1    **BreakCA, a method to discover indels using ChIP-seq and ATAC-seq reads, finds recurrent indels**

2    **in regulatory regions of neuroblastoma genomes**

3    Arko Sen[1*], Sélène T. Tyndale[1], Yi Fu[2,3], Galina Erikson[3], Graham McVicker[1*]

4    1. Integrative Biology Laboratory, Salk Institute for Biological Studies, 10010 N. Torrey Pines Road, La

5    Jolla, CA 92037, USA

6    2. UC San Diego, 9500 Gilman Dr. La Jolla, CA 92093

7    3. Razavi Newman Integrative Genomics and Bioinformatics Core, Salk Institute for Biological Studies,

8    10010 N. Torrey Pines Road, La Jolla, CA 92037, USA

9    *Corresponding Authors: AS: asen@salk.edu, GM: gmcvicker@salk.edu

10    **Abstract**

11    Most known cancer driver mutations are within protein coding regions of the genome, however, there are

12    several important examples of oncogenic non-coding regulatory mutations. We developed a method to

13    identify insertions and deletions (indels) in regulatory regions using aligned reads from chromatin

14    immunoprecipitation followed by sequencing (ChIP-seq) or the assay for transposase-accessible

15    chromatin (ATAC-seq). Our method, which we call BreakCA for Breaks in Chromatin Accessible

16    regions, allows non-coding indels to be discovered in the absence of whole genome sequencing data, out-

17    performs popular variant callers such as the GATK-HaplotypeCaller and VarScan2, and detects known

18    oncogenic regulatory mutations in T-cell acute lymphoblastic leukemia cell lines. We apply BreakCA to

19    identify indels in H3K27ac ChIP-seq peaks in 23 neuroblastoma cell lines and, after removing common

20    germline variants, we identify 23 rare germline or somatic indels that occur in multiple neuroblastoma

21    cell lines. Among them, 4 indels are candidate oncogenic drivers that are present in 4 or 5 cell lines,

22    absent from the genome aggregation database of over 15,000 whole genome sequences, and within the

23    promoters or first introns of known genes (*PHF21A*, *ADAMTS19*, *GPR85* and *RALGDS*). In addition, we

24    observe a rare 7bp germline deletion in two cell lines, which is associated with high expression of the

25    histone demethylase *KDM5B*. Overexpression of *KDM5B* is prognostic for many cancers and further

26    characterization of this indel as a potential oncogenic risk factor is therefore warranted.

27    **Introduction**

28    Several non-coding mutations are known to be important oncogenic drivers. For example, mutations

29    within the promoter of *TERT* are extremely common and cause its overexpression in numerous cancers[1-

30    3] and 2-12bp insertions create new enhancer sequences that drive overexpression of *TAL1* in 4-6% of T-

31    cell acute lymphoblastic leukemias (T-ALLs)[4]. Genome-wide scans for recurrent non-coding mutations

32    have found a handful of additional candidates including recurrent mutations in regulatory regions

33    upstream of *PLEKHS1*, *WDR74*, and *SDHD*[5]. A somatic mutation screen using WGS from chronic

34    lymphocytic leukemia patients identified recurrent T>C mutations in the 3'UTR of *NOTCH1* which cause

35    it to be aberrantly spliced, as well as clustered mutations across multiple patients in an enhancer region

36    for *PAX5*[6]. Non-coding variants can also reposition regulatory sequences so that they activate

37    oncogenes[4, 7, 8] and some non-coding germline polymorphisms that disrupt factor binding motifs are

38    associated with cancer risk. A single-nucleotide polymorphism (SNP) upstream of *MYC* impacts binding

39    of the YY1 transcription factor and is associated prostate cancer[9], SNPs in OCT1/RUNX2 and C/EBPβ

40    binding sites near *FGFR2* modulate its expression and are associated with breast cancer[10], and a SNP

41    that disrupts a GATA3 binding site in an *LMO1* enhancer is associated with neuroblastoma[11].

42    While the cost of whole-genome sequencing (WGS) has decreased dramatically, it remains expensive for

43    large panels of individuals and functional interpretation of non-coding mutations is difficult. One way to

44    overcome these challenges is to identify genetic variants using data from experiments such as chromatin

45    immunoprecipitation followed by sequencing (ChIP-seq) and the Assay for Transposase-Accessible

46    Chromatin (ATAC-seq). These experiments generate sequence reads from regulatory regions of the

47    genome, which can potentially be used to identify non-coding driver mutations in cancer samples.

2

48    Here we describe a new method to identify indels from ChIP-seq and ATAC-seq reads, which we call

49    BreakCA, for "Breaks in Chromatin Accessible" regions. We assess the performance of BreakCA on

50    ATAC-seq and ChIP-seq reads from the GM12878 lymphoblastoid cell line and the Jurkat T-ALL cell

51    line and verify that BreakCA detects known oncogenic indels that create enhancers for the *TAL1* and

52    *LMO2* genes[4, 12]. We then apply BreakCA to H3K27ac ChIP-seq data from 23 neuroblastoma cell

53    lines. After filtering the indels using a large database of known germline variants, we identify recurrent

54    rare germline or somatic indels that may be oncogenic drivers for neuroblastoma.

55    **Results**

56    **Detecting indels with BreakCA**

57    We hypothesized that aligned sequences from ChIP-seq and ATAC-seq experiments can be used to

58    identify indels in regulatory regions of the genome in the absence of whole-genome sequencing (WGS).

59    To test this hypothesis, we developed a method to detect indels by exploiting properties of mapped reads

60    such as gaps in alignments (i.e. insertions or deletions) and clipping at read ends (Fig 1a). Our method,

61    which we call BreakCA for "Breaks in Chromatin Accessible" regions, collects 16 features from mapped

62    reads and uses a random forest to identify 20bp windows that contain indels.

63

64    To train BreakCA and assess its performance, we created separate training and test datasets from 50bp

65    paired-end ATAC-seq data and 50bp single-end H3K27ac ChIP-seq data from the GM12878

66    lymphoblastoid cell line[13, 14]. To label testable windows within ChIP-seq and ATAC-seq peaks as

67    "true" or "false" we used indel calls from the Platinum Genomes (PG) project as known positives[15].

68    After training the random forest on the GM12878 training dataset, we evaluated its performance on the

69    test dataset and compared it to two popular variant callers: VarScan2[16] and the GATK-

70    HaplotypeCaller[17].

71

72    We quantified overall performance using the area under precision-recall curves (Fig 1b,c) and found that

73    BreakCA (prAUC=0.70) performs better than VarScan2 (prAUC=0.48) and comparably to the GATK-

74    HaplotypeCaller (prAUC=0.67) for paired-end ATAC-seq data. For single-end ChIP-seq data BreakCA

75    (prAUC=0.54) performs better than VarScan2 (prAUC=0.33) but worse than the GATK-HaplotypeCaller

76    (prAUC=0.60). An important advantage of BreakCA is that additional features, including output from

77    other variant callers, can be easily added to improve its performance. We added the Quality Depth (QD)

78    reported by the GATK-HaplotypeCaller as a feature for BreakCA and observed substantial improvements

79    in the prAUC for both the ChIP-seq (prAUC=0.69) and ATAC-seq (prAUC=0.79) datasets such that its

80    performance was substantially better than both GATK and VarScan2. We call this version of our method

81    BreakCA+QD.

82

83    To test the performance of BreakCA on a cancer cell line that was not used for training, we used paired-

84    end ATAC-seq and H3K27ac ChIP-seq data from the Jurkat T-ALL cell line and obtained WGS data for

85    the same cell line from a published study[18]. Since there is no gold-standard set of indel calls for this cell

86    line, we used indels identified by the GATK-HaplotypeCaller run on WGS data as our "ground truth".

87    The performance of BreakCA (prAUC= 0.58) on the Jurkat ATAC-seq data was better than both

88    VarScan2 (prAUC=0.41) and GATK-HaplotypeCaller (prAUC=0.53) and improves further when

89    information from GATK is included (prAUC=0.62). For single-end ChIP-seq, while BreakCA

90    (prAUC=0.50) out-performed VarScan2 (prAUC=0.34), its overall performance was comparable to

91    GATK-HaplotypeCaller (prAUC= 0.46) and we observed an improvement in prAUC after adding QD

92    from GATK (prAUC=0.54). While all methods appear to perform worse on the Jurkat datasets compared

93    to the GM12878 datasets (Fig 1c), it is important to note that the performance is greatly underestimated

94    due to inaccuracies in the Jurkat "ground truth" dataset (high false-negative rates) compared to the high-

95    quality platinum genomes dataset.

96

4

97    We compared BreakCA to an orthogonal method that was recently developed by Abraham et al. to

98    identify small insertions from ChIP-seq reads[12]. This method, which we refer to as Abraham's Insertion

99    Detection Pipeline (AIDP), assembles contigs from ChIP-seq reads that fail to map the reference genome.

100   AIDP was previously applied to H3K27ac ChIP-seq data from the Jurkat Cell line and we ran BreakCA

101   on the same dataset. We compared the insertion calls from AIDP and BreakCA to those from the GATK-

102   HaplotypeCaller, which was run on WGS data (Fig 1f). While BreakCA detects a much smaller number

103   of insertions (n=1372 for BreakCA compared to n=4726 for AIDP), BreakCA's overlap with the WGS-

104   identified indels is far higher (62% for BreakCA compared to 28% for AIDP). Of the 515 BreakCA

105   insertions that are not detected by GATK, most (75%) are also detected by AIDP. Only 9% of BreakCA-

106   identified indels are not called by either GATK or AIDP, suggesting that the accuracy of BreakCA is

107   high. In contrast, most of the insertions detected by AIDP (64%), are only detected by AIDP, suggesting

108   that a large proportion of them may be false-positives.

109

110   Our performance evaluations on GM12878 and Jurkat T-cells indicate that BreakCA+QD offers the best

111   balance of precision and recall for both ATAC-seq and ChIP-seq datasets and we used this approach for

112   all subsequent analyses. We selected score thresholds based on the precision-recall curves for the

113   GM12878 dataset. Specifically, we used a score threshold of ≥0.60 corresponding to ~86% precision and

114   ~75% recall for paired-end ATAC-seq and a score threshold of ≥0.33 corresponding to ~0.91% precision

115   and ~69% recall for single-end ChIP-seq data.

116

117   To test whether BreakCA detects known oncogenic mutations, we applied it to ATAC-seq from four T-

118   ALL cell lines (Jurkat, MOLT-4, CCRF-CEM and RPMI-8402) and one CML cell line (K-562). Two of

119   these cell lines (Jurkat and MOLT-4) are known to harbor oncogenic insertions 8kb upstream of the

120   TAL1 promoter[4] and BreakCA successfully detects both of them (Fig 2). In addition, we verify that

121   BreakCA detects an insertion that is known to be associated with allele-specific expression of the *LMO2*

122   oncogene in MOLT-4 cells[12] (Supplementary Fig 1). These results indicate that BreakCA can detect

123    oncogenic indels in regulatory regions of the genome using ATAC-seq or ChIP-seq reads. We next

124    applied BreakCA to characterize the non-coding regulatory landscape of neuroblastoma.

125

126    **Discovery of indels in neuroblastoma cell lines**

127    Neuroblastoma (NB) is a childhood cancer of the peripheral nervous system with low mutation rates and

128    few recurrently mutated genes[19-22]. Many neuroblastoma tumors harbor no known oncogenic

129    mutations and it has been hypothesized that many high-risk neuroblastomas are driven by rare germline

130    variants, copy number alterations or epigenetic modifications that occur during tumor evolution[19]. We

131    hypothesized that recurrent indels in regulatory regions may also be important drivers of NB. To test this

132    hypothesis, we obtained H3K27ac ChIP-seq from 26 NB cell lines and 2 normal human neural crest cells

133    (hNCCs)[23] and ran BreakCA on these samples to identify indels within regulatory regions defined by

134    H3K27ac peaks.

135

136    The SH-EP and SH-SY5Y cell lines are subclones of the SK-N-SH cell line, so we assigned variant

137    windows identified in these lines to SK-N-SH and treated the combined set of indels as a single cell line.

138    We noticed that the genotypes for the GICAN cell line are nearly identical to those of the GIMEN cell

139    line and that the *SRY* expression of the GICAN cell line is inconsistent with the reported sex (annotated

140    male with no *SRY* expression). We concluded that the GICAN cell line is probably mislabeled and

141    excluded it from further analyses. Our final set of analyzed cell lines consisted of 23 NB cell lines and 2

142    hNCCs.

143

144    We implemented a filtering pipeline to remove potential artefacts and common germline indels (Fig 3a).

145    First, we removed windows overlapping regions that were previously-identified as problematic for ChIP-

146    seq analysis based on their high ratio of multi-mapping to unique mapping reads[24]. Second, we

147    removed indel-containing windows that were detected in the two hNCCs as these are likely to be common

148    germline variants. Third, we conservatively removed 16,234 windows that contained short tandem repeat

6

149    (STR) sequences (Supplementary Fig 2). While STRs have very high mutation rates and contain many

150    true indels, they are also prone to a high rate of variant-calling artifacts caused by polymerase slippage

151    during PCR[25-27].

152

153    Since there are no matched normal tissues for the NB cell lines, germline variants cannot be distinguished

154    from somatic mutations. To eliminate common germline indels and focus on those that are either rare

155    germline variants or somatic mutations, we filtered indels based on their allele frequency in samples from

156    the Genotype-Tissue Expression Project (GTEx) and the genome aggregation database (gnomAD) of

157    15,708 whole genome sequences[28, 29]. We found 1,180 out of 17,834 windows contained indels that

158    were completely absent from GTEx and gnomAD or present with an allele frequency of less than 0.1%

159    (Fig 3a). Additionally, we ran BreakCA on WGS data from 300 GTEx samples and removed indels that

160    we detected in greater than 0.5% of samples. The 990 indel windows that remained after these filtering

161    steps contain either rare germline or somatic indels, which we refer to as RS indels.

162

163    **Recurrent rare germline or somatic indels**

164    RS indels that occur in multiple cell lines are more likely to be oncogenic drivers. To ask how many of

165    the 20bp windows contained recurrent RS indels, we focused on the 742 windows that contained at least

166    one such indel and that were testable by BreakCA (i.e. windows with at least 10 ChIP-seq reads) in 5 or

167    more cell lines. In total, 23 windows contained RS indels in two or more NB cell lines. Remarkably, 4

168    windows contained an RS indel that is present in 4 or 5 cell lines. RS indels are very unlikely to occur in

169    4 or more cell lines due to random inheritance of rare genetic variants (P < $1.8x10^{-9}$ by Poisson test; Fig

170    3c) and therefore these indels may be highly-recurrent oncogenic driver mutations. Recurrent indels can

171    also arise due to high mutation rates, however this scenario is unlikely given that (1) the total number of

172    RS indels that we observe in each cell line is low (Fig 3b) and (2) we have removed windows that overlap

173    annotated STRs, which are typically the most mutagenic sequences. Recurrent rare variants could also be

174    observed if some of the cell lines were derived from close relatives, however, none of the cell lines appear

7

175    to be closely related because the total number of RS indels shared between them is low, and the 4 highly-

176    recurrent RS indels are present in different subsets of cells (Table 1).

177

178    All of the RS indels in the highly recurrent windows (i.e. windows containing RS indels in at least 4 cell

179    lines) are not observed in gnomAD, and are close to transcription start sites (TSSs) of protein-coding

180    genes (Table 1). The first window contains a 1bp insertion that is 295bp upstream of the TSS of *PHF21A*

181    and is observed in 4 cell lines. The second window is in the first intron of *RALGDS* and contains an

182    insertion that is present in 4 cell lines. The third window is within the first intron of *ADAMTS19* and

183    contains 2 insertions that are 5bp apart: a C insertion (present in 3 cell lines) and a T insertion (present in

184    1 cell line). Finally, the fourth window is within the first intron of *GPR85* and contains a complex event

185    consisting of a 1bp insertion and multiple mismatches to the reference sequence (Supplementary Fig 5).

186    Since this complex event appears to be the same in the 5 cell lines where it is detected it may be a rare

187    germline haplotype.

188

189    **An intronic germline deletion associated with *KDM5B* expression**

190    We hypothesized that recurrent RS indels might be associated with the expression of nearby genes and we

191    therefore tested all genes located within 100kb of RS indels for differences in expression using Student's

192    t-test (assuming equal variances). The p-values from these tests show a clear departure from the null

193    expectation (Fig 3d), however our power to detect associations is limited by the fact that each indel is

194    only present in 2-5 cell lines. Under a stringent false-discovery rate threshold of 5%, a single test,

195    between an indel and the expression of the H3K4me3/me2 Lysine Demethylase 5B (*KDM5B*), is

196    significant (nominal p-value=$4.1\times10^{-4}$; Benjamini-Hochberg adjusted p-value=0.021).

197

198    The indel associated with *KDM5B* expression is a 7bp deletion (*GCCTCGG/-*), which is located in its first

199    intron and is present only in the SJNB1 and NB-EBc1 cell lines (Fig 4a & Supplementary Fig 4). This

200    deletion is a germline variant that occurs at a very low minor allele frequency in both gnomAD ($1.3\times10^{-4}$)

8

201 and GTEx ($7.9\times10^{-4}$). The expression of *KDM5B* is very high in the SJNB1 and NB-EBc1 cell lines

202 compared to cell lines that do not contain the indel (with exception of SJNB12) as well as tissues that may

203 resemble the cell-type of origin for NB including hNCCs and adrenal and spinal tissues from GTEx (Fig

204 4b). We performed a motif analysis (see methods) for the 40bp region centered on the indel and found

205 binding motifs for Transcription Factor AP-2 Beta (TFAP2B) and Gamma (TFAP2C) in the reference

206 sequence, which co-localize with the deletion. The indel disrupts the TFAP2B/2C motif and creates a

207 ZNF263 motif in its place (Fig 4c). *TFAP2B* is highly expressed in many NB cell lines (including SJNB1

208 and NB-EBc1), whereas *ZNF263* appears to be expressed across both normal hNCC and NB cells

209 (Supplementary Fig 3).

210

211 **Discussion**

212

213 We used BreakCA to identify indels in 23 neuroblastoma cell lines. One of the most interesting events we

214 detected is a 7bp deletion which replaces a *TFAP2B* binding motif with a *ZNF263* motif within the first

215 intron of *KDM5B* in the SJNB1 and NB-EBc1 cell lines (Fig 4). This indel appears to be germline, as an

216 identical event is detected in 4 out of 31,266 chromosomes surveyed by gnomAD. The presence of this

217 deletion is not associated with a difference in H3K27ac levels but is associated with overexpression of

218 *KDM5B* in these cells, most likely through disruption of the TFAP2B motif. *KDM5B* has known

219 functions in NB and its knockdown results in a 5-fold decrease in cell motility and suppresses the

220 epithelial-mesenchymal transition via downregulation of *NOTCH1* expression in NB cells[30].

221 Furthermore, NB cells that overexpress *KDM5B* form spheroids that are more resistant to in vitro

222 treatment with doxorubicin, etoposide and cisplatin[30]. Finally, *KDM5B* overexpression is associated

223 with poor outcomes in several other cancers including glioma, hepatocellular carcinoma, non-small cell

224 lung cancer, and prostate cancer[31-34].

225

226    The disruption of the TFAP2B motif is also interesting because the *TFAP2B* transcription factor is highly

227    expressed in migrating neural crest cells which are involved in the development of the sympathetic

228    nervous system and are likely cells of origin for NB[35]. In addition, low *TFAP2B* expression is

229    associated with poor survival and prevents neuronal differentiation of NB cells in vitro via

230    downregulation of *MYCN* and *REST*[36]. Our results indicate that TFAP2B may also downregulate

231    *KDM5B*.

232

233    The recurrent RS indels that are present in 4 or 5 cell lines may be oncogenic mutations or rare germline

234    predisposition variants. One of the indels is upstream of *PHF21A*, which encodes a subunit of the BRAF-

235    histone deacetylase complex that is recruited by REST to silence neuronal-specific genes[37]. Another

236    indel is in the first intron of *RALGDS*, which is a guanine exchange factor in a Ras signaling pathway[38].

237    REST is known to inhibit neuronal differentiation of neuroblastoma cells[39, 40] and members of Ras

238    signaling pathways are frequently mutated in relapsed neuroblastoma[41], so both of these indels are

239    excellent candidates for further functional characterization.

240

241    In conclusion, BreakCA allows ATAC-seq and ChIP-seq experiments to be treated as "exome capture"

242    for the regulatory genome and enables the discovery of oncogenic regulatory indels in the absence of

243    WGS data.  While we focused only on short indels in this study, future studies could combine indels

244    called by BreakCA with single nucleotide variants and larger events such as chromosome translocations

245    and copy number alterations. A caveat of BreakCA is that it cannot detect variants outside of ChIP-seq or

246    ATAC-seq peaks or variants that cause a complete loss of these peaks. However, despite this limitation,

247    BreakCA paired with rigorous filtering of common germline events, can identify potential cancer driver

248    mutations and germline risk variants that increase or at least maintain the regulatory activity of a

249    sequence. As a proof-of-principal we used BreakCA to identify recurrent indels in NB cell lines that may

250    be important for neuroblastoma progression, metastasis and drug resistance.

251

252    **Methods**

253    **ATAC-seq and ChIP-seq data**

254    ATAC-seq for the GM12878 lymphoblastoid cell was obtained from Buenrostro et al. 2013 (GEO:

255    GSE47753)[13]. H3K27ac ChIP-seq for the GM12878 lymphoblastoid cell line was obtained from Ernst

256    et al. 2011 (GEO: GSE26320)[14]. H3K27ac ChIP-seq and RNA-seq data from 26 Neuroblastoma (NB)

257    cell lines was obtained from Boeva et al. 2017 (GEO: GSE90683)[23]. H3K27ac ChIP-seq (75bp single-

258    end) for the Kelly NB cell line was obtained from Zeid et al. 2018 (GEO: GSE80151)[42]. ATAC-seq

259    experiments for the Jurkat, CCRF-CEM, RPMI-8402 and MOLT-4 T-ALL cell lines and the K-562

260    chronic myelogenous leukemia cell line were performed in our lab.

261

262    **ATAC-seq experiments**

263    ATAC-seq experiments for the Jurkat, MOLT-4, CCRF-CEM, RPMI-8402 and K-562 cell lines were

264    performed using the Omni-ATAC-seq method as described[43], with minor modifications. In each

265    experiment, $1 \times 10^5$ cells were centrifuged at $1000 \times g$ for 10 min at 4 ºC. Following aspiration, a cell

266    count of the supernatant was performed, the remaining cell number was calculated, and all further

267    reagents in the protocol were titrated to this cell number. For every $5 \times 10^4$ cells, nuclei were isolated in 50

268    μl cold ATAC-Resuspension Buffer (RSB) (10 mM Tris-HCl pH 7.4, 10 mM NaCl, 3 mM MgCl2)

269    containing 0.1% NP40, 0.1% Tween-20, and 0.01% Digitonin, and pipet-mixed up-and-down at least 5

270    times. Nuclei isolation mix was incubated on ice for 3 exactly minutes, washed in 1 ml of cold ATAC-

271    RSB containing 0.1% Tween-20 (but no NP40 or Digitonin) and centrifuged at $1000 \times g$ for 10 min at 4

272    ºC. Nuclear DNA was tagmented in 50 μl Transposition mix (25 μl $2 \times$ TD buffer, 2.5 μl transposase

273    (100nM final), 0.5 μl 1% digitonin, 0.5 μl 10% Tween-20, 16.5 μl PBS and 5 μl diH2O), and incubated in

274    a thermomixer at 37 °C, $1000 \times g$ for 30 min. Tagmented DNA was purified with Zymo DNA Clean and

275    Concentrator-5 Kit (cat# D4014). Library amplification was performed using custom indexing Nextera

276    primers from IDT in a 50 μl Kapa Hi Fi Hot Start PCR reaction (cat# KK2602). Following 3 initial

11

277      cycles, 1 µl of PCR reaction was used in a quantitative PCR (Kapa qPCR Library Quantitation Kit cat#

278      KK4824) to calculate the optimum number of final amplification cycles. Library amplification was

279      followed by SPRI size selection with Kapa Pure Beads (cat# KK8002) to retain only fragments between

280      80-1,200bp. Library size was obtained on an Aglient Bio-Analyzer or TapeStation using a High

281      Sensitivity DNA kit and factored into final Kapa qPCR results to calculate the final size-adjusted molarity

282      of each library. Libraries were pooled and sequenced on an Illumina NextSeq500 in Paired-End 42 base

283      pair configuration at the Salk Next Generation Sequencing Core. ATAC-seq data quality was assessed

284      using the fraction of reads within peaks (FRiP) and fraction of mitochondrial reads (Fmito) metrics

285      (Supplementary Table 1).

286

287      **Aligning reads to the genome and calling peaks**

288      ChIP-seq and ATAC-seq reads were aligned to the hg19/GRCh37 reference genome using BWA-MEM

289      (version 0.7.15-r1140)[44] with default parameters. BWA-MEM performs local alignment and retains

290      reads that only partially map to the genome as soft-clipped alignments, which are useful for identifying

291      indels. Reads were filtered using samtools (version=0.1.19)[45] and only non-duplicated reads with

292      mapping quality (MAPQ $\geq$ 30) were kept for downstream analysis. ChIP-seq and ATAC-seq peak regions

293      were identified using MACS2 (version 2.1.1)[46] with default parameters.

294

295      **RNA-seq and gene expression**

296      RNA-seq reads were aligned to the hg19/GRCh37 reference genome using STAR (version

297      STAR_2.5.3a). Mapped reads were filtered using samtools (version 0.1.19) and only non-duplicated reads

298      with mapping quality (MAPQ $\geq$ 20) were kept for expression analysis. Read counts per gene were

299      calculated with featureCounts (version 1.6.3)[47] using the GENCODE v19 GTF file in paired-end mode

300      and converted to RPKM using edgeR[48].

301

302      **Obtaining features for indel prediction**

303    We divided the genome into non-overlapping 20bp windows and collected 16 features from each window

304    that could help predict the presence of an indel (Supplementary Table 2). Many of the features are based on

305    the proportion of read alignments that contain characteristics such as insertions, deletions, or clipping. For

306    example, one of the features that we consider is the proportion of reads that contain an insertion starting at

307    one base position. We obtain posterior estimates of these proportions using an empirical Bayesian approach,

308    which prevents over-estimation of the proportion when the number of aligned reads at a position is small

309    (i.e. this is a Bayesian alternative to using a pseudocount). We assume that the count of reads with a given

310    characteristic (e.g. insertion) at genomic position $i$, is a binomially-distributed random variable, $X_i$ with

311    proportion parameter $p_i$:

312    $$X_i \sim \text{Binom}(n_i, p_i)$$

313    where, $n_i$ is the total number of reads overlapping genomic position $i$. We place a Beta prior on $p_i$:

314    $$p_i \sim \text{Beta}(\alpha, \beta)$$

315    with $\alpha$ and $\beta$ hyperparameters that describe the shape of the distribution. We estimate $\alpha$ and $\beta$ empirically

316    using the estimated proportion mean ($\hat{\mu}$) and variance ($\hat{\sigma}^2$) computed across all positions within peaks:

317    $$\hat{\mu} = \frac{1}{N} \sum_{i=1}^{N} \frac{x_i}{n_i}$$

318    $$\hat{\sigma}^2 = \frac{1}{N} \sum_{i=1}^{N} \left(\frac{x_i}{n_i} - \hat{\mu}\right)^2$$

319    where $N$ is the total number of positions, $n_i$ is the number of reads overlapping position $i$, and $x_i$ is the

320    number of reads with a characteristic at that position (e.g. insertion). The $\alpha$ and $\beta$ hyperparameters are

321    then calculated as:

322    $$\alpha = \left(\frac{1 - \hat{\mu}}{\hat{\delta}^2} - \frac{1}{\hat{\mu}}\right) \hat{\mu}^2$$

323    $$\beta = \alpha \left(\frac{1}{\hat{\mu}} - 1\right)$$

324    The Beta prior distribution is conjugate with the Binomial likelihood and the corresponding posterior

325    distribution for the proportion is:

13

326
$$p_i \sim \text{Beta}(x_i + \alpha, n_i - x_i + \beta)$$

327    Using this distribution, we calculate the posterior mean and posterior standard deviation at each site, *i,* for

328    the following proportions: (1) insertion reads, (2) deletion reads, (3) reads with clipping on the left side,

329    (4) reads with clipping on the right side. We only perform posterior calculations for base positions

330    covered by at least 10 reads.

331

332    For positions with left-clipped or right-clipped reads we collect additional clipping features including (1)

333    the mean clipping length, (2) the standard deviation in clipping length, and (3) the information content of

334    the clipped sequences (see below).

335

336    We compute the above features for each site, but we perform predictions on windows containing 20 sites.

337    To assign features to 20bp windows we use the sites with the highest posterior means for each of the 4

338    proportion types. Additional features related to clipped reads, such as the mean right clipping length, are

339    taken from the sites with the highest posterior mean clipping proportions.

340

341    **Information content of clipped reads**

342    We compute an information content (*IC*) for clipped reads, which is defined so that highly similar clipped

343    sequences have high *IC*, and dissimilar clipped sequences (perhaps arising from multiple genomic

344    locations) have low *IC*. To compute the *IC* of overlapping clipped reads that start clipping at the same

345    position (*i*=1), we define $s_{i,j}$ as the nucleotide at position *i* in clipped sequence *j*. We also define $T_i$ as the

346    total number of overlapping clipped sequences at position *i*. We assume the first clipped position is *i*=1,

347    and the last position with at least two clipped reads ($T_i > 1$) is *i=S*. We first compute the proportion of

348    each nucleotide $m \in$ (A, C, T, G) at each position *i* as:

349
$$p_{i,m} = \frac{1}{T_i + cT_i} \sum_{j=1}^{T_i} g(s_{i,j}, m) + c$$

14

350 where $c = 0.1$ is a small pseudocount and $g$ returns 1 if two nucleotides are equal:

351
$$g(m,n) = \begin{cases} 1 \text{ if } m == n \\ 0 \text{ otherwise} \end{cases}$$

352 The Shannon entropy of a set of overlapping clipped sequences is then:

353
$$H = \sum_{i=1}^{S} \sum_{m \in (A,C,T,G)} -p_{i,m} \log_2(p_{i,m})$$

354 We define the information content of a set of clipped sequences as the difference between the observed

355 entropy and the maximum possible entropy:

356
$$IC = H_{\max} - H$$

357 The maximum possible entropy, $H_{\max}$, is computed assuming that the nucleotides are evenly distributed

358 across the four possible nucleotides at each site, taking into account that the number of sequences may not

359 be divisible by 4. For example, if there are 6 clipped reads overlapping a position, the maximum possible

360 entropy occurs when the four nucleotide proportions are 2/6, 2/6, 1/6, and 1/6.

361

362 **Training and testing BreakCA.**

363 BreakCA uses a random forest to predict whether a given 20bp window contains an indel using 16 features

364 described in Supplementary Table 2. We created training and test datasets for the GM12878 lymphoblastoid

365 cell line using ATAC-seq (50bp paired-end reads) data[13], H3K27ac ChIP-seq (50bp single-end reads)

366 data[14] and indel genotypes from the Platinum Genomes Project[15]. We labeled 20bp windows centered

367 around the start and end positions of indels within ChIP-seq or ATAC-seq peaks as "true" windows and

368 20bp windows located within peaks and not co-localizing with indels as "false" windows. We used 50%

369 of the windows to create a training dataset and the remaining to create a test dataset. In total there were

370 2980 true (indel-containing) and 842,739 false (non indel-containing) windows in the training dataset and

371 2980 true and 842,738 false windows in the test dataset for paired-end ATAC-seq. For single-end ChIP-

372 seq, there were 623 true and 254,328 false windows in the training dataset and 623 true and 254,327 false

373 windows in the test dataset.

374

375    To implement the random forest model, we used the mlr package in R and tuned three hyperparameters by

376    performing a grid search of reasonable hyperparameter values and choosing the values that yielded the

377    highest accuracy (defined as mean (response == truth)) in 5-fold cross-validation. The 3 hyperparameters

378    were ntree (number of trees to grow), mtry (number of predictors to use for node-split) and node-size

379    (number of observations in the terminal node which is associated with the depth of the decision trees). After

380    choosing hyperparameter values, we trained the random forest on the complete training dataset and applied

381    it to the test dataset, using the fraction of true votes from the decision trees as the prediction score.

382

383    We compared the performance of BreakCA to two popular variant callers: VarScan2[16] and the GATK-

384    HaplotypeCaller[17]. For VarScan2 and the GATK-HaplotypeCaller we extracted the start and end position

385    of the indels using VariantAnnotation[49] and overlapped them with true and false windows in the test

386    dataset. For VarScan2, we used 1.0 - *p-value* as the score and assigned the highest score to each window.

387    For the GATK-HaplotypeCaller we used QD (phred-scaled variant call confidence normalized by allele

388    depth) as the score and assigned the highest score to overlapping windows. For true and false windows

389    with insufficient coverage to call variants, we set the prediction output to 0 for BreakCA, VarScan2, and

390    GATK. To draw precision-recall curves and compute the area under them, we interpolated between

391    datapoints by setting the precision at each recall level, $r$, to $p_{interp}(r) = \max(p(r'))$, where $r' \geq r$ (see

392    https://nlp.stanford.edu/IR-book/html/htmledition/evaluation-of-ranked-retrieval-results-1.html). We did

393    not evaluate the performance of popular indel callers MANTA[50], pindel[51], or Delly[52] because all

394    of these methods require paired-end reads and/or rely on the distribution of insert sizes from mapped

395    reads, which limits the number of ChIP-seq and ATAC-seq datasets these callers can be applied to. In

396    addition, the use of insert sizes is unlikely to work well for ATAC-seq reads, which have a multimodal

397    distribution that changes depending on whether the reads are from nucleosome-free or nucleosome-

398    containing genomic regions[13].

399

400    We further evaluated the performance of BreakCA using ATAC-seq from the Jurkat T-ALL cancer cell line.

401    For this cell line we used indel genotypes from the GATK-HaplotypeCaller applied to WGS as "ground

402    truth"[18]. In total, there were 5020 true and 1,023,159 false windows in the Jurkat ATAC-seq dataset and

403    6403 true and 1,267,161 false windows for the Jurkat H3K27ac ChIP-seq dataset. For these windows, we

404    made predictions using the model trained on the GM12878 ATAC-seq training set and evaluated the

405    precision-recall as described above.

406

407    **Overlapping variant windows with repeat regions.**

408    Genomic positions of known repeat sequences in the hg19 genome were downloaded from the UCSC

409    Genome Browser's RepeatMasker track. Short Tandem Repeat (STR) locations were obtained from

410    Willems et al. 2017[53]. We used GenomicRanges (version 1.24.3)[54] to overlap indel windows with

411    known repeat positions in the human genome and removed indel windows which overlapped STRs.

412

413    **Frequency of indels in gnomAD and GTEx**

414    To filter indels that are likely to be common germline variants we overlapped our 20bp variant windows

415    with indel calls from two large datasets: (1) 15,708 whole genomes from v2.1 of the genome aggregation

416    database (gnomAD)[55] and (2) 635 whole genomes from V7 of the genotype-tissue expression project

417    (GTEx)[29]. We used GATK-SelectVariants to identify gnomAD and GTEx indels overlapping our 20bp

418    variant windows and used SNPSift to extract genomic position, reference and alternate allele and allele

419    frequency (AF) fields from the VCF file[56]. We also required indels to have high coverage across

420    gnomAD samples, removing those with a median coverage of less than 20 reads. For our analysis of rare

421    germline/somatic indels, we retained variant windows with MAF $\leq 5 \times 10^{-3}$.

422

423    **Running BreakCA on GTEx samples**

17

424    In addition to filtering based on indels identified by the gnomAD and GTEx variant calling pipelines, we

425    added a filter for indels identified by BreakCA on 300 GTEx WGS samples[29]. This purpose of this

426    filter is to remove common germline variants (or artefacts) that are detected by BreakCA, but that are not

427    detected by GATK. We ran BreakCA on 300 GTEX WGS samples, using the same $\alpha$ and $\beta$ prior

428    distribution hyperparameters values that we estimated from the GM12878 ATAC-seq dataset. After

429    running BreakCA, we calculated the fraction of samples (SF) with a predicted indel in each window (only

430    considering samples with $\geq$10 reads overlapping a window as 'testable'). (We use SF rather than MAF

431    here because BreakCA calls indels as present/absent and does not distinguish between heterozygotes and

432    homozygotes). For our analysis of rare germline/somatic indels, we retained variant windows with SF $\leq$

433    $5\times10^{-3}$.

434

435    **Estimating the expected frequency of recurrent indels**

436    To create Fig 3C, we calculate the expected number of windows, $T_x$, that would contain RS indels in $x$

437    cell lines, conditional on seeing an indel in one cell line:

438    $$X \sim \text{Binom}(n, p)$$

439    $$E[T_x|\boldsymbol{p}] = \sum_{i=1}^{W} \left(1 + \Pr(X = x|n_i - 1, p_i)\right)$$

440    where $W$ is the total number of windows, and $\boldsymbol{p}$ is a vector of length $W$, with elements $p_i$ that give the

441    expected proportion of cells with an indel. We define $n_i$ is the total number of testable cell lines for

442    window $i$ (i.e. those with sufficient read depth), and subtract 1 from $n_i$ to account for the fact that we have

443    conditioned on seeing an indel in one cell line already. We assume Hardy-Weinberg equilibrium and set

444    $p_i=2f_i(1-f_i) + f_i^2$, where $f_i$ is the allele frequency of the indel in window $i$. We set the allele frequency to

445    either $f_i=2\times10^{-3}$ for all windows, $f_i=1\times10^{-3}$ for all windows, or $f_i=\max(g_i, 1\times10^{-4})$, where $g_i$ is the observed

446    allele frequency of the indel in gnomAD. For windows that contained multiple gnomAD indels, we used

447    the one with the highest allele frequency. We use $f_i=1\times10^{-3}$, because this is the gnomAD allele frequency

448    cutoff used to identify RS indels. We use $f_i=2\times10^{-3}$ as a conservative assumption that gnomAD

449    underestimates some allele frequencies. We use $f_i=\max(gi, 1\times10^{-4})$ to match observed allele frequencies in

450    gnomAD (which are typically much lower than $1\times10^{-3}$).

451

452    To calculate a p-value for the number of observed windows with RS indels in 4 or more cell lines we

453    calculate the probability of observing 4 or more indels using the Poisson cumulative distribution function:

454    $$Z \sim \text{Pois}(\lambda)$$

455    $$pval = \Pr(Z \geq 4 | \lambda = \sum_{x=4}^{M=23} \text{E}[T_x|\boldsymbol{p}])$$

456    Where $\lambda$ is the expected rate of windows containing 4 or more RS indels and M=23 is the number of cell

457    lines in our study.

458

459    RS indels could also be observed due to elevated mutation rates within some windows or due to false

460    positive indel calls. The former possibility is only likely if the window mutation rate is substantially

461    exceeds the allele frequencies that we assume above. Even if we assume an allele frequency of 0.05%

462    (equivalent to a window indel mutation rate of approximately 2p=0.001), the expected number of

463    windows with 4 or more RS indels is 0.38, far lower than the observed 4 (p=$6.4\times10^{-4}$ by Poisson test).

464    This number of windows with RS indels is also unlikely to result from indel call errors (assuming the

465    errors occur independently) because we estimate the per-window false discovery rate for BreakCA on

466    ChIP-seq data to be $1.7\times10^{-4}$—an order of magnitude lower than the allele frequencies we assume above.

467

468    **Testing recurrent indels for association with gene expression**

469    We used GenomicRanges (version 1.34.0) to find promoters of genes located within 100kb of recurrent

470    rare germline or somatic (RS) indels and used Student's t-test to test for differences in mean expression

471    between the indel and non-indel groups. To test if the association with expression is expected by chance

472    we permuted the sample labels for each test and compared the signals with quantile-quantile plots.

473

**Transcription factor motif discovery**

474

475     We obtained the reference sequence (hg19/GRCh37) for 40bp regions centered around predicted indels

476     and introduced the indel to create a non-reference sequence. We used TFBSTools[57] to search the

477     JASPAR2016 database[58] for known transcription factor binding motifs located within the reference and

478     the non-reference sequences. We filtered the motifs using p-value $\leq 0.001$ (computed by TFMPvalue[59])

479     and kept only the top 10% of the motifs found uniquely in either reference or non-reference sequences as

480     our most-reliable hits.

481

**Data and source code availability**

482

483     ATAC-seq data from the Jurkat, RPMI-8402, MOLT-4, CCRF-CEM and K-562 cell lines has been

484     submitted to GEO under accession GSE129086. The BreakCA source code is available from

485     https://github.com/Arkosen/BreakCA.

486

**Author Contributions**

487

488     GM and AS conceived of the project. GM supervised the project. AS developed the BreakCA software

489     and analyzed the data. STT performed the ATAC-seq experiments on the T-ALL cell-lines. YF performed

490     bioinformatic processing and QC of the ATAC-seq data under the supervision of GE. GM and AS wrote

491     and edited the manuscript.

492

498    for technical support. The Razavi Newman Integrative Genomics and Bioinformatics Core Facility of the

499    Salk Institute is funded through the NIH-NCI CCSG: P30 014195 and the Helmsley Charitable Trust.

500

## 501    **References**

502    1.    Borah S, Xi L, Zaug AJ, Powell NM, Dancik GM, Cohen SB, Costello JC, Theodorescu D, Cech TR:
503          **Cancer. TERT promoter mutations and telomerase reactivation in urothelial cancer.** *Science*
504          2015, **347:**1006-1010.
505    2.    Huang FW, Hodis E, Xu MJ, Kryukov GV, Chin L, Garraway LA: **Highly recurrent TERT promoter**
506          **mutations in human melanoma.** *Science* 2013, **339:**957-959.
507    3.    Vinagre J, Almeida A, Pópulo H, Batista R, Lyra J, Pinto V, Coelho R, Celestino R, Prazeres H, Lima
508          L, et al: **Frequency of TERT promoter mutations in human cancers.** *Nat Commun* 2013, **4:**2185.
509    4.    Mansour MR, Abraham BJ, Anders L, Berezovskaya A, Gutierrez A, Durbin AD, Etchin J, Lawton L,
510          Sallan SE, Silverman LB, et al: **Oncogene regulation. An oncogenic super-enhancer formed**
511          **through somatic mutation of a noncoding intergenic element.** *Science* 2014, **346:**1373-1377.
512    5.    Weinhold N, Jacobsen A, Schultz N, Sander C, Lee W: **Genome-wide analysis of noncoding**
513          **regulatory mutations in cancer.** *Nat Genet* 2014, **46:**1160-1165.
514    6.    Puente XS, Bea S, Valdes-Mas R, Villamor N, Gutierrez-Abril J, Martin-Subero JI, Munar M, Rubio-
515          Perez C, Jares P, Aymerich M, et al: **Non-coding recurrent mutations in chronic lymphocytic**
516          **leukaemia.** *Nature* 2015, **526:**519-524.
517    7.    Dalla-Favera R, Bregni M, Erikson J, Patterson D, Gallo RC, Croce CM: **Human c-myc onc gene is**
518          **located on the region of chromosome 8 that is translocated in Burkitt lymphoma cells.** *Proc*
519          *Natl Acad Sci U S A* 1982, **79:**7824-7827.
520    8.    Johnson NA, Savage KJ, Ludkovski O, Ben-Neriah S, Woods R, Steidl C, Dyer MJS, Siebert R,
521          Kuruvilla J, Klasa R, et al: **Lymphomas with concurrent BCL2 and MYC translocations: the critical**
522          **factors associated with survival.** *Blood* 2009, **114:**2273-2279.
523    9.    Meyer KB, Maia A-T, O'Reilly M, Ghoussaini M, Prathalingam R, Porter-Gill P, Ambs S, Prokunina-
524          Olsson L, Carroll J, Ponder BAJ: **A functional variant at a prostate cancer predisposition locus at**
525          **8q24 is associated with PVT1 expression.** *PLoS Genet* 2011, **7:**e1002165.
526    10.   Meyer KB, Maia A-T, O'Reilly M, Teschendorff AE, Chin S-F, Caldas C, Ponder BAJ: **Allele-specific**
527          **up-regulation of FGFR2 increases susceptibility to breast cancer.** *PLoS Biol* 2008, **6:**e108.
528    11.   Oldridge DA, Wood AC, Weichert-Leahey N, Crimmins I, Sussman R, Winter C, McDaniel LD,
529          Diamond M, Hart LS, Zhu S, et al: **Genetic predisposition to neuroblastoma mediated by a**
530          **LMO1 super-enhancer polymorphism.** *Nature* 2015, **528:**418-421.
531    12.   Abraham BJ, Hnisz D, Weintraub AS, Kwiatkowski N, Li CH, Li Z, Weichert-Leahey N, Rahman S,
532          Liu Y, Etchin J, et al: **Small genomic insertions form enhancers that misregulate oncogenes.** *Nat*
533          *Commun* 2017, **8:**14385.
534    13.   Buenrostro JD, Giresi PG, Zaba LC, Chang HY, Greenleaf WJ: **Transposition of native chromatin**
535          **for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and**
536          **nucleosome position.** *Nat Methods* 2013, **10:**1213-1218.
537    14.   Ernst J, Kheradpour P, Mikkelsen TS, Shoresh N, Ward LD, Epstein CB, Zhang X, Wang L, Issner R,
538          Coyne M, et al: **Mapping and analysis of chromatin state dynamics in nine human cell types.**
539          *Nature* 2011, **473:**43-49.
540    15.   Eberle MA, Fritzilas E, Krusche P, Kallberg M, Moore BL, Bekritsky MA, Iqbal Z, Chuang HY,
541          Humphray SJ, Halpern AL, et al: **A reference data set of 5.4 million phased human variants**

542      **validated by genetic inheritance from sequencing a three-generation 17-member pedigree.**
543      *Genome Res* 2017, **27:**157-164.

544   16.   Koboldt DC, Zhang Q, Larson DE, Shen D, McLellan MD, Lin L, Miller CA, Mardis ER, Ding L,
545      Wilson RK: **VarScan 2: somatic mutation and copy number alteration discovery in cancer by**
546      **exome sequencing.** *Genome Res* 2012, **22:**568-576.

547   17.   De Summa S, Malerba G, Pinto R, Mori A, Mijatovic V, Tommasi S: **GATK hard filtering: tunable**
548      **parameters to improve variant calling for next generation sequencing targeted gene panel**
549      **data.** *BMC Bioinformatics* 2017, **18:**119.

550   18.   Gioia L, Siddique A, Head SR, Salomon DR, Su AI: **A genome-wide survey of mutations in the**
551      **Jurkat cell line.** *BMC Genomics* 2018, **19:**334.

552   19.   Pugh TJ, Morozova O, Attiyeh EF, Asgharzadeh S, Wei JS, Auclair D, Carter SL, Cibulskis K, Hanna
553      M, Kiezun A, et al: **The genetic landscape of high-risk neuroblastoma.** *Nat Genet* 2013, **45:**279-
554      284.

555   20.   Pugh TJ, Weeraratne SD, Archer TC, Pomeranz Krummel DA, Auclair D, Bochicchio J, Carneiro
556      MO, Carter SL, Cibulskis K, Erlich RL, et al: **Medulloblastoma exome sequencing uncovers**
557      **subtype-specific somatic mutations.** *Nature* 2012, **488:**106-110.

558   21.   Cancer Genome Atlas Research Network: **Comprehensive genomic characterization of**
559      **squamous cell lung cancers.** *Nature* 2012, **489:**519-525.

560   22.   Banerji S, Cibulskis K, Rangel-Escareno C, Brown KK, Carter SL, Frederick AM, Lawrence MS,
561      Sivachenko AY, Sougnez C, Zou L, et al: **Sequence analysis of mutations and translocations**
562      **across breast cancer subtypes.** *Nature* 2012, **486:**405-409.

563   23.   Boeva V, Louis-Brennetot C, Peltier A, Durand S, Pierre-Eugene C, Raynal V, Etchevers HC,
564      Thomas S, Lermine A, Daudigeos-Dubus E, et al: **Heterogeneity of neuroblastoma cell identity**
565      **defined by transcriptional circuitries.** *Nat Genet* 2017, **49:**1408-1413.

566   24.   ENCODE Project Consortium: **An integrated encyclopedia of DNA elements in the human**
567      **genome.** *Nature* 2012, **489:**57-74.

568   25.   Huang QY, Xu FH, Shen H, Deng HY, Liu YJ, Liu YZ, Li JL, Recker RR, Deng HW: **Mutation patterns**
569      **at dinucleotide microsatellite loci in humans.** *Am J Hum Genet* 2002, **70:**625-634.

570   26.   Kayser M, Roewer L, Hedman M, Henke L, Henke J, Brauer S, Kruger C, Krawczak M, Nagy M,
571      Dobosz T, et al: **Characteristics and frequency of germline mutations at microsatellite loci from**
572      **the human Y chromosome, as revealed by direct observation in father/son pairs.** *Am J Hum*
573      *Genet* 2000, **66:**1580-1588.

574   27.   Chakraborty R, Kimmel M, Stivers DN, Davison LJ, Deka R: **Relative mutation rates at di-, tri-,**
575      **and tetranucleotide microsatellite loci.** *Proc Natl Acad Sci U S A* 1997, **94:**1041-1046.

576   28.   Lek M, Karczewski KJ, Minikel EV, Samocha KE, Banks E, Fennell T, O'Donnell-Luria AH, Ware JS,
577      Hill AJ, Cummings BB, et al: **Analysis of protein-coding genetic variation in 60,706 humans.**
578      *Nature* 2016, **536:**285.

579   29.   GTEx Consortium: **The Genotype-Tissue Expression (GTEx) project.** *Nat Genet* 2013, **45:**580-
580      585.

581   30.   Kuo YT, Liu YL, Adebayo BO, Shih PH, Lee WH, Wang LS, Liao YF, Hsu WM, Yeh CT, Lin CM:
582      **JARID1B Expression Plays a Critical Role in Chemoresistance and Stem Cell-Like Phenotype of**
583      **Neuroblastoma Cells.** *PLoS One* 2015, **10:**e0125343.

584   31.   Kuo KT, Huang WC, Bamodu OA, Lee WH, Wang CH, Hsiao M, Wang LS, Yeh CT: **Histone**
585      **demethylase JARID1B/KDM5B promotes aggressiveness of non-small cell lung cancer and**
586      **serves as a good prognostic predictor.** *Clin Epigenetics* 2018, **10:**107.

587   32.   Shigekawa Y, Hayami S, Ueno M, Miyamoto A, Suzaki N, Kawai M, Hirono S, Okada KI,
588      Hamamoto R, Yamaue H: **Overexpression of KDM5B/JARID1B is associated with poor prognosis**
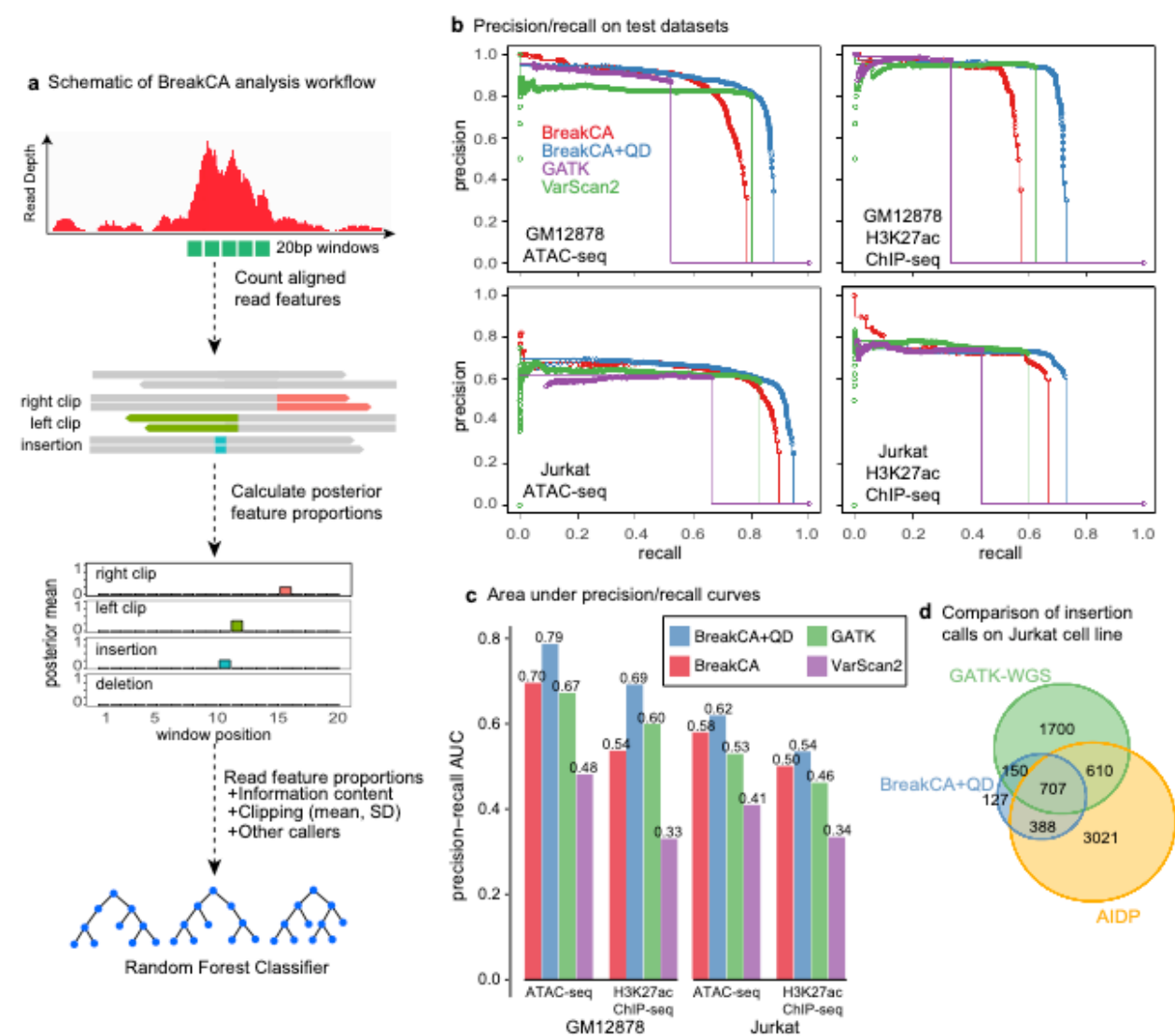589      **in hepatocellular carcinoma.** *Oncotarget* 2018, **9:**34320-34335.

590  33.  Xiang Y, Zhu Z, Han G, Ye X, Xu B, Peng Z, Ma Y, Yu Y, Lin H, Chen AP, Chen CD: **JARID1B is a**
591      **histone H3 lysine 4 demethylase up-regulated in prostate cancer.** *Proc Natl Acad Sci U S A*
592      2007, **104:**19226-19231.
593  34.  Fang L, Zhao J, Wang D, Zhu L, Wang J, Jiang K: **Jumonji AT-rich interactive domain 1B**
594      **overexpression is associated with the development and progression of glioma.** *Int J Mol Med*
595      2016, **38:**172-182.
596  35.  Schmidt M, Huber L, Majdazari A, Schutz G, Williams T, Rohrer H: **The transcription factors AP-**
597      **2beta and AP-2alpha are required for survival of sympathetic progenitors and differentiated**
598      **sympathetic neurons.** *Dev Biol* 2011, **355:**89-100.
599  36.  Ikram F, Ackermann S, Kahlert Y, Volland R, Roels F, Engesser A, Hertwig F, Kocak H, Hero B,
600      Dreidax D, et al: **Transcription factor activating protein 2 beta (TFAP2B) mediates**
601      **noradrenergic neuronal differentiation in neuroblastoma.** *Mol Oncol* 2016, **10:**344-359.
602  37.  Hakimi MA, Bochar DA, Chenoweth J, Lane WS, Mandel G, Shiekhattar R: **A core-BRAF35**
603      **complex containing histone deacetylase mediates repression of neuronal-specific genes.** *Proc*
604      *Natl Acad Sci U S A* 2002, **99:**7420-7425.
605  38.  Urano T, Emkey R, Feig LA: **Ral-GTPases mediate a distinct downstream signaling pathway from**
606      **Ras that facilitates cellular transformation.** *EMBO J* 1996, **15:**810-816.
607  39.  Holzel M, Huang S, Koster J, Ora I, Lakeman A, Caron H, Nijkamp W, Xie J, Callens T, Asgharzadeh
608      S, et al: **NF1 is a tumor suppressor in neuroblastoma that determines retinoic acid response**
609      **and disease outcome.** *Cell* 2010, **142:**218-229.
610  40.  Singh A, Rokes C, Gireud M, Fletcher S, Baumgartner J, Fuller G, Stewart J, Zage P,
611      Gopalakrishnan V: **Retinoic acid induces REST degradation and neuronal differentiation by**
612      **modulating the expression of SCF(beta-TRCP) in neuroblastoma cells.** *Cancer* 2011, **117:**5189-
613      5202.
614  41.  Eleveld TF, Oldridge DA, Bernard V, Koster J, Colmet Daage L, Diskin SJ, Schild L, Bentahar NB,
615      Bellini A, Chicard M, et al: **Relapsed neuroblastomas show frequent RAS-MAPK pathway**
616      **mutations.** *Nat Genet* 2015, **47:**864-871.
617  42.  Zeid R, Lawlor MA, Poon E, Reyes JM, Fulciniti M, Lopez MA, Scott TG, Nabet B, Erb MA, Winter
618      GE, et al: **Enhancer invasion shapes MYCN-dependent transcriptional amplification in**
619      **neuroblastoma.** *Nat Genet* 2018, **50:**515-523.
620  43.  Corces MR, Trevino AE, Hamilton EG, Greenside PG, Sinnott-Armstrong NA, Vesuna S, Satpathy
621      AT, Rubin AJ, Montine KS, Wu B, et al: **An improved ATAC-seq protocol reduces background**
622      **and enables interrogation of frozen tissues.** *Nature methods* 2017, **14:**959-962.
623  44.  Li H: **Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.** *arXiv [q-*
624      *bioGN]* 2013.
625  45.  Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R: **The**
626      **Sequence Alignment/Map format and SAMtools.** *Bioinformatics* 2009, **25:**2078-2079.
627  46.  Zhang Y, Liu T, Meyer CA, Eeckhoute J, Johnson DS, Bernstein BE, Nusbaum C, Myers RM, Brown
628      M, Li W, Liu XS: **Model-based Analysis of ChIP-Seq (MACS).** *Genome Biology* 2008, **9:**R137.
629  47.  Liao Y, Smyth GK, Shi W: **featureCounts: an efficient general purpose program for assigning**
630      **sequence reads to genomic features.** *Bioinformatics* 2014, **30:**923-930.
631  48.  Robinson MD, McCarthy DJ, Smyth GK: **edgeR: a Bioconductor package for differential**
632      **expression analysis of digital gene expression data.** *Bioinformatics* 2010, **26:**139-140.
633  49.  Obenchain V, Lawrence M, Carey V, Gogarten S, Shannon P, Morgan M: **VariantAnnotation : a**
634      **Bioconductor package for exploration and annotation of genetic variants.** *Bioinformatics* 2014,
635      **30:**2076-2078.

636 50. Chen X, Schulz-Trieglaff O, Shaw R, Barnes B, Schlesinger F, Kallberg M, Cox AJ, Kruglyak S,
637        Saunders CT: **Manta: rapid detection of structural variants and indels for germline and cancer**
638        **sequencing applications.** *Bioinformatics* 2016, **32:**1220-1222.
639 51. Ye K, Guo L, Yang X, Lamijer EW, Raine K, Ning Z: **Split-Read Indel and Structural Variant Calling**
640        **Using PINDEL.** *Methods Mol Biol* 2018, **1833:**95-105.
641 52. Rausch T, Zichner T, Schlattl A, Stütz AM, Benes V, Korbel JO: **DELLY: structural variant discovery**
642        **by integrated paired-end and split-read analysis.** *Bioinformatics (Oxford, England)* 2012,
643        **28:**i333-i339.
644 53. Willems T, Zielinski D, Yuan J, Gordon A, Gymrek M, Erlich Y: **Genome-wide profiling of**
645        **heritable and de novo STR variations.** *Nat Methods* 2017, **14:**590-592.
646 54. Lawrence M, Huber W, Pages H, Aboyoun P, Carlson M, Gentleman R, Morgan MT, Carey VJ:
647        **Software for computing and annotating genomic ranges.** *PLoS Comput Biol* 2013, **9:**e1003118.
648 55. Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alföldi J, Wang Q, Collins RL, Laricchia KM,
649        Ganna A, Birnbaum DP, et al: **Variation across 141,456 human exomes and genomes reveals**
650        **the spectrum of loss-of-function intolerance across human protein-coding genes.** *bioRxiv*
651        2019**:**531210.
652 56. Cingolani P, Platts A, Wang le L, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden DM: **A program**
653        **for annotating and predicting the effects of single nucleotide polymorphisms, SnpEff: SNPs in**
654        **the genome of Drosophila melanogaster strain w1118; iso-2; iso-3.** *Fly (Austin)* 2012, **6:**80-92.
655 57. Tan G, Lenhard B: **TFBSTools: an R/bioconductor package for transcription factor binding site**
656        **analysis.** *Bioinformatics* 2016, **32:**1555-1556.
657 58. Mathelier A, Fornes O, Arenillas DJ, Chen C-Y, Denay G, Lee J, Shi W, Shyr C, Tan G, Worsley-
658        Hunt R, et al: **JASPAR 2016: a major expansion and update of the open-access database of**
659        **transcription factor binding profiles.** *Nucleic acids research* 2016, **44:**D110-D115.
660 59. Touzet H, Varre JS: **Efficient and accurate P-value computation for Position Weight Matrices.**
661        *Algorithms Mol Biol* 2007, **2:**15.
662 60. Robinson JT, Thorvaldsdottir H, Winckler W, Guttman M, Lander ES, Getz G, Mesirov JP:
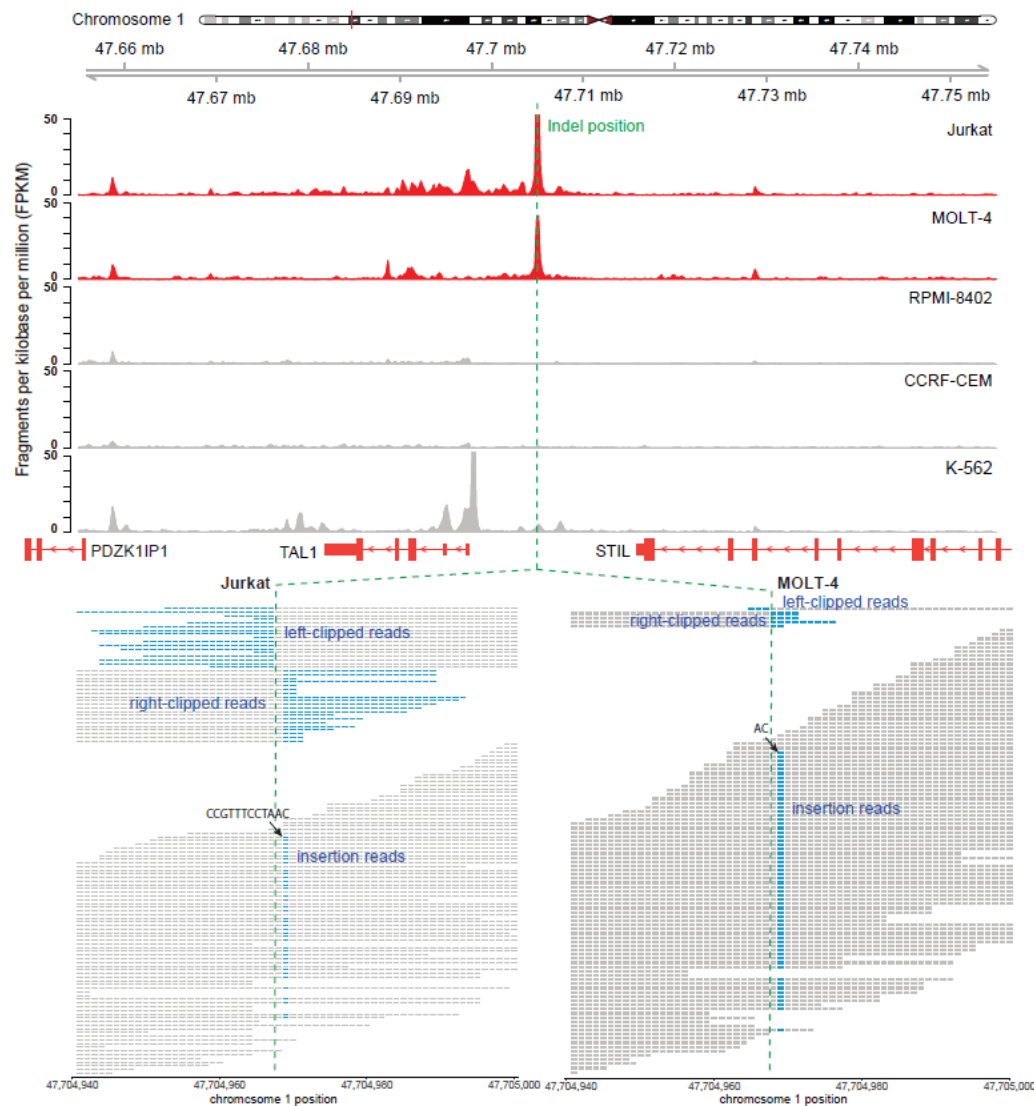663        **Integrative genomics viewer.** *Nat Biotechnol* 2011, **29:**24-26.
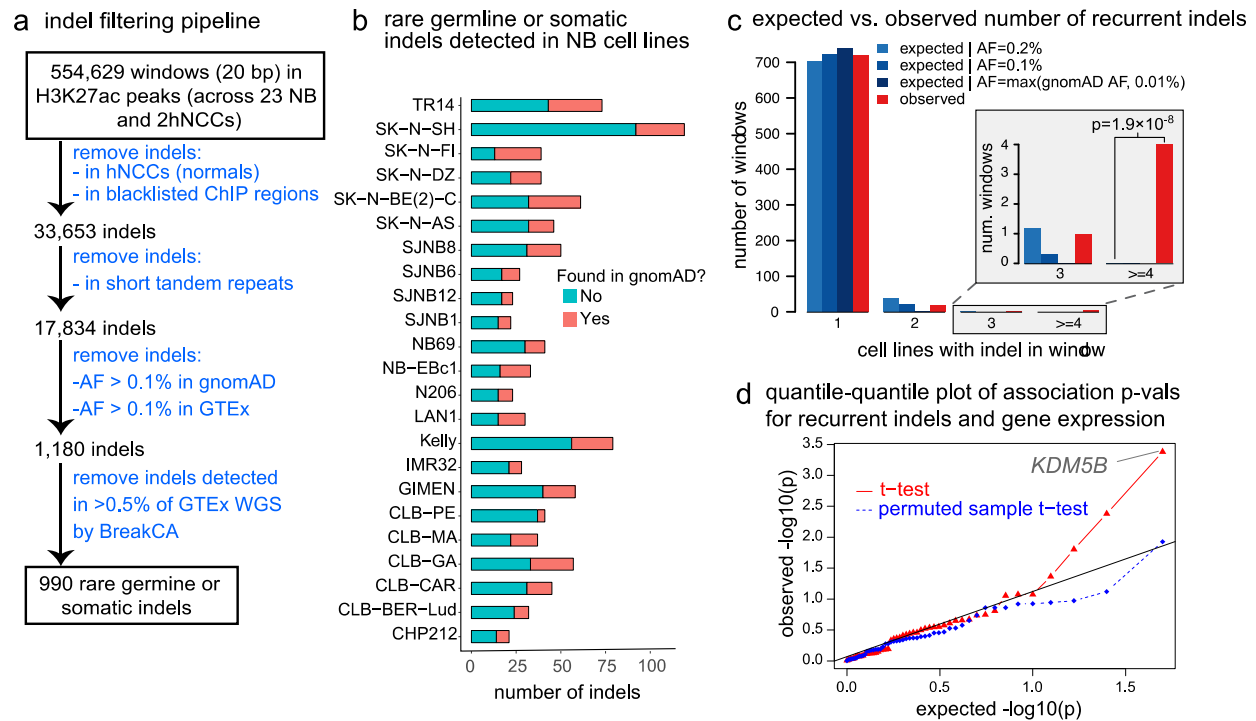
664

665 **Figures**



666

667 **Figure 1: a)** BreakCA detects insertions and deletions (indels) in cancer genomes using mapped ATAC-

668 seq or ChIP-seq reads. **b)** Precision-recall curves for indels scored by four callers: BreakCA,

669 BreakCA+QD (BreakCA with GATK Quality Depth), the GATK-HaplotypeCaller and VarScan2. Each

670 panel shows precision recall curves for a different test dataset: 50bp paired-end ATAC-seq from

671 GM12878 (2980 true and 842,738 false windows); 50bp single-end H3K27ac ChIP-seq from GM12878

672 (623 true and 254,328 false windows); 42bp paired-end ATAC-seq from Jurkat (5020 true and 1,023,159

673 false windows) and 40bp single-end H3K27ac from Jurkat (6403 true and 1,267,161 false windows). Indels

674 called by Platinum Genomes are used as the "ground truth" for GM12878 and indels called by the

25

675    GATK-HaplotypeCaller applied to whole-genome sequence (WGS) data are used as the "ground truth"

676    for Jurkat. **c)** Area under the precision recall curves for each of the test datasets and indel callers. **d)**

677    Comparison of insertions identified in the Jurkat cell line from: GATK-HaplotypeCaller applied to WGS

678    and then filtered for ChIP-seq peaks; BreakCA applied to H3K27ac ChIP-seq; Abraham's Insertion

679    Detection Pipeline (AIDP) applied to H3K27ac ChIP-seq.

680

681



682

**Figure 2:** BreakCA detects known oncogenic insertions upstream of *TAL1* in the Jurkat and MOLT-4 cell lines. The insertion is well-covered by both clipped and insertion-containing reads and inspection of the read pileup reveals a *CCGTTTCCTAAC* insertion in Jurkat and an *AC* insertion in MOLT-4 cell lines (bottom panel). Mapped reads per base position are in grey and soft-clipped bases and insertion positions are colored in blue.

688



689

**Figure 3: a)** Filtering pipeline for identifying rare germline or somatic (RS) indels in 23 neuroblastoma cell-lines. **b)** The number of rare germline or somatic indel windows detected in each neuroblastoma cell-line after filtering, divided into those that are present/absent in gnomAD+GTEx  **c)** The expected number of windows containing RS indels in one or more cell lines (blue bars), assuming that they are inherited germline variants with given allele frequencies (AFs). The plots are conditional on seeing the indel in at least one cell line, and on the number of testable cell lines being at least 5. The red bars are the observed number of windows with one or more RS indels. **d)** Quantile-quantile plot of expected and observed -log10 p-values for RS indel-gene pairs. Each RS indel was tested for association with the expression of all genes within 100kb.

**Figure 4: a)** H3K27ac read depth for two indel-containing neuroblastoma cell lines (in red), 2 human neural crest cell lines, and two non-indel neuroblastoma cell lines. The germline deletion is located within the first intron of *KDM5B* and is covered by both soft-clipped and deletion reads in both cell lines where it is detected. Soft-clipped and deletion base positions are colored in blue. **b)** *KDM5B* gene expression in Reads Per Kilobase Per Million mapped reads (RPKM) in neuroblastoma (NB) cell-lines (from Boeva et al. 2017), human neural crest cells (hNCCs), and in adrenal and spinal tissues from the GTEx project **c)** The *GCCTCGG/-* 7bp deletion disrupts a TFAP2B motif (top panel) and creates a new ZNF263 motif.

708    **Tables**

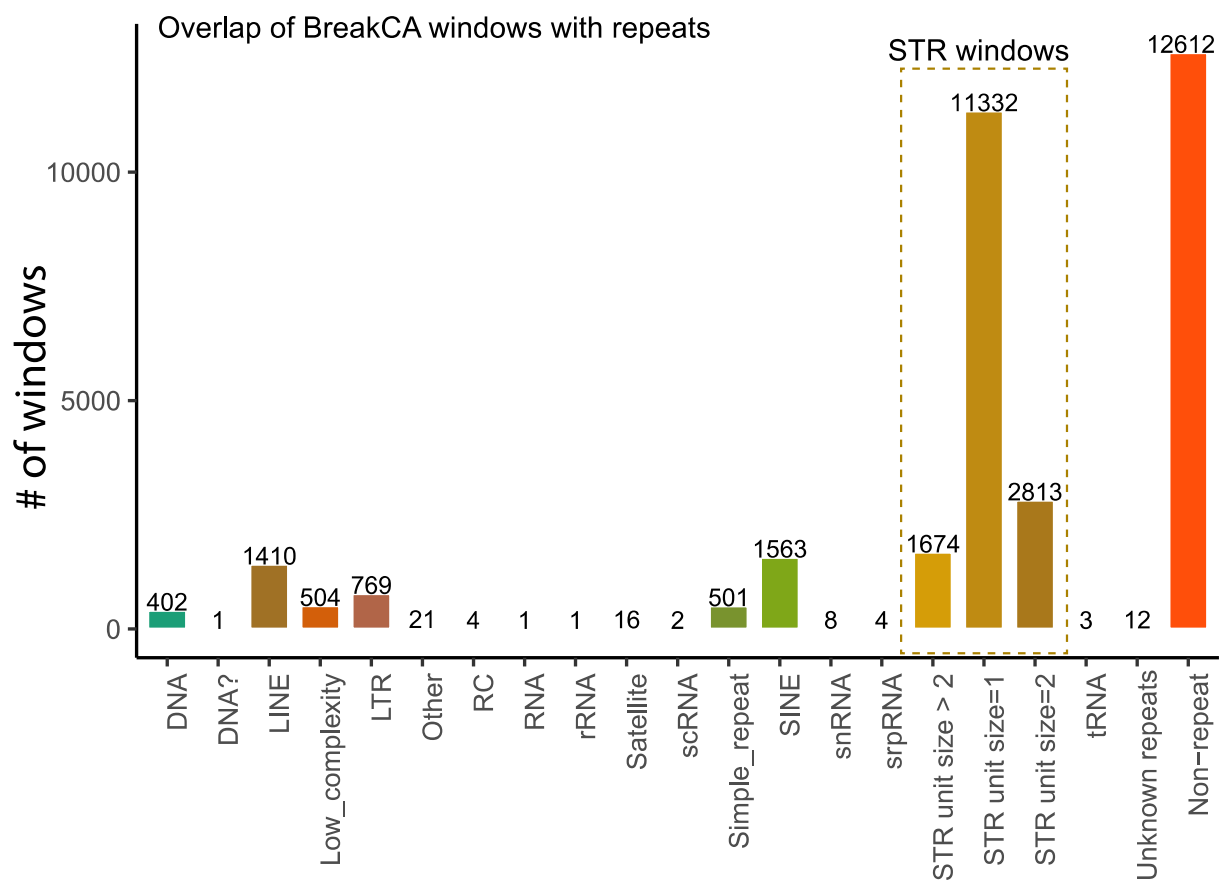| Window | Testable cell lines | Cell lines with indel | gnomAD AF | Gene | Location | Alleles | Notes |
|---|---|---|---|---|---|---|---|
| chr11: 46,144,281-46,144,300 | 21 | CLB-GA, SJNB12, SJNB6, SK-N-BE (2)-C | 0 | PHF21A | Promoter | C/CG | Insertion in all cell-lines |
| chr5: 128,796,561-128,796,580 | 16 | CLB-CAR, CLB-GA, CLB-PE, SJNB6 | 0 | ADAMTS19 | First Intron | T/TC (CLB-CAR, CLB-GA,CLB-PE), C/CT (SJNB6) | Insertion in all cell lines but insertion position differs |
| chr7: 112,726,881-112,726,900 | 19 | CLB-GA, CLB-PE, SJNB6, SK-N-AS, TR14 | 0 | GPR85 | First Intron | Complex event | 1bp insertion and multiple mismatches to reference |
| chr9: 135,994,681-135,994,700 | 22 | CLB-PE, SJNB1, SK-N-AS, TR14 | 0 | RALGDS | First Intron | T/TC | Insertion in all cell lines |

709    **Table 1:** Rare germline or somatic (RS) indels that are observed in at least 4 cell lines.

710
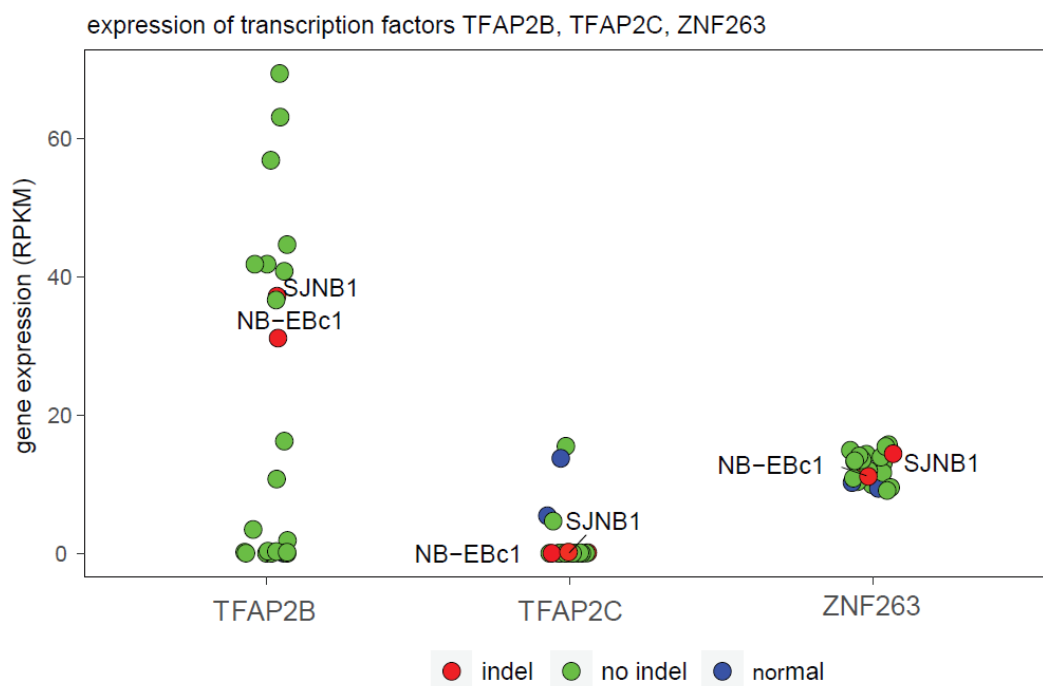
711     **Supplementary Figures**



712

713     **Supplementary Figure 1:** BreakCA detects a known 8bp CGGTTTAA insertion upstream of the LMO2

714     oncogene in the MOLT-4 cell line. Mapped read alignments are grey with soft-clipped bases and insertion

715     positions indicated in blue.

716

**Supplementary Figure 2:** Overlap of indel windows with repeats regions from RepeatMasker or short

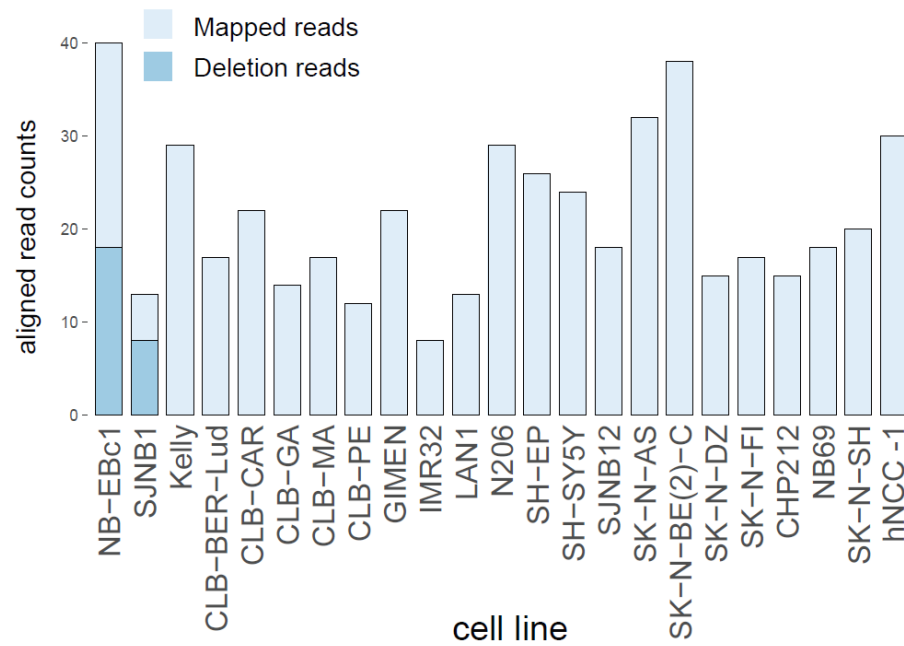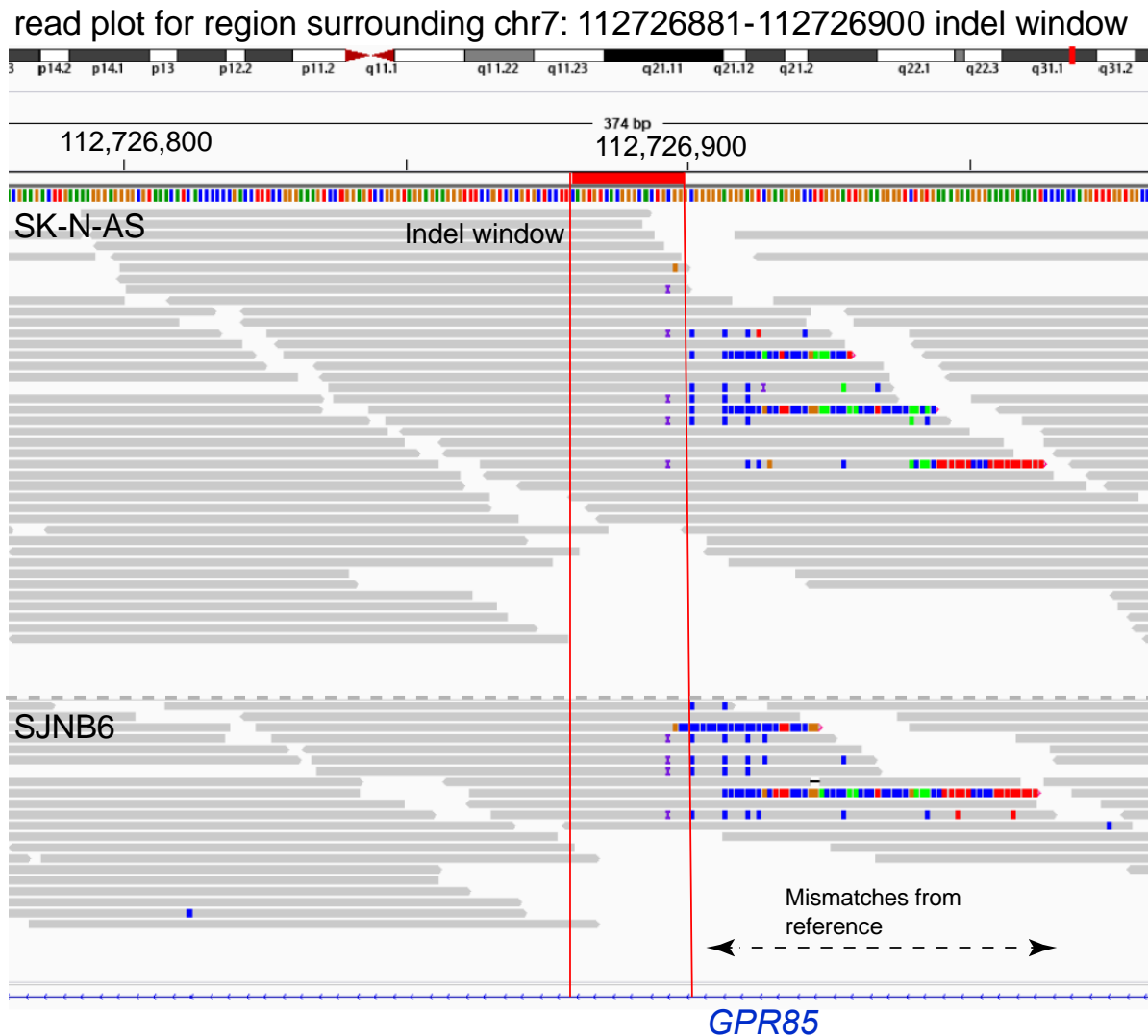tandem repeats (STRs). Indel windows overlapping STRs are not included in our analyses.

719

720    **Supplementary Figure 3**: Gene expression of the *TFAP2B*, *TFAP2C* and *ZNF263* transcription factors in

721    neuroblastoma cell lines and human neural crest cells (hNCCs) demonstrating that *TFAP2B* and *ZNF263*

722    are expressed in the SJNB1 and NB-EBc1 cell lines, while *TFAP2C* is not expressed in NB cell-lines.

aligned reads support KDM5B deletion in 2 cell lines

723

**Supplementary Figure 4:** The count of aligned reads that contain deletions at position chr1:202777149.

Only the SJNB1 and NB-EBc1 cell lines have any deletion reads at this position.

726

34

727

**Supplementary Figure 5:** Read plots from the Integrative Genomics Viewer[60] for two of the

neuroblastoma cell lines containing the complex event in the GPR85 intron. Reads are colored when they

are clipped or have a base mismatch from the reference. Insertions are indicated with a small purple "I".

**Supplementary Tables**

**Supplementary Table 1:** Read statistics including the fraction of reads in peaks (FRiP) for ATAC-seq

datasets generated from the Jurkat, CCRF-CEM, RPMI-8402, MOLT-4, and K-562 cell lines.

**Supplementary Table 2:** Description of prediction features used by the BreakCA random forest.

735    **Supplementary Table 3:** List of rare-germline and somatic indels detected in 23 NB cell-lines with

736    features used for BreakCA prediction.