# LinkedSV for detection of mosaic structural variants from linked-read exome and genome sequencing data

4   Li Fang[1], Charlly Kao[2], Michael V Gonzalez[2], Fernanda A Mafra[2], Renata Pellegrino da Silva[2],

5   Mingyao Li[3], Sören Wenzel[4], Katharina Wimmer[4], Hakon Hakonarson[2,5], Kai Wang[1,6*]

7   [1] Raymond G. Perelman Center for Cellular and Molecular Therapeutics, Children's Hospital of

8   Philadelphia, PA 19104, USA

9   [2] Center for Applied Genomics, Children's Hospital of Philadelphia, PA 19104, USA

10   [3] Department of Biostatistics, University of Pennsylvania, Philadelphia, PA 19104, USA

11   [4] Section for Human Genetics, Department of Medical Genetics, Molecular and Clinical

12   Pharmacology, Medical University Innsbruck, Innsbruck, Austria

13   [5] Department of Pediatrics, University of Pennsylvania, Philadelphia, PA 19104, USA

14   [6] Department of Pathology and Laboratory Medicine, University of Pennsylvania, Philadelphia,

15   PA 19104, USA

16   * Email: wangk@email.chop.edu

## Abstract

18    Linked-read sequencing provides long-range information on short-read sequencing data by

19    barcoding reads originating from the same DNA molecule, and can improve the detection and

20    breakpoint identification for structural variants (SVs). We present LinkedSV for SV detection on

21    linked-read sequencing data. LinkedSV considers barcode overlapping and enriched fragment

22    endpoints as signals to detect large SVs, while it leverages read depth, paired-end signals and

23    local assembly to detect small SVs. Benchmarking studies demonstrates that LinkedSV

24    outperforms existing tools, especially on exome data and on somatic SVs with low variant allele

25    frequencies. We demonstrate clinical cases where LinkedSV identifies disease causal SVs from

26    linked-read exome sequencing data missed by conventional exome sequencing, and show

27    examples where LinkedSV identifies SVs missed by high-coverage long-read sequencing. In

28    summary, LinkedSV can detect SVs missed by conventional short-read and long-read

29    sequencing approaches, and may resolve negative cases from clinical genome/exome sequencing

30    studies.

31

32

## Introduction

Genomic structural variants (SVs) have been implicated in a variety of phenotypic diversity and human diseases[1]. Several approaches such as split-reads [2, 3], discordant read-pairs [3, 4], and assembly-based methods [5, 6] have been developed for SV discovery from short reads. However, reliable detection of SVs from these approaches remains challenging. The split-reads and discordant read-pairs approaches require that the breakpoint-spanning reads/read-pairs are sequenced and confidently mapped. Genomic rearrangements are often mediated by repeats and thus breakpoint junctions of SVs are very likely to reside in repetitive regions [7, 8, 9]. Therefore, the breakpoint-spanning reads/read-pairs may be multi-mapped and have low mapping qualities. It is also difficult to perform assembly at repeat regions. Long-read sequencing such as SMRT sequencing and Nanopore sequencing are better for SV detection [10, 11], but their application is limited by the higher cost and per-base error rate.

Linked-read sequencing technology developed by 10X Genomics combines the throughput and accuracy of short-read sequencing with the long-range information. In this approach, nanogram amounts of high-molecular weight (HMW) DNA molecules are dispersed into more than 1 million droplet partitions with different barcodes by a microfluidic system [12]. Thus, only a small number of HMW DNA molecules (~10) are loaded per partition [13]. The HMW DNA molecules can be up to several hundred kilobases in size and have a length-weighted mean DNA molecule length of about 50 kb. Within an individual droplet partition, HMW DNA molecules are primed

52    and amplified by primers with a partition-specific barcode. The barcoded DNA molecules are

53    released from the droplets and sequenced by standard Illumina paired-end sequencing [12]. The

54    sequenced short reads derived from the same HMW DNA molecule can be linked together,

55    providing long-range information for mapping, phasing and SV calling. In addition, linked-read

56    whole exome sequencing (WES) has also been developed [12], which provides an attractive and

57    efficient option for clinical genetic testing.

58    In linked-read sequencing data, barcode similarities between any two nearby genome locations

59    are very high, because the reads tend to originate from the same sets of HMW DNA molecules.

60    In contrast, barcode similarities between any two distant genome locations are very low, because

61    the reads of the two genome locations originate from two different sets of HMW DNA molecules

62    and it is highly unlikely that two different sets of HMW DNA molecules share multiple barcodes.

63    Thus, the presence of multiple shared barcodes between two distant locations indicates that the

64    two distant locations are close to each other in the alternative genome [14]. A few pipelines and

65    software tools have adopted this principle to call SVs from linked-read sequencing data, such as

66    Longranger [12], GROC-SVs [14], NAIBR [15]. Longranger is the official pipeline developed by 10X

67    Genomics. Longranger bins the genome into 10 kb windows and finds the barcodes of high

68    mapping quality reads within each window. A binomial test is used to find all pairs of regions

69    that are distant and share more barcodes than what would be expected by chance. A sophisticated

70    probabilistic model is used to assign a likelihood and remove low quality events[12]. GROC-SVs

71  uses a similar method to find candidate SV loci but performed assembly to identify precise

72  breakpoint locations. GROC-SVs also provides functionality to interpret complex SVs [14].

73  NAIBR detects structural variants using a probabilistic model that incorporates signals from both

74  linked-reads and paired-end reads and into a unified model [15].

75  However, SV detection from linked-read datasets is still in the early stage. The available SV

76  callers face challenges if we want to detect: i) SVs from targeted region sequencing (e.g.WES); ii)

77  somatic SVs in cancer or somatic mosaic SVs that have low variant allele frequencies (VAFs,

78  also known as variant allele fractions); iii) SVs of which the exact breakpoints have no coverage

79  or are located in repeat regions. In this study, we introduce LinkedSV, a novel computational

80  method and software tool for linked-read sequencing, which aims to address all the above

81  challenges. LinkedSV detects large SVs using two types of evidence and quantifies the evidence

82  using a novel probabilistic model. It also leverages read depth, paired-end signals and local

83  assembly to detect small deletions. We evaluated the performance of LinkedSV on both whole-

84  genome and whole-exome sequencing data sets. In each case, LinkedSV outperformed other

85  existing tools, including Longranger, GROC-SVs and NAIBR, especially on exome data and on

86  somatic SVs with low variant allele frequencies. We additionally demonstrated clinical cases

87  where LinkedSV identified disease causal SVs from linked-read exome sequencing data missed

88  by conventional exome sequencing, and showed examples where LinkedSV identifies SVs

89  missed by high-coverage long-read sequencing.

90

## **Results**

**Illustration of two types of evidence near SV breakpoints**

93    Two types of evidence may be introduced while a genomic rearrangement happens: 1) reads

94    from one HMW DNA molecule which spans the breakpoint being mapped to two genomic

95    locations and 2) reads from two distant genome locations that get mapped to adjacent positions.

96    Both types of evidence can be used for SV detection.

97    First, we describe the signals of type 1 evidence. After reads mapping, the original HMW DNA

98    molecules can be computationally reconstructed from the sequenced short reads using their

99    barcodes and mapping positions. In order to distinguish them from the physical DNA molecules,

100   we use fragments to refer to the computationally reconstructed DNA molecules. A fragment has

101   a left-most mapping position, which we call L-endpoint, and a right-most mapping position,

102   which we call R-endpoint. As a result of genomic rearrangement, reads from one breakpoint-

103   spanning HMW DNA molecule would be mapped to two different genome loci on the reference

104   genome. This split-molecule event has two consequences: 1) observing two separate fragments

105   sharing the same barcode; 2) each of the two fragment has one endpoint close to the true

106   breakpoints. Therefore, in a typical linked-read WGS data set, multiple split-molecule events

107  could be captured and we would usually observe multiple shared barcodes between two distant

108  genome loci and multiple fragment endpoints near the breakpoints.

109  To illustrate this, Figure 1a shows the split-molecule events of a deletion, where breakpoints 1

110  and 2 are marked by red arrows. Multiple fragment endpoints are enriched near the two

111  breakpoints of a large deletion. This can be observed in deletions with minimal size of about 5-

112  10 kb. Figure 1b and Supplementary Figures 1-3 show the patterns of enriched fragment

113  endpoints that are introduced by different types of SVs. As an example, Figure 1c shows the

114  number of fragment endpoints in a 5-kb sliding window near two deletion breakpoints, based on

115  a 35X coverage linked-read WGS data generated from the NA12878 genome (genome of a

116  female individual extensively sequenced by multiple platforms). At the breakpoints, the number

117  of fragment endpoints in the 5-kb sliding window is more than 100 and is five times more than

118  normal regions, forming peaks in the figure.

119  Since the fragments can be paired according to their barcodes, we can also observe fragment

120  endpoints of this deletion in a two-dimensional view. As shown in Figure 1d, each dot indicates

121  two endpoints from a pair of fragments which share the same barcode. The x-value of the dot is

122  the position of the first fragment's R-endpoint and the y-value of the dot is the position of the

123  second fragment's L-endpoint. The bottom panel and right panel in Figure 1d shows number of

124  dots that are projected to the x-axis and y-axis. Similar with the one-dimensional plot (Figure 1c),

125  a peak is formed near each breakpoint, which is marked by the red arrow. The background noise

126    of the two-dimensional plot is cleaner than the one-dimensional plot since the fragments that do

127    not share barcodes are excluded. Therefore, the two-dimensional plot is more useful when the

128    variant allele frequency (VAF) is very low and there are only a few supporting fragments.

129    Next, we describe the signals of type 2 evidence. The barcodes between two nearby genome

130    locations is highly similar because the two locations are spanned by almost the same set of input

131    HMW DNA molecules. However, due to the genome rearrangement, the reads mapped to the left

132    side and right side of a breakpoint may originate from different locations of the alternative

133    genome and thus have different barcodes (Figure 1e). Dropped barcode similarity between two

134    nearby loci therefore indicates an SV breakpoint. LinkedSV detects this type of evidence by a

135    twin-window method, which uses two adjacent sliding windows to scan the genome and find

136    regions where the barcode similarity between the two nearby window regions is significantly

137    decreased. Figure 1f illustrates an inversion breakpoint detected by LinkedSV from the

138    NA12878 genome. The change of barcode similarity was plotted and a peak was formed at the

139    breakpoint. After searching for the two types of evidence, LinkedSV combines the candidate SV

140    regions and quantifies the evidence using a novel probabilistic model. The breakpoints are

141    further refined using short-read information, including discordant read pairs and split-reads.

142

143    **Performance evaluation on simulated WGS data**

144    To assess LinkedSV's performance, we simulated a 35X linked-read WGS data set with 1,175

145    SVs inserted using LRSIM[13] (see Methods for details). The breakpoints of the simulated SVs

146    were designed to be located in repeat regions, since we found that LinkedSV and other available

147    SV callers performed very well when the breakpoints were located in non-repeat regions, and

148    thus we set to test the performances of all the SV callers under more challenging situations. This

149    makes sense because SV breakpoints are more likely to be in repeat regions [7, 8, 9], and because

150    these situations represent those that are difficult to be addressed by conventional short-read

151    sequencing approaches.

152    The simulated reads were aligned to the reference genome using the Longranger [12] package

153    provided by 10X Genomics. The Longranger pipeline internally uses the Lariat aligner [16], which

154    was designed for the alignment of linked reads. SV calling was performed using LinkedSV as

155    well as three other available SV callers designed for linked-read sequencing: Longranger,

156    GROC-SVs [14] and NAIBR [15]. Two widely used short-read SV callers (Delly[3] and Lumpy[17]) were

157    also used.

158    We used recalls, precisions and F1 scores to evaluate the performance of the six SV callers on

159    this data set. As shown in Figure 2a, the four linked-read SV callers showed higher F1 scores

160    than the two short-read SV callers. LinkedSV achieved the highest recall and F1 score among all

161    methods. GROC-SVs had a good precision but its recall was lower than LinkedSV, so we further

162    analyzed the false negative calls of GROC-SVs to understand the underlying reason. A major

163    portion of the false negative calls by GROC-SVs represents duplications that are smaller than

164    twice the fragment length. For large duplications, the reads of the alternative allele are separated

165    by a large gap so that we can observe two sets of fragments with the same set of barcodes, which

166    indicate an SV (Supplementary Figure 4a). If the duplication is not large enough, the reads will

167    be probably clustered into one fragment (Supplementary Figure 4b). Even in this case, we can

168    observe enriched fragment endpoints near the duplication breakpoints in LinkedSV. As an

169    example, Figure 2b shows the endpoint signal of a missed duplication call by GROC-SVs. The

170    supporting fragments of this duplication is shown in Figure 2c. A detailed explanation of the

171    pattern of duplication can be found in Supplementary Figure 1 and Supplementary Movie 1.

172    Figure 2d showed the extra read depth in this region. We also evaluated the breakpoint precision

173    of LinkedSV. Most of breakpoints predicted by fragment endpoints are within 20 bp (Figure 2e)

174    and refined breakpoints using discordant read-pairs and split-reads have base-pair resolution

175    (Figure 2f).

176

177    **Benchmarking on WGS data with somatic SVs of low VAF**

178    Somatic SVs are commonly found in cancer genomes [18, 19, 20]. However, due to the high

179    heterogeneity of genomic alteration in cancer genomes, somatic SVs often have low (as opposed

180    to ~50% in a germline genome) VAF and thus are more difficult to detect by SV callers designed

181    for germline SVs. We simulated two WGS data sets with VAF of 10% and 20%, respectively.

182    Recalls, precisions and F1 values of the six SV callers were evaluated on both data sets (Figure

183    3a, Figure 3b). When the VAF was 20%, the recall of LinkedSV (0.803) was much higher than

184    that of Longranger (0.306), GROC-SVs (0.324) and NAIBR (0.679) The F1 score of LinkedSV

185    (0.855) was also the highest among all the SV callers. When the VAF was 10%, LinkedSV still

186    had a recall of 0.761, which was 72% higher than the second best SV caller NAIBR. Longranger

187    detected 17% of the SVs while GROC-SVs almost completely failed to detect the SVs. The

188    recall rates of Delly and Lumpy were 0.28 and 0.72, respectively, indicating that some of the

189    SVs can be detected even without barcode information. These observations confirmed that other

190    SV callers were mainly designed for germline genomes and had substantial difficulty in

191    detecting SVs with somatic mosaicism. However, due to the combination of barcode overlapping

192    and enriched fragment endpoints in our statistical model (see Methods for details), LinkedSV

193    was able to achieve a good performance even when VAF was very low. We manually checked

194    the barcode overlapping evidence of some SV calls using the Loupe software developed by 10X

195    Genomics. Figure 3c shows an inversion that was missed by Longranger, and NAIBR but

196    detected by LinkedSV (at VAF of 10%). Although the variant frequency is low, the overlapped

197    barcodes between the two inversion breakpoints can be clearly visualized (in the black circles) in

198    the figure. Figure 3d shows the supporting fragments of the inversion detected by LinkedSV.

199    Each horizontal line represent two fragments that share the same barcode and support the SV.

200    These results suggest that the manufacturer-provided software tool has limitations for SV

201    detection, despite its strong functionality in visualization.

202    To test the performance of LinkedSV on the detection of disease casual SVs, we simulated one

203    germline and two somatic (VAF = 10% and 20%) linked-read WGS data sets with 51

204    deletions/duplications that were known to cause CNV (copy number variation) syndromes

205    involved in developmental disorders (see Method for details). The size distribution of the events

206    was shown in Supplementary Figure 5. The performances of LinkedSV as well as 5 other SV

207    callers were shown in Supplementary Figure 6. The results were similar to those of the above

208    simulations. LinkedSV had the highest F1 score on both germline and mosaic data sets, followed

209    by NAIBR.

210

211    **Benchmarking of deletion detection on the HG002 genome**

212    Recently, the Genome in a Bottle (GIAB) Consortium released a benchmark call set for the

213    evaluation of germline SV detection[21]. The benchmark set was based on the HG002 genome and

214    was generated from integrating multiple SV calling methods from multiple sequencing platforms

215    including 10X Genomics sequencing and PacBio long-read sequencing. The current GIAB call

216    set only contains insertions and deletions. Since LinkedSV and the other three linked-read SV

217    callers cannot detect insertions, we only benchmarked the performance to detect deletions using

218    this benchmarking data set.

219    LinkedSV uses different strategies to detect deletions of different sizes. For deletions that are

220    more than 10 kb, LinkedSV uses the two types of evidence from barcode signals as described

221    above; for deletions that are within 1 - 10 kb, LinkedSV uses a combination of read depth and

222    paired-end signals, with additional consideration of local haplotypes; for detection of SVs that

223    are less than 1kb, LinkedSV uses a local assembly-based method. Specifically, we modified the

224    FermiKit[22] *de novo* assembly pipeline to be a local assembler to improve speed and reduce the

225    complexity of the assembly graph (see Method for details).

226    Supplementary Figure 7a showed the performance on detection of deletions that were more than

227    10 kb.  The recall and F1 score of LinkedSV was the highest among the seven methods. The four

228    linked-read SV callers performed better than the three short-read SV callers in terms of F1 score.

229    The performance on detection of deletions that within 1 - 10 kb were shown in Supplementary

230    Figure 7b. The performance of LinkedSV was similar to Longranger, which also provided an

231    algorithm to detect small deletions. NAIBR and GROC-SVs did not perform well because they

232    were not designed to detect small events including small deletions. For deletions that were less

233    than 1 kb, LinkedSV (with modified FermiKit) performed best (recall = 0.48, F1 score = 0.64), it

234    detected more calls than the original *de novo* assembly version (recall = 0.43, F1 score = 0.60),

235    indicating that local assembly reduced the complexity of the assembly graph and improved the

236    performance. NAIBR, GROC-SVs and Lumpy did not perform well on deletions of this scale.

237    (Supplementary Figure 7c). Size distribution of SV events (including deletions, duplications and

238    inversions) detected by LinkedSV was shown in Supplementary Figure 8.

239

240    **Performance evaluation on simulated WES data**

241    Compared with WGS, WES is currently widely used in clinical settings to identify disease causal

242    variants on patients with suspected genetic diseases, partly due to the lower cost of WES. Since

243    WES only covers a small portion of regions in the whole genome, it is far more challenging to

244    detect SVs from WES data, especially when the SV breakpoints are not in the capture regions.

245    However, by combining linked-read sequencing with WES capture platforms, it is possible to

246    alleviate this problem, and significantly improve the sensitivity of SV detection using WES.

247    To evaluate SV detection on linked-read WES data, we simulated a 40X coverage linked-read

248    WES data set with 1,160 heterozygous SVs (see Methods for details). 44.3% of the breakpoints

249    were not in exon regions. SV calling was performed using the six SV callers. As shown in Figure

250    4a, LinkedSV had the highest recall (0.79) and highest F1 score (0.86). In terms of the balanced

251    accuracy (F1 score), NAIBR was the second best caller, followed by GROC-SVs.

252    We analyzed false negative calls of the second best SV caller NAIBR. NAIBR tends to miss

253    some SV events that have shared barcodes but lack short-read support.  For example, Figure 4b

254    showed a deletion between chr1:172545561-173504265. Both breakpoints were located outside

255    of capture regions. Breakpoint 1 (chr1:172545561) was 768 bp away from the nearest capture

256    region and breakpoint 2 (chr1: 173504265) was 392 bp away from the nearest capture region.

257    Unfortunately, no discordant read pairs that support the deletion could be found. However,

258    shared barcodes between the two breakpoints were clearly indicated by the Loupe software

259    (Figure 4b). In addition, LinkedSV also detected 28 pairs of fragments that share the same

260    barcodes and support the SV. These fragments were plotted in Figure 4c. Although no short-read

261    support was found, the SV type could be determined using the pattern of enriched fragment

262    endpoints shown in Figure 1b. In this SV event, R-endpoints were highly enriched for the first

263    set of fragments and L-endpoints were highly enriched for second set of fragments. Thus, the SV

264    type was predicted as deletion.

265

266    **Detection of *F8* inversion from clinical WES data**

267    We also tested the performance of LinkedSV on several clinical samples with linked-read WES

268    data. First, we applied LinkedSV on a WES sample of a male individual with Hemophilia A.

269    Previous experiments had shown that the patient had type I inversion of the *F8* gene, where the

270    two breakpoints resided in intronic/intergenic regions, thus the inversion and its breakpoints

271    cannot be inferred from conventional WES. The *F8* gene is located in Xq28. The intron 22 of *F8*

272    gene contains a GC-rich sequence (named int22h-1) that is duplicated at two positions towards

273    the Xq-telomere (int22h-2 and int22h-3). Int22h-2 has the same direction with int22h-1 while

274    int22h-3 has the inverted direction. The type I inversion is induced by the recombination

275    between int22h-1 and int22h-3 [23, 24] (Figure 5a). BLAST alignment of int22h-1 and int22h-3

276    showed that the two sequences had 99.88% identity. Since the breakpoints were located in two

277    segmental duplications with nearly identical sequences, the inversion is undetectable by

278    conventional short-read sequencing. Delly [3] and Lumpy [17] failed to detect the inversion from the

279    linked-read WES data (results were shown in Supplementary Tables 1-2).

280    Longranger, GROC-SVs, NAIBR and LinkedSV were also used to detect SVs from this sample.

281    None of the first three methods detected this inversion (results were shown in Supplementary

282    Tables 3-5), although the overlapping barcodes can be visualized using the Loupe software

283    (Figure 5b). However, LinkedSV successfully detected this inversion by combining two types of

284    evidence. As described above, barcode similarity between two nearby regions are very high but

285    drops suddenly at the breakpoints. Figure 5c shows the suddenly drop of barcode similarity at the

286    two breakpoints. Each dot in the figure represents the reciprocal of the barcode similarity

287    between its left 40 kb window and right 40 kb window thus the Y value of the dots are reversely

288    related to the barcode similarity and positively related to the probability of being a breakpoint.

289    The barcode similarities are lowest at the two breakpoints and thus form two peaks in the figure

290    (marked by red arrow). In addition, LinkedSV also identified the supporting fragments of the SV

291    using type 1 evidence (Figure 5d). The predicted breakpoint positions are consistent with the

292    genomic positions of int22h-1 and int22h-3.

293

**Detection of mosaic *NF1* deletion from clinical WES data**

295 Another linked-read WES sample was from an individual who was clinically diagnosed with

296 Neurofibromatosis type 1. Neurofibromatosis type 1 is caused by mutations in the *NF1* gene on

297 chromosome 17q11.2, which encodes neurofibromin, a GTPase activating protein that has a role

298 in the regulation of RAS signaling. Since standard genetic testing techniques including cDNA

299 sequencing and multiplex ligation-dependent probe amplification (MLPA) revealed no

300 constitutional or mosaic pathogenic mutation in this patient, we hypothesized that this patient

301 may carry an SV affecting the *NF1* gene that escapes the detection by the applied standard

302 techniques. To evaluate LinkedSV, we utilized the 10X Genomics Chromium platform to

303 generate linked-read WES data to confirm and resolve the mutation. SV detection was conducted

304 using the four linked-read SV callers as well as Delly and Lumpy. Longranger detected

305 overlapped barcodes between exon 54 of the *NF1* gene and intron 3 of *RAB11FIP4*. However,

306 the SV type was unknown and no supporting read pairs or split-reads were found. GROC-SVs,

307 NAIBR, Delly and Lumpy failed to detect this SV (Supplementary Tables 6-9). As shown in

308 Figure 6a, LinkedSV detected 16 fragment pairs that may support a deletion spanning the region

309 of chr17:29684175-29822527. In addition, a discordant read pair spanning the two breakpoints

310 were found (Figure 6b), which gave further evidence supporting the deletion. The breakpoints

311 were estimated from this discordant read pair and thus the resolution is a few hundred base pairs.

312    In Figure 6, each colored line represent a reconstructed fragment, and ~13% of the fragments

313    belong to the variant allele, indicating the somatic mosaicism of this deletion. The right

314    breakpoint was within an AluJr sequence masked by repeat masker, which may explain why the

315    deletion was difficult to be detected by conventional methods.

316    In comparison, the clinical lab used massive parallel sequencing (TruSightCancer panel on a

317    MiSeq platform (Illumina)) and successfully revealed in exon 54 a transition of *NF1* sequences

318    into a non-*NF1* derived sequence. This sequence transition at *NF1* position c.7886_7887 was

319    present in 8% of the reads covering this site in germline DNA of the patient. Analysis of the

320    reads displaying the aberrant sequence in exon 54 showed that the non-*NF1* derived sequence

321    was part of an Alu element that matched best a sequence in intron 3 of the *RAB11FIP4* gene

322    located 138 kb downstream of *NF1* exon 54. These results suggested a low-level (~8%)

323    mosaicism of a deletion encompassing the region intervening between *NF1*:c.7886 and

324    *RAB11FIP4*:c.337-22216, so that the true deletion spans chr17:29684367-29822453, which is

325    very similar to our estimated breakpoints from LinkedSV above. In summary, our analysis on

326    two clinical samples with *F8* inversion and *NF1* deletion demonstrated the unique advantage of

327    linked-read sequencing in confirming and resolving structural variants in repetitive regions and

328    challenging situations.

329

330    **Comparison with SVs detected from long-read sequencing**

331    We previously reported the *de novo* genome assembly of a Chinese individual (HX1) [25]. This

332    genome was sequenced deeply at 103X coverage using PacBio long-read sequencing. Recently,

333    the developers of SMRT-SV [10, 26] reported the SV calls of HX1 detected from the PacBio data.

334    Additionally, we have also generated a 37X linked-read WGS dataset on HX1. Therefore, in the

335    current study, we detected SVs from the linked-read data using LinkedSV and compared the SV

336    calls detected by LinkedSV and SMRT-SV. The SMRT-SV call set has 17 large deletions

337    (≥10kb), all of which were detected by LinkedSV. In addition, LinkedSV detected another 46

338    large deletions, which were missed by SMRT-SV. To validate these deletion calls, we mapped

339    the PacBio reads of HX1 to GRCh38 reference genome using minimap2[27], and manually

340    examined all the SV affected regions in both PacBio data and linked-read data, using the

341    Integrative Genomics Viewer (IGV) [28] and the Loupe software tool. We classify a deletion as a

342    true deletion if there are decreased read depth in the deletion region and clear boundaries at the

343    breakpoints. After the manual inspection, we found that among the 46 deletions that are only

344    detected by LinkedSV, 34 of them have clear evidence of deletion in the WGS data; 10 of them

345    are complex SV events that need to be fully resolved; and 2 of them are false positive events.

346    Figure 7a-c showed an example of a deletion that were detected by LinkedSV but missed by

347    SMRT-SV. This is a 45 kb deletion located in chr2:110395971-110441346. A deletion pattern

348    was clearly indicated by the Loupe software tool (Figure 7a). After examine the PacBio reads,

349    we were able to found clipped reads at the breakpoint positions (Figure 7b, 7c). However, for

350    most of the clipped reads, the clipped sequences were aligned to the hs38d1 decoy sequence,

351   except for 5 reads with clipped sequence > 7 kb. Analysis of the 5 reads revealed that the two

352   breakpoints in chr2 were not directly joined. There was a 6 kb insertion in between. The inserted

353   sequence was from hs38d1 (coordinates: 1381394-1387327). The proposed variant allele was

354   shown in Supplementary Figure 9a. To validate this deletion/insertion event, we aligned all the

355   PacBio reads to a new reference genome with all sequences of GRCh38 plus hs38d1 and the

356   sequence of the proposed variant allele. The reads aligned to the proposed variant allele were

357   shown in Supplementary Figure 9b. There were 33 reads spanning the chr2-hs38d1 junction, 48

358   reads spanning the hs38d1-chr2 junction and 13 reads spanning both junctions. *De novo*

359   assembly of all the reads aligned to the proposed variant allele generated a single contig of 42.7

360   kb, which also spanned both junctions (Supplementary Figure 9b, bottom track). These analysis

361   showed that the large deletion event detected by LinkedSV is true and with PacBio long reads

362   the details of complex SV events could be resolved.

363   We also compared the duplication calls of LinkedSV and SMRT-SV and manually examined

364   discordant SV calls. LinkedSV reported 6 large duplications ($\geq$ 10kb), 5 of which were not

365   reported by SMRT-SV.   Figure 7d-f showed the evidence of a 61 kb duplication call

366   (chr19:27338390-27399298), which was only reported by LinkedSV. A two-fold increase of

367   read depth could be observed in the duplication region (Figure 7d), and the breakpoints were also

368   clearly indicated in the alignments of PacBio long reads, as shown in IGV (Figure 7e,f). The read

369   depths of PacBio raw reads and error-corrected reads were shown in Supplementary Figure 10.

370   The increase of read depth in the duplication region can also be observed. After the manual

371   inspection of the left duplication breakpoint, a small duplication event was found next to the

372   main event. The boundaries of the small duplication can be observed in the alignments of linked

373   reads and error-corrected PacBio reads, but not in the alignments of PacBio raw reads,

374   potentially because of mapping errors (Supplementary Figure 11). SMRT-SV reported 194 large

375   duplications ($\geq$ 10 kb). Unexpectedly, 193 of the duplication calls were not detected by

376   LinkedSV. In addition, none of these duplications could be detected by Sniffles[11], another widely

377   used long-read SV caller. After comparing with the segmental duplication database [29], we found

378   that 182 of the 193 duplication calls (94.3%) were located in large segmental duplication regions.

379   Both long reads and linked reads could not be reliably mapped in these regions. As an example,

380   we plotted the read depth distribution in the region around a 25 kb duplication call of SMRT-SV.

381   Neither long reads (Supplementary Figure 12a) nor linked reads (Supplementary Figure 12b) had

382   mapping quality >20 in the duplication region. Therefore, SV detection in the super-large

383   segmental duplication regions is still very challenging. In summary, our comparative analysis

384   demonstrated unique advantages of linked-read WGS in resolving large SVs that may be failed

385   by even long-read sequencing platform with very deep coverage.

386

## Discussion

387   **Discussion**

388     In this study, we present LinkedSV, a novel open source algorithm for structural variant

389     detection from linked-read sequencing data. We assessed the performance of LinkedSV on three

390     simulated data sets and two real data sets. By incorporating the two types of evidence as outlined

391     below, LinkedSV outperforms all existing linked-read SV callers including Longranger, GROC-

392     SVs and NAIBR on both WGS and WES data sets.

393     Type 1 evidence gives information about which two genomic positions are connected in the

394     alternative genome. It has two observations: 1) fragments with shared barcodes between two

395     genomic locations and 2) enriched fragment endpoints near breakpoints. Current existing linked-

396     read SV callers only use the first observation to detect SVs while LinkedSV incorporates both

397     observations in the statistical model and is therefore more sensitive and can detect SVs with

398     lower allele frequencies, such as somatic SVs in cancer genomes and mosaic structural variations.

399     Type 2 evidence gives information about which genomic position is interrupted with the

400     observation that the reads on the left side and right side of a genomic position have different

401     barcodes and should be derived from different HMW DNA molecules. LinkedSV is the only SV

402     caller that use type 2 evidence to detect breakpoints. Type 2 evidence is independent to type 1

403     evidence, and gives additional confidence to identify the breakpoints. In addition, type 2

404     evidence can be detected locally, which means we can detect a weird genomic location without

405     looking at the barcodes of the other genomic locations. This is particularly useful in two

406     situations: 1) novel sequence insertions where there is only one breakpoint; 2) only one

407    breakpoint is detectable and the other breakpoint located in a region where there is little coverage

408    within 50 kb, which is often the case in target region sequencing. As LinkedSV incorporates two

409    types of evidence from barcodes, and performs local assembly to detect small deletions, the

410    computation time of LinkedSV is longer than NAIBR, but shorter than GROC-SVs and

411    Longranger (Supplementary Figure 13).

412    In recent years, WES has been widely used to identify disease causal variants for patients with

413    suspected genetic diseases in clinical settings. Identification of SVs from WES data sets are more

414    challenging because the SV breakpoints may not be in the capture regions and thus there would

415    be little coverage at the breakpoints. Linked-read sequencing increases the chance of resolving

416    such type of SVs by providing long-range information. As well as there are a few capture regions

417    nearby, the fragments can still be reconstructed and type 1 and type 2 evidence can still be

418    observed. Our statistical models for both type 1 evidence and type 2 evidence were designed to

419    handle both WGS and WES data sets. GROC-SVs uses a local-assembly method to verify the SV

420    call, which requires sufficient coverage at the breakpoints. By using these two types of evidence,

421    LinkedSV can be less relied on short-read information (e.g. pair-end reads and split-reads). We

422    demonstrated that LinkedSV has better recall and balanced accuracy (F1 score) on the simulated

423    WES data set and can detect SVs even when the breakpoints were not located in capture regions

424    and have no short-read support. In addition, LinkedSV is also the only SV caller that clearly

425    detected the *F8* intron 22 inversion and *NF1* deletion from the clinical WES data sets.

426    Linked-read sequencing has several advantages over traditional short-read sequencing on the

427    purpose of SV detection. First, the human genome is highly repetitive. Previous studies have

428    shown that SVs are closely related to repeats and many SVs are directly mediated by

429    homologous recombination between repeats [30]. In traditional short-read sequencing, if the

430    breakpoint falls in a repeat region, the supporting reads would be multi-mapped and thus the SV

431    cannot be confidently identified. However, this type of SVs are detectable by linked-read

432    sequencing when the HMW DNA molecules span the repeat region. We can observe type 1 and

433    type 2 evidence in the non-repeat region nearby. In our benchmarking, LinkedSV detected more

434    SVs than Delly and Lumpy, especially when the VAF is low. Secondly, SVs are undetectable

435    from traditional short-read sequencing if there is little coverage at the breakpoints, which is often

436    the case in WES data sets. As described above, this type of SV can also be resolved by linked-

437    read sequencing and LinkedSV. Third, linked-read sequencing requires less coverage for

438    detection of SVs with low variant allele frequencies. In linked-read sequencing data, short read

439    pairs are sparsely and randomly distributed along the HMW DNA molecule. In a typical linked-

440    read WGS data set, the average distance between two read pairs derived from the same HMW

441    DNA molecule is about 1000 bp and each HMW DNA molecule only has a short-read coverage

442    of about 0.2X. Therefore, there are about 150 HMW DNA molecules (reconstructed fragments)

443    covering a genomic location of 30X depth. An SV of 10% VAF will has15 supporting fragment

444    pairs in a 30X depth location in linked-read WGS data set, which is sufficient to be detected by

445    LinkedSV. However, an SV of 10% VAF will only has 3 supporting read pairs in a 30X depth

446    location in traditional short read WGS, which makes the detection more challenging.

447    Linked-read sequencing also has several advantages over long-read sequencing in terms of SV

448    detection. The fragment length of linked-reads (typically 50-100 kb) is longer than the read

449    length of regular long-read sequencing (typically 20-30 kb). Therefore, linked-read sequencing

450    has unique advantages for detection of large SVs. In our study, LinkedSV detected several large

451    SVs that were missed in the long-read SV call set. We also showed that the sequencing error (13-

452    15%) of long-read sequencing technologies potentially had a negative effect on reads mapping

453    and subsequent SV calling (Supplementary Figure 11). In terms of library preparation, Linked-

454    read sequencing only requires 1 ng input DNA, which is two orders of magnitude smaller than

455    what is needed by long-read sequencing. Therefore, disease samples of very low DNA amount

456    can be easily sequenced by linked-read sequencing. In addition, SNPs, indels and SVs can be

457    detected from linked-read sequencing simultaneously.

458    LinkedSV may have limitations on detection of SVs in large segmental duplication regions,

459    where the linked reads have low mapping qualities. SMRT-SV was able to find 194 large

460    duplications in the HX1 genome, which were not detected by LinkedSV and Sniffles, two

461    alignment-based SV callers. SMRT-SV detects SVs using an assembly-based approach. During

462    the assembly process, the assembly contigs were error corrected and polished by the PacBio

463    reads. Therefore, the assembly contigs are potentially more accurate and longer than each of the

464    raw reads. Thus, it is possible for SMRT-SV to detect SVs in these large segmental duplication

465    regions.

466    The linked-read technology provides strong evidence to detect large SVs, but it provides little

467    additional evidence to detect small SVs. Therefore, LinkedSV has limited power to detect small

468    SVs such as small duplications and inversions. However, based on our analysis of SV size

469    distribution, large SVs are associated with diseases such as cancers and CNV syndromes

470    (Supplementary Note 2). Therefore, we expect that linked-read technology can help resolve

471    disease associated SVs. Similar to the existing linked-read SV callers, LinkedSV currently does

472    not handle insertions and repeat expansions. As a future direction, we plan to detect novel

473    sequence insertions using type 2 evidence, since this type of SV also cause a decrease of barcode

474    similarity between nearby regions and can be detected by the twin-window method. The exact

475    insertion sequence may then be inferred from the assembly of all the reads that share barcodes

476    with the candidate breakpoint. LinkedSV currently already supports local assembly to detect

477    deletions, but it has not been parameterized and optimized to be combined with type 2 evidence

478    for detection of insertions.

479    In summary, we present LinkedSV, a novel SV caller for linked-read sequencing. LinkedSV

480    outperformed current existing SV callers, especially for identifying SVs with low allele

481    frequency or identifying SVs from target region sequencing such as linked-read WES. We expect

482    that LinkedSV will facilitate the detection of SVs from linked-read sequencing data and help

483    solve negative cases from conventional short-read sequencing.


484    **Methods**

485    **Breakpoint detection from type 1 evidence**

486    First, LinkedSV reconstructs the original long DNA fragments from the reads using mapping

487    positions and barcode information. All mapped reads are partitioned according to the barcode

488    and sorted by mapping position. We define gap distance as the distance between two nearest

489    reads with the same barcode. Two nearby reads are considered from the same long DNA

490    fragment if they have the same barcode and their gap distance is less than a certain distance $G$. $G$

491    is determined using two steps. First, we use $G$ = 50 kb (the same as Zheng et.al [12]) to group the

492    reads into fragments. This value is suitable for detection of large SVs. However, it may be too

493    large for detection of SVs that are smaller than 50 kb. Therefore, we calculate the empirical

494    distribution of intra-fragment gap distance, which is the distance of two nearby reads that are

495    grouped in one fragment. The empirical distribution of intra-fragment gap distance is calculated

496    from all the fragments, and we assign $G$ as the $99^{th}$ percentile of this distribution. $G$ is a fixed

497    number for all fragments and is usually between 5-15 kb, depending on the data set. Fragments

498    with a gap distance larger than $G$ potentially span a breakpoint and will be separated to two

499    fragments.

500

501 In non-SV regions, all the reads from the same HMW DNA molecule would be reconstructed

502 into a single DNA fragment. The reads from the breakpoint-spanning HMW DNA molecule will

503 be mapped to two different positions in the genome. As illustrated in the Result section, this

504 split-molecule event has two consequences: 1) observing two fragments sharing the same

505 barcode; 2) each of the two fragment has one endpoint close to the breakpoints. Therefore, we

506 could observe enriched fragment endpoints near the breakpoints, in both one-dimensional view

507 (Figure 1c) and two-dimensional view (Figure 1d). The type of the endpoints (L-endpoint or R-

508 endpoint) that enriched near the breakpoints depends on the type of SV (Figure 1b). The two-

509 dimensional view has less background noise because the fragments that do not share barcodes

510 and thus do not support the SVs are excluded. Therefore, we detect the enriched endpoints in the

511 two-dimensional view.

512 We now describe how we detect the type 1 evidence of deletion calls, but the method can be

513 applied to other types of SVs. We define fragment pair to be two fragments sharing the same

514 barcode. Let $b_1$, $b_2$ be the positions of the two breakpoint candidates (assuming $b_1 < b_2$). Let $n$ be

515 the number of fragment pairs that may support the SV between $b_1$ and $b_2$. Let $F_{i1}$, $F_{i2}$ denote the

516 $i^{th}$ fragment pair that support the SV. Let $B(F)$ denote the barcode of fragment $F$. Therefore, we

517 have:

518
$$B(F_{i1}) = B(F_{i2}), i = 1, 2, 3, \ldots, n. \quad (1)$$

519    Let $L(F)$ denote the L-endpoint position (i.e., left-most position) of fragment $F$, $R(F)$ denote the

520    R-endpoint position (i.e., right-most position) of fragment $F$. Since this is a deletion and $b_1 < b_2$,

521    $R(F_{i1})$ is the position on $F_{i1}$ that is closest to $b_1$ and $L(F_{i2})$ is the position on $F_{i2}$ that is closest to

522    $b_2$ (Supplementary Figure 14a). The distance between the fragment endpoint and its

523    corresponding breakpoint should be within gap distance distribution (explained in

524    Supplementary Figure 15). Therefore, for almost all (99% × 99%) of the fragment pairs, we have:

525 $$b_1 - G \leq R(F_{i1}) \leq b_1;\ b_2 \leq L(F_{i2}) \leq b_2 + G.\ \ (2)$$

526    As described above, $G$ is the $99^{th}$ percentile of the empirical distribution of intra-fragment gap

527    distance.

528    If we regard $(R(F_{i1}), L(F_{i2}))$ as a point in a two-dimensional plane, according to equation (2), for

529    almost all (98.01%) of the fragment pairs $(F_{i1}, F_{i2})$, $((R(F_{i1}), L(F_{i2}))$ is restricted in a $G \times G$

530    square region with the point $(b_1, b_2)$ being a vertex (Supplementary Figure 14b).

531    We used a graph-based method to fast group the points into clusters and find square regions

532    where the numbers of points were more than expected. First, every possible pair of endpoints

533    $(R(F_1), L(F_2))$ meeting $B(F_1) = B(F_2)$ formed a point in the two-dimensional plane. Each point

534    indicated a pair of fragments that share the same barcode. For example, if 10 fragments share the

535    same barcode, $C_{10}^2$ pairs of endpoints will be generated. A point/pair of endpoints may or may

536    not support an SV because there are two possible reasons for observing two fragments sharing

537    the same barcode: 1) the two fragments originated from two different HMW DNA molecules but

538    were dispersed into the same droplet partition and received the same barcode; 2) the two

539    fragments originated from the same HMW DNA molecule but the reads were reconstructed into

540    two fragments due to an SV. The points are sparsely distributed in the two-dimensional plane

541    and it is highly unlikely to observe multiple points in a specific region. Next, a k-d tree ($k = 2$)

542    was constructed, of which each node stores the ($X$, $Y$) coordinates of one point. A k-d tree is a

543    binary tree that enable fast query of nearby nodes. Therefore, we could quickly find all pairs of

544    points within a certain distance. Any two points ($x_1$, $y_1$) and ($x_2$, $y_2$) were grouped into one cluster

545    if $|x_1 - x_2| < G$ and $|y_1 - y_2| < G$. For each cluster, if the number of points in the cluster was more

546    than a user-defined threshold (default: 5), it was considered as a potential region of enriched

547    fragment endpoints. If the points in the cluster were not within a $G \times G$ square region, we used a

548    $G \times G$ moving square to find a square region where the points are best enriched. Theoretically,

549    the best enriched square region should contain 98.01% ($0.99 \times 0.99$) of the points, according to

550    equation (2). The predicted breakpoints were the X and Y coordinates of the right-bottom vertex

551    of the square. The points in the square region were subjected to a statistical test describe below.


552


553    **Quantification of type 1 evidence**


554    Let $n$ be the number of points in the square region. Each point corresponds to a pair of fragment

555    $F_{i1}$, $F_{i2}$, ($i = 1, 2, 3, \ldots, n$) that may support the SV. Let $b_1$ and $b_2$ be the coordinates of the

556    predicted breakpoint. Equation (1) and (2) hold for all the fragment pairs $F_{i1}$, $F_{i2}$ ($i$ = 1, 2, 3, …,

557    $n$). We then test the null hypothesis that there is no SV between $b_1$ and $b_2$.

558    First, we test the hypothesis that the $n$ fragment pairs $F_{i1}$, $F_{i2}$ have originated from different DNA

559    molecules, but coincidently received the same barcode. Here we define two fragments $F_a$ and $F_b$

560    as an independent fragment pair if $F_a$ and $F_b$ share the same barcode but have originated from

561    different DNA molecules. Thus, $R(F_a)$ and $L(F_b)$ are independent variables. All the fragment

562    pairs that do not support SVs are independent fragment pairs. It is reasonable to assume the

563    generation of HMW DNA molecules from chromosomal DNA is a random process thus both

564    $R(F_a)$ and $L(F_b)$ are uniformly distributed across the chromosome. Therefore, the point $((R(F_a),$

565    $L(F_b))$ is equal likely to be in any place in the two-dimensional plane. Technically, we connect

566    all the chromosomes in a head-to-tail order so that both intra-chromosomal events and inter-

567    chromosomal can be analyzed at the same time. Observing at least $n$ independent fragment pairs

568    meeting equation (2) is equivalent to the event that observing at least $n$ points $((R(F_{i1}), L(F_{i2}))$

569    located in a squared region with an area of $G^2$ on the two-dimensional plane. The probability of

570    this event is:

571
$$p_1 = \sum_{j=n}^{N} \text{Binomial\_pmf}\left(n, N_{\text{ifp}}, \frac{G^2}{L^2}\right), (3)$$

572    where Binomial_pmf is the probability mass function of binomial distribution; $L$ is the total

573    length of the genome (also the side length of the two-dimensional plane); $N_{\text{ifp}}$ is the total number

574    of independent fragment pairs.

575  Since we are doing multiple hypothesis testing in the data set, the probability need to be adjusted.

576
$$p_{\text{adjusted1}} = p_1 \frac{G^2}{L^2}. \quad (4)$$

577  We reject the hypothesis if $p_{\text{adjusted1}} < p_{\text{threshold}}$. $p_{\text{threshold}}$ is $10^{-5}$ by default.

578  Next, we test the hypothesis that fragment pairs $F_{i1}$, $F_{i2}$ ($i = 1, 2, 3, \ldots, n$) have originated from

579  the same DNA molecule, but no reads were sequenced in the gap between $R(F_{i1})$ and $L(F_{i2})$. Let

580  $g_i$ denote the length of the gap between $F_{i1}$ and $F_{i2}$, $\overline{g}$ denote the mean of $g_i$, and we have:

581
$$g_i = L(F_{i2}) - R(F_{i1}), \quad (5)$$

582
$$\overline{g} = \frac{1}{n} \sum_{i=0}^{n} g_i. \quad (6)$$

583  If $\overline{g}$ is too large such that the probability of no reads being generated is smaller than a threshold,

584  we can reject this hypothesis.

585  Similar to the model described by 10X Genomics[12], we assume the read generation on a DNA

586  molecule is a Poisson process with constant rate $\lambda$ across the genome. Let $r$ be the number of

587  reads generated in a region of length $g$, then $r \sim \text{Pois}(\lambda g)$. Let $P_{\text{gap}}(g)$ denote the probability of no

588  read being generated in length $g$, we have:

589
$$P_{\text{gap}}(g) = P(r = 0 \mid \lambda g) = \frac{e^{-\lambda g}(\lambda g)^0}{0!} = e^{-\lambda g} \quad (7)$$

590  Therefore, the gap length $g_i$ follows Exponential distribution: $g_i \sim \text{Exp}(\lambda)$. Recalling that 1) the

591  Exponential distribution with rate parameter $\lambda$ is a Gamma distribution with shape parameter 1

592    and rate parameter $\lambda$; 2) the sum of $n$ independent random variables from Gamma $(1, \lambda)$ is a

593    Gamma random variable from Gamma $(n, \lambda)$, we have:

594    
$$\sum_{i=0}^{n} g_i \sim \text{Gamma } (n, \lambda), \text{ (8)}$$

595    
$$\bar{g} = \frac{\sum_{i=0}^{n} g_i}{n} \sim \text{Gamma } (n, n\lambda), \text{ (9)}$$

596    Therefore, the probability that observing $n$ gap regions with mean length equal to or larger than $\bar{g}$

597    is:

598    
$$p_2 = 1 - \text{Gamma\_cdf } (n, n\lambda), \text{ (10)}$$

599    where Gamma_cdf is the cumulative distribution function of Gamma distribution.

600    Since we are doing multiple hypothesis testing in the data set, the probability need to be adjusted.

601    
$$p_{\text{adjusted2}} = p_2 \frac{N_{\text{rp}}}{n}, \text{ (11)}$$

602    where $N_{\text{rp}}$ is the total number of read pairs.

603    We reject the hypothesis if $p_{\text{adjusted2}} < p_{\text{threshold}}$. $p_{\text{threshold}}$ is set as $10^{-5}$ by default. If both $p_{\text{adjusted1}}$

604    and $p_{\text{adjusted2}}$ are less than $p_{\text{threshold}}$, we accept the hypothesis that the SV is true. For each

605    candidate SV, we report a confidence score for type 1 evidence as:

606    
$$\text{Confidence score } 1 = -\log_{10} (\max (p_{\text{adjusted1}}, p_{\text{adjusted2}})). \text{ (12)}$$

607

**Breakpoint detection from type 2 evidence**

608

609    Barcode similarity between two nearby regions is very high because the reads originate from

610    almost the same set of HMW DNA molecules. However, at the SV breakpoint, the aligned reads

611    from the left side and right side may have originated from different locations in the alternative

612    genome. Thus, the barcode similarity between the left side and right side of the breakpoint are

613    dramatically reduced (as described in the Result section and shown in Figure 1e-f). To detect this,

614    LinkedSV uses two adjacent sliding windows (twin windows, moving 100 bp) to scan the

615    genome and calculate the barcode similarity between the twin windows. The window length can

616    be specified by user. By default, it is $G$ for WGS data sets and 40 kb for WES data sets.

617    The barcode similarity can be simply calculated as the fraction of shared barcodes. This method

618    is suitable for WGS, where the coverage is continuous and uniform. But it does not perform well

619    for WES, where the numbers of reads in the sliding windows vary a lot due to capture bias and

620    the length of capture regions. Therefore, we use a model that considering the variation of

621    sequencing depth and capture region positions. The barcode similarity is calculated as:

622
$$S = \frac{x}{m_1^a m_2^b} n e^{-\alpha d} \quad (13)$$

623    where:

624    $m_1$ is the number of barcodes in window 1,

625    $m_2$ is the number of barcodes in window 2,

626     $x$ is the number of barcodes in both windows,

627     $d$ is the weight distance between reads of the left window and the right window,

628     $n$ is a constant representing the characteristic of the library,

629     $\alpha$ is a parameter of fragment length distribution,

630     a and b are two parameters between 0 and 1,

631     $n$, $\alpha$, $a$ and $b$ are estimated from the data using regression. Detailed explanation of this model is

632     in Supplementary Note 1.

633

634     Next, we calculate the empirical distribution of barcode similarity. Regions where the barcode

635     similarity less than a threshold ($5^{th}$ percentile of the empirical distribution by default) were

636     regarded as breakpoint candidates. If a set of consecutive regions have barcode similarity lower

637     than the threshold, we only retain the region that has the lowest barcode similarity. If the barcode

638     similarity of a breakpoint candidate is $S_0$, the empirical $p$-value is calculated as:

639     $$p_{empirical} = \frac{\text{number of twin windows with } S \leq S_0}{\text{total number of twin windows}}. \quad (14)$$

640     The confidence score of type 2 evidence is:

641     $$\text{Confidence score 2} = -\log_{10}(p_{empirical}). \quad (15)$$

642

### Combination of both types of evidence

644    Type 1 evidence gives pairs of endpoints that indicate two genomic positions are joined in the

645    alternative genome. Type 2 evidence gives genomic positions where the barcodes suddenly

646    changed, regardless of which genomic position can be joined. Therefore, type 1 and type 2

647    evidence are independent. The candidate breakpoints detected from type 2 evidence were

648    searched against the candidate breakpoint pairs detected from type 1 evidence so that the calls

649    were merged. The combined confidence score is:

650    Combined score = Confidence score 1 + Confidence score 2a + Confidence score 2b, (16)

651    where Confidence score 1 is the confidence score calculated from type 1 evidence (equation 12);

652    Confidence score 2a and Confidence score 2b are the confidence scores of the two breakpoints

653    calculated from type 2 evidence (equation 14).

654

### Refining breakpoints using short-read information

656    For large SV events, we search for discordant read-pairs and clipped reads that are within 10 kb

657    to the predicted breakpoint pairs by the above approach. We use a graph-based approach that is

658    similar to DELLY[3] to cluster the discordant read-pairs. We define the supporting split-reads as

659    the clipped reads that can be mapped to the both breakpoints, and the map direction matches the

660    SV type. If both discordant read-pairs and split-reads are found to support the SV, we use the

661    breakpoints inferred by split-reads as the final breakpoint position.

662

663    **Detection of small deletions that are within 50 bp -10 kb**

664    We use a 1Mb moving window (with 0.1 Mb overlapping) to scan the genome. For each window,

665    all the aligned reads (including phased and un-phased reads) were extracted and were assembled

666    by the FermiKit pipeline. Regions with extreme high coverage (more than 20-fold of average

667    coverage) were skipped. The resulting contigs were mapped back to the 1 Mb reference sequence

668    of the moving window using bwa-mem and deletions were called from the aligned contigs if the

669    alignments were unique within the 1 Mb moving window. The local assembly based process

670    mainly contribute to the detection of deletions within 50-1000 bp. To detect deletions that are

671    larger than 1 kb and might be missed by the assembly-based process, we use a 500 bp moving

672    window (with no overlapping) to find candidate regions where the read depth of either haplotype

673    is less than 10% of the average depth of the haplotype. Next we extract all the read pairs of this

674    haplotype and test if the mean insert size of these read pairs is significantly larger than the mean

675    value of the whole genome, assuming the average insert size of $n$ read pairs follows normal

676    distribution: $N(\mu, \sigma^2/n)$, where $\mu$ and $\sigma$ are the mean and standard deviation of the insert size of

677    the whole genome.

678    We use a read depth based method to detect deletions that are larger than 1 kb and lack read pair

679    support. If there are $m$ consecutive windows where the read depths are less than 10% of the

680    average depth, we assume the read depth of each window is independent, and calculate the $p$

681    value using the simple equation: $p = (a/b)^m$, where $b$ is the total number moving windows and $a$

682    is the total number of moving windows where the read depths are less than 10% of the average

683    depth. A deletion is called if $p < 10^{-10}$.

684

685    **Generation of simulated linked-read WGS data set**

686    The linked reads were simulated by LRSIM, which can generate linked-reads from a given

687    FASTA file containing the genome sequences. We generated a diploid FASTA file based on

688    hg19 reference genome with SNPs and SVs inserted. The purpose of inserting SNPs was to

689    mimic real data. The generation of the diploid FASTA file is described below. First, we inserted

690    SNPs to hg19 using vcf2diploid [31]. The inserted SNPs were from the gold standard SNP call set

691    (v.3.3.2) of NA12878 genome [32]. The vcf2diploid software generated two FASTA files, each of

692    which was a pseudohaplotype (paternal or maternal) with the phased SNPs inserted. Next, we

693    insert SVs into the paternal FASTA file using our custom script. The breakpoints were located in

694    the repetitive regions in hg19 and the distance between the two breakpoints were in the range of

695    50 kb to 10 Mb. In total, we simulated 351 deletions, 386 duplications, 353 inversions and 85

696    translocations, all of which were in the paternal copy and were heterozygous SVs. We then

697     concatenate the paternal and maternal FAST file into a single FASTA file and simulated linked-

698     reads using LRSIM. To mimic real data, the barcode sequences and molecule length distribution

699     used for simulation were from the NA12878 whole-genome data set released by 10X Genomics.

700     The number of read pairs was set to 360 million so that a 35X coverage data set was generated.

701     The genome coordinates of simulated SVs was shown in Supplementary Data 1. The size

702     distribution of the simulated SVs was shown in Supplementary Figure 16a.

703

704     **Generation of WGS data set with low VAF**

705     In cancer samples or mosaic samples, the total DNA is a mix of a small portion of variant alleles

706     and a large portion of normal alleles. To simulate the WGS data sets with low variant

707     frequencies, we used the same paternal and maternal FASTA file described above but the

708     combined FASTA file contained multiple copies of the normal allele (the maternal FASTA) and

709     only one copy of the variant allele (the paternal FASTA). For example, to simulate a WGS data

710     set with VAF of 20%, four copies of the maternal FASTA and one copy of the paternal FASTA

711     were combined. The linked reads were simulated using LRSIM with the same parameters and a

712     35X coverage data set was generated.

713

714     **Simulation of deletions and duplications that cause diseases**

715  To test the performance of LinkedSV on the detection of disease casual SVs, we downloaded a

716  list of expert-curated deletions and duplications that were known to cause CNV syndromes

717  involved in developmental disorders. This list was downloaded from the DECIPHER database,

718  and contained 67 CNV syndromes. Some syndromes were affected by CNV events in the same

719  region. After removing redundant syndromes, we got 51 CNV events (Supplementary Table 10).

720  Based on the 51 CNV events we simulated a germline WGS data set and two mosaic WGS data

721  sets (VAF = 10% and 20%) using the same method described above.

722

723  **Generation of simulated linked-read WES data set**

724  To generate the linked-read WES data set, we first generate a 100X linked-read WGS data set

725  and then down-sample it to be a WES data set. Generation of the simulated linked-read WGS

726  data set with SNPs and SVs inserted was similar to the method described above. In total, we

727  inserted 1160 heterozygous SVs. The SV breakpoints were randomly selected from regions that

728  were within 2000 bp of an exon. Among the 2320 breakpoints (two breakpoints per each SV),

729  1028 breakpoints (44.3%) were in intronic or intergenic regions. The SV sizes are in the range of

730  50 kb to 10 Mb (Supplementary Figure 16b). Supplementary Data 2 showed the list of simulated

731  SVs. The number of inserted SVs in the simulated WES data set was slightly smaller than that in

732  the simulated WGS data set because the SV breakpoints were designed to reside within 2000 bp

733  of an exon. The simulated reads were generated using LRSIM and were mapped to hg19

734   reference genome using the Longranger pipeline (default settings). The phased bam generated by

735   Longranger was down-sampled to be a simulated WES data set. To mimic real WES data set, we

736   used the coverage distribution of the linked-read WES data set of NA12878 genome (released by

737   10X Genomics) to guide the down-sampling process. We bin the genome into 10 bp windows

738   and calculate number of reads mapped to each window (left mapping positions were used) in

739   NA12878 linked-read WES data. The simulated WES data set was generated by sampling reads

740   from the 100X WGS data according to number of reads mapped to the same 10 bp window in the

741   NA12878 WES. The down sampling was at read pair level, if the one read is retained, the paired

742   read would also be retained.

743

744   **Benchmarking of deletion detection on the HG002 genome**

745   The HG002 benchmark set (version 0.6) was downloaded from the FTP site: ftp://ftp-

746   trace.ncbi.nlm.nih.gov/giab/ftp/data/AshkenazimTrio/analysis/NIST_SVs_Integration_v0.6/. The

747   benchmarking process was performed according to the authors' suggestions[21]. The benchmark

748   set contains a Tier 1 benchmark regions, where all the insertions/deletions are resolved and any

749   extra calls were putative false positives. This region covers 2.66 Gbp of the human genome. A

750   deletion call was considered to be a true positive call if it had at least 50% reciprocal overlap (the

751   overlapped region was more than 50% of both calls) with a deletion call with filter = PASS in the

752     Tier 1 vcf file. Otherwise, it was considered to be a false positive call. This 50% reciprocal

753     overlapping criterion was chosen to follow what was done by a previous study [33].

754     Recall, precision and F1 score were calculated as follows.

755     $$\text{Recall} = \frac{\text{Number of true positive calls}}{\text{Total number of deletion calls with filter=PASS in the Tier 1 vcf file}}; (17)$$

756     $$\text{Precision} = \frac{\text{Number of true positive calls}}{\text{Total number of deletion calls of the query set}}; (18)$$

757     $$\text{F1 score} = \frac{2*\text{Recall}*\text{Precision}}{\text{Recall}+\text{Precision}}. (19)$$

758

## Competing Interests

759

760     The authors declare no competing interests.

761

## Acknowledgments

762

## Data Availability

767

768 The 10X Genomics sequencing data of the HX1 genome was generated in this study and can be

769 obtained from the NCBI SRA database with the accession code SRX5781869

770 [https://www.ncbi.nlm.nih.gov/sra/?term=SRX5781869].

771 The PacBio sequencing data of the HX1 genome was previously published and can be obtained

772 from the NCBI SRA database with the accession code SRX1424851

773 [https://www.ncbi.nlm.nih.gov/sra?term=SRX1424851]. The 10X Genomics sequencing data of

774 the HG002 genome was released by 10X Genomics and can be downloaded from

775 https://support.10xgenomics.com/de-novo-assembly/datasets/2.1.0/ash.

776 Due to potential compromise of individual privacy, full datasets of the clinical samples (*F8* and

777 *NF1*) are available from the authors on reasonable request and institutional data use agreement.

778 All other relevant data is available upon request.

## Code Availability

779

780 The source code of LinkedSV is publicly available on GitHub

781 (https://github.com/WGLab/LinkedSV). A detailed description of how to use LinkedSV is also

782 provided in the GitHub repository.

783

## Author contributions

784

785    L.F. and K.W. (Kai Wang) designed the study. L.F. implemented the tool and performed the

786    analysis. F.A.M. and R.P.S. generated the 10X Genomics sequencing data of the *F8* inversion

787    sample and the *NF1* deletion sample and C. K. and M.V.G analyzed the data. S.W. and K.W.

788    (Katharina Wimmer) performed targeted Illumina MiSeq sequencing of the *NF1* deletion sample

789    and analyzed the data. M. L. and H.H. guided on method development and data analysis. L.F.

790    drafted the manuscript. All authors read, revised, and approved the manuscript.

## References

791

792    1.    Weischenfeldt J, Symmons O, Spitz F, Korbel JO. Phenotypic impact of genomic structural
793        variation: insights from and for human disease. *Nat Rev Genet* **14**, 125-138 (2013).

794

795    2.    Ye K, Schulz MH, Long Q, Apweiler R, Ning Z. Pindel: a pattern growth approach to detect
796        break points of large deletions and medium sized insertions from paired-end short reads.
797        *Bioinformatics* **25**, 2865-2871 (2009).

798

799    3.    Rausch T, Zichner T, Schlattl A, Stutz AM, Benes V, Korbel JO. DELLY: structural variant
800        discovery by integrated paired-end and split-read analysis. *Bioinformatics* **28**, i333-i339 (2012).

801

802    4.    Chen K*, et al.* BreakDancer: an algorithm for high-resolution mapping of genomic structural
803        variation. *Nat Methods* **6**, 677-681 (2009).

804

805    5.    Chong Z*, et al.* novoBreak: local assembly for breakpoint detection in cancer genomes. *Nat*
806        *Methods* **14**, 65-67 (2017).

807

808    6.    Wala JA*, et al.* SvABA: genome-wide detection of structural variants and indels by local
809          assembly. *Genome Res* **28**, 581-591 (2018).

810
811    7.    Carvalho CM, Lupski JR. Mechanisms underlying structural variant formation in genomic
812          disorders. *Nat Rev Genet* **17**, 224-238 (2016).

813
814    8.    Payer LM*, et al.* Structural variants caused by Alu insertions are associated with risks for many
815          human diseases. *Proc Natl Acad Sci U S A* **114**, E3984-E3992 (2017).

816
817    9.    Sharp AJ*, et al.* Segmental duplications and copy-number variation in the human genome. *Am J*
818          *Hum Genet* **77**, 78-88 (2005).

819
820    10.   Chaisson MJ*, et al.* Resolving the complexity of the human genome using single-molecule
821          sequencing. *Nature* **517**, 608-611 (2015).

822
823    11.   Sedlazeck FJ*, et al.* Accurate detection of complex structural variations using single-molecule
824          sequencing. *Nat Methods*,  (2018).

825
826    12.   Zheng GX*, et al.* Haplotyping germline and cancer genomes with high-throughput linked-read
827          sequencing. *Nat Biotechnol* **34**, 303-311 (2016).

828
829    13.   Luo R, Sedlazeck FJ, Darby CA, Kelly SM, Schatz MC. LRSim: A Linked-Reads Simulator
830          Generating Insights for Better Genome Partitioning. *Comput Struct Biotechnol J* **15**, 478-484
831          (2017).

832
833    14.   Spies N*, et al.* Genome-wide reconstruction of complex structural variants using read clouds. *Nat*
834          *Methods*,  (2017).

835
836    15.   Elyanow R, Wu HT, Raphael BJ. Identifying structural variants using linked-read sequencing
837          data. *Bioinformatics*,  (2017).

838
839    16.   Bishara A*, et al.* Read clouds uncover variation in complex regions of the human genome.
840          *Genome Res* **25**, 1570-1580 (2015).

841

842    17.    Layer RM, Chiang C, Quinlan AR, Hall IM. LUMPY: a probabilistic framework for structural
843           variant discovery. *Genome Biol* **15**, R84 (2014).

844

845    18.    Campbell PJ*, et al.* Identification of somatically acquired rearrangements in cancer using
846           genome-wide massively parallel paired-end sequencing. *Nat Genet* **40**, 722-729 (2008).

847

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

850

851    20.    Stephens PJ*, et al.* Complex landscapes of somatic rearrangement in human breast cancer
852           genomes. *Nature* **462**, 1005-1010 (2009).

853

854    21.    Zook JM*, et al.* A robust benchmark for germline structural variant detection. Preprint at
855           https://www.biorxiv.org/content/10.1101/664623v3, (2019).

856

857    22.    Li H. FermiKit: assembly-based variant calling for Illumina resequencing data. *Bioinformatics* **31**,
858           3694-3696 (2015).

859

860    23.    Lakich D, Kazazian HH, Jr., Antonarakis SE, Gitschier J. Inversions disrupting the factor VIII
861           gene are a common cause of severe haemophilia A. *Nat Genet* **5**, 236-241 (1993).

862

863    24.    De Brasi CD, Bowen DJ. Molecular characteristics of the intron 22 homologs of the coagulation
864           factor VIII gene: an update. *J Thromb Haemost* **6**, 1822-1824 (2008).

865

866    25.    Shi L*, et al.* Long-read sequencing and de novo assembly of a Chinese genome. *Nat Commun* **7**,
867           12065 (2016).

868

869    26.    Huddleston J*, et al.* Discovery and genotyping of structural variation from long-read haploid
870           genome sequence data. *Genome Res* **27**, 677-685 (2017).

871

872    27.    Li H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094-3100
873           (2018).

874
875    28.    Robinson JT*, et al.* Integrative genomics viewer. *Nat Biotechnol* **29**, 24-26 (2011).

876
877    29.    She X*, et al.* Shotgun sequence assembly and recent segmental duplications within the human
878           genome. *Nature* **431**, 927-930 (2004).

879
880    30.    Startek M*, et al.* Genome-wide analyses of LINE-LINE-mediated nonallelic homologous
881           recombination. *Nucleic Acids Res* **43**, 2188-2198 (2015).

882
883    31.    Rozowsky J*, et al.* AlleleSeq: analysis of allele-specific expression and binding in a network
884           framework. *Mol Syst Biol* **7**, 522 (2011).

885
886    32.    Zook JM*, et al.* Integrating human sequence data sets provides a resource of benchmark SNP and
887           indel genotype calls. *Nat Biotechnol* **32**, 246-251 (2014).

888
889    33.    Pendleton M*, et al.* Assembly and diploid architecture of an individual human genome via single-
890           molecule technologies. *Nat Methods* **12**, 780-786 (2015).

891
892    34.    Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM. Canu: scalable and
893           accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome Res* **27**,
894           722-736 (2017).

895
896    35.    Li H*, et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-2079
897           (2009).

898
899    36.    Ruan J, Li H. Fast and accurate long-read assembly with wtdbg2. *bioRxiv*, (2019).

900

901

## Figure Legends

**Figure 1**

**Two types of evidence near SV breakpoints. a**) Type 1 evidence. Reads from HMW DNA molecules that span the breakpoints of a deletion are mapped to two genomic locations, resulting in two sets of observed fragments and two sets of newly introduced fragment endpoints (large dots). **b**) The patterns of enriched fragment endpoints indicate the SV types. Please refer to Supplementary Figures 1-3 and Supplementary Movie 1 for detailed explanations of how the patterns are formed. **c**) Enriched fragment endpoints detected near two breakpoints of a deletion on NA12878 genome. L-endpoints and R-endpoints are plotted separately. The breakpoint positions are marked by red arrows. **d**) Two-dimensional view of enriched endpoints near the two breakpoints of the deletion. Each dot indicates a pair of fragments which share the same barcode and thus may support the SV. The x-value of the dot is the position of the first fragment's R-endpoint and the y-value of the dot is the position of the second fragment's L-endpoint. The background of the 2D plot is cleaner than the 1D plot (panel c) since the fragments that do not share barcodes are excluded. **e**) Type 2 evidence. Reads from two breakpoints of an inversion being mapped to nearby positions (in the grey rectangles), resulting in decreased barcode similarity between the two nearby positions. **f**) Decreased barcode similarity near the breakpoints of an inversion on NA12878 genome. The reciprocal of barcode similarity is shown in the figure. The peaks indicate the positions of the breakpoint.

**Figure 2**

922    **Performance of LinkedSV on the simulated WGS data set. a**) Recalls, precisions and F1

923    scores of six SV callers on the simulated WGS data set. **b**) Fragment endpoint signals of a small

924    duplication that was missed by GROC-SVs. The peaks indicate the approximate breakpoint

925    positions. **c**) Supporting fragments of the tandem duplication. These are fragments span the

926    junction of the first copy and the second copy. Please refer to Supplementary Figure 1 and

927    Supplementary Movie 1 for detailed explanations of how the patterns are formed. Horizontal

928    lines represent linked reads with the same barcode; dots represent reads; colors indicate barcodes;

929    dashed vertical grey lines represent breakpoint positions. **d**) Read depth distribution near the

930    duplication region. The black lines showed the depth of reads with mapping quality $\geq 20$ while

931    the grey lines showed the depth of reads with mapping quality $\geq 0$ (e.g. all reads). Red lines

932    indicate breakpoints predicted by LinkedSV and the blue line indicate the average depth of the

933    whole genome. **e**) Precision of breakpoints predicted by LinkedSV without checking short-read

934    information. **f**) Precision of LinkedSV refined breakpoints using discordant read-pairs and split-

935    reads. Source data is provided as a Source Data file.

936

937    **Figure 3**

938    **Performance of LinkedSV on the simulated WGS data with low variant allele frequencies.**

939    **a, b**) Recalls, precisions and F1 scores of six SV callers on the simulated WGS data set with

940    VAF of 10% and 20%. **c**) Heap map of overlapping barcodes in chr1:193412560-194518464

941    (hg19 coordinates) showing an inversion that was missed by Longranger, and NAIBR (VAF =

942    10%). The overlapping barcodes between the two inversion breakpoints can be clearly visualized

943    (in the black circles). The heat map was plotted by the Loupe software (10X Genomics). Dots

944    represent overlapping barcodes. **d**) Supporting fragments of the inversion detected by LinkedSV.

945    Horizontal lines represent linked reads with the same barcode; dots represent reads; colors
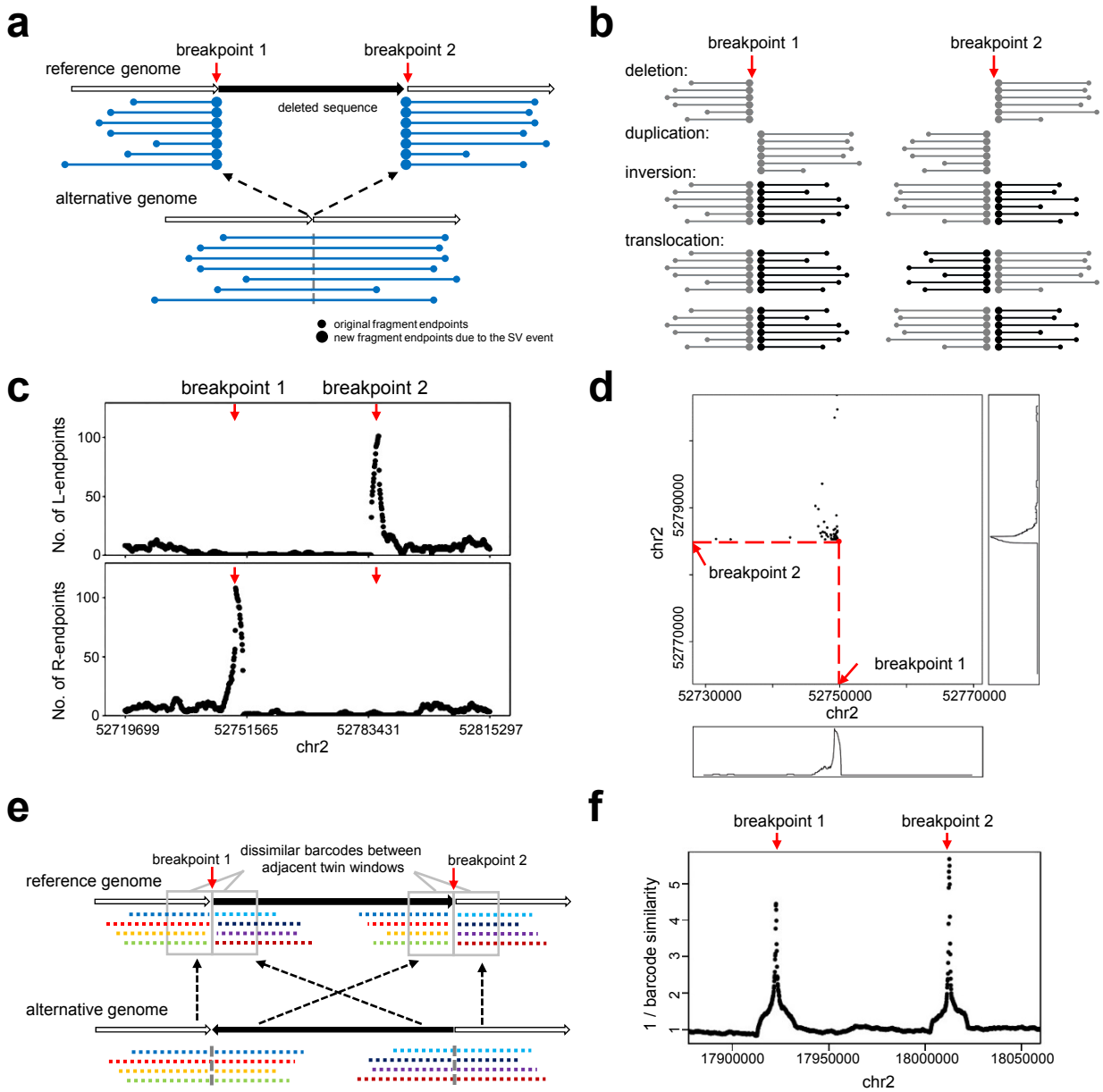
946    indicate barcodes. Predicted breakpoint positions are marked by red arrows. Source data is

947    provided as a Source Data file.

948

949    **Figure 4**

950    **Performance of LinkedSV on the simulated WES data set. a**) Recalls, precisions and F1

951    scores of six linked-read SV callers on the simulated WES data set. **b**) Heat map showing a

952    deletion that was missed by NAIBR. The overlapping barcodes between the two breakpoints can

953    be clearly visualized (in the black circles). The heat map was plotted by the Loupe software.

954    Dots represent overlapping barcodes. **c**) Supporting fragments of the deletion detected by

955    LinkedSV. Horizontal lines represent linked reads with the same barcode; dots represent reads;

956    colors indicate barcodes. Predicted breakpoint positions are marked by vertical grey lines.

957    Capture regions were shown as vertical bars in the bottom. Source data is provided as a Source

958    Data file.

959

960 **Figure 5**

961 **Detection of *F8* inversion from clinical exome sequencing data.** (**a**) Illustration of type I

962 inversion of *F8* gene. A portion of intron 22 has three copies in chrX (int22h-1, int22h-2, int22h-

963 3). The inversion is induced by the homologous recombination between two inverted copies

964 int22h-1 and int22h-3. Int22h-1 is located in intron 22 of *F8* gene and int22h-3 is located in the

965 intergenic regions. (**b**) Heat map of overlapping barcodes in chrX:153916335-154862316 (hg19

966 coordinates), plotted by the Loupe software tool. Black circles indicate overlapping barcodes

967 near the inversion breakpoints. Dots represent overlapping barcodes. (**c**) Decreased barcode

968 similarity at breakpoints detected by the twin window method of LinkedSV. Window size = 40

969 kb (**d**) Supporting fragments detected by LinkedSV. Horizontal lines represent linked reads with

970 the same barcode; dots represent reads; colors indicate barcodes. Dashed vertical grey lines

971 represent breakpoints. Capture regions were shown as vertical bars in the bottom.

972

973 **Figure 6**

974 **Detection of *NF1* deletion from clinical exome sequencing data.** (**a**) Plot of linked-reads for

975 *NF1* WES sample spanning chr17:29645000-29855000. In the normal allele (top), there are 71

976 fragments crossing over the left breakpoint and 38 fragments crossing over the right breakpoint.

977 In the variant allele (bottom), the linked reads are separated by a large gap. Horizontal lines

978 represent linked reads with the same barcode; dots represent reads; colors indicate barcodes.

979    Dashed vertical red lines represent breakpoints. (**b**) Zoom-in plot of supporting fragments for the

980    deletion. One read pair was found to support the deletion.
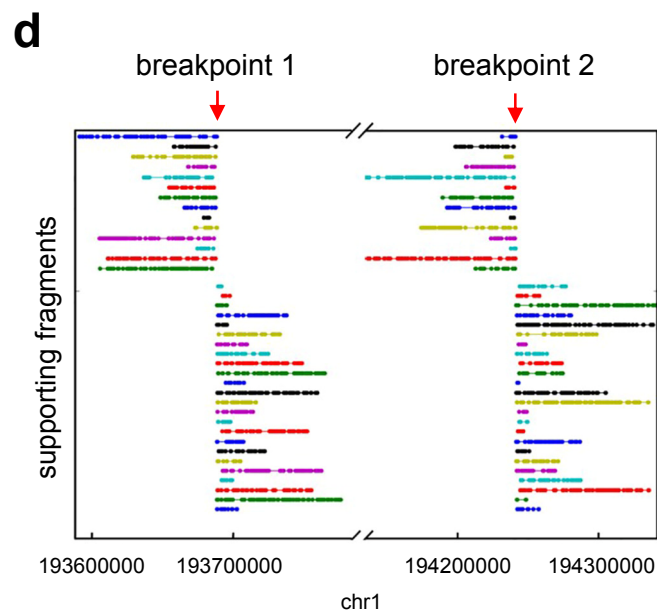
981
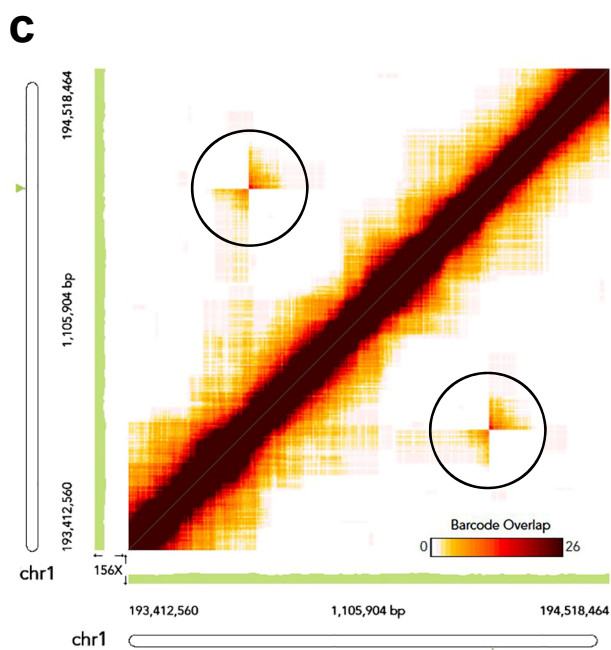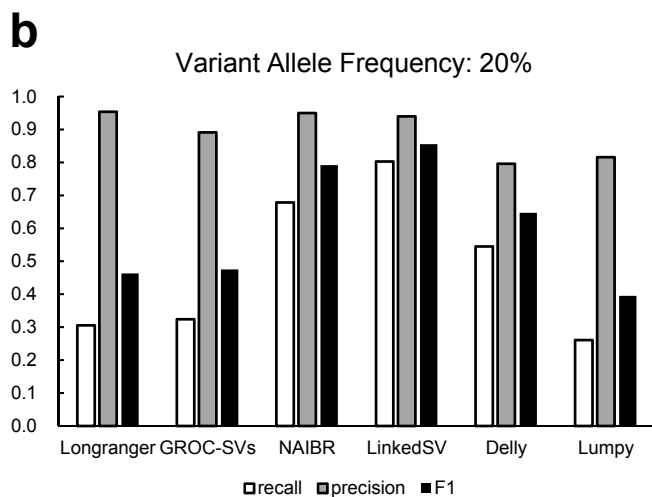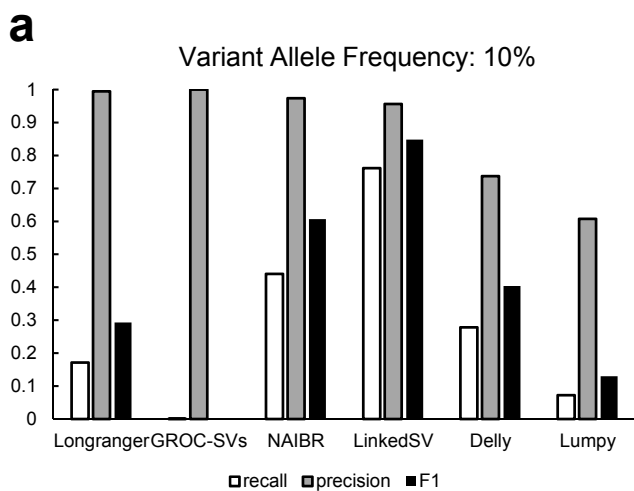
982    **Figure 7**

983    **Structural variants detected from linked-read WGS data of the HX1 genome. a**) Heatmap of

984    overlapping barcodes for a 45 kb deletion on chromosome 2 (chr2:110395971-110441346, hg38

985    coordinates), plotted by the Loupe software tool. Black circles indicate overlapping barcodes

986    near the breakpoints. The deletion was not detected by SMRT-SV from PacBio long reads; **b**)

987    and **c**) Alignments of PacBio reads near the breakpoints of the 45 kb deletion in chr2 in the HX1

988    genome. Clipped reads were marked by vertical pink lines (5'-clipping) or pink arrows (3'-

989    clipping). The figures were generated by IGV. Reads with mapping qualities equal to 0 were in

990    white color. **d**) Read depth distribution near a 61 kb duplication region on chromosome 19

991    (chr19:27338390-27399298, hg38 coordinates). The calculation was based on the bam file of

992    linked-reads. Only reads with mapping quality >= 20 were counted. The dotted blue line showed

993    the average depth across the whole genome. The predicted breakpoints were indicated by vertical

994    red lines. The duplication was not detected by SMRT-SV using PacBio long reads; **e**, **f**) Aligned

995    PacBio raw reads near the two breakpoints of the duplication, as shown in IGV. Increased

996    alignment mismatches due to the SV were observed in **e**) (black rectangles). A clear duplication

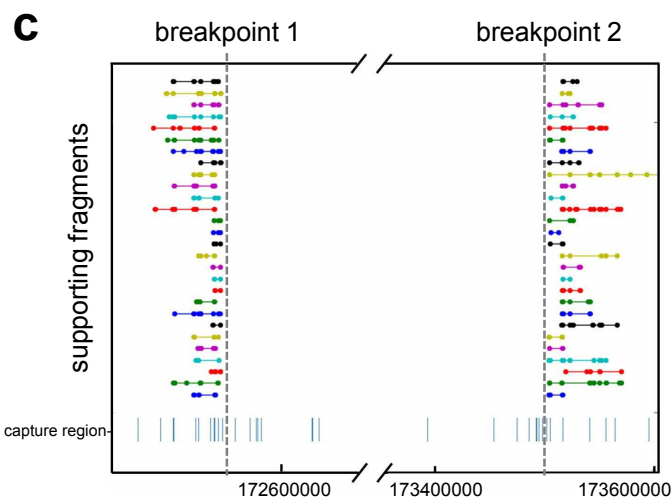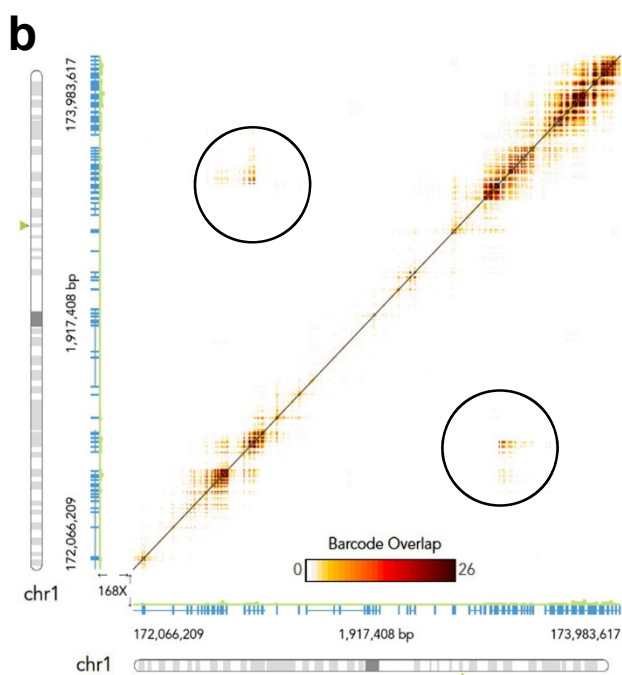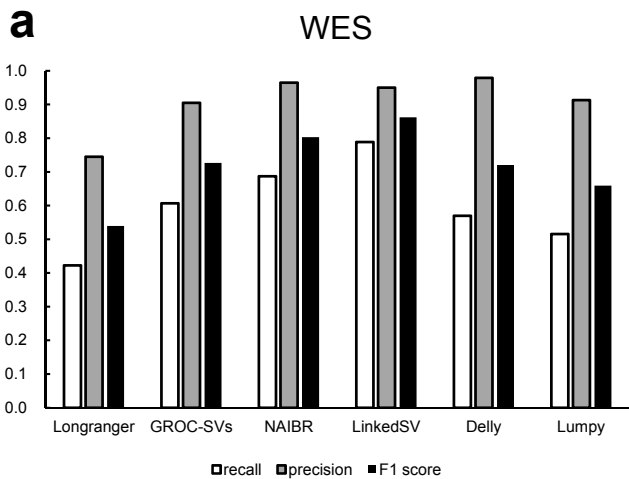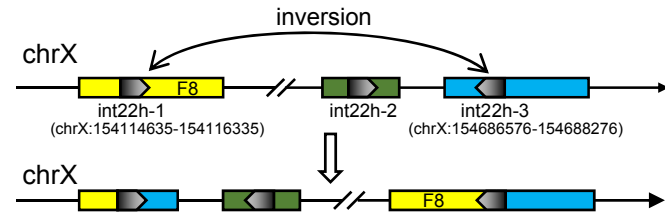997    breakpoint was observed in **f**).

**Figure 1**

**a**

**b** breakpoint 1    breakpoint 2

**c**    breakpoint 1    breakpoint 2

supporting fragments

chr1

**d** read depth

MapQ ≥ 20
MapQ ≥ 0

read depth

**e** breakpoints predicted by endpoints

breakpoint shift (bp)

**f** refined breakpoints

breakpoint shift (bp)

**Figure 2**

**a**

Variant Allele Frequency: 10%

**b**

Variant Allele Frequency: 20%

□recall  ▢precision  ■F1

**c**

Barcode Overlap

0 ▭ 26

chr1

156X

193,412,560    1,105,904 bp    194,518,464

chr1

**d**

breakpoint 1    breakpoint 2

supporting fragments

193600000    193700000    194200000    194300000

chr1

# Figure 3

**a** WES

**b**

Barcode Overlap
0 — 26

**c** breakpoint 1 breakpoint 2

supporting fragments

capture region

**Figure 4**

**Figure 5**

**a**

normal allele

variant allele

predicted deletion region

chr17 position

**b**

one read pair

chr17 position

**Figure 6**

**a**

110,486,721

136,125 bp

110,350,596

128X

110,350,596          136,125 bp          110,486,721

deletion

Barcode Overlap

0                              26

**b**

Alignments of PacBio reads near chr2:110395942

**c**

Alignments of PacBio reads near chr2:110441266

**d**

linked reads

chr19:27338390-27399298 DUP

read depth

120
100
80
60
40
20
0

27277000          27337908          27398816

**e**

**f**

# Figure 7