU times for SV prediction using ETCHING, with reads filtered by PGKN or with unfiltered reads,
408 and using novoBreak. **b.** Algorithmic differences between ETCHING, novoBreak, and others
409 (DELLY, LUMPY, Manta, and SvABA). **c.** RAM usage by the SV callers on TCGA-A2-A04P (20X
410 tumor, 37X normal) and TCGA-A1-A0SM (61X tumor, 31X normal) datasets with 60 threads.

411

412 **METHODS**

413 ***k*-mer counting**

414 An efficient *k*-mer counting tool, KMC, was applied to count all possible *k*-mers (31-mers) from
415 tumor and normal reads. *k*-mer counting can be done with multi-process (MP) computation. The
416 results of *k*-mer counting are summarized in a histogram (Supplementary Fig. 1b) showing the *k*-mer
417 depth (count) on the x-axis and the number of *k*-mers on the y-axis; the *k*-mer frequency shows a
418 bimodal distribution for WGS data. A histogram of error-free *k*-mers is known to be close to a normal
419 (Poisson) distribution, whereas rare *k*-mers, considered to be those with sequencing errors, show an
420 exponentially decreasing curve over low depths. Hence, the local minimum was generally determined
421 to be between *k*-mer depth 3–10, varying with the sequencing depth, quality, and tumor heterogeneity.
422 Therefore, tumor *k*-mers with depth below the local minimum (the cutoff for erroneous *k*-mers) were
423 removed, and the remaining error-free *k*-mers were subjected to the following steps (Supplementary
424 Fig. 1a). For normal *k*-mers, those below *k*-mer depth 2 were removed and the remainder were added
425 to the *k*-mer set (PGKN).

426 For targeted gene panel sequencing data, the local minimum is usually not presented as in
427 WGS data. As the local minimum *k*-mer depth in WGS data is generally observed at a point about 10%
428 of the distribution value at *k*-mer depth 2, we used the point as the local minimum *k*-mer depth in
429 panel sequencing data.

430

431 **Roll-encoding**

432 To efficiently process *k*-mers, we introduced a roll-encoding strategy, which encodes a *k*-mer to a
433 series of 2-bit numbers by our encoding rules: A to 00, C to 01, G to 11, and T to 10. Because the *k*-1
434 nucleotides of the *i*-th and (*i*+1)th *k*-mers overlap, we can obtain the (*i*+1)th encoded *k*-mer simply by
435 sliding a 2-bit number. This approach means that a new 2-bit number is added to the last nucleotide of
436 the (*i*+1)th *k*-mer while the first 2-bit number is removed from the *i*-th encoded *k*-mer (Supplementary
437 Fig. 16a). This procedure is repeated until the end of a read. Our roll-encoding also simultaneously
438 encodes *k*-mer reverse complements. The smaller of the forward- and reverse-encoded values was

439 stored as a canonical encoded k -mer. This roll-encoding method appeared to be faster than methods
440 with conventionally encoded and ordinary (not encoded) k -mers (Supplementary Fig. 16b).

441

442 **The reference and normal k -mer sets**

443 The reference k -mer set, PGK, is a unique set of k -mers from references (10 human genome
444 assemblies; Supplementary Table 1) and those embedding common non-medical (nonpathogenic)
445 SNPs in hg19 (GRCh37.p13) from dbSNP (release number 150). The normal k -mer set is from
446 matched normal input reads. PGKN is a unique set of the PGK and the normal k -mers. For the
447 YH_1.0 genome assembly, which includes uncertain bases, all possible nucleotides were assigned to
448 generate the k -mer set. The reference k -mer set (PGK) is stored as a binary database file for reuse. The
449 PGK binary file can be downloaded from our website (<http://big.hanyang.ac.kr/ETCHING>).

450

451 **Filter module**

452 The saved reference k -mer set (PGK) is loaded to a hash table in the Filter module. If there is a
453 matched normal sample as input, then normal k -mers are added to the k -mer set (PGK + Normal).
454 When tumor sequencing data are used as the input, they are decomposed into tumor k -mers. The
455 tumor k -mers are then searched in the reference k -mer sets (PGK or PGKN). The tumor k -mers
456 present in the reference k -mer set are regarded as reference k -mers; otherwise, they are regarded as TS
457 k -mers and subjected to the following read-collection step. The read-collection step collects TS reads
458 embedding a TS k -mer. To speed up the read-collection step, a multi-processing procedure for
459 simultaneously treating reading, collecting, and writing substeps was implemented (Supplementary
460 Fig. 16c, d).

461

462 **Reduced read mapping**

463 From the total input tumor reads, only TS reads collected through the Filter were mapped to the
464 reference genome (hg19) using BWA-MEM with default parameters. We also used default parameters
465 in read mapping for benchmarking tools.

466

467 **Caller module**

468 After the TS reads are mapped, the Caller module finds BND candidates (BP pairs) by analyzing split
469 reads with supplementary alignment (SA) tags, as follows (Supplementary Fig. 3). We focused on
470 simply clipped pairs only, not on complex or double clipped reads, to reduce FP calls (Supplementary
471 Fig. 3a). First, we defined a BP by its vector or chromosome (or contig/scaffold) name, its clipped
472 position on the chromosome, and its clipped direction (Supplementary Fig. 3b). If a read was clipped
473 in a region that is downstream of the BP, its clipped direction s is denoted as “+”. If a read was
474 clipped in a region that is upstream of the BP, its clipped direction s is denoted as “-”. Thus, reads
475 clipped at a locus can define a BND with a BP pair. A lack of SA tags in a clipped read indicates that
476 there is a single BP that we called as a single-breakend (SND). Once all of the BNDs and SNDs are
477 defined, BNDs are then classified by SV type (such as DEL, DUP, INV, or TRA), with their
478 chromosome, BP position, and clipped direction information (Supplementary Fig. 3c).

479

480 **Sorter module**

481 The Sorter module is a machine learning classifier that removes FP SVs from the Caller module
482 outputs. Because ensemble machines usually show optimal performance in diverse problems, we
483 applied RF (<https://github.com/crfllynn/skranger>), and extreme gradient boosting (XGB,
484 <https://github.com/dmlc/xgboost>) models to this study. To train the models, we randomly selected 31
485 training and 9 hold-out test samples from 55 BRCA samples (Supplementary Table 2; Supplementary
486 Fig. 4a) as follows. We first predicted all possible SVs using five benchmarking SV callers and
487 summarized tumor purities and sequencing depths for all 55 samples. Based on this information, we
488 excluded (1) nine samples that had a low number of predicted SVs (<100) for at least one caller, (2)
489 four samples with too many predicted SVs (>50,000 on average), and (3) two additional samples, one
490 with the highest tumor read depth (93X) and one with the lowest tumor purity (0.474), to avoid
491 extreme cases. From the remaining 40 samples, we randomly selected 31 and 9 samples so that there
492 would be about a 3:1 ratio of SV candidates in the training and hold-out test datasets, respectively

493 (Supplementary Fig. 4a). There were 894,333 and 278,627 SV candidates in 31 training and 9 hold-
494 out samples. For training data, we selected 315,949 SV candidates detected by the Caller module,
495 which were subjected to the training step of the Sorter module.

496 There is no ground truth exhaustively validated by experiments for the TCGA dataset. Thus,
497 we used silver standard SVs detected by multiple SV callers. Of 315,949 SV candidates predicted by
498 the Caller module in 31 training samples, 10,736 SVs were simultaneously predicted by at least three
499 SV callers (Supplementary Fig. 4b). We regarded them as silver standard SVs and the remainder
500 (314,507 SVs) as false (Supplementary Fig. 4b, c). With the true and false SVs, we trained the models
501 with six different features – clipped-read count (*CR*), split-reads count (*SR*), supporting paired-end
502 read count (*PE*), average mapping quality (*MQ*), depth difference (*DD*), and total length of clipped
503 bases (*TC*) (see Supplementary Note for more details). Our training procedure consists of an outer 10-
504 fold cross-validation (*CV*) loop for training and an inner 10-fold *CV* loop for model selection
505 (Supplementary Fig. 4d). The SVs in the training samples were split evenly into eleven sets, including
506 ten outer-training sets (*TR_out*) and one validation set (*VA*). During the outer 10-fold *CV*, a test set is
507 selected (*TE*) from *TR_out*, and the remaining nine sets were subjected to inner-training (*TR_in*). The
508 model selection process was done by inner 10-fold *CV* using *TR_in*, which was evaluated on *TE*. The
509 procedure was iteratively performed through an outer 10-fold *CV* loop. A final model was obtained
510 by averaging ten trained models. We validated the final model on the *VA*.

511 We then searched the optimal classification cutoffs of RF and XGB scores using the *VA* set
512 (Supplementary Fig. 4f). F1-scores of RF (or XGB) showed robust performances in the range from
513 0.2 to 0.8 (from 0.05 to 0.95 for XGB). We used RF as default ML module in this study.

514

515 **Parameter optimization for benchmarking SV callers**

516 ETCHING was benchmarked to the popular, high performing SV callers DELLY, LUMPY, Manta,
517 SvABA, and novoBreak over WGS data, cfDNA reference materials, and targeted gene panel
518 sequencing data from tumor samples. We also benchmarked Breakmer for cfDNA reference
519 materials and targeted gene panel sequencing data.

520 For a fair comparison on WGS datasets, we searched optimal parameters of benchmarking
521 tools corresponding to the nearest points to the perfect performance (where precision and recall rates
522 are 100%) over PR-curves on HCC1395 data (Fig. 2a). The point minimizes the distance,
523 $\sqrt{(1-P)^2 + (1-R)^2}$, to (1,1) on given PR-curve, where P and R refer to precision and recall,
524 respectively. DELLY's optimal parameter was near its default parameter (-a 0.2), LUMPY was -m 12
525 option, and Manta was *minEdgeObservations* = 12 and *minCandidateSpanningCount* = 12. For
526 SvABA, log-odd ratios of real and artifact variants ≥ 32 was the optimal one. novoBreak's PR curve
527 was closest to the corner for its statistical quality score ≥ 40 . The statistical quality score is defined
528 as $-10 \log_{10} \frac{P(D|reference\ alleles\ or\ germline\ variations)}{P(D|somatic\ variations)}$, where D is the number of read counts
529 supporting each variation or reference allele.

530 For cfDNA reference materials and targeted gene panel sequencing data, all tools were
531 applied with default parameter sets. Manta was run with --tumorBam --exome options.

532

533 Evaluation metrics

- 534 True positive (TP): Predicting true SVs (or biomarkers) as positive.
535 False negative (FN): Predicting true SVs (or biomarkers) as negative.
536 False positive (FP): Predicting false SVs (or biomarkers) as positive.
537 True negative (TN): Predicting false SVs (or biomarkers) as negative.

538 Given the TP, FN, FP, and TN metrics, the recall, sensitivity, precision, specificity, F1-score, and
539 accuracy are estimated as follows:

540 • $Recall = Sens = \frac{TP}{TP+FN}$

541 • $Precision = \frac{TP}{TP+FP}$

542 • $Specificity = \frac{TN}{TN+FP}$

543 • $F1 = \frac{2\ Recall\ Precision}{Recall + Precision}$

544 • $Accuracy = \frac{TP+TN}{TP+TN+FP+FN}$

545

546 **Public WGS datasets**

547 55 BRCA WGS datasets and 4 THCA WGS datasets were downloaded from TCGA
548 (<https://cancergenome.nih.gov>).

549

550 **MM WGS data**

551 Tumor cells were collected from bone marrow using CD138+ MACS sorting (Miltenyi Biotec,
552 Auburn, CA) and DNA was extracted from the tumor cells for WGS library preparation. For matched
553 normal samples, DNA was extracted from patients' saliva with RNase treatment. Sequencing libraries
554 were generated using a TruSeq nano DNA library prep kit (Illumina, San Diego, CA) following the
555 manufacturer's recommendations and sheared DNA fragments were end-repaired and size-selected to
556 obtain DNA fragments around 350bp. Following PCR amplification, the DNA libraries were
557 sequenced using the HiSeq™ X platform (Illumina). The 26 MM WGS datasets were produced and
558 deposited in the CODA (registration number: R002594) of the Korean National Institute of Health.
559 The study was approved by the Internal Review Board of Seoul National University Hospital (H-
560 1103-004-353).

561

562 **FISH and karyotyping**

563 Cytogenetic studies were performed at SNUH. Unstimulated bone marrow cells obtained at MM
564 diagnosis were cultured for 24h; then, karyotypes were analyzed using the standard G-banding
565 technique. The karyotypes were constructed and chromosomal abnormalities were reported according
566 to the International System for Human Cytogenetic Nomenclature³⁷. Interphase FISH was performed
567 on myeloma cells from the bone marrow samples obtained at diagnosis according to the probe
568 manufacturer's instructions. Seven commercially available FISH probe sets were used. These
569 included *IGH* dual-color, break-apart rearrangement probe; *TP53* SpectrumOrange probe; *RBI*
570 D13S25 (13q14.3) SpectrumOrange probe; *IGH-FGFR3* dual-color, dual-fusion translocation probe;
571 1q21 SpectrumGreen probe; and *p16* (9p21, *CDKN2A*), SpectrumOrange/CEP9 SpectrumGreen probe

572 (Abbott Diagnostics, Abbott Park, IL). The FISH experiments were performed on 26 MM specimens.

573 The FISH probe sequences are summarized in Supplementary Table 6.

574

575 **PCR validation of actionable targets**

576 PCR amplification was performed using the primer sets listed in Supplementary Table 9. Targets were

577 amplified using primers designed in the flanking region of the junction. GAPDH was used as a

578 control for assessing the PCR efficiency and for subsequent analysis by agarose gel electrophoresis.

579

580 **Manual curation of biomarkers and actionable targets in MM samples**

581 The SV biomarkers and actionable targets were manually curated with all mapped reads. The

582 candidate DELs and DUPs were checked by considering minor allele frequencies and read depth

583 changes across chromosomes (Supplementary Fig. 11a,b), remaining focal DELs and DUPs. For *IGH*-

584 associated TRA, candidate TRAs with which >10 paired-reads (mapping quality ≥ 20) are connected

585 between *IGH* locus (14q32) and other loci in tumor but not in normal were selected as true somatic

586 TRAs (Supplementary Fig. 11c,e). The candidates with the connection both in tumor and normal were

587 considered as germline TRAs (Supplementary Fig. 11d). The read depth, minor allele frequency,

588 discordant paired-read data to inspect true SVs during manual curation were summarized in

589 Supplementary Material.

590

591 **Cell-free DNA reference materials**

592 Targeted sequencing data from cfDNA reference materials (SeraCare, Milford, MA) were generated.

593 DNA libraries were prepared using a KAPA Hyper Prep kit (Kapa Biosystems, Woburn, MA) as

594 described previously. Hybrid selection for target enrichment was performed using customized baits

595 targeting 38 cancer-related genes. After hybrid selection, the libraries were pooled, amplified, purified,

596 quantified, and then subjected to cluster amplification according to the manufacturer's protocol

597 (Illumina). Flow cells were sequenced in the 150bp paired-end mode using a NextSeq 500/550 High

598 Output Kit v2.5 (Illumina). The mean target coverage was 2023X. Two kinds of DNA mixtures, with

599 the frequency of variant alleles ranging from 0.5–5.0% (CMM and MMv2), and a plasma-like DNA
600 mixture, with the frequency of variant alleles ranging from 0.5–2.5% (CR), were generated along with
601 WT DNA (Supplementary Table 10). The WT material was used as the matched normal. Note that
602 DELLY displayed a low recall in normal-matched case, since it excessively removed SV calls using
603 matched-normal data in the filtration step (Supplementary Fig. 12). The BreaKmer tool was excluded
604 from this analysis because it failed to call variants from any sample, presumably because its approach
605 is not feasible for such low allele frequencies.

606

607 **Cancer panel datasets**

608 For cancer panel data of BreaKmer, we downloaded hybrid capture targeted gene panel data (110
609 replicates from 38 cancer samples). Because the normal samples that were provided are not matched
610 to the cancer samples, they were excluded from the analysis. One sample with three replicates was
611 also excluded from this analysis, since it was marked as non-cancer sample rather than diagnosed
612 cancer type (SRR1304190-2). Two datasets (SRR1304204, SRR1304210) failed to run in at least one
613 benchmarking tool, so the remaining 105 replicates from 37 sample (216X mean coverage of the
614 targets) were analyzed. Because the sample labels in SRA are inconsistent with those in BreaKmer
615 paper, we used ones described in the paper.

616 To reproduce the results of previous BreaKmer study, we needed to install the same version
617 of BreaKmer with detailed information of target bait. However, we failed to install the same version
618 of BreaKmer in their publication, and the bait information was also unavailable. Hence, we tested two
619 other releases, v0.0.4 and v0.0.6. The version v0.0.6 found 78 true SVs out of 105 and only 487
620 additional calls, while v0.0.4 found 70 true SVs with 17,738 additional calls. Thus, we selected v0.0.6
621 for comparison. To substitute the missing target bait information, we used the genomic coordinates of
622 target gene regions.

623 In case of novoBreak, it requires normal sequencing data. However, there is no matched
624 normal samples in the panel data. For the reasons, we simulated WGS reads (30X coverage) from
625 hg19 using an in-house script for novoBreak.

626 *FLT3* indel (30–73bp) and *KIT* deletion (48bp) were included in the list of known target
627 alterations. As SvABA separately reports indels as output, we used SvABA high-confidence indel
628 report along with its SVs. However, although Manta also reports indels, we did not use them because
629 they are unfiltered candidates.

630

631 **Supplementary information**

632 **Supplementary Note:**

633 Graph theory presentation for SV analysis

634 Six features for machine learning

635 Commands for benchmarking tools

636 High quality SVs

637

638 **Supplementary Fig. 1. a.** Detailed workflow of the ETCHING pipeline. MP and SP indicate multiple
639 and single processing, respectively. **b.** A representative k -mer distribution of WGS data.

640

641 **Supplementary Fig. 2. a.** The size of the unique set containing the hg19 and PGK k -mers. **b.** The size
642 of the unique set containing the k -mers not present in hg19.

643

644 **Supplementary Fig. 3. a.** The Caller module uses simply clipped reads (left side) but excludes reads
645 that make complex clipped pairs (right side). **b.** Each BND is a pair of BPs, *i.e.* (BP_i, BP_j) . An SND is
646 a single-BND consisting of one BP with a dangling point, *i.e.* (BP_i, ϕ) . If a BP displays reads clipped
647 in a direction at position x on chromosome c , we define that BP as a node (c, x, s) , where s indicates
648 its clipped direction (+1 or -1). **c.** Classification of SV types. For a BND (BP_i, BP_j) , x_i and x_j are the
649 positions on chromosome c , and s_i and s_j are the clipped directions of each BP. The table on the right
650 side shows the classification criteria for SVs.

651

652 **Supplementary Fig. 4. a.** A flowchart for selecting training data from BRCA samples. **b.** A Venn
653 diagram of SVs predicted by ETCHING and other tools in the training set. **c.** The numbers of SVs
654 predicted by multiple callers were tallied in a histogram. The number of SVs are indicated on the y-
655 axis and the number of tools that predicted the corresponding SVs are indicated on the x-axis. The
656 vertical line denotes the cutoff for selecting silver standard SVs. **d.** A schematic workflow for training
657 machine learning modules. **e.** Training and validation results of machine learning. **f.** Optimized
658 cutoffs of machine learning methods. We set the optimized cutoff to 0.4 for RF and XGB.

659

660 **Supplementary Fig. 5.** The wall-clock times used by ETCHING and other tools on nine hold-out
661 BRCA samples, which were measured using 30 threads.

662

663 **Supplementary Fig. 6.** The effectiveness of the Filter module on HCC1395 data. **a.** The percentage
664 of TS reads that passed the PGK, Normal, and PGKN filters. **b.** The precision, recall, and F1-scores of
665 the ETCHING results from total reads (unfiltered) and TS reads collected by the PGK, Normal, and
666 PGKN filters.

667

668 **Supplementary Fig. 7.** Benchmarking results on THCA samples by SV type.

669

670 **Supplementary Fig. 8.** Strategies for HQ SV detection using tumor (HCC1395) and normal
671 sequencing data (HCC1395 BL). **a.** The landscape of the depth change within the HQ DELs and HQ
672 DUPs. **b.** The landscape of the discordant read-pairs connecting BPs within the HQ INVs and HQ
673 TRAs. (**c** and **d**). The ROC curve (left) and the density distribution (right) for setting the cutoff using
674 the depth difference score (DS) of DELs (**c**) and DUPs (**d**). **e.** The ROC curve for setting the cutoff
675 using the connected-pair score (CS) of the HQ INVs and TRAs. **f.** Bar plot showing the count of HQ
676 SVs and all SVs. DS and CS are defined in Supplementary Note.

677

678 **Supplementary Fig. 9.** Benchmarking results on BRCA, MM, and THCA samples with HQ SVs.

679

680 **Supplementary Fig. 10.** Benchmarking results on MM samples by SV type.

681

682 **Supplementary Fig. 11.** Rational for SV manual curation **a.** *RBI* biomarker shown with unbalanced
683 minor allele frequency (MAF) and read depth change on Chr13 in MM20. **b.** *AKT1* DUP locus shown
684 with unbalanced MAF and read depth change on Chr14 in MM22. **c.** Discordant paired-reads
685 connected between *IGH* locus and Chr11 of MM6. The blue dot near 69M indicates a TRA,
686 t(11;14)(q13;q32), including *CCND1* gene in tumor. **d.** Germline TRA with discordant paired-reads
687 both in tumor and normal. **e.** An instance of manually curated TRAs is shown in IGV.

688

689 **Supplementary Fig. 12.** SV and FG prediction on targeted gene panel sequencing data paired with
690 sequencing data from WT alleles (regarded as matched normal). Otherwise, as in Fig. 4a–c.

691

692 **Supplementary Fig. 13.** Indels associated with *FLT3* in eight different samples. The index numbers
693 are the same ones in Fig. 4d.

694

695 **Supplementary Fig. 14.** The wall-clock times used by ETCHING, ETCHING without filter
696 (unfiltered), and novoBreak on HCC1395 data on 30 threads.

697

698 **Supplementary Fig. 15.** The wall-clock times used by ETCHING with different thread numbers
699 ranging from 5 to 50.

700

701 **Supplementary Fig. 16. a.** Schematic of the roll-encoding algorithm for processing *k*-mers. As a *k*-
702 mer window slides, it updates an encoded value using our encoding rule. **b.** The computing costs of
703 the Read-collector using the conventional encoding method, ordinary *k*-mers, and roll-encoding
704 methods on tumor (46X) and normal (31X) WGS data with 30 threads. **c.** A schematic workflow of

705 parallel computing for read collection. **d.** Data from the read collection step are processed by parallel
706 computing.

707

708 **Supplementary Tables**

709 **Supplementary Table 1.** Reference genomes and dbSNP used in PGK

710 **Supplementary Table 2.** BRCA samples

711 **Supplementary Table 3.** THCA samples

712 **Supplementary Table 4.** MM samples

713 **Supplementary Table 5.** FISH and karyotype in MM samples

714 **Supplementary Table 6.** FISH probe sets

715 **Supplementary Table 7.** Manually curated partner BPs of *IGH* TRAs.

716 **Supplementary Table 8.** BreaKmer panel data

717 **Supplementary Table 9.** PCR primer sets

718 **Supplementary Table 10.** Reference material

719

720 **References**

721 1. Chiang, C. et al. The impact of structural variation on human gene expression. *Nat Genet* **49**,
722 692-699 (2017).

723 2. Sharp, A.J., Cheng, Z. & Eichler, E.E. Structural variation of the human genome. *Annu Rev*
724 *Genomics Hum Genet* **7**, 407-442 (2006).

725 3. Mitelman, F., Johansson, B. & Mertens, F. The impact of translocations and gene fusions on
726 cancer causation. *Nat Rev Cancer* **7**, 233-245 (2007).

727 4. Stratton, M.R., Campbell, P.J. & Futreal, P.A. The cancer genome. *Nature* **458**, 719-724
728 (2009).

729 5. Meyerson, M., Gabriel, S. & Getz, G. Advances in understanding cancer genomes through
730 second-generation sequencing. *Nature Reviews Genetics* **11**, 685-696 (2010).

- 731 6. Stankiewicz, P. & Lupski, J.R. Structural variation in the human genome and its role in
732 disease. *Annu Rev Med* **61**, 437-455 (2010).
- 733 7. Macintyre, G., Ylstra, B. & Brenton, J.D. Sequencing Structural Variants in Cancer for
734 Precision Therapeutics. *Trends Genet* **32**, 530-542 (2016).
- 735 8. Di Fiore, P.P. et al. erbB-2 is a potent oncogene when overexpressed in NIH/3T3 cells.
736 *Science* **237**, 178-182 (1987).
- 737 9. Slamon, D.J. et al. Human breast cancer: correlation of relapse and survival with
738 amplification of the HER-2/neu oncogene. *Science* **235**, 177-182 (1987).
- 739 10. Soda, M. et al. Identification of the transforming EML4-ALK fusion gene in non-small-cell
740 lung cancer. *Nature* **448**, 561-566 (2007).
- 741 11. Lugo, T.G., Pendergast, A.M., Muller, A.J. & Witte, O.N. Tyrosine kinase activity and
742 transformation potency of bcr-abl oncogene products. *Science* **247**, 1079-1082 (1990).
- 743 12. Abyzov, A., Urban, A.E., Snyder, M. & Gerstein, M. CNVnator: an approach to discover,
744 genotype, and characterize typical and atypical CNVs from family and population genome
745 sequencing. *Genome Res* **21**, 974-984 (2011).
- 746 13. Chen, K. et al. BreakDancer: an algorithm for high-resolution mapping of genomic structural
747 variation. *Nat Methods* **6**, 677-681 (2009).
- 748 14. Ye, K., Schulz, M.H., Long, Q., Apweiler, R. & Ning, Z. Pindel: a pattern growth approach to
749 detect break points of large deletions and medium sized insertions from paired-end short
750 reads. *Bioinformatics* **25**, 2865-2871 (2009).
- 751 15. Wang, J. et al. CREST maps somatic structural variation in cancer genomes with base-pair
752 resolution. *Nat Methods* **8**, 652-654 (2011).
- 753 16. Schroder, J. et al. Socrates: identification of genomic rearrangements in tumour genomes by
754 re-aligning soft clipped reads. *Bioinformatics* **30**, 1064-1072 (2014).
- 755 17. Rausch, T. et al. DELLY: structural variant discovery by integrated paired-end and split-read
756 analysis. *Bioinformatics* **28**, i333-i339 (2012).

- 757 18. Yang, L. et al. Diverse mechanisms of somatic structural variations in human cancer genomes.
758 *Cell* **153**, 919-929 (2013).
- 759 19. Layer, R.M., Chiang, C., Quinlan, A.R. & Hall, I.M. LUMPY: a probabilistic framework for
760 structural variant discovery. *Genome Biol* **15**, R84 (2014).
- 761 20. Chen, X. et al. Manta: rapid detection of structural variants and indels for germline and cancer
762 sequencing applications. *Bioinformatics* **32**, 1220-1222 (2016).
- 763 21. Wala, J.A. et al. SvABA: genome-wide detection of structural variants and indels by local
764 assembly. *Genome Res* **28**, 581-591 (2018).
- 765 22. Chong, Z. et al. novoBreak: local assembly for breakpoint detection in cancer genomes. *Nat*
766 *Methods* **14**, 65-67 (2017).
- 767 23. Consortium, I.T.P.-C.A.o.W.G. Pan-cancer analysis of whole genomes. *Nature* **578**, 82-93
768 (2020).
- 769 24. Collins, R.L. et al. A structural variation reference for medical and population genetics.
770 *Nature* **581**, 444-451 (2020).
- 771 25. Cameron, D.L., Di Stefano, L. & Papenfuss, A.T. Comprehensive evaluation and
772 characterisation of short read general-purpose structural variant calling software. *Nat*
773 *Commun* **10**, 3240 (2019).
- 774 26. Gong, T., Hayes, V.M. & Chan, E.K.F. Detection of somatic structural variants from short-
775 read next-generation sequencing data. *Brief Bioinform* (2020).
- 776 27. Zhang, J. et al. INTEGRATE: gene fusion discovery using whole genome and transcriptome
777 data. *Genome Res* **26**, 108-118 (2016).
- 778 28. Li, H. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
779 *arXiv:1303.3997* (2013).
- 780 29. Wright, M.N. & Ziegler, A. ranger: A Fast Implementation of Random Forests for High
781 Dimensional Data in C++ and R. *Journal of Statistical Software* **77** (2017).
- 782 30. Kosugi, S. et al. Comprehensive evaluation of structural variation detection algorithms for
783 whole genome sequencing. *Genome Biol* **20**, 117 (2019).

- 784 31. Avet-Loiseau, H. et al. High incidence of translocations t(11;14)(q13;q32) and
785 t(4;14)(p16;q32) in patients with plasma cell malignancies. *Cancer Res* **58**, 5640-5645 (1998).
- 786 32. Chakravarty, D. et al. OncoKB: A Precision Oncology Knowledge Base. *JCO Precis Oncol*
787 **2017** (2017).
- 788 33. Abo, R.P. et al. BreaKmer: detection of structural variation in targeted massively parallel
789 sequencing data using kmers. *Nucleic Acids Res* **43**, e19 (2015).
- 790 34. Kokot, M., Dlugosz, M. & Deorowicz, S. KMC 3: counting and manipulating k-mer statistics.
791 *Bioinformatics* **33**, 2759-2761 (2017).
- 792 35. Zito Marino, F. et al. A new look at the ALK gene in cancer: copy number gain and
793 amplification. *Expert Rev Anticancer Ther* **16**, 493-502 (2016).
- 794 36. Pasini, L. et al. TrkA is amplified in malignant melanoma patients and induces an anti-
795 proliferative response in cell lines. *BMC Cancer* **15**, 777 (2015).
- 796 37. Slovak, M. & Campbell, L. International System of Human Cytogenetic Nomenclature. *ISCN*,
797 *S Karger AG, Basel, Switzerland* (2009).
- 798

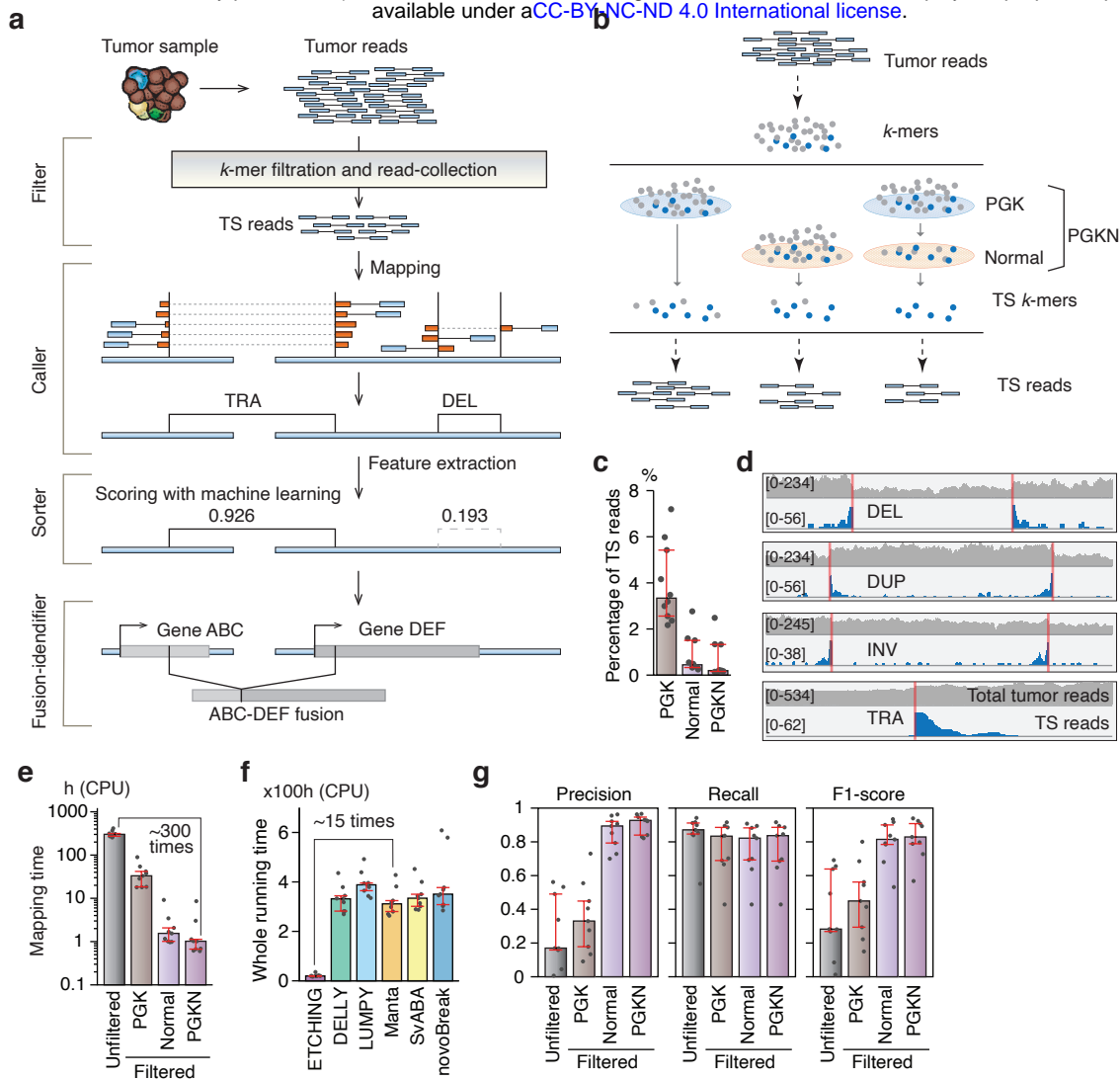


Fig. 2

bioRxiv preprint doi: <https://doi.org/10.1101/2020.10.25.354456>; this version posted October 26, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

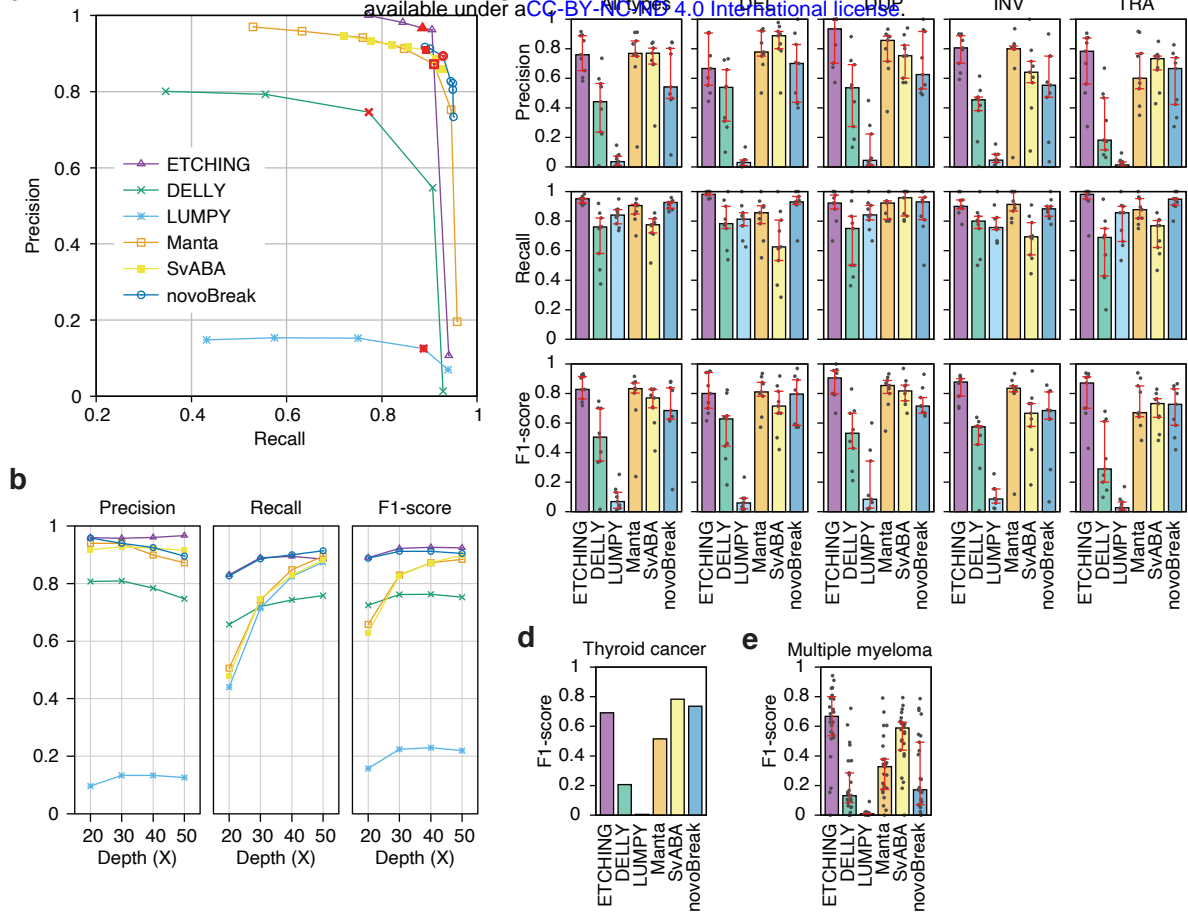


Fig. 3

bioRxiv preprint doi: <https://doi.org/10.1101/2020.10.25.354456>; this version posted October 26, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

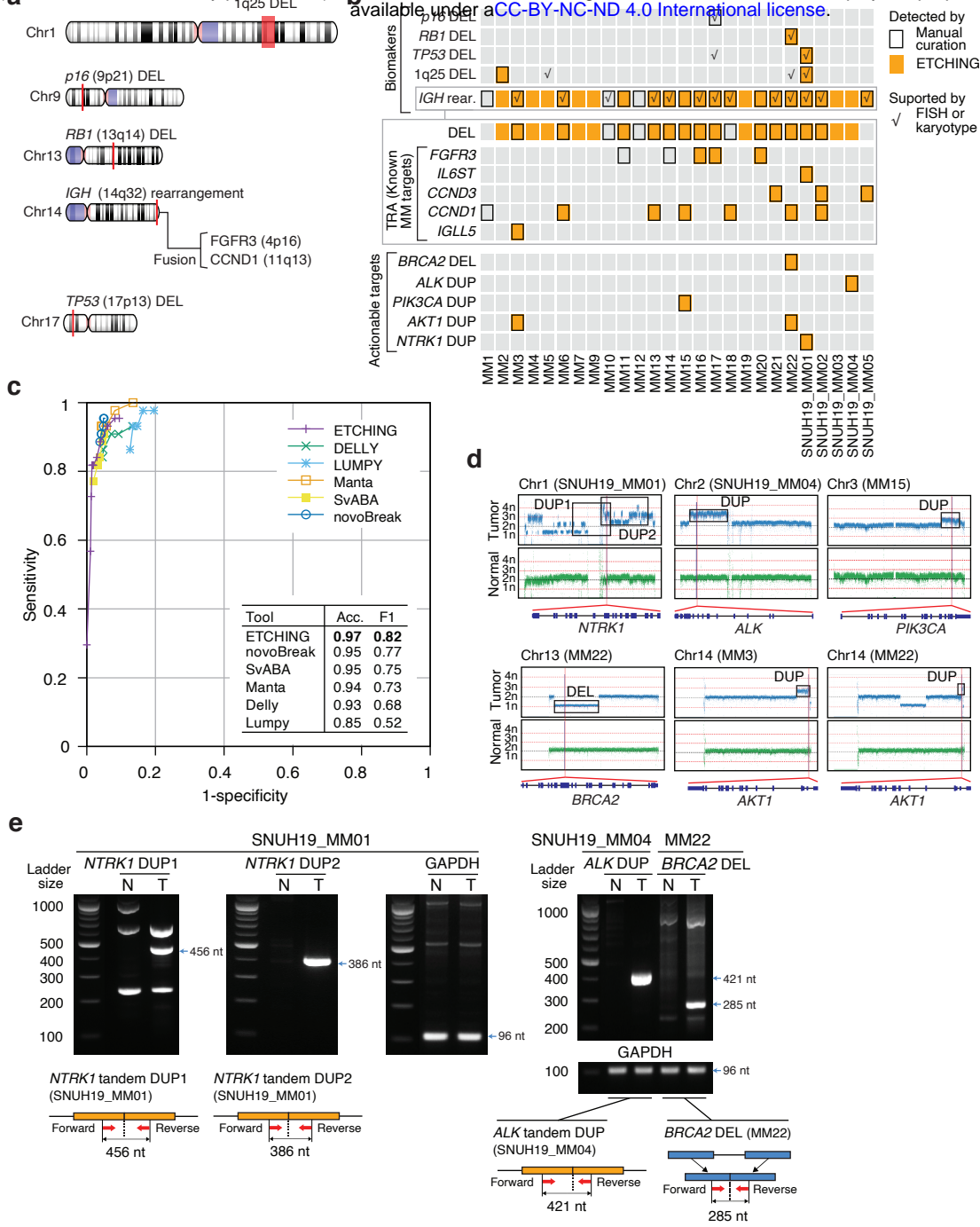


Fig 4 preprint doi: <https://doi.org/10.1101/2020.10.25.354456>; this version posted October 26, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](#).

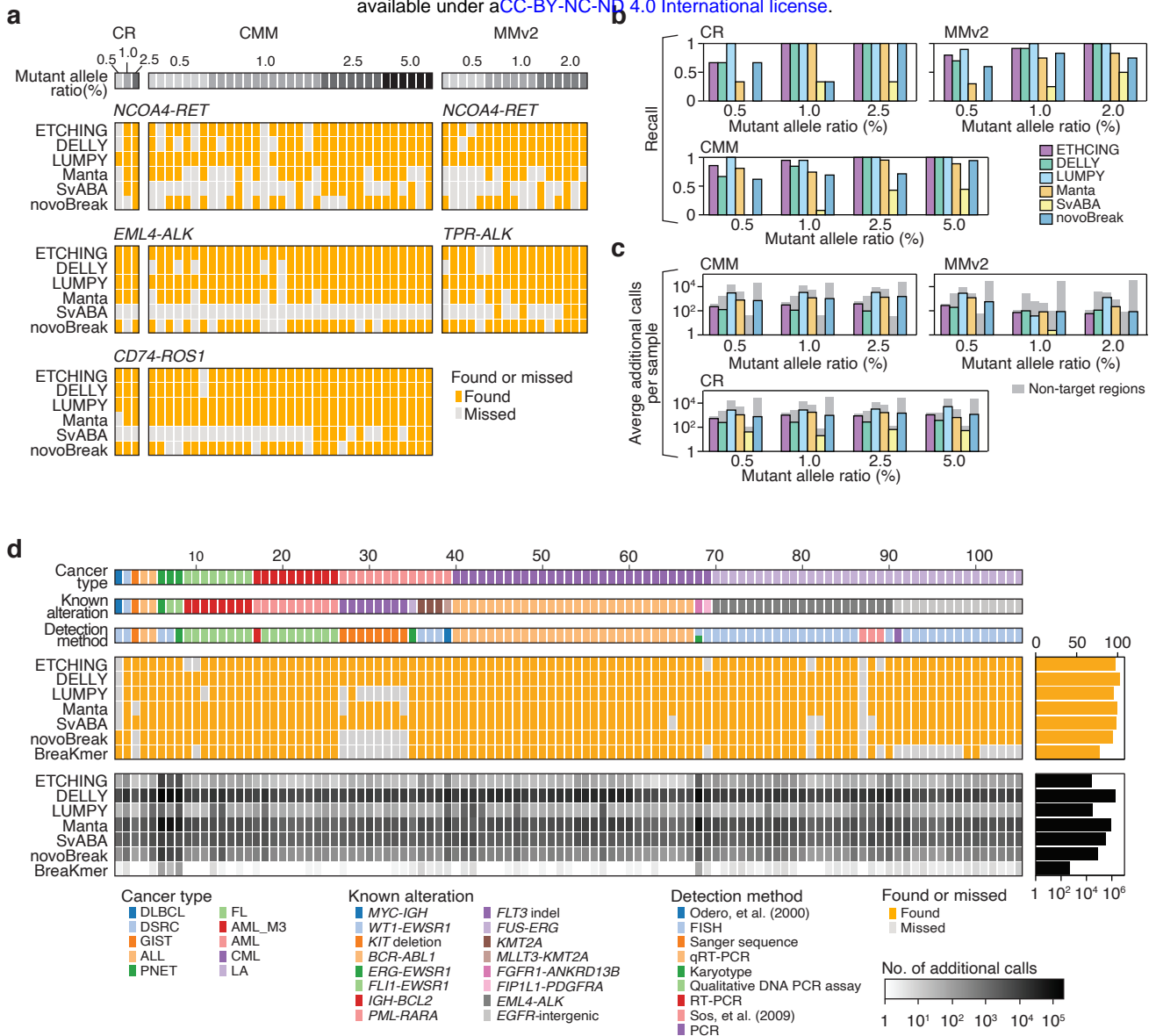


Fig. 5

