

1

2 **Title:**

3 **A systematic screen of breast cancer patients' exomes**
4 **for retrotransposon insertions reveals disease**
5 **associated genes**

6

7

8 **Authors:**

9 **Sylvia De Brakeleer, Jacques De Grève and Erik Teugels***

10

11

12 **Affiliation:**

13 **Laboratory of Molecular and Medical Oncology**

14 **Vrije Universiteit Brussel**

15 **Laarbeeklaan 103, 1090 Brussels**

16 **Belgium**

17

18

19 **Email addresses:**

20 **Sylvia De Brakeleer** **sylvia.debrakeleer@gmail.com**

21 **Jacques De Grève** **Jacques.DeGreve@uzbrussel.be**

22 **Erik Teugels** **eteugels@vub.be**

23

24

25 ***corresponding author**

26

27 **ABSTRACT**

28 **Background:** Retrotransposons are genetic elements that jump within the genome via an RNA
29 intermediate. Although they had a strong impact on human genome evolution, only a very tiny
30 fraction of them can be reactivated nowadays, most often with neutral or detrimental
31 consequences. The pathological outcomes associated with such genetic alterations are poorly
32 investigated in the clinic, merely due to their difficult detection.

33 **Results:** We developed a strategy to detect rare retrotransposon mediated insertions in Whole
34 Exome Sequencing data from 65 familial breast cancer patients. When restricting our search to
35 high confidence retrotransposition events occurring in less than 10% of the samples, we
36 identified only ten different Alu elements, two L1 elements, one SVA and two processed
37 pseudogenes. Only two of these insertions occurred within protein coding sequences and
38 interestingly, several of the targeted genes have been previously linked to cancer, in three cases
39 even to increased breast cancer risk (*GHR*, *DMBT1* and *NEK10*). When investigating the
40 molecular consequences of four Alu insertions at the mRNA level, we found that the element
41 present in the 3'UTR of *GHR* repressed expression of the corresponding allele. Moreover, the
42 analysis of a near exonic Alu insertion in *PTPN14* (a mediator of *P53* tumor suppressor activity)
43 revealed that this gene was imprinted and that the presence of an intronic Alu element can lead to
44 loss of imprinting.

45 **Conclusions:** Our data underline the relevance of incorporating the search for uncommon
46 retrotransposition events in Next Generation Sequencing pipelines when analyzing patients with
47 a suspected genetic disease.

48

49 **KEYWORDS:** retrotransposition, Alu, L1, pseudogene, breast cancer predisposition, exome,
50 *GHR, PTPN14, ZNF442, Clorf194*

51

52

53 **BACKGROUND**

54 More than 50% of the human genome is built up with sequences that originated from the activity
55 of transposable elements[1], mainly retrotransposons which are able to move to new locations in
56 the genome via an RNA intermediate using the copy/paste principle[2-4]. L1 retrotransposons
57 belong to the LINE (long interspersed element) superfamily and are the only transposons still
58 active in human. Full length L1 elements (~6kb long) code for the proteins necessary for
59 retrotransposition, including a reverse transcriptase. Although L1 sequences represent about 17%
60 of the human genome, only a very small fraction of them (80-100 copies) retained their
61 transposition capacity[5, 6]. Alu and SVA retrotransposons (about 300 and 1400bp long
62 respectively) are SINEs (short interspersed elements) and do not code for proteins. Rather, they
63 are fully dependent on L1 retrotransposition and parasitize its propagation system. Accidentally,
64 cellular RNA can also misuse the L1 retrotransposition machinery resulting in the generation of
65 processed pseudogenes[7]. The CpG rich sequences present in L1 promoters, Alu and SVA
66 elements are sites for DNA methylation and heterochromatin formation, causing epigenetic
67 silencing[8]. Interestingly, the vast majority of evolutionary stabilized Alu insertions are
68 located in gene-rich regions. They are often embedded in sequences encoding pre-mRNAs or
69 mature mRNAs, usually as part of their introns or UTRs where they can potentially contribute to
70 transcriptome variation[9].

71 In 1985 the first de novo Alu element insertion was reported in a B cell lymphoma[10]. Three
72 years later, Kazazian et al.[11] reported a Haemophilia A causing L1 element insertion in
73 germline DNA. In a review article published in 2016[6], the list of retrotransposons associated
74 with human diseases counts 124 entries. Obviously, the number of disease causing
75 retrotransposition events identified in human[12] is increasing steadily, but the attention these
76 mutations receive in the clinic is still very poor when compared to classical mutational events
77 (nucleotide substitutions and small insertions/deletions). When searching for breast cancer (BC)
78 predisposing mutations in high risk families, we previously identified one Alu insertion in exon
79 11 of *BRCA1* and one Alu insertion in exon 3 of the *BRCA2* gene[13], this last one appeared later
80 on to be a recurrent founder mutation restricted to the Portuguese population[14]. When Next
81 Generation Sequencing tools were introduced in the diagnostic field, only few labs adapted their
82 IN SILICO pipeline to allow the identification of pathogenic retrotransposition events in the
83 genes they analyze[15, 16]. In the present work, we developed a strategy allowing the detection
84 of uncommon retrotransposition events in Whole Exome Sequencing (WES) data using BAM
85 files generated from a cohort of familial breast cancer patients previously used to identify more
86 classical cancer predisposing mutations[17]. In a second step, we investigated the molecular
87 consequences generated by 5 candidate pathogenic Alu insertions revealed during our screen,
88 and highlight the relevance of screening for rare retrotransposition events in the clinical context.

89

90

91 **RESULTS**

92

93 **IN SILICO screen for retrotransposition mediated insertions**

94 Because of their shared mechanisms of transposition using an RNA molecule as intermediate,
95 retrotransposons all bear a long polyA tail at their 3' extremity. Also, their specific mechanism
96 for insertion into the target DNA results in the duplication of a short sequence at the position
97 where the retrotransposon has integrated. This duplicated sequence (called Target Site
98 Duplication or TSD) is usually between 4 and 20 nucleotides long and flanks the retrotransposon
99 specific sequences. We exploited those two characteristics to identify retrotransposon insertions
100 not present in the human reference genome (hg19) using WES data generated with the SeqCap
101 EZ Exome v3.0 kit from Roche. Paired-end sequencing usually generates read pairs that
102 perfectly map in close proximity on a reference genome. However, the presence of a non-
103 reference retrotransposon insertion will result in read pairs that do not map closely (Figure1). We
104 selected these discordantly mapped read pairs by applying the open-source RetroSeq
105 software[18] on the BAM files previously obtained from 65 BC patients[17] (positive *BRCAl*
106 control included), with a first additional restriction that one read of the discordantly mapped read
107 pairs must harbor a long polyA stretch. On average, we obtained 987 candidate insertions for
108 each patient. When introducing additional filtering steps (using an Excel program) to remove
109 false positives and frequently occurring non reference polymorphic insertions, we ended with a
110 list counting only 32 different candidate polyA containing insertion sites for the pooled 65
111 patients (Additional file 1). A detailed description of our detection strategy and its application on
112 a positive control sample is presented in the Method section, wherein we also compare two
113 widely used exome enrichment kits.

114

115 **High confident retrotransposition mediated insertions**

116 When performing a manual selection (using the Integrative Genomics Viewer: IGV[19]) for the
117 presence of a TSD on the 32 polyA containing insertions we had identified, we ended up with 18
118 high confidence candidates retrotransposon insertions (Table 1). The further characterization of
119 these retrotransposition mediated insertions (12 Alu, 2 L1, 1 SVA and 3 pseudogenes, see Table
120 1) was made possible by the presence of a stretch of nucleotides (maximum ~50bp long) in a
121 fraction of the aligned read pairs that do not match with the reference sequence, but
122 corresponded the 5' extremity of the inserted sequence. These sequences could be identified
123 using available web tools (see Methods). A representative screenshot showing the outcome of an
124 IGV analysis of the genomic region wherein the Alu element integrated into the BRCA1 gene of
125 our positive control sample is presented in Figure 2. Only 2 out of the 18 retrotransposon
126 mediated events occurred within protein coding regions (Alu element in *ZNF442* and *UQCR10*
127 transcript in *Clorf194*, see Table 1) which might be surprising since WES probes are expected to
128 be specifically designed towards exonic sequences. Indeed, according to the kit provider (Roche)
129 only protein coding parts of the transcripts are targeted, but for exons that are smaller than 100
130 bp, the target region is extended to 100 bp. Nevertheless, we expect that the large majority of
131 retrotransposition mediated insertions occurring in the well covered protein coding regions will
132 be picked up with our screen, their low incidence being in agreement with the literature[5, 20,
133 21]. Conversely, we are convinced that a significant fraction of the polymorphic retrotransposons
134 located in the vicinity of protein coding regions (regions often poorly or not covered at all by
135 WES) is missed with our WES based approach, although the large majority of identified
136 insertions (15/17, positive control excluded) are located in intronic, 3'UTR, promoter or even
137 intergenic sequences. As a sufficient coverage at sequencing is essential to allow the
138 identification of polymorphic retrotransposons in the non-coding regions, we wanted to estimate

139 the number of false negative carriers of such highly confident retrotransposition mediated
140 insertions among the 65 samples we investigated. Therefore, all samples were rescreened for the
141 presence of the 18 IN SILICO identified insertion types using IGV (see Table 1). All insertions
142 identified with the IN SILICO pipeline could be confirmed with IGV analysis (no false
143 positives). As expected, IGV inspection revealed higher rates of carriers for 9 insertions (Table
144 1). In two cases, insertion carriership was even observed in 50% of the samples, with IGV and
145 IN SILICO data discrepancies being well explainable by the low coverage at sequencing. Of
146 note, about half of the high confident retrotransposon insertions we identified are reported in the
147 database of retrotransposon insertion polymorphisms in human (dbRIP, see Table 1).

148

149 **Retrotransposon insertions in non-coding regions can affect allelic expression**

150 To investigate the potential consequences of Alu insertions located in non-coding gene regions at
151 the molecular level[22], we collected blood samples from available carriers (and non-carriers for
152 control) and extracted the corresponding RNA material. Since allele specific expression levels
153 can be monitored by Sanger sequencing only when the Alu carrier is also heterozygous for a
154 coding SNP in the same gene, the number of suitable insAlu carriers among the 65 investigated
155 patients was not always sufficient and additional familial BC patients were genotyped using
156 mutation specific hemi-nested PCRs (Table 2). Sufficient blood samples could be collected to
157 initiate the analysis of 5 Alu targeted genes (*GHR*, *GSTA5*, *NEK10*, *PTPN14* and *UPF2*).

158 Sequence analyses of the cDNA regions flanking the allele discriminating SNPs revealed that the
159 intronic Alu insertions in *UPF2* and *NEK10* do not result in an obvious decrease or increase of
160 the mRNA levels generated from the retrotransposon targeted allele (Figure 3). In sharp contrast,
161 the Alu insertion in the 3'UTR of *GHR* clearly resulted in allele silencing (Figure 4A) suggesting

162 that this Alu sequence may contain sites for microRNA directed degradation, a mechanism
163 previously proposed to explain the occurrence of evolutionary stabilized Alu derived microRNA
164 binding sites in the 3'UTRs of specific gene[23]. Conversely, the two available carriers of the
165 intronic Alu insertion in *PTPN14* seem to express both alleles equally well but among the three
166 non-carriers, one expresses virtually only one allele while the two others show a reduced
167 expression of this same allele (Figure 4B). Partial or near absolute mono-allelic expression of
168 *PTPN14* is in agreement with a study listing this gene as a high-confidence imprinted human
169 gene candidate with maternal expression[24]. On the other hand, the apparent loss of imprinting
170 (LOI) observed in the Alu carriers suggests that the integration of one single Alu element into the
171 gene sequence of the imprinted allele would be sufficient to fully reactivate this allele, which
172 matches with earlier observations reporting a sharp inflection in SINE content at transitions from
173 imprinted to non-imprinted genomic regions[25]. Unfortunately, differential expression of the
174 *GSTA5* Alu containing allele (insertion occurred in the promoter region) could not be
175 investigated due to the very low expression level of *GSTA5* in leucocytes compared to *GSTA1*
176 and *GSTA2* (see GeneCards.org), three homologues that are hard to discriminate when using
177 PCR based assays (only *GSTA1* and/or *GSTA2* sequences could be detected after Sanger
178 sequencing even by using gene discriminating primers for PCR, results not shown).

179

180 **PolyA containing insertions without identifiable TSD**

181 In addition to the 18 high confident retrotransposition mediated insertions for which we could
182 identify a TSD (described above), our IN SILICO pipeline also identified 14 candidate insertion
183 sites for which a TSD could not be deduced using IGV (see Additional file 1 and Table 3). A
184 summary of each of the targeted gene's function is presented in Additional file 2 (see also Table

185 4 for PubMed hits). A lower coverage at sequencing most probably explains the problematic
186 discovery and characterization of these polyA containing insertions. However, for eight of them
187 we found parts of retrotransposon sequences (7x Alu, 1x L1) by analyzing the mates of the
188 discordantly mapped mate pairs (Table 3). In addition, 7 out of these 14 expected
189 retrotransposition events map at the same genomic position as do elements reported in the
190 database for Retrotransposition Insertion Polymorphisms in Humans (dbRIP[26]). Consequently
191 10 out of the 14 polyA containing insertions can be linked to a retrotransposon, making them still
192 strong candidate retrotransposon insertion sites with potential pathological consequences (Table
193 3). The high rate of confirmed retrotransposon mediated insertions observed among the
194 candidate polyA containing insertions we identified (28 out of 32) is a good indicator for the
195 high specificity of the followed screening strategy.

196

197

198 **DISCUSSION**

199 De novo retrotransposition events occurring within protein coding regions are expected to be
200 very rare as they would strongly affect protein integrity, most often in a detrimental manner [5,
201 20, 21]. Our results fully match with these earlier observations as the only rare polymorphic
202 retrotransposon that inserted exome wide into a coding sequence among our 65 patients (positive
203 BRCA1 control excluded) was an Alu element into the *ZNF442* gene (a short summary of the
204 function of all genes targeted with a high confident retransposon mediated insertion is presented
205 in Additional file 3). A mutation specific PCR screen (see Table 2) detected this variant in only
206 one out of 710 familial BC patients. Since this patient has North African roots, we cannot

207 exclude that the mutation is more prevalent in that region. *ZNF442* has been poorly investigated
208 till now (see outcomes of PubMed searches for the different targeted genes in Table 5) although
209 the gene has been reported as a good candidate driver gene for the neoplastic process of breast
210 and colorectal cancer[27]. The second identified retrotransposition mediated event that destroys
211 protein integrity is the *UQCR10* transcript inserting into the *Clorf194* coding sequence (creating
212 a *UQCR10* processed pseudogene). This polymorphic insertion has been previously observed[7,
213 28] and was present in 7 of our 65 samples. The mutant allele is unable to produce normal
214 *Clorf194* protein, but we cannot exclude that the *UQCR10* pseudogene is expressed as a
215 recombinant mRNA leading to overproduction of *UQCR10* protein since a near full length
216 *UQCR10* transcript (from c.-15 till polyA tail) inserted in the 5' coding region (c.58_70) of the
217 *Clorf194* gene in the same orientation (Table 1). Little is known about these two genes, but none
218 has been linked to breast cancer yet. Interestingly, two seemingly dominant heterozygous
219 missense mutations in *Clorf194* were recently associated with Charcot-Marie-Tooth disease[29],
220 the most common form of inherited peripheral neuropathy. The authors reported a Ca^{2+}
221 regulatory function for *Clorf194* and therefore suggested that this gene (and consequently any
222 mutated form at it) may also be associated with other neurodegenerative disorders characterized
223 by altered Ca^{2+} homeostasis.

224

225 The large majority of high confident retrotransposon insertions we could identify were located in
226 non-coding regions (15 out of 17). However, such regions are not the primary targets when
227 performing WES and consequently their coverage at sequencing will often be (very) low, or they
228 will not be covered at all. Therefore, we believe much more intronic (near exonic) and UTR
229 located polymorphic retrotransposons escaped our attention, the 15 ones reported herein

230 representing the tip of an iceberg. The majority of the identified retrotransposon insertions
231 incorporated into intronic sequences, close to the intron-exon boundaries. The genes targeted by
232 these retrotransposition events can be linked to several biochemical processes (Additional file 3)
233 that might be involved in breast cancer development: they can have a tyrosine kinase activity
234 (*NEK10*[30, 31]), a Tyrosine phosphatase activity (*PTPNI4*[32]), being involved in the
235 ubiquitination process (*UBA6*[33]), in the posttranscriptional methylation of internal adenosine
236 residues in eukaryotic mRNAs (*METTL3*[34, 35]), in the regulation of the nonsense mediated
237 decay pathway (*UPF2*[36]), have a suppressive role in osteosarcoma progression (*TMIGD3*[37])
238 or are considered as candidate tumor suppressor gene for different cancer types (*DMBT1*[38-
239 40]). Interestingly, four of the five intronic Alu insertions occurred in antisense orientation,
240 suggesting that these insertions might contribute to the generation of alternative transcript
241 forms[41]. We also identified two retrotransposon insertions in 3'UTR gene regions (in *GHR* and
242 *HSD17B12*), one in the promoter of *GSTA5* (a glutathione S transferase catalyzing the
243 conjugation of reduced form of glutathione to xenobiotic substances for the purpose of
244 detoxification) and three in intergenic regions (in one case with nearest gene linked to BC: *MIF*-
245 *ASI*[42]). More convincing for their role in breast cancer predisposition is that 3 retrotransposons
246 inserted into a gene that had been previously associated to breast cancer risk (*GHR*[43-45],
247 *DMBT1*[46-48] and *NEK10*[49, 50]). The potential cancer risk associated to these particular
248 mutations should be confirmed by investigating a much higher number of BC cases and controls.
249 Furthermore, we could successfully investigate the consequences of four different Alu element
250 insertions in non-coding gene regions at the molecular level by determining the relative
251 expression level of the wild type versus mutant allele in leucocytes. Analysis of the Alu
252 insertions in *NEK10* and *UPF2* did not reveal a marked decrease or increase of the mRNA levels

253 generated from the retrotransposon targeted alleles. Nevertheless, a pathogenic (or protective)
254 contribution by these mutated alleles cannot be excluded yet, as they may for instance generate
255 splicing alterations that cannot be detected with the performed test or may result in cell type
256 specific effects not observed in leucocytes. Moreover, the recently revealed mechanism of onco-
257 exaptation[51, 52] whereby cryptic regulatory elements within transposons can be epigenetically
258 reactivated in precancerous cells may also drive cancer risk. In sharp contrast, the Alu insertion
259 in the 3'UTR of *GHR* clearly resulted in allele silencing. Since several association studies
260 indicated that *GHR* is implicated in breast cancer predisposition[43-45], it is tempting to
261 designate the Alu insertion in *GHR* as (one of) the causal mutation(s) for the reported cancer risk
262 modulation. The interpretation of the data obtained when investigating the Alu insertion in
263 *PTPN14* was more challenging since (near) mono allelic expression was already observed in the
264 three patients that did not carry the germline Alu insertion in *PTPN14*. Fortunately, literature
265 digging revealed a table wherein *PTPN14* was listed as a high-confidence imprinted human gene
266 candidate with maternal expression[24]. The LOI observed in our two Alu carriers suggests that
267 the integration of one single Alu element into the imprinted *PTPN14* gene allele would be
268 sufficient to fully reactivate this allele[25]. The observed LOI also suggest that both patients
269 inherited the mutated *PTPN14* allele from their father (could not be verified) and that in case the
270 mutated allele is maternally inherited, LOI will not be observed. Since *PTPN14* was recently
271 identified as a mediator of the tumor suppressor activity of *p53* and regularly mutated in
272 cancer[32], its imprinted nature and the loss of imprinting (LOI) induced by particular mutations
273 are two elements that contribute additional levels of complexity to the molecular mechanism
274 leading to cancer when the *PTPN14* gene is involved.

275

276 CONCLUSIONS

277 Our study indicates that raw WES data obtained from clinical samples (whose availability is
278 exponentially growing) hide a manageable number of retrotransposition mediated polymorphic
279 mutations that can be dug up when using appropriate IN SILICO tools. A significant fraction of
280 these mutations will affect gene function, even when located outside protein coding regions, and
281 consequently may be (one of) the pathogenic factor(s) responsible for a patient's disease. Our
282 investigation was restricted to 65 non BRCA1/2 familial BC patients, but much more patients
283 should be analyzed (in future genome wide) to obtain a broader view of all possible
284 retrotransposon mediated insertions involved in the disease, and to determine their respective
285 molecular and pathogenic consequences. Moreover, transposon mediated insertions resulting in
286 modified breast cancer risk may also generate other clinical phenotypes. For example, the growth
287 hormone receptor (which is encoded by the *GHR* gene) is expressed in a broad range of tissues
288 and involved in fundamental biological processes such as growth promotion, metabolism, cell
289 division and regeneration[53], the most typical clinical syndrome associated to GHR deficiency
290 being dwarfism (Laron syndrome). Accordingly, the clinical consequences associated to the
291 presence of an Alu insertion in the 3'UTR of *GHR* (which leads to allele silencing) are most
292 probably not limited to breast cancer risk. As pathogenic retrotransposon insertions are not
293 limited to the familial BC syndrome, the described mutation screen (or an alternative version of
294 it, e.g. for the analysis of WES data obtained with the exome capture kit from Agilent) should be
295 applied for all diseases with a suspected genetic etiology. Indeed, screening patients with a
296 specific genetic disease will enrich for insertions in genes involved in that disease (in
297 preparation). Centralized databases registering all identified polymorphic retrotransposition
298 events should be further expanded[7, 26, 54], with inclusion of population specific allele

299 frequencies, as is the case for classical variations. In order to identify good candidate pathogenic
300 Alu element insertions, Payer et al.[55] used the outputs of GWAS studies to restrict their search
301 to genomic loci previously implicated in disease risk. Conversely, our data indicate that
302 microchips should be developed allowing genome wide genotyping for the presence of this
303 particular type of polymorphic insertions, in order to perform retrotransposition specific
304 GWAS's for a multitude of diseases or traits. Finally, the genome wide identification and
305 investigation of polymorphic retrotransposon insertions in clinical samples will not only lead to a
306 better understanding of diseases, but will also contribute to elucidate more basic genetic
307 mechanisms such as gene imprinting, and will help to evaluate the impact that de novo
308 retrotransposition mediated insertions still have on human genome evolution.

309

310

311 **METHODS**

312

313 **Patient Material**

314 A WES study using the SeqCap EZ Exome v3.0 kit from Roche has previously been performed
315 using blood samples from probands of 65 unrelated high risk BC families with the aim to
316 identify classical cancer predisposing mutations (protein truncating SNPs and Indels, and splice
317 site mutations)[17]. All patients were recruited at the UZ-Brussel hospital and met the
318 requirements for a diagnostic *BRCA1* and *BRCA2* mutation analysis. All except one (the positive
319 control: c.1739_1740insAlu in *BRCA1*[13]) were negative for a pathogenic *BRCA1* or *BRCA2*
320 mutation. A large subset of these 65 patients (57) belongs to families with at least two first

321 degree blood relatives with BC. The FastQ files generated during this previous study were reused
322 for the present study. The genomic DNA from the positive control was resubmitted for WES
323 analysis using the SureSelect Human All Exon V6 kit from Agilent (performed by BGI
324 Genomics Co, China) to compare both WES approaches. In addition, blood samples obtained
325 from 710 probands from independent non BRCA1/2 BC families (including the 65 probands
326 used for the WES analysis) were used for a PCR based genomic screen to identify additional
327 carriers of Alu element insertions in case insufficient samples were available for RNA analysis.

328

329 **Strategy for retrotransposon detection using WES-BAM files**

330 Binary alignment files (BAM files) are outputs (binary versions) of a short read aligner (e.g.
331 BWA[56]) that can map the read pairs generated during the sequencing process (paired-end
332 sequencing) towards a reference sequence, for instance the reference human genome (hg19). In
333 particular circumstances, these read pairs are discordantly mapped while generated from the
334 same small genomic fragment, meaning that the mate reads map at very different positions on the
335 reference genome.

336 To trace retrotransposon insertions not present in the human reference sequence we used, in a
337 first try-out, the open-source RetroSeq software described by Keane et al.[18] as this software
338 can be run on virtually any computer[57]. This software needs as input a BAM file (we used the
339 WES BAM file generated with a positive control sample harboring the c.1739_1740insAlu
340 mutation in *BRCA1*; the library was prepared using the exome capture kit SeqCap EZ Exome
341 v3.0 from Roche), a reference genome (hg19) and a library of mobile element sequences.
342 RetroSeq operates in two phases, the first being the discovery phase where discordantly mapped
343 (and one end mapped) mate pairs are screened for the presence of a mobile element and in the

344 affirmative assigned to a particular type of mobile element. In the second phase (calling phase),
345 the sequence of the anchoring mate read is used to localize the polymorphic mobile element on
346 the reference genome. Unfortunately, when applying RetroSeq on our positive control sample we
347 could not detect the Alu insertion in *BRCAl*. RetroSeq is primarily designed to identify non
348 reference mobile elements in Whole Genome Sequencing data, while companies that design kits
349 for exome enrichment try to avoid the capture of nucleotide fragment harboring repetitive
350 sequences. To circumvent this obstacle, we decide in a first step to search not for mobile
351 elements but for polyA stretches (settings: minimum 80 nucleotides long; 90% ID) within the
352 discordantly mapped reads of each patient data set using the first module of the RetroSeq
353 software (discovery phase). Among the discordantly mapped read pairs traceable in BAM files
354 we expect to find those generated from genomic DNA fragments that contain both non reference
355 retrotransposon sequences and exon (flanking) sequences (see Figure 1). By selecting
356 discordantly mapped reads containing a long polyA stretch in one of the paired reads, we will
357 only retain the junctions at the 3' end of the retrotransposon. The location of the mate of each
358 polyA containing read identified during the discovery phase is recorded in a BED file by the
359 RetroSeq program (together with additional information). This file can be recovered as a .tab file
360 for further analysis (see Additional file 4 for the .tab file generated with the positive control
361 sample saved in Excel).

362 If coverage is sufficient, several anchor reads (the mates of the polyA containing reads) are
363 expected to align in the same genomic region. Clusters of such reads are subsequently generated
364 when they occur within a 300 bp long genomic interval and a maximum inter read distance of
365 200 bp. In a next step, each cluster is represented by the anchor read expected to be closest to the
366 integration site (TSD in case of retrotransposons) while the number of reads within each cluster

367 is recorded in a separate column (duplicate reads are removed). To perform these cluster
368 calculation steps, the patient specific .tab file generated by RetroSeq is saved in Excel format,
369 non-relevant information is first deleted and the remaining data are pasted in the first sheet of a
370 preformatted Excel file (provided as Additional file 5). The detailed procedure for cluster
371 calculation is described in next section. We considered clusters with three or more anchor reads
372 as potential candidate insertion sites and retained them for further investigation. When applying
373 this detection strategy to our 65 BC patients, we obtained on average 987 such clusters for each
374 individual. The positive control sample (see Additional file 6) shows 641 clusters with three or
375 more units, with a clear cluster of 17 units at the level of the Alu insertion site in *BRCA1* (TSD is
376 at chr17:41245809-41245825).

377 To identify and in a second step discard the clusters (genomic regions) picked up in a large
378 fraction of individuals (e.g. false positives resulting from technical artefacts, or highly recurrent
379 polymorphic insertions not present in the reference genome), the Excel outputs obtained from
380 each BC patient were pooled (keeping patient ID tracked) and clusters were again generated.
381 Candidate genomic insertion sites detected in more than 10% of the samples (corresponding to
382 clusters with more than 6 units for the BC cohort) were identified and subtracted from the patient
383 specific output data. After this filtering step, on average 9 clusters with 5 or more reads are
384 obtained for each individual. Visual inspection using the IGV software revealed that the majority
385 of the remaining candidate insertion sites are generated by the presence of a polyA stretch in the
386 reference sequence that resulted in DNA polymerase slippage during the NGS process.
387 Consequently, these candidate genomic insertion sites were manually tracked and listed, and
388 used for subsequent filtering. Additional file 7 provides a preformatted Excel file to perform the
389 combined filtering steps. After filtering, only half of the BC patient samples presented a

390 candidate insertion event (31/65). One patient presented 4 such events (Additional file 1). To
391 minimize the possibility that these candidate insertion sites would be missed in a subset of
392 individuals because of poorer coverage, all genomic positions identified as insertional target sites
393 during the first screening round in any individual (positions deduced from the presence of a
394 cluster with minimum 5 anchor reads) were re-screened for the presence of a smaller cluster (3 or
395 4 reads) in all individuals (Additional file 1, columns M & N). Further IGV inspection allowed
396 the identification of a TSD in a significant fraction (18/32) of the obtained candidate insertion
397 sites, strongly suggesting that a retrotransposition event had occurred at those genomic locations
398 (Additional file 1, column K). For validation purposes, the presence or absence of each of these
399 18 insertions was verified in each of the 65 patients using IGV (Table 1). The recovery rates
400 obtained IN SILICO and manually using IGV matched very well except for two insertions (the
401 *SET* transcript in *DPP10* and an Alu element in *HSD17B12*), the discrepancies being explainable
402 by the poor coverage at sequencing. For the identification of the inserted sequences, we used the
403 Dfam database of repetitive DNA families[58] (<https://dfam.org/home>) and the Basic Local
404 Alignment Search Tool (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>) (Additional file 1, column I).
405 As only small stretches of nucleotides (about 50 bp long) from the 5' end of the inserted
406 sequences can be deduced from the WES generated data, it was not possible to further define the
407 sub-family of the identified retrotransposons.

408 A screenshot showing read visualization by IGV of the region where the Alu insertion occurred in
409 the positive control sample (*BRCA1*) when using the SeqCap EZ Exome v3.0 kit from Roche for
410 exome capture is presented in Figure 2. A corresponding screenshot showing read visualization
411 when using the SureSelect Human All Exon V6 kit from Agilent is presented in Figure 5. Note
412 that very few discordantly mapped reads are observed when using the Agilent kit, indicating that

413 the detection strategy we described here above is unable to detect retrotransposon insertions
414 when using the WES BAM files generated with this kit, although several other identifiers of an
415 Alu insertion are clearly present and useable for the development of an adapted IN SILICO
416 pipeline.

417

418 **Procedure for cluster calculation starting with the data obtained from RetroSeq:**

419 Save the RetroSeq generated tab. file in Excel format and delete the 2 last columns (the .tab file
420 generated with the positive control sample is provided as “Additional file 4”).

421 Remove all lines not referring to an autosome or to chX (all samples are from females).

422 Sort (A>Z) according to sign in last (6th) column.

423 Select/copy all filled cells and paste in Sheet 1 (“calculate cluster”) of the preformatted Excel file
424 for cluster calculation (Additional file 5).

425 In this Sheet 1, select for “TRUE” in column I (OR) and select/copy all filled cells in columns A
426 to J.

427 Paste the data in Sheet 2 (“simplify cluster”), and in column A, order by “color”. Delete the
428 entire rows corresponding to the cells with pink background (this will remove duplicates).

429 In column O (chr), deselect “blank” and select/copy all informative cells in columns O to U.

430 Paste these data in Sheet 3 (“all clusters”), and deselect “2” in column G (to keep clusters with
431 minimum 3 reads).

432 Select/copy all filled cells and paste in Sheet 4 (annotation list, use the first column for sample
433 ID). The positive control sample will generate 641 clusters (see Additional file 6).

434

435 **Procedure to filter out the clusters of secondary interest (false positives or clusters highly**
436 **recurrent among patients):**

437 Based on the data we obtained from our 65 samples, we generated a table (first sheet of
438 Additional file 7) listing (a) all clusters with minimal 3 reads that were observed in 6 or more of
439 our patients (= clusters present ~10% or more of the samples) and (b) all clusters that resulted
440 from the presence of a polyA stretch in the reference sequence. This table was subsequently used
441 to identify and discard the less relevant clusters from the patient specific cluster lists (as
442 described above). This list is provided with the only intension to familiarize with the described
443 procedure. The content of this list can depend on the type and number of samples used for its
444 generation (type of disease, ethnicity...) but also on the followed wet and dry bench procedures.

445

446 Select/copy all filled cells of Sheet 4 (annotation list) obtained from a patient specific cluster
447 calculation (see above).

448 Open Sheet 1 (comm65>=6 + sample) of the preformatted Excel file for cluster filtering
449 (Additional file 7) and paste the data (in column A) down the preexisting list.

450 Apply Sort (A>Z) on column C (start) and thereafter on column B (chr).

451 Select/copy the generated dataset and paste in Sheet 2 (calculate cluster_com65>=6).

|

452 Select for “0” in column L (“cluster”), for “FALSE” in column K (“OR”) and for “sampleID” in
453 column A.

454 Select/copy/paste clusters with five or more reads to the “results” Sheet 3.

455 The obtained candidate insertion sites (represented by each of the generated clusters in Sheet 3)
456 can be manually validated by IGV inspection of the corresponding BAM files at the indicated
457 positions.

458

459 **Hemi-nested PCR’s for genomic validation and screen**

460 All uncommon Alu element insertions identified with high confidence (presence of a TSD)
461 within or close to exonic gene sequences or promoter sequences (7 in total, see Table 1) were
462 validated by a 2 step PCR. During the first step, and starting with an input of 6.25 ng genomic
463 DNA, primers flanking the suspected integration site are used to generate a 200-300 bp long
464 PCR fragment. This fragment is used as template (after 2000x dilution) for the second step of the
465 hemi-nested PCR wherein one primer of the first step PCR is used in combination with an Alu
466 specific primer. The first step of the nested PCR should work for all genomic DNA samples. The
467 second step will work only when the patient is carrier of the targeted Alu insertion and the
468 appropriate primer combination is used. Both PCR steps were run on a real-time PCR instrument
469 (LightCycler 480 II from Roche). DNA samples from Alu carriers (and negative controls)
470 according to our WES screen were used for validation. During the first step, an amplification
471 with reproducible Ct and melting curve was observed with all samples while during the second
472 step only samples with an Alu insertion gave an amplification signal (with reproducible Ct and
473 melting curve). The obtained PCR fragments were further evaluated by agarose gel

474 electrophoresis and all showed the expected size. The primers used for PCR amplification are
475 shown in Table 6. The Alu specific primer (Alu/Rev) is located in a well conserved region and
476 points to the 5' extremity of the transposon. Consequently, this primer will never allow
477 amplification of the 3' extremity (polyA tail) of the transposon. The choice of the second primer
478 for the second step of the PCR will depend on the orientation of the Alu element compared to the
479 orientation of the targeted gene (see Table1, 2nd column). In order to identify additional carriers
480 of the Alu insertions characterized during the raw WES data screen, the same nested PCRs were
481 performed on genomic DNA from 710 familial BC patients (Table 2).

482

483 **Genotyping at polymorphic sites**

484 To investigate whether an Alu element that integrated into the non-coding regions of a gene
485 affects expression levels of the targeted allele, a polymorphic site allowing allelic discrimination
486 at the cDNA level must first be identified. This was done by IGV inspection of the BAM files
487 from patients that had revealed a polymorphic Alu insertion during our screen (Additional file 1).
488 To increase the number of samples that could be enrolled in the allele expression study, the
489 patients whose genomic DNA revealed an interesting Alu insertion upon mutation specific PCR
490 analysis were genotyped for the corresponding polymorphic site. PCR primers are presented in
491 Table 7. Heterozygosity was determined by High Resolution Melting Curve Analysis.

492

493 **RNA isolation, cDNA synthesis and RT-PCR**

494 RNA isolation was performed using the RNeasy Mini kit (Qiagen) with additional DNase step
495 using manufacturer's protocol. cDNAs were prepared using Superscript II reverse transcriptase

496 (Life Sciences) using hexamer random primers and qPCR was performed using the SybrGreen

497 Master kit (Roche) on a LightCycler 480 II. Primers are presented in Table 8.

498

499 **Sequencing**

500 Purified PCR fragments were Sanger sequenced by the VIB Genetic Service Facility, from the

501 University of Antwerp.

502

503 **DECLARATIONS**

504 Ethics approval and consent to participate: Patient recruitment and blood sampling were
505 performed according to the ethical procedures approved by the institutional ethics committee of
506 the UZ Brussel. For the concerned cases, peripheral blood was collected after obtaining a written
507 informed consent for a broad genomic analysis covering also incidental findings in genes
508 predictive for other diseases.

509

510 Consent for publication: Not applicable

511

512 Availability of data and materials: a dataset allowing the reproduction of the outputs for a
513 positive control is provided as “additional file”. Additional datasets are available on reasonable
514 request

515

516 Competing interests: the authors declare that they have no competing interests

517

518 Funding: financial support from “Kom Op Tegen Kanker” is acknowledged

519

520 Authors’ contribution: SDB and ET conceived the project, designed and performed the
521 experiments, interpreted the data and wrote the manuscript. JDG is responsible for patient
522 recruitment and counselling. All authors reviewed and approved the manuscript.

|

523 Acknowledgements: not applicable

524

525 Authors' information (optional): not applicable

526

527

528

529

530 REFERENCES

- 531 1. Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, Devon K, Dewar K, Doyle M,
532 FitzHugh W, et al: **Initial sequencing and analysis of the human genome.** *Nature* 2001, **409**:860-
533 921.
- 534 2. Jurka J: **Evolutionary impact of human Alu repetitive elements.** *Curr Opin Genet Dev* 2004,
535 **14**:603-608.
- 536 3. Cordaux R, Batzer MA: **The impact of retrotransposons on human genome evolution.** *Nat Rev*
537 *Genet* 2009, **10**:691-703.
- 538 4. Kazazian HH, Moran JV: **Mobile DNA in Health and Disease.** *N Engl J Med* 2017, **377**:361-370.
- 539 5. Rishishwar L, Wang L, Wang J, Yi SV, Lachance J, Jordan IK: **Evidence for positive selection on**
540 **recent human transposable element insertions.** *Gene* 2018, **675**:69-79.
- 541 6. Hancks DC, Kazazian HH: **Roles for retrotransposon insertions in human disease.** *Mob DNA*
542 2016, **7**:9.
- 543 7. Ewing AD, Ballinger TJ, Earl D, Harris CC, Ding L, Wilson RK, Haussler D, Platform BIGSAPA:
544 **Retrotransposition of gene transcripts leads to structural variation in mammalian genomes.**
545 *Genome Biol* 2013, **14**:R22.
- 546 8. Yoder JA, Walsh CP, Bestor TH: **Cytosine methylation and the ecology of intragenomic**
547 **parasites.** *Trends Genet* 1997, **13**:335-340.
- 548 9. Daniel C, Behm M, Öhman M: **The role of Alu elements in the cis-regulation of RNA processing.**
549 *Cell Mol Life Sci* 2015, **72**:4063-4076.
- 550 10. Economou-Pachnis A, Tschlis PN: **Insertion of an Alu SINE in the human homologue of the**
551 **MLVI-2 locus.** *Nucleic Acids Res* 1985, **13**:8379-8387.
- 552 11. Kazazian HH, Wong C, Youssoufian H, Scott AF, Phillips DG, Antonarakis SE: **Haemophilia A**
553 **resulting from de novo insertion of L1 sequences represents a novel mechanism for mutation**
554 **in man.** *Nature* 1988, **332**:164-166.
- 555 12. Payer LM, Burns KH: **Transposable elements in human genetic disease.** *Nat Rev Genet* 2019,
556 **20**:760-772.
- 557 13. Teugels E, De Brakeleer S, Goelen G, Lissens W, Sermijn E, De Grève J: **De novo Alu element**
558 **insertions targeted to a sequence common to the BRCA1 and BRCA2 genes.** *Hum Mutat* 2005,
559 **26**:284.
- 560 14. Peixoto A, Santos C, Pinheiro M, Pinto P, Soares MJ, Rocha P, Gusmão L, Amorim A, van der Hout
561 A, Gerdes AM, et al: **International distribution and age estimation of the Portuguese BRCA2**
562 **c.156_157insAlu founder mutation.** *Breast Cancer Res Treat* 2011, **127**:671-679.
- 563 15. De Brakeleer S, De Grève J, Lissens W, Teugels E: **Systematic detection of pathogenic alu**
564 **element insertions in NGS-based diagnostic screens: the BRCA1/BRCA2 example.** *Hum Mutat*
565 2013, **34**:785-791.
- 566 16. Qian Y, Mancini-DiNardo D, Judkins T, Cox HC, Brown K, Elias M, Singh N, Daniels C, Holladay J,
567 Coffee B, et al: **Identification of pathogenic retrotransposon insertions in cancer predisposition**
568 **genes.** *Cancer Genet* 2017, **216-217**:159-169.
- 569 17. Shahi RB, De Brakeleer S, Caljon B, Pauwels I, Bonduelle M, Joris S, Fontaine C, Vanhoeij M, Van
570 Dooren S, Teugels E, De Grève J: **Identification of candidate cancer predisposing variants by**
571 **performing whole-exome sequencing on index patients from BRCA1 and BRCA2-negative**
572 **breast cancer families.** *BMC Cancer* 2019, **19**:313.
- 573 18. Keane TM, Wong K, Adams DJ: **RetroSeq: transposable element discovery from next-**
574 **generation sequencing data.** *Bioinformatics* 2013, **29**:389-390.

- 575 19. Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, Mesirov JP:
576 **Integrative genomics viewer.** *Nat Biotechnol* 2011, **29**:24-26.
- 577 20. Stewart C, Kural D, Strömberg MP, Walker JA, Konkel MK, Stütz AM, Urban AE, Grubert F, Lam
578 HY, Lee WP, et al: **A comprehensive map of mobile element insertion polymorphisms in**
579 **humans.** *PLoS Genet* 2011, **7**:e1002236.
- 580 21. Witherspoon DJ, Zhang Y, Xing J, Watkins WS, Ha H, Batzer MA, Jorde LB: **Mobile element**
581 **scanning (ME-Scan) identifies thousands of novel Alu insertions in diverse human populations.**
582 *Genome Res* 2013, **23**:1170-1181.
- 583 22. Chen LL, Yang L: **ALU alternative Regulation for Gene Expression.** *Trends Cell Biol* 2017, **27**:480-
584 490.
- 585 23. Spengler RM, Oakley CK, Davidson BL: **Functional microRNAs and target sites are created by**
586 **lineage-specific transposition.** *Hum Mol Genet* 2014, **23**:1783-1793.
- 587 24. Luedi PP, Dietrich FS, Weidman JR, Bosko JM, Jirtle RL, Hartemink AJ: **Computational and**
588 **experimental identification of novel human imprinted genes.** *Genome Res* 2007, **17**:1723-1730.
- 589 25. Greally JM: **Short interspersed transposable elements (SINEs) are excluded from imprinted**
590 **regions in the human genome.** *Proc Natl Acad Sci U S A* 2002, **99**:327-332.
- 591 26. Wang J, Song L, Grover D, Azrak S, Batzer MA, Liang P: **dbRIP: a highly integrated database of**
592 **retrotransposon insertion polymorphisms in humans.** *Hum Mutat* 2006, **27**:323-329.
- 593 27. Sjöblom T, Jones S, Wood LD, Parsons DW, Lin J, Barber TD, Mandelker D, Leary RJ, Ptak J,
594 Silliman N, et al: **The consensus coding sequences of human breast and colorectal cancers.**
595 *Science* 2006, **314**:268-274.
- 596 28. Olsen TK, Panagopoulos I, Gorunova L, Micci F, Andersen K, Kilen Andersen H, Meling TR, Due-
597 Tønnessen B, Scheie D, Heim S, Brandal P: **Novel fusion genes and chimeric transcripts in**
598 **ependymal tumors.** *Genes Chromosomes Cancer* 2016, **55**:944-953.
- 599 29. Sun SC, Ma D, Li MY, Zhang RX, Huang C, Huang HJ, Xie YZ, Wang ZJ, Liu J, Cai DC, et al:
600 **Mutations in C1orf194, encoding a calcium regulator, cause dominant Charcot-Marie-Tooth**
601 **disease.** *Brain* 2019, **142**:2215-2229.
- 602 30. van de Kooij B, Creixell P, van Vlimmeren A, Joughin BA, Miller CJ, Haider N, Simpson CD, Linding
603 R, Stambolic V, Turk BE, Yaffe MB: **Comprehensive substrate specificity profiling of the human**
604 **Nek kinome reveals unexpected signaling outputs.** *Elife* 2019, **8**.
- 605 31. Chivukula RR, Montoro DT, Leung HM, Yang J, Shamseldin HE, Taylor MS, Dougherty GW,
606 Zariwala MA, Carson J, Daniels MLA, et al: **A human ciliopathy reveals essential functions for**
607 **NEK10 in airway mucociliary clearance.** *Nat Med* 2020, **26**:244-251.
- 608 32. Mello SS, Valente LJ, Raj N, Seoane JA, Flowers BM, McClendon J, Bieging-Rolett KT, Lee J,
609 Ivanochko D, Kozak MM, et al: **A p53 Super-tumor Suppressor Reveals a Tumor Suppressive**
610 **p53-Ptpn14-Yap Axis in Pancreatic Cancer.** *Cancer Cell* 2017, **32**:460-473.e466.
- 611 33. Liu X, Sun L, Gursel DB, Cheng C, Huang S, Rademaker AW, Khan SA, Yin J, Kiyokawa H: **The non-**
612 **canonical ubiquitin activating enzyme UBA6 suppresses epithelial-mesenchymal transition of**
613 **mammary epithelial cells.** *Oncotarget* 2017, **8**:87480-87493.
- 614 34. Cai X, Wang X, Cao C, Gao Y, Zhang S, Yang Z, Liu Y, Zhang X, Zhang W, Ye L: **HBXIP-elevated**
615 **methyltransferase METTL3 promotes the progression of breast cancer via inhibiting tumor**
616 **suppressor let-7g.** *Cancer Lett* 2018, **415**:11-19.
- 617 35. Lan Q, Liu PY, Haase J, Bell JL, Hüttelmaier S, Liu T: **The Critical Role of RNA m.** *Cancer Res* 2019,
618 **79**:1285-1292.
- 619 36. Gupta P, Li YR: **Upf proteins: highly conserved factors involved in nonsense mRNA mediated**
620 **decay.** *Mol Biol Rep* 2018, **45**:39-55.
- 621 37. Ranjan A, Iyer SV, Iwakuma T: **Suppressive roles of A3AR and TMIGD3 i1 in osteosarcoma**
622 **maligancy.** *Cell Cycle* 2017, **16**:903-904.

- 623 38. Mollenhauer J, Wiemann S, Scheurlen W, Korn B, Hayashi Y, Wilgenbus KK, von Deimling A,
624 Poustka A: **DMBT1, a new member of the SRCR superfamily, on chromosome 10q25.3-26.1 is**
625 **deleted in malignant brain tumours.** *Nat Genet* 1997, **17**:32-39.
- 626 39. Wu W, Kemp BL, Proctor ML, Gazdar AF, Minna JD, Hong WK, Mao L: **Expression of DMBT1, a**
627 **candidate tumor suppressor gene, is frequently lost in lung cancer.** *Cancer Res* 1999, **59**:1846-
628 1851.
- 629 40. Braidotti P, Nuciforo PG, Mollenhauer J, Poustka A, Pellegrini C, Moro A, Bulfamante G, Coggi G,
630 Bosari S, Pietra GG: **DMBT1 expression is down-regulated in breast cancer.** *BMC Cancer* 2004,
631 **4**:46.
- 632 41. Sorek R, Ast G, Graur D: **Alu-containing exons are alternatively spliced.** *Genome Res* 2002,
633 **12**:1060-1067.
- 634 42. Ding J, Wu W, Yang J, Wu M: **Long non-coding RNA MIF-AS1 promotes breast cancer cell**
635 **proliferation, migration and EMT process through regulating miR-1249-3p/HOXB8 axis.** *Pathol*
636 *Res Pract* 2019, **215**:152376.
- 637 43. Wagner K, Hemminki K, Grzybowska E, Bermejo JL, Butkiewicz D, Pamula J, Pekala W, Försti A:
638 **Polymorphisms in the growth hormone receptor: a case-control study in breast cancer.** *Int J*
639 *Cancer* 2006, **118**:2903-2906.
- 640 44. Canzian F, Cox DG, Setiawan VW, Stram DO, Ziegler RG, Dossus L, Beckmann L, Blanché H,
641 Barricarte A, Berg CD, et al: **Comprehensive analysis of common genetic variation in 61 genes**
642 **related to steroid hormone and insulin-like growth factor-I metabolism and breast cancer risk**
643 **in the NCI breast and prostate cancer cohort consortium.** *Hum Mol Genet* 2010, **19**:3873-3884.
- 644 45. Rudd MF, Webb EL, Matakidou A, Sellick GS, Williams RD, Bridle H, Eisen T, Houlston RS,
645 Consortium G: **Variants in the GH-IGF axis confer susceptibility to lung cancer.** *Genome Res*
646 2006, **16**:693-701.
- 647 46. Tchatchou S, Riedel A, Lyer S, Schmutzhard J, Strobel-Freidekind O, Gronert-Sum S, Mietag C,
648 D'Amato M, Schlehe B, Hemminki K, et al: **Identification of a DMBT1 polymorphism associated**
649 **with increased breast cancer risk and decreased promoter activity.** *Hum Mutat* 2010, **31**:60-66.
- 650 47. Gao C, Devarajan K, Zhou Y, Slater CM, Daly MB, Chen X: **Identifying breast cancer risk loci by**
651 **global differential allele-specific expression (DASE) analysis in mammary epithelial**
652 **transcriptome.** *BMC Genomics* 2012, **13**:570.
- 653 48. Blackburn AC, Hill LZ, Roberts AL, Wang J, Aud D, Jung J, Nikolcheva T, Allard J, Peltz G, Otis CN,
654 et al: **Genetic mapping in mice identifies DMBT1 as a candidate modifier of mammary tumors**
655 **and breast cancer risk.** *Am J Pathol* 2007, **170**:2030-2041.
- 656 49. Ahmed S, Thomas G, Ghossaini M, Healey CS, Humphreys MK, Platte R, Morrison J, Maranian
657 M, Pooley KA, Luben R, et al: **Newly discovered breast cancer susceptibility loci on 3p24 and**
658 **17q23.2.** *Nat Genet* 2009, **41**:585-590.
- 659 50. Decker B, Allen J, Luccarini C, Pooley KA, Shah M, Bolla MK, Wang Q, Ahmed S, Baynes C, Conroy
660 DM, et al: **Targeted Resequencing of the Coding Sequence of 38 Genes Near Breast Cancer**
661 **GWAS Loci in a Large Case-Control Study.** *Cancer Epidemiol Biomarkers Prev* 2019, **28**:822-825.
- 662 51. Babaian A, Mager DL: **Endogenous retroviral promoter exaptation in human cancer.** *Mob DNA*
663 2016, **7**:24.
- 664 52. Jang HS, Shah NM, Du AY, Dailey ZZ, Pehrsson EC, Godoy PM, Zhang D, Li D, Xing X, Kim S, et al:
665 **Transposable elements drive widespread expression of oncogenes in human cancers.** *Nat*
666 *Genet* 2019, **51**:611-617.
- 667 53. Guevara-Aguirre J, Guevara A, Palacios I, Pérez M, Prócel P, Terán E: **GH and GHR signaling in**
668 **human disease.** *Growth Horm IGF Res* 2018, **38**:34-38.
- 669 54. Mir AA, Philippe C, Cristofari G: **euL1db: the European database of L1HS retrotransposon**
670 **insertions in humans.** *Nucleic Acids Res* 2015, **43**:D43-47.

- 671 55. Payer LM, Steranka JP, Yang WR, Kryatova M, Medabalimi S, Ardeljan D, Liu C, Boeke JD,
672 Avramopoulos D, Burns KH: **Structural variants caused by *Alu* insertions are associated with**
673 **risks for many human diseases.** *Proc Natl Acad Sci U S A* 2017, **114**:E3984-E3992.
- 674 56. Li H: **Toward better understanding of artifacts in variant calling from high-coverage samples.**
675 *Bioinformatics* 2014, **30**:2843-2851.
- 676 57. Rishishwar L, Mariño-Ramírez L, Jordan IK: **Benchmarking computational tools for polymorphic**
677 **transposable element detection.** *Brief Bioinform* 2017, **18**:908-918.
- 678 58. Hubley R, Finn RD, Clements J, Eddy SR, Jones TA, Bao W, Smit AF, Wheeler TJ: **The Dfam**
679 **database of repetitive DNA families.** *Nucleic Acids Res* 2016, **44**:D81-89.

680

681

682 to be added to ref35 (bold): ...role of RNA m6A **Methylation in Cancer.**

683 (not accepted by EndNote)

684

685

686

687

688 **FIGURES (titles and legends)**

689

690 Figure 1

691 Title: Schematic representation of the “discordantly mapped read pairs” used for the detection of
692 polymorphic retrotransposons

693 Legend: After genome fragmentation (broken arrows), size selected DNA fragments are
694 sequenced from both ends (paired-end sequencing). DNA fragments having one breakpoint in a
695 gene and the other breakpoint in a retrotransposon not included in the reference human genome
696 will generate read pairs (grey filled arrows) that do not map at nearby positions on the reference
697 genome after alignment. The discordantly mapped read pairs obtained by analysis of the full
698 exome of a patient (traceable in the corresponding BAM file) are further selected for the
699 presence of a long polyA stretch in one of the mate reads. An additional landmark for
700 retrotransposition is the presence of a target site duplication (TSD) flanking the transposon.

701

702

703 Figure 2

704 Title: Screenshots of IGV outputs visualizing the Alu element insertion in the positive control
705 sample

706 Legend: The screenshots are restricted to the genomic region where the Alu element insertion
707 occurred (c.1739_1740insAlu in *BRCAl*) in the positive control sample. The BAM file used for
708 IGV visualization was generated from WES data obtained with the exome capture kit SeqCap EZ

|

709 Exome v3.0 from Roche. Group alignment is by chromosome of mate. Note the presence of a
710 high number of reads (22) whose mate reads are located on another chromosome (encircled and
711 labelled “A”). All those BRCA1 specific reads point in the same direction, the 3’ extremity of
712 the Alu insertion that ends with a long polyA stretch (3 of these 22 reads already end with a
713 small polyA stretch). Consequently, the mates of these 22 reads have a high probability to
714 contain a long polyA stretch and to be picked up with our detection strategy. Five such
715 discordantly mapped read pairs are sufficient to retain this particular genomic region as a
716 candidate insertion site. Note also that a number of correctly mapped read pairs (labelled “B” and
717 “C”) start with a sequence that do not correspond to the *BRCA1* sequence. The reads labelled
718 “B” confirm the presence (and position) of a polyA stretch adjacent to the *BRCA1* sequence,
719 while the reads labelled “C” allow the identification of the inserted sequence most upstream of
720 the polyA stretch, in this case an Alu element. Inserted retrotransposon sequences can be
721 identified making use of the Dfam [58] database (<https://dfam.org/home>), processed pseudogenes
722 with BLAST (<https://blast.ncbi.nlm.nih.gov/Blast.cgi>). Finally, and characteristic for
723 retrotransposition events, the inserted element is flanked by a short repeated sequence (TSD,
724 17bp long in this case) as is indicated with label “D”. Often, a sharp increase of the coverage at
725 the level of the TSD is observed (less pronounced in this example, see label “E”).

726

727

728 Figure 3

729 Title: Alu element insertions in *UPF2* and *NEK10* do not induce allelic up or down regulation

730 Legend: Sanger sequencing experiments were performed on RNA material obtained from blood
731 leucocytes. The chromatograms show a nucleotide region containing a polymorphic position (a
732 SNP, indicated by an arrow) for which all tested samples were heterozygous at the genomic
733 level. (A) In addition to be heterozygous for c.1539A>C (rs11595168), all 5 samples were also
734 heterozygous for the Alu insertion in the intronic region of *UPF2* (c.3408+98_+111insAlu). The
735 c.1539C allele is very rare and was not detected among available control samples. (B) In addition
736 to be heterozygous for c.1674A>G (rs674303), 1 of the 4 samples was also heterozygous for the
737 Alu insertion in the intronic region of *NEK10* (c.26+891_908insAlu).

738

739

740 Figure 4

741 Title: Alu element insertions in *GHR* and *PTPN10* modulate allelic expression

742 Legend: Sanger sequencing experiments performed on RNA material obtained from blood
743 leucocytes. The chromatograms show a nucleotide region containing a polymorphic position
744 (arrow) for which the tested samples were heterozygous at the genomic level. (A) In addition to
745 be heterozygous for c.558A>G (rs2397118), 3 of the 5 samples were also heterozygous for the
746 Alu insertion in the 3'UTR of *GHR* (c.*+413_+422insAlu). Allele silencing is observed in all
747 Alu mutation carriers. (B) In addition to be heterozygous for c.978A>G (rs7550799), 2 of the 5
748 samples were also heterozygous for the Alu insertion in the intronic region of *PTPN14*
749 (c.10066+35_+44insAlu). Note that the two Alu mutation carriers express both alleles well,
750 while one control sample (sample *PTPN14*_ctrl1) expresses almost exclusively the G allele, the 2
751 other control samples showing a deficit in the G allele.

752

753

754 Figure 5

755 Title: IGV visualization of the control Alu insertion using BAM file generated with the Agilent
756 kit

757 Legend: Screenshots of IGV outputs restricted to the genomic region where the Alu element
758 insertion occurred (c.1739_1740insAlu in *BRCAL*) in the positive control sample. The BAM file
759 used for IGV visualization was generated from WES data obtained with the SureSelect Human
760 All Exon V6 kit (exome capture kit) from Agilent. Group alignment is by chromosome of mate.
761 Note that this exome capture kit generates much less discordantly mapped read pairs than the kit
762 from Roche (2 versus 22, compare with Figure 2). Accordingly, the required number (five) of
763 discordantly mapped read pairs containing a long polyA track will not be reached and the
764 insertion will not be revealed with our detection strategy. Note that all other characteristics
765 observed for a retrotransposon insertion when using the exome capture kit from Roche (see
766 Figure 2) are also observed when using the kit from Agilent.

767

768 Table 1: List with the 18 high confidence retrotransposition mediated insertions detected in 65
769 familial BC patients

Targeted gene (region)	Location TSD (RefSeq); orientation of the insertion compared to the targeted gene	Transposed sequence	5' extremity of the transposed sequence ^a	# carriers (cl_5) ^b	# carriers (cl_3) ^c	# carriers IGV ^d	reported in dbRIP?
BRCA1 (exon)	c.1723_1739 (NM_007294); sense	Alu	positive control	1	1	1	yes
GHR (3'UTR)	c.*+413_+422 (NM_000163.5); antisense	Alu ^e	1_42 of AluYe6	3	3	5	yes
GSTA5 (promoter)	1,3 Kb up transcription start (NM_153699.1); antisense	Alu ^e	1_59 of AluYb9	4	4	4	yes
METTL3 (intron)	c.1456+69_+84 (NM_019852.5); antisense	Alu ^e	1_58 of AluYd8	1	1	1	no
NEK10 (intron)	c.26+891_+908 (CCDS77713); antisense	Alu ^e	2_53 of AluYe6	1	1	3	no
PTPN14 (intron)	c.10066+35_+44 (NM_005401.5);sense	Alu ^e	1_57 of AluYd8	3	4	4	no
UPF2 (intron)	c.3408+98_+111 (NM_080599.2); antisense	Alu ^e	2_53 of AluYe6	3	3	3	no
ZNF442 (exon)	c.1402_1413 (NM_030824.3); antisense	Alu ^e	1_47 of AluSx3	1	1	1	no
TMIGD3 (intron)	c.458-623_-638 (NM_020683.7); antisense	Alu	2_55 of AluYe6	1	4	7	yes
intergenic	closest gene at 10kb (GPR42)	Alu	1_57 of AluYd8	2	2	2	no
intergenic	2 genes at 30Kb (MIF-AS1 & GSTT2B)	Alu	1_57 of AluYb9	1	5	5	yes
HSD17B12 (3'UTR)	c.*+642_+650 (NM_016142.3); antisense	Alu	1_46 of AluYe6	1	4	29	yes
intergenic	closest gene at 85kb (CLVS1)	L1	2247_2296 of L1P1_orf2	1	4	5	no
DMBT1 (intron)	c.1459+48_+63 (NM_007329.2); antisense	L1	4_54 of L1HS	1	1	1	no
UBA6 (intron)	c.1097+17_+25 (NM_018227.6); sense	SVA	731_778 of SVA_F	1	1	1	no
C1orf194 (exon)	c.58_70 (NM_001122961.1); sense	UQCR10	starts at c.-15 (NM_013387.4)	1	4	7	pseudogene
PSMA7 (intron)	c.654+95_+107 (NM_002792.4); antisense	MRPL50	starts at c.-7 (NM_019051.3)	1	1	1	pseudogene
DPP10 (intron)	c.442-70581_-70595 (NM_001321907.1); antisense	SET	starts at c.*+199 (NM_003011)	1	4	35	pseudogene

770

771 Strong candidate integration sites are retained when generated by clusters counting at least 5
772 polyA containing discordantly mapped read pairs, when present in less than 10% of the samples,
773 and when a TSD is deducible upon IGV inspection in at least one patient. ^a indicates the number
774 of nucleotides available to identify retrotransposon type (using the Dfam database [58]). As only
775 about 50 nucleotides are available for transposon identification, it was not possible to define sub-
776 classes. ^b indicates the number of patients (out of 65) who carry the mutation based on clusters
777 counting at least 5 discordantly mapped read pairs. ^c indicates the number of patients who carry
778 the same mutation based on clusters counting at least 3 discordantly mapped read pairs. ^d
779 indicates the number of patients harboring the same mutation based upon IGV inspection. Note
780 that for two insertions (SET in *DPP10* and Alu in *HSD17B12*) a much higher number of carriers
781 is revealed upon IGV inspection. This discrepancy results from the poor coverage at sequencing.
782 ^e indicates the Alu insertions for which a genomic hemi-nested PCR has been developed and
783 applied for validation (see Methods).

784

785 Table 2: Results from insAlu genotyping experiments performed on 710 non related familial BC
786 patients

Alu targeted gene	# carriers out of 710
UPF2	107
PTPN14	10
GHR	17
ZNF442	1
METTL3	1

787 The primers used for mutation specific PCR amplification are presented in Table 6

788

789

790

791

792

793 Table 3: List with the 14 polyA associated insertions for which a TSD could not be defined.

Targeted gene (region)	Insertion type	Position of insertion in cDNA (RefSeq)	Position of insertion in hg19	Retrotransposon sequence in mate?	# carriers cl_5	# carriers cl_3	reported in dbRIP?
ANAPC16 (3'UTR)	polyA	c.*+397 (NM_001242546.1)	chr10:73993271	no	2	2	yes (Alu)
ARHGAP1 (intron)	polyA	c.317+71 (NM_004308)	chr11:46709653	yes (Alu)	1	1	no
AUH (intron)	polyA	c.599-122 (NM_001698)	chr9:94058481	yes (L1)	1	4	yes (L1)
CLEC1B (intron)	polyA	c.438+306 (NM_016509)	chr12:10149139	no	1	2	no
DHRS1 (intron)	polyA	c.374+111 (NM_138452)	chr14:24765604	yes (Alu)	1	1	no
DYNC2H1 (intron)	polyA	c.6139+43 (NM_001377)	chr11:103048592	yes (Alu)	1	1	yes (Alu)
ETF1 (3'UTR)	polyA	c.*+2210 (NM_004730)	chr5:137841784	no	1	1	no
GCNT1 (3'UTR)	polyA	c.*+2545 (NM_001490)	chr9:79121129	yes (Alu)	1	1	no
GSAP (intron)	polyA	c.1491+495 (NM_017439)	chr7:76980766	yes (Alu)	1	1	yes (Alu)
MEOX2 (intron)	polyA	c.690+95 (NM_005924)	chr7:15666276	yes (Alu)	1	1	yes (Alu)
MIER1 (3'UTR)	polyA	c.*+1193 (NM_020948)	chr1:67453349	no	1	1	no
USP38 (intron)	polyA	c.1210-89 (NM_032557)	chr4:144127097	yes (Alu)	1	1	yes (Alu)
intergenic	polyA	LINC0251 at 46kb	chr6:114017537	no	3	4	no
intergenic	polyA	SPIN1 at 6,1kb	chr9:91099726	no	1	3	yes (Alu)

794 Note that for 8 of the 14 polyA containing insertion sites, a partial retrotransposon sequence

795 could be traced in the mate of at least one discordantly mapped read upon IGV inspection (5th

796 column). Screening the database for Retrotransposon Insertion Polymorphisms in Human

797 (dbRIP) at the corresponding positions confirmed the presence of a retrotransposon at 7 of the 14

798 positions (8th column).

799

800 Table 4: Number of hits obtained in PubMed when searching for the polyA targeted genes.

PubMed domain	title	title/abstract	title/abstract	title/abstract
Search terms used	"all gene names"	"all gene names"	"all gene names" AND "cancer"	"all gene names" AND "breast cancer"
gene (synonyms)^a				
ANAPC16 (APC16, C10orf104, CENP-27, FLJ33728, bA570G20.3)	5	5	0	0
ARHGAP1 (CDC42GAP, RhoGAP , p50rhoGAP)	31	93	15	4
AUH	-	-	-	-
CLEC1B (CLEC2)	9	49	6	0
DHRS1 (FLJ25430, MGC20204, SDR19C1)	1	4	1	0
DYNC2H1 (DHC1b, DHC2, DNCH2, DYH1B, hdhc11)	26	86	4	0
ETF1 (ERF , ERF1, RF1, SUP45L1, TB3-1)	185	845	21	5
GCNT1 (C2GNT, NACGT2, NAGCT2)	8	124	35	1
GSAP (LOC54103, PION)	6	31	1	0
MEOX2 (GAX, MOX2)	85	254	22	1
MIER1 (KIAA1610, MI-ER1, hMI-ER1)	10	15	2	2
USP38 (HP43.8KD, KIAA1891)	6	10	3	1
intergenic: LINC0251 at 46kb	0	0	0	0
intergenic: SPIN1 at 6kb (SPIN , TDRD24)	18	34	14	4

801 This table shows the number of hits obtained in PubMed when using as search term a specific
802 gene or its synonym (= "all gene names") in combination with "cancer" or "breast cancer" and
803 when restricting the search to particular domains (title or abstract). ^agene names generating
804 excessive number of false positive hits were removed (strikethrough). Search was performed on
805 14/01/2020

806

807 Table 5: Number of hits obtained in PubMed when searching for high confidence targeted genes.

PubMed domain Search terms used	title	title/abstract	title/abstract	title/abstract
	"all gene names"	"all gene names"	"all gene names" AND "cancer"	"all gene names" AND "breast cancer"
gene (synonyms)^a				
BRCA1 (BRCC1, FANCS, PPP1R53, RNF53); positive control	6503	14198	11505	7402
GHR (GHBP)	208	2416	129	35
GSTA5	2	29	7	2
METTL3 (MTA , MT-A70, Spo8)	80	251	65	6
NEK10 (FLJ32685)	1	17	11	8
PTPN14 (PEZ)	33	98	37	6
UPF2 (DKFZP434D222, KIAA1408, RENT2, smg-3)	32	189	5	0
ZNF442 (FLJ14356)	0	0	0	0
TMIGD3 (AD026)	2	2	1	0
intergenic: GPR42 at 10 kb (FFAR3L, GPR41L, GPR42P)	2	4	0	0
intergenic: MIF-AS1 at 30 kb (LOC284889, MIF-AS)	4	5	2	1
intergenic: GSTT2B at 30kb (GSTT2P)	2	8	2	0
intergenic: CLVS1 at 85kb (CGorf212L, CRALBPL, MGC34646, RLBP1L1)	2	7	1	0
DMBT1 (GP340, Gp-340, SALSA , hensin, muclin, vomeroglandin)	166	389	68	7
UBA6 (FLJ10808, UBE1L2)	14	32	5	1
C1orf194	1	2	0	0
PSMA7 (C6 , HSPC , RC6-1, XAPC7)	19	60	18	0
UQCR10 (HSPC051, QCR9, UCCR7.2, UCRC) transcript in C1orf194	3	35	4	0
MRPL50 (FLJ20493, MRP-L50) transcript in PSMA7	0	2	0	0

808 This table shows the number of hits obtained in PubMed when using as search term a specific
809 gene or its synonym (= "all gene names") in combination with "cancer" or "breast cancer" and
810 when restricting the search to specific domains (title or abstract). ^agene names generating
811 excessive number of false positive hits were removed (strikethrough). Search was performed on
812 13/01/2020.

813

814

815 Table 6: Primers used for the detection of 7 Alu element insertions in genomic DNA

Primer name	Primer sequence	wt fragment size (bp)
PTPN14/Gen/F1	GTGTGGTGAGCACTACTCGG	205
PTPN14/Gen/R1	CTGACAGTCTAGGCCTCCAC	
GHR/Gen/F1	AAATCAGGTGGCTTTTGGCGG	275
GHR/Gen/R1	AGAGGGGTATACCAACTGC	
GSTA5/Gen/F1	TGGTATAAACGGTGGTGGCA	286
GSTA5/Gen/R1	TAGCTTCAACAGGCACAA	
NEK10/DNA/F1	AATGGCAGCTTACTTCACGG	240
NEK10/DNA/R1	AGGAATTGTCTTCTCACCTTCTTT	
UPF2/E17/F	GGAAAATGAAACCGATGAAGA	272
UPF2/I17/R	GGCAAAGCCTTTTAGTATTGA	
METTL3/Gen/F1	CCTGTCTGCTCAGAAAACCTCG	249
METTL3/Gen/R1	AGCGGATATCACACAGATCCA	
ZNF442/Gen/F1	GCCTCCGTATTTCTAGTTCCC	220
ZNF442/Gen/R1	TTTCAGCCATGTGAGTCCTTCT	
Alu/Rev	TTTTTAGTAGAGACGGGGTTTC	

816

817

818 Table 7: Primers used for genotyping at polymorphic positions (for RNA study)

Primer name	Primer sequence	fragment size (bp)	SNP
PTPN14/ex11/F	CCTCCCTGTCTATCTGGTCTTT	298	rs7550799
PTPN14/ex11/R	CAAGGGGGACAAAACGGGTG		c.978A>G; p.Arg326=
GHR/ex6/F	CCACCATTGCCCTCAACTG	187	rs6179
GHR/ex6/R	AGGTGTAGCAACTCTTACCATT		c.558A>G; p.Gly186=
GSTA5/i2/F	GGGTGGGTTTCATAGACACTTCATA	183	rs2397118
GSTA5/RNA/R1	AGGGCTCTCTCCTTCATGTCT		c.163G>A; p.Val55Ile
NEK10/ex20/F1	AGGTTAGAAAGCATAGTGGTCAAA	222	rs674303
NEK10/ex20/R1	TGGACAAGAGCACCACAACCT		c.1674A>G; p.Gly558=
UPF2/ex6/F1	ACTTTTGTGACGAGGATTTTCA	274	rs11595168
UPF2/ex6/R1	GGCAAGAGCTCTGTGTGTA		c.1539A>C; p.Ser513=

819

820

821 Table 8: Primers used for cDNA amplification (and Sanger sequencing)

Primer name	Primer sequence	fragment size (bp)
PTPN14/RNA/F1	GTTGTTTGCCACACGACAC	308
PTPN14/RNA/R1	AGTTGACGCTGTGAACGCTA	
GHR/RNA/S1	ATTGCCCTCAACTGGACTTTA	241
GHR/RNA/R1	TGGATCTCACACGCACTTCA	
GSTA5/RNA/F1	GCTGCAGCTGGAGTAGAGTT	203
GSTA5/RNA/R1	AGGGCTCTCTCCTTCATGTCT	
NEK10/RNA/F2	CAGAGCCTTGAGATTTCTCTCAG	386
NEK10/RNA/R2	TTACGCTGCTGTCTCGATCTTT	
UPF2/RNA/F1	TTTGTCCAGCCATCTTGTT	250
UPF2/RNA/R1	GTGCTGGCTTCTCATCTTC	

822

823

824 **ADDITIONAL MATERIAL > INFORMATION LIST**

825

826 Additional file 1 (Excel file)

827 Title: Strong candidates retrotransposon mediated insertion sites

828 Description: This table lists all strong candidate insertion sites detected with the RetroSeq
829 software (minimum 5 discordantly mapped read pairs pointing to the same genomic location)
830 and further filtered using Excel software to remove false positives and insertions occurring in
831 more than 10% of the investigated patient samples.

832

833 Additional file 2 (Word file)

834 Title: Retrotransposon targeted genes with no TSD identified

835 Description: For each gene presumably targeted by a retrotransposition event that could not be
836 confirmed by the presence of a TSD, a short description is given. The gene/protein descriptions
837 provided by “Entrez Gene Summary” (<https://www.ncbi.nlm.nih.gov/gene>), by “GeneCards
838 Summary” (<https://www.genecards.org>) and by “UniProtKB/Swiss-Prot Summary”
839 (<https://uniprot.org>) are presented, when available.

840

841 Additional file 3 (Word file)

842 Title: Retrotransposon targeted genes with TSD identified

|

843 Description: For each gene targeted by a retrotransposition event that could be confirmed by the
844 presence of a TSD, a short description is given. The gene/protein descriptions provided by
845 “Entrez Gene Summary” (<https://www.ncbi.nlm.nih.gov/gene>), by “GeneCards Summary”
846 (<https://www.genecards.org>) and by “UniProtKB/Swiss-Prot Summary” (<https://uniprot.org>) are
847 presented, when available.

848

849 Additional file 4 (Excel file)

850 Title: RetroSeq generated list of anchor reads when using the positive control sample

851 Description: This table lists all discordantly mapped read pairs (19.674 in total) selected by the
852 RetroSeq software when analyzing the positive control sample (patient harboring the
853 c.1739_1740insAlu mutation in BRCA1). The input BAM file was generated using the exome
854 capture kit SeqCap EZ Exome v3.0 from Roche for library preparation. To be selected, one read
855 of the discordantly mapped read pair must contain a long polyA stretch, while the position of the
856 mate of this read on the reference genome is used for ordering the RefSeq selected read pairs.

857

858 Additional file 5 (Excel file)

859 Title: Excel file for cluster calculation

860 Description: This pre-formated Excel file can be used to extract potential retrotransposon
861 insertion sites from lists generated with the RetroSeq software (e.g. the list presented in
862 Additional file 4)

863

|

864 Additional file 6 (Excel file)

865 Title: Output of cluster calculation for the positive control sample

866 Description: This table lists all potential retrotransposon integration sites observed in the positive
867 control sample (641 in total). Integration sites are retained only when 3 or more anchor reads
868 cluster in the same genomic interval (max 300bp long).

869

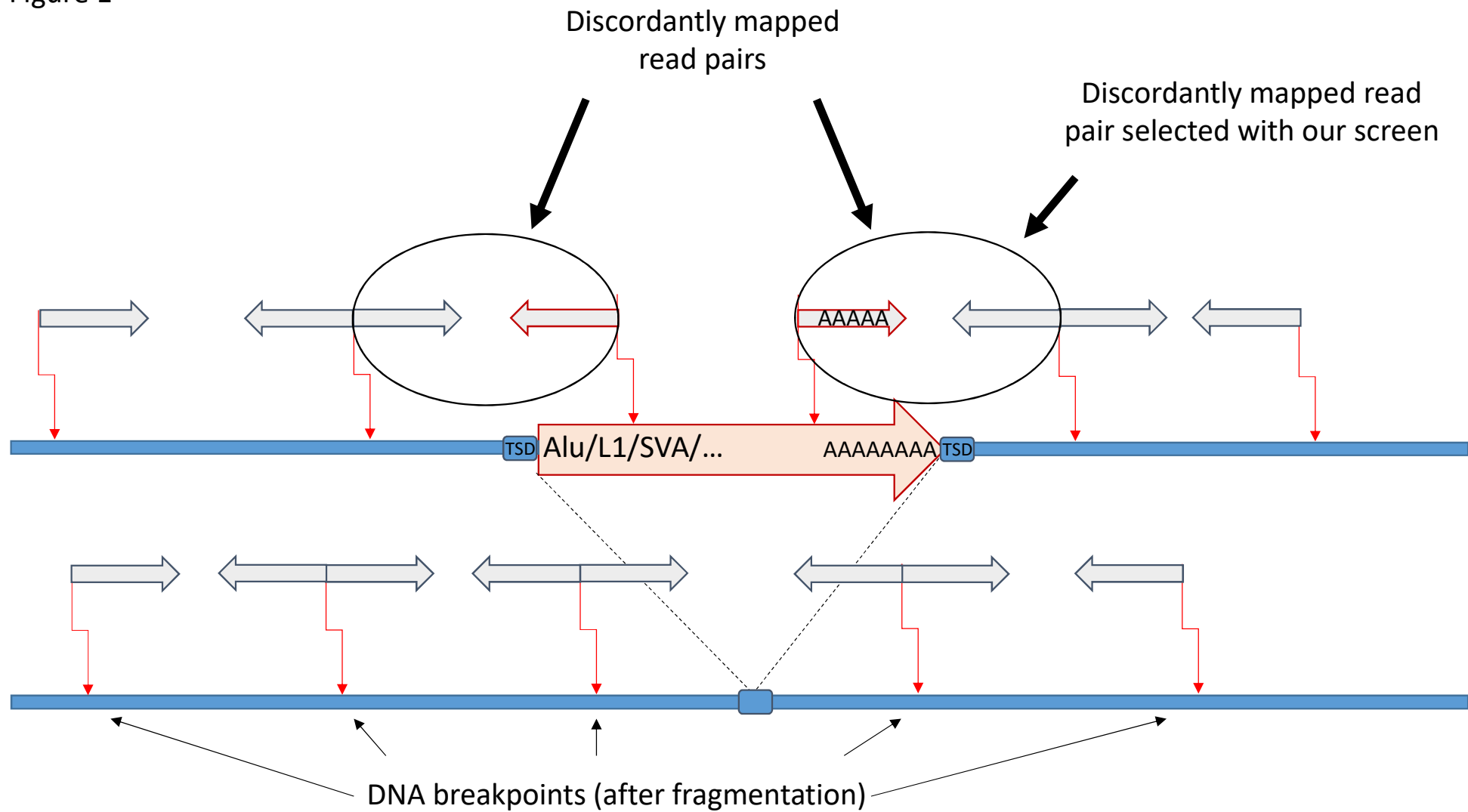
870 Additional file 7 (Excel file)

871 Title: Excel file for cluster filtering

872 Description: This pre-formated Excel file allows further filtering of the potential retrotransposon
873 integration sites (output of Additional file 5). Highly recurrent integration sites (occurring in
874 more than 10% of the investigated population) or false positives are removed during this filtering
875 step.

876

Figure 1



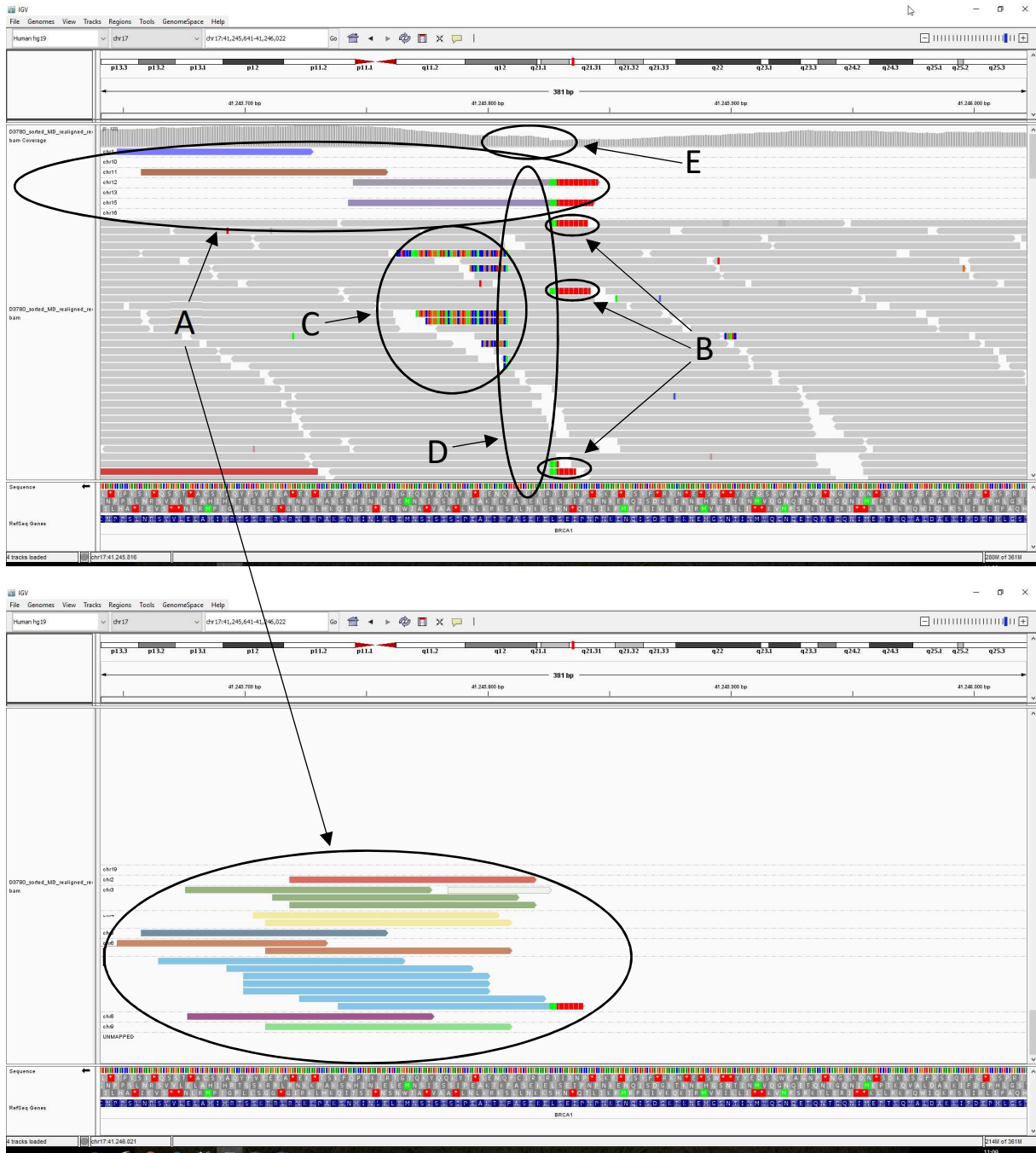


Figure 2

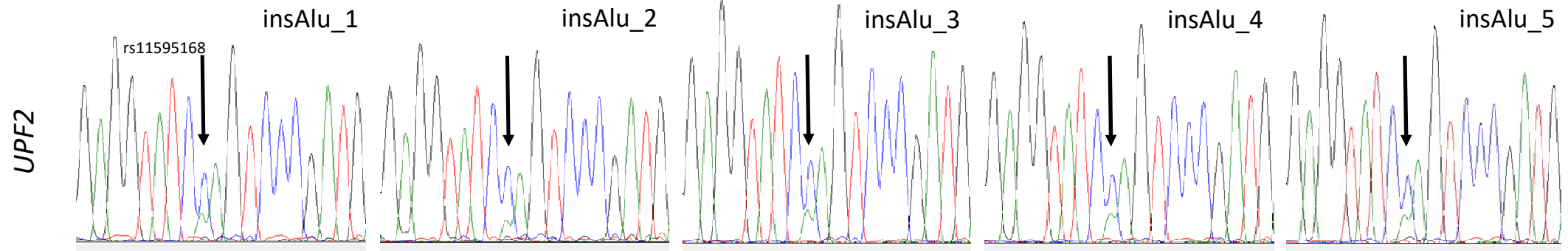
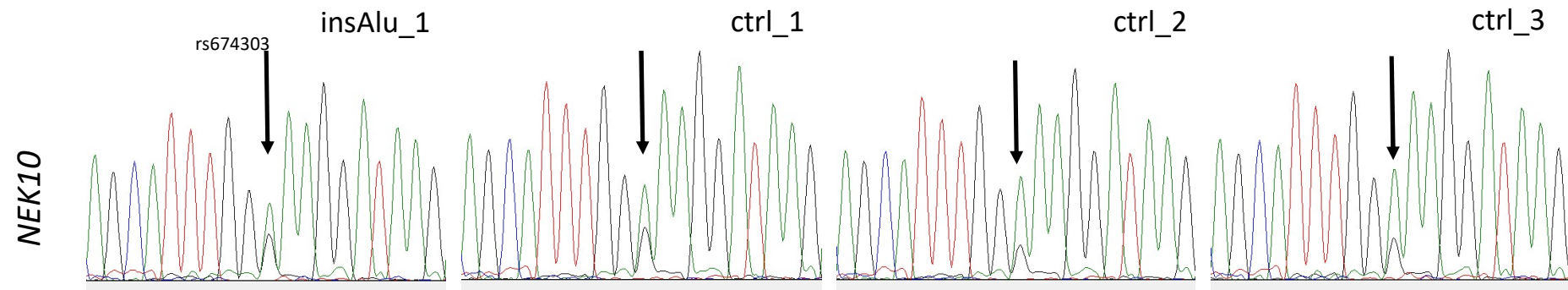
A**B**

Figure 3

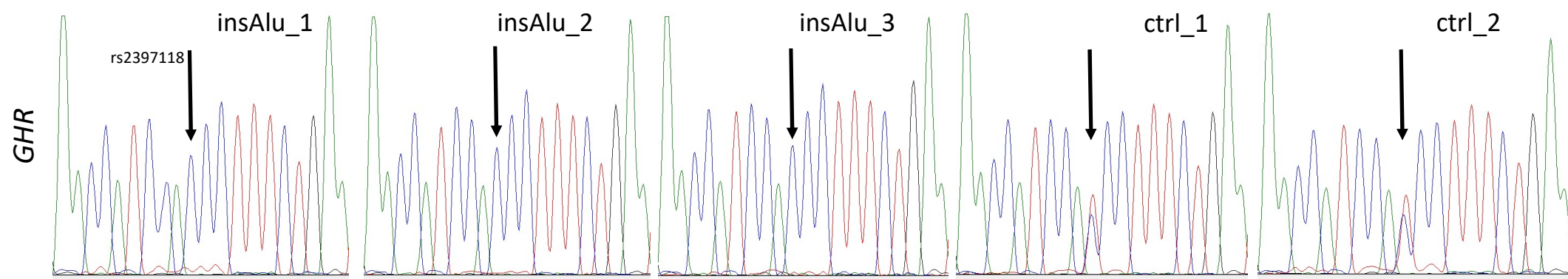
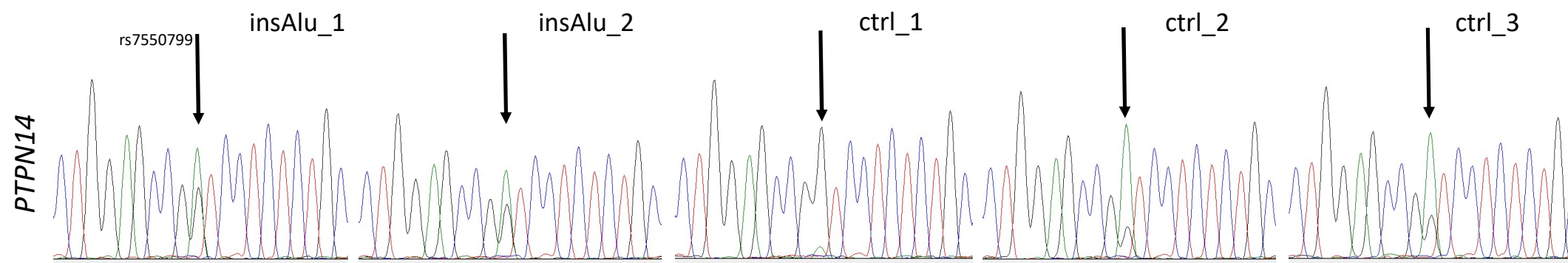
A**B**

Figure 4

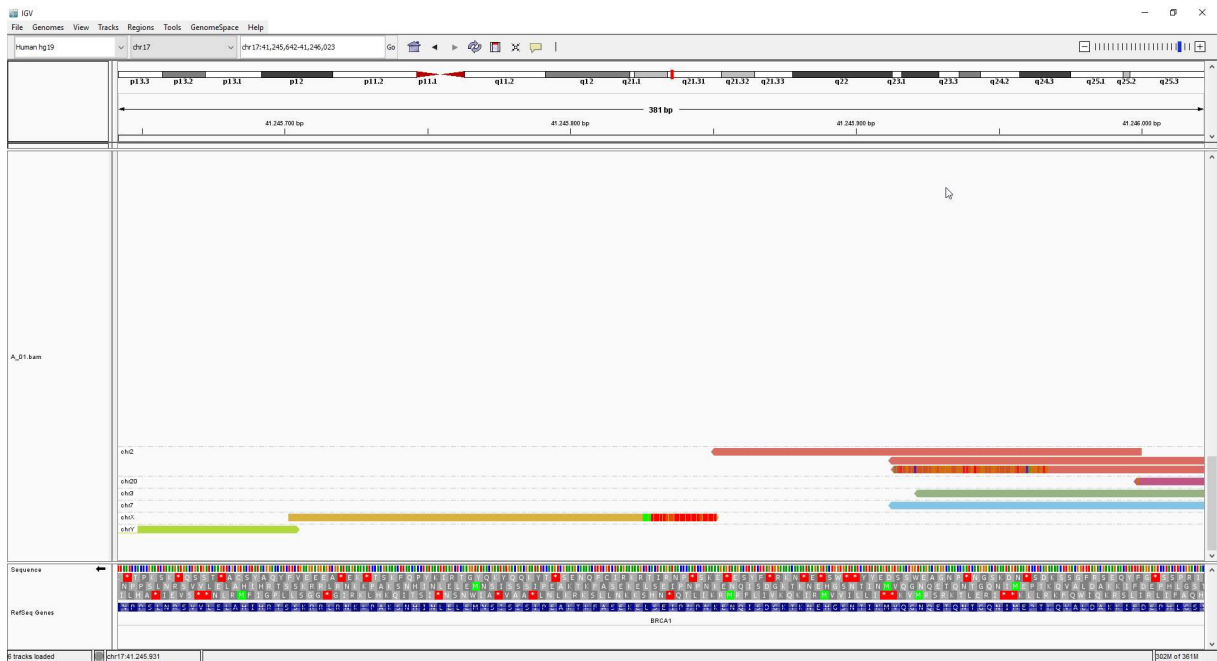


Figure 5