

# **LinkedSV: Detection of mosaic structural variants from linked-read exome and genome sequencing data**

**Li Fang<sup>1</sup>, Charly Kao<sup>2</sup>, Michael V Gonzalez<sup>2</sup>, Renata Pellegrino da Silva<sup>2</sup>, Mingyao Li<sup>3</sup>, Hakon Hakonarson<sup>2,4</sup>, Kai Wang<sup>1,5\*</sup>**

<sup>1</sup> Raymond G. Perelman Center for Cellular and Molecular Therapeutics, Children's Hospital of Philadelphia, PA 19104, USA

<sup>2</sup> Center for Applied Genomics, Children's Hospital of Philadelphia, PA 19104, USA

<sup>3</sup> Department of Biostatistics, University of Pennsylvania, Philadelphia, PA 19104, USA

<sup>4</sup> Department of Pediatrics, University of Pennsylvania, Philadelphia, PA 19104, USA

<sup>5</sup> Department of Pathology and Laboratory Medicine, University of Pennsylvania, Philadelphia, PA 19104, USA

\* Email: wangk@email.chop.edu

## Abstract

Reliable detection of structural variants (SVs) from short-read sequencing remains challenging, mainly due to the presence of repetitive DNA elements that are longer than typical short reads (~100-150bp). Linked-read sequencing provides long-range information from short-read sequencing data by linking reads originating from the same HMW DNA molecule, and thus has the potential to improve the sensitivity of SV detection and accuracy of breakpoint identification for certain classes of SVs. We present LinkedSV (<https://github.com/WGLab/LinkedSV>), a novel SV detection algorithm which combines two types of evidence. Simulation and real data analysis demonstrated that LinkedSV outperforms several existing tools including Longranger, GROCC-SVs and NAIBR. LinkedSV works particularly well on exome sequencing data and on SVs with low variant allele frequencies due to somatic mosaicism. Our results support the use of linked-read sequencing to detect hidden SVs missed by conventional short-read sequencing approaches and helps resolve negative cases from clinical genome or exome sequencing.

## 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 still remains challenging. The split-reads and discordant read-pairs approaches require the breakpoint-spanning reads/read-pairs being sequenced and confidently mapped. Genomic rearrangements are often mediated by repeats and thus breakpoint junctions of SVs are highly likely to reside in repetitive regions [7-9]. 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 high 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 ten 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 and amplified by primers with a partition-specific barcode. The barcoded DNA molecules are released from the droplets and sequenced by standard Illumina paired-end sequencing [12]. The sequenced short reads derived from the same HMW DNA molecule can be linked together, providing long-range information for mapping, phasing and SV calling. In addition, linked-read whole exome sequencing (WES) has also been developed [12], which provides a more cost-effective way for clinical genetic testing.

In linked-read sequencing data, barcode similarities between any two nearby genome locations are very high, because the reads tend to originate from the same sets of HMW DNA molecules. In contrast, barcode similarities between any two distant genome locations are very low, because the reads of the two genome locations originate from two different sets of HMW DNA molecules and it is highly unlikely that two different sets of HMW DNA molecules share multiple barcodes. Thus, the presence of multiple shared barcodes between two distant locations indicates that the two distant locations are close to each other in the alternative genome [14]. A few pipelines and software tools have adopted this principle to call SVs from linked-read sequencing data, such as Longranger [12], GROC-SVs [14], NAIBR [15]. Longranger is the official pipeline developed by 10X genomics. Longranger bins the genome into 10 kb windows and finds the barcodes of high mapping quality reads within each window. A binomial test is used to find all pairs of regions that are distant and share more barcodes than what would be expected by chance. A probabilistic approach is used to clean up this initial candidate list [12]. GROC-SVs uses a similar method to find candidate SV loci but performed assembly to identify precise breakpoint locations. GROC-SVs also provides functionality to interpret complex SVs [14]. NAIBR detects structural variants using a probabilistic model that incorporates signals from both linked-reads and paired-end reads and into a unified model [15].

However, SV detection from linked-read datasets is still in the early stage. The available SV callers face challenges if we want to detect: i) SVs from target region sequencing (e.g. Whole-exome sequencing); ii) somatic SVs in cancer or somatic mosaic SVs that have low allele frequencies; iii) SVs of which the exact breakpoints have no coverage or located in repeat regions. In this study, we introduce LinkedSV, a novel open source SV caller for linked-read sequencing, which aims to address all the above challenges. LinkedSV detects candidate breakpoints using two types of evidence and quantified the evidence using a novel probabilistic model. We evaluated the performance of LinkedSV on both whole-genome and whole-exome sequencing data sets. In each case, LinkedSV performs better than other existing tools including Longranger, GROC-SVs and NAIBR.

## Results

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

Two types of evidence may be introduced while a genomic rearrangement happens: 1) reads from one HMW DNA molecule which spans the breakpoint being mapped to two genomic locations and 2) reads from two distant genome locations being mapped to adjacent positions. Both types of evidence can be used for SV detection.

First, we describe the signals of type 1 evidence. After reads mapping, the original HMW DNA molecules can be computationally reconstructed from the sequenced short reads using their barcodes and mapping positions. In order to distinguish from the physical DNA molecules, we use “fragments” to refer to the computationally reconstructed DNA molecules. A fragment has a left-most mapping position, which we call 5'-endpoint, and a right-most mapping position, which we call 3'-endpoint. As a result of genomic rearrangement, reads from one breakpoint-spanning HMW DNA molecule would be mapped to two different genome loci on the reference genome. This split-molecule event has two consequences: 1) observing two fragments sharing the same barcode; 2) each of the two fragment has one endpoint close to the true breakpoints. Therefore, in a typical linked-read WGS data set, multiple split-molecule event could be captured and we could usually observe: multiple share barcodes between two distant genome loci and multiple fragment endpoints near the breakpoints.

To illustrate this, Figure 1a shows the split-molecule events of a deletion, where breakpoints 1 and 2 are marked by red arrows. Multiple fragment endpoints enriched near the two breakpoints of a large deletion, which is typically longer than 5 kb for whole-genome sequencing (WGS) data sets. Figure 1b shows the patterns of “enriched fragment endpoints” that are introduced by different types of SVs. As an example, Figure 1c shows the number of fragment endpoints in a 5-kb sliding window near two deletion breakpoints, based on a 35X coverage linked-read WGS data on the NA12878 genome (genome of a female individual extensively sequenced by multiple platforms). At the breakpoints, the number of fragment endpoints in the 5-kb sliding window is more than 100 and is five times more than normal regions, forming “peaks” in the figure.

Since the fragments can be paired according to their barcodes, we can also observe fragment endpoints of this deletion in a two-dimensional view. As shown in Figure 1d, each dot indicates two endpoints from a pair of fragments which share the same barcode. The x-value of the dot is the position of the first fragment's 3'-endpoint and the y-value of the dot is the position of the second fragment's 5'-endpoint. The bottom panel and right panel in Figure 1d shows number of dots that are projected to the x-axis and

y-axis. Similar with the one-dimensional plot (Figure 1c), a peak is formed near each breakpoint, which is marked by the red arrow. The background noise of the two-dimensional plot is cleaner than the one-dimensional plot since the fragments that do not share barcodes are excluded. Therefore, the two-dimensional plot is more useful when the variant allele frequency (VAF) is very low and there are only a few supporting fragments.

Next, we describe the signals of type 2 evidence. The barcodes between two nearby genome locations is highly similar because the two locations are spanned by almost the same set of input HMW DNA molecules. However, due to the genome rearrangement, the reads mapped to the left side and right side of a breakpoint may originate from different locations of the alternative genome and thus have different barcodes (Figure 1e). Dropped barcode similarity between two nearby loci therefore indicates a SV breakpoint. LinkedSV detects this type of evidence by a twin-window method, which uses two adjacent sliding windows to scan the genome and find regions where the barcode similarity between the two nearby window regions are significantly decreased. Figure 1f illustrates an inversion breakpoint detected by LinkedSV on the NA12878 genome. The change of barcode similarity was plotted and a peak was formed at the breakpoint. After searching for the two types of evidence, LinkedSV combines the candidate SV regions and quantified the evidence using a novel probabilistic model. The breakpoints are further refined using short-read information, including discordant read pairs and split-reads.

### **Performance evaluation on simulated whole-genome sequencing data set**

To assess LinkedSV's performance, we simulated a 35X linked-read WGS data set with 1175 SVs inserted using LRSIM[13] (see Methods for details). The inserted SVs includes 351 deletions, 386 duplications, 353 inversions, and 85 translocations. The breakpoints of the simulated SVs are designed to be located in repeat regions, since we found that the LinkedSV and other available SV callers perform very well when the breakpoints are located in non-repeat regions, and thus we set to test the performances of all the SV callers under more challenging situations. This makes sense because SV breakpoints are more likely to be in repeat regions [7-9], and because these situations represent those that are impractical to be addressed by conventional short-read sequencing approaches.

The simulated reads were aligned to the reference genome using the Longranger [12] package provided by 10X genomics. The Longranger pipeline internally uses the Lariat aligner [16], which was designed for the alignment of linked reads. SV calling was performed using LinkedSV as well as three other available SV callers designed for linked-read sequencing: Longranger, GROCSVs [14] and NAIBR [15].

We used recalls, precisions and F1 scores to evaluate the four SV callers on this data set. As shown in Figure 2a, LinkedSV has the highest recall and F1 score among all methods. NABIR has a good recall (0.81), but relatively low precision (0.42). GROC-SVs has a good precision but its recall is lower than LinkedSV, so we further analyzed the false negative calls of GROC-SVs to understand this underlying reasons. A major portion of the false negative calls by GROC-SVs represents duplications that are smaller than twice of the fragment length. For large duplications, the reads of the alternative allele are separated by a large gap so that we can observe two sets of fragments with the same set of barcodes, which indicate a SV (Supplemental Figure 5a). If the duplication is not large enough, the reads will be probably clustered into one fragment (Supplemental Figure 5b). Even in this case, we can observe enriched fragment endpoints near the duplication breakpoints in LinkedSV. As an example, Figure 2b shows the endpoint signals of a missed duplication call of GROC-SVs. The supporting fragments of this duplication is shown in Figure 2c. We also evaluated the breakpoint precision of LinkedSV. Most of breakpoints predicted by fragment endpoints are within 20 bp (Figure 2d) and refined breakpoints using discordant read-pairs and split-reads have base-pair resolution (Figure 2e).

## **Performance evaluation on WGS data set with somatic SVs with low variant allele frequencies**

Somatic SVs are often found in cancer genomes [17-19]. However, due to the high heterogeneity of genomic alteration in cancer genomes, somatic SVs often have low allele frequencies (as opposed to ~50% in a germline genome) and thus are more difficult to detect by typical SV callers designed for germline SVs. We simulated two WGS data sets with VAF of 10% and 20%. Recalls, precisions and F1 values of the four linked-read SV callers were evaluated on the two data sets (Figure 3a, Figure 3b). When the VAF was 20%, the recall of LinkedSV (0.803) was much higher than that of Longranger (0.238), GROC-SVs (0.402) and NAIBR (0.679). The precisions of the four SV callers range from 0.87-0.95. The F1 score of LinkedSV (0.855) was also the highest among all the SV callers. When the VAF was 10%, LinkedSV still had a recall of 0.761, which was 72% higher than the second best SV caller NAIBR. Longranger and GROC-SVs almost completely failed to detect the SVs. These observations confirmed that other SV callers were mainly designed for germline genomes and had substantial difficulty in detecting SVs with somatic mosaicism. However, due to the combination of barcode overlapping and enriched fragment endpoints in our statistical model (see Methods for details), LinkedSV was able to achieve a good performance even when VAF was very low. We manually checked the barcode overlap evidence of some SV calls using the Loupe software developed by 10X Genomics. Figure 3c showed an inversion that was missed by Longranger, and NAIBR but detected by LinkedSV (at VAF of 10%). Although the variant

frequency is low, the overlapped barcodes between the two inversion breakpoints be clearly visualized (in the black circle) in the figure. Figure 3d showed the supporting fragments of the inversion detected by LinkedSV. Each horizontal line represent two fragments which shared the same barcode and support the SV. This suggests that the manufacturer-provided software tool has limitations for SV detection, despite its strong functionality in visualization.

### **Performance evaluation on simulated whole-exome sequencing data set**

Compared with WGS, whole-exome sequencing (WES) is currently widely used in clinical settings to identify disease causal variants on patients with suspected genetic diseases, partly due to the slightly lower cost of WES and due to reimbursement reasons. Since the WES only cover a small portion of regions in the whole genome, it is far more challenging to detect SVs from WES data, especially when the SV breakpoints are not in the capture regions. However, by combining linked-read sequencing with WES capture platforms, it is possible to alleviate this problem, and significantly improve the sensitivity of SV detection using WES.

To evaluate SV detection on linked-read WES data, we simulated a 40X coverage linked-read WES data set with 1160 heterozygous SVs (see Methods for details). 44.3% of the breakpoints are not in exon regions. SV calling was performed using LinkedSV, Longranger, GROC-SVs and NAIBR. As shown in Figure 4a, LinkedSV has the highest recall (0.79) and highest F1 score (0.86). NAIBR has the highest precision (0.97) but its recall is lower (0.69). GROC-SVs has a good precision (0.90) but the recall (0.61) is not high.

We analyzed false negative calls of the second best SV caller NAIBR. NAIBR tends to miss some SV events that have shared barcodes but lack short-read support. For example, Figure 4b shows a deletion between chr1:172545561-173504265. Both breakpoints are not located in capture regions. Breakpoint 1 (chr1:172545561) is 768 bp away from the nearest capture region and breakpoint 2 (chr1: 173504265) is 392 bp away from the nearest capture region. Unfortunately, no discordant read pairs that support the deletion can be found. However, shared barcodes between the two breakpoints can be clearly seen using the loupe software (Figure 4b). In addition, LinkedSV also detected 28 pairs of fragments that share the same barcodes and support the SV. These fragments were plotted in Figure 4c. Although no short-read support was found, the SV type could be determined using the pattern of “enriched fragment endpoints” shown in Figure 1b. In this SV event, 3'-endpoints were highly enriched for the first set of fragments and

5'-endpoints were highly enriched for second set of fragments. Thus, the SV type was predicted as "deletion".

## **Detection of F8 inversion from clinical exome sequencing data on a patient with hemophilia A**

We also tested the performance of LinkedSV in a WES sample with known SV breakpoints in introns from a male patient with Hemophilia A. Previous experiments had shown that the patient had type I inversion in intron 22 of F8 gene. The F8 gene is located in Xq28. The intron 22 of F8 gene contains a GC-rich sequence (named int22h-1) that is duplicated at two positions towards the Xq-telomere (int22h-2 and int22h-3). Int22h-2 has the same direction with int22h-1 while int22h-3 has the inverted direction. The type I inversion is induced by the recombination between int22h-1 and int22h-3 [20, 21] (Figure 5a). BLAST alignment of int22h-1 and int22h-3 showed that the two sequences had 99.88% identity. Since the breakpoints were located in two segmental duplications with nearly identical sequences, the inversion is undetectable by conventional short-read sequencing. Two popular short-read SV callers Delly [3] and Lumpy [22] failed to detect the inversion from the linked-read WES data (results were shown in Supplemental Table 1 and Supplemental Table 2).

Longranger, GROCSVs, NAIBR and LinkedSV were also used to detect SVs from this sample. None of Longranger, GROCSVs and NAIBR detected this inversion (results were shown in Supplemental Table 3, Supplemental Table 4, Supplemental Table 5), although the overlapped barcodes can be visualized using the Loupe software (Figure 5b). However, LinkedSV successfully detected this inversion, by combining two types of evidence. As described above, barcode similarity between two nearby regions are very high but drops suddenly at the breakpoints. Figure 5c shows the suddenly drop of barcode similarity at the two breakpoints. Each dot in the figure represents the reciprocal of the barcode similarity between its left 40 kb window and right 40 kb window thus the Y value of the dots were reversely related to the barcode similarity and positively related to the probability of being a breakpoint. The barcode similarities are lowest at the two breakpoints and thus form two peaks in the figure (marked by red arrow). In addition, LinkedSV also identified the supporting fragments of the SV using type 1 evidence (Figure 5d). The predicted breakpoint positions are consistent with the genomic positions of int22h-1 and int22h-3.

## **Methods**

### **Breakpoint detection from type 1 evidence**



First, LinkedSV reconstructs the original long DNA fragments from the reads using mapping positions and barcode information. All mapped reads are partitioned according to the barcode and sorted by mapping position. We define gap distance as the distance between two nearest reads with the same barcode. Two nearby reads are considered from the same long DNA fragment if they have the same barcode and their gap distance is less than a certain distance  $G$ .  $G$  is determined using two steps. First, we use  $G = 50$  kb (the same as Zheng et.al [12]) to group the reads into fragments. This value is suitable for detection of large SVs. However, it may be too large for detection of SVs that are smaller than 50 kb. Therefore, we calculate the empirical distribution of intra-fragment gap distance, which is the distance of two nearby reads that are grouped in one fragment. We assign  $G$  as the 99<sup>th</sup> percentile of the empirical distribution of intra-fragment gap distance.  $G$  is usually between 8~15kb, depending on the data set. Fragments with a gap distance larger than  $G$  will be separated to two fragments.

In non-SV regions, all the reads from the same HMW DNA molecule would be reconstructed into a single DNA fragment. The reads from the breakpoint-spanning HMW DNA molecule will be mapped to two different positions in the genome. As illustrated in the Result section, this split-molecule event has two consequences: 1) observing two fragments sharing the same barcode; 2) each of the two fragment has one endpoint close to the breakpoints. Therefore, we could observe enriched fragment endpoints near the breakpoints, in both one-dimensional view (Figure 1c) and two-dimensional view (Figure 1d). The type of the endpoints (5' or 3') that enriched near the breakpoints depends on the type of SV (Figure 1b). The two-dimensional view has less background noise because the fragments that do not share barcodes and thus do not support the SVs are excluded. Thus, we detect the enriched endpoints in the two-dimensional view.

We now describe how we use detect the type 1 evidence of deletion calls, but the method can be applied to other types of SVs. We define *fragment pair* to be two fragments sharing the same barcode. Let  $b_1, b_2$  be the positions of the two breakpoint candidates (assuming  $b_1 < b_2$ ). Let  $n$  be the number of fragment pairs that may support the SV between  $b_1$  and  $b_2$ . Let  $F_{i1}, F_{i2}$  denote the  $i^{\text{th}}$  fragment pair that support the SV. Let  $B(F)$  denote the barcode of fragment  $F$ . Therefore, we have:

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

Let  $L(F)$  denote the 5'-endpoint position (i.e., left-most position) of fragment  $F$ ,  $R(F)$  denote the 3'-endpoint position (i.e., right-most position) of fragment  $F$ . Since this is a deletion and  $b_1 < b_2$ ,  $R(F_{i1})$  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  $b_2$  (Supplemental Figure 6a). The distance between the fragment endpoint and its corresponding breakpoint should be within gap distance distribution (explained in Supplemental Figure 4). Therefore, we have:

$$b_1 - G \leq R(F_{i1}) \leq b_1, b_2 \leq L(F_{i2}) \leq b_2 + G \quad (2)$$

As described above,  $G$  is the 99<sup>th</sup> percentile of the empirical distribution of intra-fragment gap distance.

If we regard  $(R(F_{i1}), L(F_{i2}))$  as a point in a two-dimensional plane, according to equation (2),  $((R(F_{i1}), L(F_{i2})))$  is restricted in a square region with a side length of  $G$  and the point  $(b_1, b_2)$  being a vertex (Supplemental Figure 6b).

We used a graph-based method to fast group the points into clusters and find square regions where the numbers of points were more than expected. First, every possible pair of endpoints  $(R(F_1), L(F_2))$  meeting  $B(F_1) = B(F_2)$  formed a point in the two-dimensional plane. Each point indicated a pair of fragments that share the same barcode. For example, if 10 fragments share the same barcode,  $C_{10}^2$  pairs of endpoints will be generated. A point/pair of endpoints may or may not support an SV because there are two possible reasons for observing two fragments sharing the same barcode: 1) the two fragments originated from two different HMW DNA molecules but were dispersed into the same droplet partition and received the same barcode; 2) the two fragments originated from the same HMW DNA molecule but the reads were reconstructed into two fragments due to an SV. The points are sparsely distributed in the two-dimensional plane and it is highly unlikely to observe multiple points in a specific region. Next, a k-d tree ( $k = 2$ ) was constructed, of which each node stores the  $(X, Y)$  coordinates of one point. K-d tree is binary tree that enable fast query of nearby nodes. Therefore, we could quickly find all pairs of points within a certain distance. Any two points  $(x_1, y_1)$  and  $(x_2, y_2)$  were grouped into one cluster 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 than a user-defined threshold (default: 5), it was considered as a potential region of enriched fragment endpoints. If the points in the cluster were not within a  $G \times G$  square region, we used a  $G \times G$  moving square to find a square region which contained the most of points. The predicted breakpoints were the X and Y coordinates of the right-bottom vertex of the square. The points in the square region were subjected to a statistical test describe below.

## Quantification of type 1 evidence

Let  $n$  be the number of points in the square region. Each point corresponds to a pair of fragment  $F_{i1}, F_{i2}$ , ( $i = 1, 2, 3, \dots, n$ ) that may support the SV. Let  $b_1$  and  $b_2$  be the coordinates of the predicted breakpoint. Equation (1) and (2) hold for all the fragment pairs  $F_{i1}, F_{i2}$  ( $i = 1, 2, 3, \dots, n$ ). We then test the null hypothesis that there is no SV between  $b_1$  and  $b_2$ .

First, we test the hypothesis that the  $n$  fragment pairs  $F_{i1}, F_{i2}$  have originated from different DNA molecules, but coincidentally received the same barcode. Here we define two fragments  $F_a$  and  $F_b$  as an *independent fragment pair* if  $F_a$  and  $F_b$  share the same barcode but have originated from different DNA molecules. Thus,  $R(F_a)$  and  $L(F_b)$  are independent variables. All the fragment pairs that do not support SVs are independent fragment pairs. It is reasonable to assume the generation of HMW DNA molecules from chromosomal DNA is a random process thus both  $R(F_a)$  and  $L(F_b)$  are uniformly distributed across the chromosome. Therefore, the point  $((R(F_a), L(F_b)))$  is equal likely to be in any place in the two-dimensional plane. Technically, we connect all the chromosomes in a head-to-tail order so that both intra-chromosomal events and inter-chromosomal can be analyzed at the same time. Observing at least  $n$  independent fragment pairs meeting equation (2) is equivalent to the event that observing at least  $n$  points  $((R(F_{i1}), L(F_{i2})))$  located in a squared region with an area of  $G^2$  on the two-dimensional plane. The probability of this event is:

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

where Binomial\_pmf is the probability mass function of binomial distribution;  $L$  is the total length of the genome (also the side length of the two-dimensional plane);  $N_{ifp}$  is the total number of independent fragment pairs.

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

$$p_{adjusted1} = p_1 \frac{G^2}{L^2}$$

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

Next, we test the hypothesis that fragment pairs  $F_{i1}, F_{i2}$  ( $i = 1, 2, 3, \dots, n$ ) have originated from the same DNA molecule, but no reads were sequenced in the gap between  $R(F_{i1})$  and  $L(F_{i2})$ . Let  $g_i$  denote the length of the gap between  $F_{i1}$  and  $F_{i2}$ ,  $\bar{g}$  denote the mean of  $g_i$ , we have:

$$g_i = L(F_{i2}) - R(F_{i1}) \quad (4)$$

$$\bar{g} = \frac{1}{n} \sum_{i=0}^n g_i \quad (5)$$

If  $\bar{g}$  is too large such that the probability of no reads being generated is smaller than a threshold, we can reject this hypothesis.

Similar to the model described by 10X Genomics[12], we assume the read generation on a DNA molecule is a Poisson process with constant rate  $\lambda$  across the genome. Let  $r$  be the number of reads generated in a region of length  $g$ , then  $r \sim \text{Pois}(\lambda g)$ . Let  $P_{gap}(g)$  denote the probability of no read being generated in length  $g$ , we have:

$$P_{gap}(g) = P(r = 0 | \lambda g) = \frac{e^{-\lambda g} (\lambda g)^0}{0!} = e^{-\lambda g} \quad (6)$$

Therefore, the gap length  $g_i$  follows Exponential distribution:  $g_i \sim \text{Exp}(\lambda)$ . Recalling that 1) the Exponential distribution with rate parameter  $\lambda$  is a Gamma distribution with shape parameter 1 and rate parameter  $\lambda$ ; 2) the sum of  $n$  independent random variables from Gamma (1,  $\lambda$ ) is a Gamma random variable from Gamma ( $n$ ,  $\lambda$ ), we have:

$$\sum_{i=0}^n g_i \sim \text{Gamma}(n, \lambda) \quad (7)$$

$$\bar{g} = \frac{\sum_{i=0}^n g_i}{n} \sim \text{Gamma}(n, n\lambda) \quad (8)$$

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

$$p_2 = 1 - \text{Gamma\_cdf}(n, n\lambda) \quad (9)$$

where Gamma\_cdf is the cumulative distribution function of Gamma distribution.

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

$$p_{adjusted2} = p_2 \frac{N_{rp}}{n} \quad (10)$$

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

We reject the hypothesis if  $p_{adjusted2} < p_{threshold}$ .  $p_{threshold}$  is set as  $10^{-5}$  by default. If both  $p_{adjusted1}$  and  $p_{adjusted2}$  are less than  $p_{threshold}$ , we accept the hypothesis that the SV is true. For each candidate SV, we report a confidence score for type 1 evidence as:

$$\text{Confidence score1} = -\log_{10}(\max(p_{adjusted1}, p_{adjusted2})) \quad (11)$$

## Breakpoint detection from type 2 evidence

Barcode similarity between two nearby regions is very high because the reads originate from almost the same set of HMW DNA molecules. However, at the SV breakpoint, the aligned reads from the left side and right side may have originated from different locations in the alternative genome. Thus, the barcode similarity between the left side and right side of the breakpoint are dramatically reduced (as described in the Result section and shown in Figure 1e-f). To detect this evidence, LinkedSV uses two adjacent sliding windows (twin windows, moving 100 bp) to scan the genome and calculate the barcode similarity between the twin windows. The window length can be specified by user. By default, it is  $G$  for WGS data sets and 40 kb for WES data sets.

The barcode similarity can be simply calculated as the fraction of shared barcodes. This method is suitable for whole-genome sequencing (WGS), where the coverage is continuous and uniform. But it does not perform well for whole exome sequencing (WES), where the numbers of reads in the sliding windows vary a lot due to capture bias and the length of capture regions. Therefore, we use a model that considering the variation of sequencing depth and capture region positions. The barcode similarity is calculated as:

$$S = \frac{x}{m_1^a m_2^b} n e^{-\alpha d} \quad (12)$$

where:

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

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

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

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

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

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

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

$n$ ,  $\alpha$ ,  $a$  and  $b$  are estimated from the data using regression. Detailed explanation of this model is in supplemental note 1.

Next, we calculate the empirical distribution of barcode similarity. Regions where the barcode similarity less than a threshold (5<sup>th</sup> percentile of the empirical distribution by default) were regarded as breakpoint candidates. If a set of consecutive regions have barcode similarity lower than the threshold, we only retain

the region that has the lowest barcode similarity. If the barcode similarity of a breakpoint candidate is  $S_0$ , the empirical  $p$ -value is calculated as:

$$\text{empirical } p = \frac{\text{number of twin windows with } S \leq S_0}{\text{total number of twin windows}} \quad (13)$$

The confidence score of type 2 evidence is:

$$\text{Confidence score2} = -\log_{10}(\text{empirical } p) \quad (14)$$

### **Combination of both types of evidence**

Type 1 evidence gives pairs of endpoints that indicate two genomic positions are joined in the alternative genome. Type 2 evidence gives genomic positions where the barcodes suddenly changed, regardless of which genomic position can be joined. Therefore, type 1 and type 2 evidence are independent. The candidate breakpoints detected from type 2 evidence were searched against the candidate breakpoint pairs detected from type 1 evidence so that the calls were merged. The combined confidence score is:

$$\text{Combined score} = \text{Confidence score1} + \text{Confidence score2a} + \text{Confidence score2b}$$

Where Confidence score1 is the confidence score calculated from type 1 evidence (equation 11); Confidence score2a and Confidence score2b are the confidence scores of the two breakpoints calculated from type 2 evidence (equation 14).

### **Refining breakpoints using discordant read-pairs and split-reads**

We search for discordant read-pairs and clipped reads that are within 10 kb to the predicted breakpoint pairs by the above approach. We use a graph-based approach that is similar to DELLY[3] to cluster the discordant read-pairs. We define the supporting split-reads as the clipped reads that can be mapped to the both breakpoints, and the map direction matches the SV type. If both discordant read-pairs and split-reads are found to support the SV, we use the breakpoints inferred by split-reads as the final breakpoint position.

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

The linked reads were simulated by LRSIM, which can generate linked-reads from a given FASTA file containing the genome sequences. We generated a diploid FASTA file based on hg19 reference genome with SNPs and SVs inserted. The purpose of inserting SNPs was to mimic real data. The generation of the diploid FASTA file is described below. First, we inserted SNPs to hg19 using vcf2diploid [23]. The inserted SNPs were from the gold standard SNP call set (v.3.3.2) of NA12878 genome [24]. The vcf2diploid software generated two FASTA files, each of which was a pseudohaplotype (paternal or maternal) with the phased SNPs inserted. Next, we insert SVs into the paternal FASTA file using our custom script. The breakpoints were located in the repetitive regions in hg19 and the distance between the two breakpoints were in the range of 50 kb to 1 Mb. In total, we simulated 351 deletions, 386 duplications, 353 inversions and 85 translocations, all of which were in the paternal copy and were heterozygous SVs. We then concatenate the paternal and maternal FAST file into a single FASTA file and simulated linked-reads using LRSIM. To mimic real data, the barcode sequences and molecule length distribution used for simulation were from the NA12878 whole-genome data set released by 10X genomics. The number of read pairs was set to 360 million so that a 35X coverage data set was generated.

### **Generation of WGS data set with low variant allele frequencies**

In cancer samples or mosaic samples, the total DNA is a mix of a small portion of variant alleles and a large portion of normal alleles. To simulate the WGS data sets with low variant frequencies, we used the same paternal and maternal FASTA file described above but the combined FASTA file contained multiple copies of the normal allele (the maternal FASTA) and only one copy of the variant allele (the paternal FASTA). For example, to simulate a WGS data set with VAF of 20%, four copies of the maternal FASTA and one copy of the paternal FASTA were combined. The linked reads were simulated using LRSIM with the same parameters and a 35X coverage data set was generated.

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

To generate the linked-read WES data set, we first generate a 100X linked-read WGS data set and then down-sample it to be a WES data set. Generation of the simulated linked-read WGS data set with SNPs and SVs inserted was similar to the method described above. In total, we inserted 1160 heterozygous SVs. The SV breakpoints were randomly selected from regions that were within 2000 bp of an exon. Among the 2320 breakpoints (two breakpoints per each SV), 1028 breakpoints (44.3%) were in intronic or intergenic regions. The SV sizes are in the range of 50 kb to 1 Mb. The number of inserted SVs in the

simulated WES data set was slightly smaller than that in the simulated WGS data set because the SV breakpoints were designed to reside within 2000 bp of an exon. The simulated reads were generated using LRSIM and were mapped to hg19 reference genome using the Longranger pipeline (default settings). The phased bam generated by Longranger was down-sampled to be a simulated WES data set. To mimic real WES data set, we used the coverage distribution of the linked-read WES data set of NA12878 genome (released by 10X genomics) to guide the down-sampling process. We bin the genome into 10 bp windows and calculate number of reads mapped to each window (left mapping positions were used) in NA12878 linked-read WES data. The simulated WES data set was generated by sampling reads from the 100X WGS data according to number of reads mapped to the same 10 bp window in the NA12878 WES. In this process, if the one read is retained, the paired read will also be retained.

## Discussion

In this study, we present LinkedSV, a novel open source algorithm for structural variants detection from linked-read sequencing data. We assessed the performance of LinkedSV on three simulated data sets and one real data set. By incorporating two types of evidence, LinkedSV outperforms all existing linked-read SV callers including Longranger, GROCSV and NAIBR on both WGS and WES data sets.

Type 1 evidence gives information about which two genomic positions are connected in the alternative genome. It has two observations: 1) fragments with shared barcodes between two genomic locations and 2) enriched fragment endpoints near breakpoints. Current existing linked-read SV callers only use the first observation to detect SVs while LinkedSV incorporates both observations in the statistical model and is therefore more sensitive and can detect SVs with lower allele frequencies, such as somatic SVs in cancer genomes and mosaic structural variations.

Type 2 evidence gives information about which genomic position is “interrupted” with the observation that the reads on the left side and right side of a genomic position have different barcodes and should be derived from different HMW DNA molecules. LinkedSV is the only SV caller that use type 2 evidence to detect breakpoints. Type 2 evidence is independent to type 1 evidence, and gives additional confidence to identify the breakpoints. In addition, type 2 evidence can be detected locally, which means we can detect a weird genomic location without looking at the barcodes of the other genomic locations. This is particularly useful in two situations: 1) novel sequence insertions where there is only one breakpoint; 2) only one breakpoint is detectable and the other breakpoint located in a region where there is little coverage within 50 kb, which is often the case in target region sequencing.



In recent years, WES has been widely used to identify disease causal variants for patients with suspected genetic diseases in clinical settings. Identification of SVs from WES data sets are more challenging because the SV breakpoints may not be in the capture regions and thus there would be little coverage at the breakpoints. Linked-read sequencing increases the chance of resolving such type of SVs by providing long-range information. As well as there are a few capture regions nearby, the fragments can still be reconstructed and type 1 and type 2 evidence can still be observed. Our statistical models for both type 1 evidence and type 2 evidence were designed to handle both WGS and WES data sets. GROC-SVs uses a local-assembly method to verify the SV call, which requires sufficient coverage at the breakpoints. By using two types of evidence, LinkedSV can be less relied on short-read information (e.g. pair-end reads and split-reads). We demonstrated that LinkedSV has better recall and balanced accuracy (F1 score) on the simulated WES data set and can detect SVs even when the breakpoints were not located in capture regions and have no short-read support. In addition, LinkedSV is also the only SV caller that clearly detected the F8 intron 22 inversion from the WES data set.

Linked-read sequencing has several advantages over traditional short-read sequencing on the purpose of SV detection. First, the human genome is highly repetitive. Previous studies have shown that SVs are closely related to repeats and many SVs are directly mediated by homologous recombination between repeats [25]. In traditional short-read sequencing, if the breakpoint falls in a repeat region, the supporting reads would be multi-mapped and thus the SV cannot be confidently identified. However, this type of SVs are detectable by Linked-read sequencing as well as the HMW DNA molecules span the repeat region. We can observe type 1 and type 2 evidence in the non-repeat region nearby. In our benchmarking, LinkedSV resolved 95% SVs of which breakpoints located in repeat regions. Secondly, SVs are undetectable from traditional short-read sequencing if there is little coverage at the breakpoints, which is often the case in WES data sets. As described above, this type of SV can also be resolved by linked-read sequencing and LinkedSV. Third, linked-read sequencing requires less coverage for detection of SVs with low variant allele frequencies. In linked-read sequencing data, short read pairs are sparsely and randomly distributed along the HMW DNA molecule. In a typical linked-read WGS data set, the average distance between two read pairs derived from the same HMW DNA molecule is about 1000 bp and each HMW DNA molecule only has a short-read coverage of about 0.2X. Therefore, there are about 150 HMW DNA molecules (reconstructed fragments) covering a genomic location of 30X depth. A SV of 10% VAF will have 15 supporting fragment pairs in a 30X depth location in linked reads WGS data set, which is sufficient to be detected by LinkedSV. However, a SV of 10% VAF will only 3 supporting read pairs in a 30X depth location in traditional short read WGS, which makes the detection more challenging.

Linked-read sequencing also has several advantages over long-read sequencing for SV detection. Linked-read sequencing has a much lower cost compared with long-read sequencing technologies. Thus, it is possible to sequencing a sample at very high converge and find SVs with very low VAF. It is also possible to do population-based large cohort studies. In addition, long-read sequencing technologies tend to have higher error rates (13-15%) and there may be some false positive calls due to sequencing errors [11].

The linked-read technology provides strong evidence to detect large SVs, but it provides little additional evidence to detect small SVs. Therefore, LinkedSV has limited power to detect SVs that are shorter than 10 kb. However, large SVs are more likely to be harmful and to cause diseases. Therefore, we still expect a strong ability of linked-read technology to resolve harmful and disease causal SVs. Like the existing linked-read SV callers, LinkedSV currently does not handle novel sequence insertions. As future work, we plan to detect novel sequence insertions using type 2 evidence, since this type of SV also cause a decrease of barcode similarity between nearby regions and can be detected by the twin-window method.

In summary, we present LinkedSV, a novel SV caller for linked-read sequencing. LinkedSV outperformed current existing SV caller, especially for identifying SVs with low allele frequency or identifying SVs from target region sequencing such as linked-read WES. We expect LinkedSV will facilitate the detection of structural variants from linked-read sequencing data and help solve negative cases of conventional short-read sequencing.

## **Competing Interests**

The authors declare no competing interests.

## **Acknowledgments**

The authors would like to thank members of the Center for Applied Genomics for generating the linked read exome sequencing data on the patient with F8 inversion, and thank 10X Genomics for making the NA12878 linked read genome sequencing data publicly available for benchmarking software tools. We would like to thank the authors of the simulation software LRSIM to provide tools that facilitated our benchmarking study.

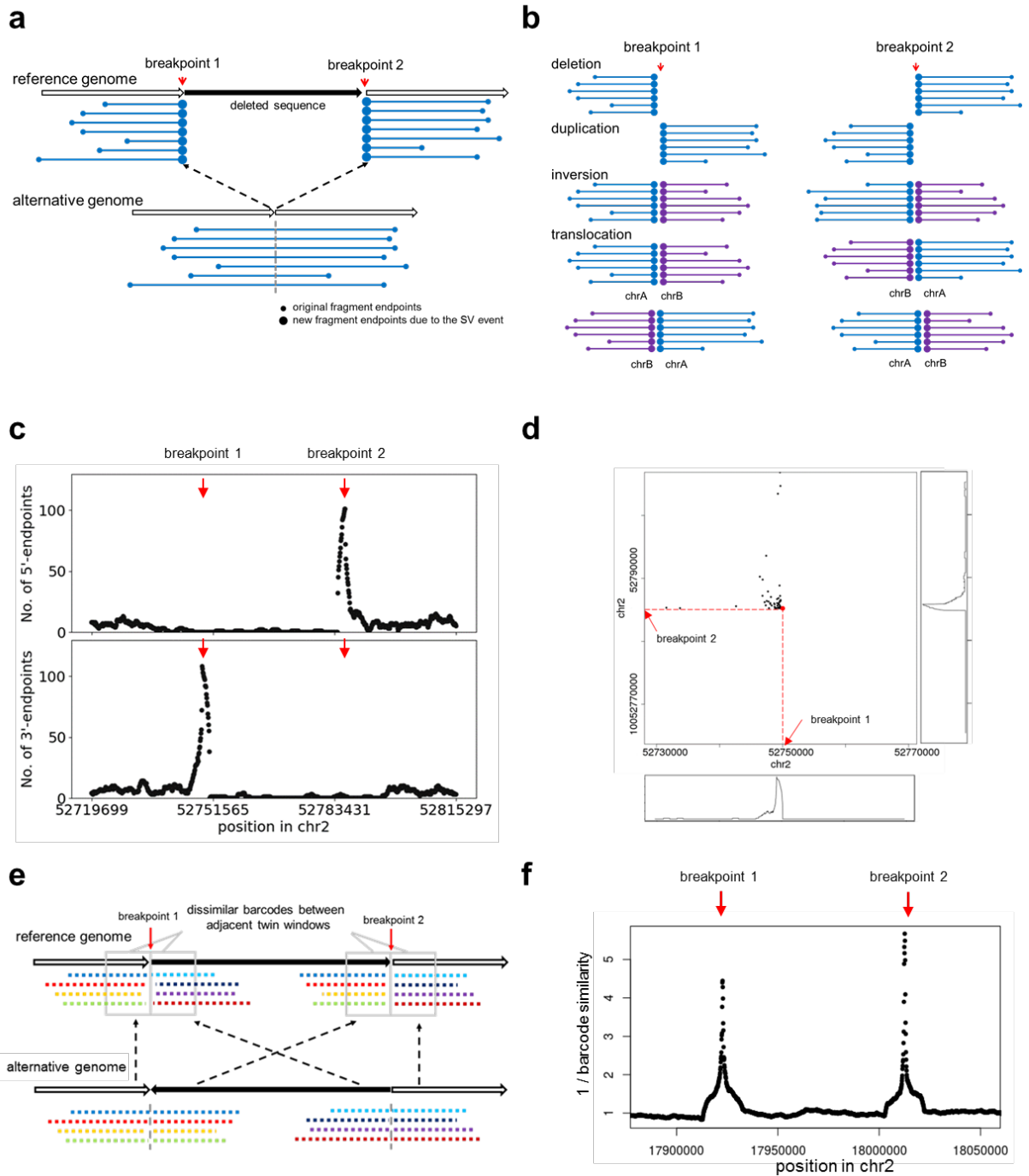
## References

1. Weischenfeldt J, Symmons O, Spitz F, Korbel JO. Phenotypic impact of genomic structural variation: insights from and for human disease. *Nat Rev Genet* 2013, 14(2):125-138.
2. Ye K, Schulz MH, Long Q, Apweiler R, Ning Z. Pindel: a pattern growth approach to detect break points of large deletions and medium sized insertions from paired-end short reads. *Bioinformatics* 2009, 25(21):2865-2871.
3. Rausch T, Zichner T, Schlattl A, Stutz AM, Benes V, Korbel JO. DELLY: structural variant discovery by integrated paired-end and split-read analysis. *Bioinformatics* 2012, 28(18):i333-i339.
4. Chen K, Wallis JW, McLellan MD, Larson DE, Kalicki JM, Pohl CS, McGrath SD, Wendl MC, Zhang Q, Locke DP, Shi X, Fulton RS, Ley TJ, Wilson RK, Ding L, Mardis ER. BreakDancer: an algorithm for high-resolution mapping of genomic structural variation. *Nat Methods* 2009, 6(9):677-681.
5. Chong Z, Ruan J, Gao M, Zhou W, Chen T, Fan X, Ding L, Lee AY, Boutros P, Chen J, Chen K. novoBreak: local assembly for breakpoint detection in cancer genomes. *Nat Methods* 2017, 14(1):65-67.
6. Wala JA, Bandopadhyay P, Greenwald NF, O'Rourke R, Sharpe T, Stewart C, Schumacher S, Li Y, Weischenfeldt J, Yao X, Nusbaum C, Campbell P, Getz G, Meyerson M, Zhang CZ, Imielinski M, Beroukhi R. SvABA: genome-wide detection of structural variants and indels by local assembly. *Genome Res* 2018, 28(4):581-591.
7. Carvalho CM, Lupski JR. Mechanisms underlying structural variant formation in genomic disorders. *Nat Rev Genet* 2016, 17(4):224-238.
8. Payer LM, Steranka JP, Yang WR, Kryatova M, Medabalimi S, Ardeljan D, Liu C, Boeke JD, Avramopoulos D, Burns KH. Structural variants caused by Alu insertions are associated with risks for many human diseases. *Proc Natl Acad Sci U S A* 2017, 114(20):E3984-E3992.
9. Sharp AJ, Locke DP, McGrath SD, Cheng Z, Bailey JA, Vallente RU, Pertz LM, Clark RA, Schwartz S, Segreaves R, Oseroff VV, Albertson DG, Pinkel D, Eichler EE. Segmental duplications and copy-number variation in the human genome. *Am J Hum Genet* 2005, 77(1):78-88.
10. Chaisson MJ, Huddleston J, Dennis MY, Sudmant PH, Malig M, Hormozdiari F, Antonacci F, Surti U, Sandstrom R, Boitano M, Landolin JM, Stamatoyannopoulos JA, Hunkapiller MW, Korlach J, Eichler EE. Resolving the complexity of the human genome using single-molecule sequencing. *Nature* 2015, 517(7536):608-611.
11. Sedlazeck FJ, Rescheneder P, Smolka M, Fang H, Nattestad M, von Haeseler A, Schatz MC. Accurate detection of complex structural variations using single-molecule sequencing. *Nat Methods* 2018.

12. Zheng GX, Lau BT, Schnall-Levin M, Jarosz M, Bell JM, Hindson CM, Kyriazopoulou-Panagiotopoulou S, Masquelier DA, Merrill L, Terry JM, Mudivarti PA, Wyatt PW, Bharadwaj R, Makarewicz AJ, Li Y, Belgrader P, Price AD, Lowe AJ, Marks P, Vurens GM, Hardenbol P, Montesclaros L, Luo M, Greenfield L, Wong A, Birch DE, Short SW, Bjornson KP, Patel P, Hopmans ES, Wood C, Kaur S, Lockwood GK, Stafford D, Delaney JP, Wu I, Ordonez HS, Grimes SM, Greer S, Lee JY, Belhocine K, Giorda KM, Heaton WH, McDermott GP, Bent ZW, Meschi F, Kondov NO, Wilson R, Bernate JA, Gauby S, Kindwall A, Bermejo C, Fehr AN, Chan A, Saxonov S, Ness KD, Hindson BJ, Ji HP. Haplotyping germline and cancer genomes with high-throughput linked-read sequencing. *Nat Biotechnol* 2016, 34(3):303-311.
13. Luo R, Sedlazeck FJ, Darby CA, Kelly SM, Schatz MC. LRSim: A Linked-Reads Simulator Generating Insights for Better Genome Partitioning. *Comput Struct Biotechnol J* 2017, 15:478-484.
14. Spies N, Weng Z, Bishara A, McDaniel J, Catoe D, Zook JM, Salit M, West RB, Batzoglou S, Sidow A. Genome-wide reconstruction of complex structural variants using read clouds. *Nat Methods* 2017.
15. Elyanow R, Wu HT, Raphael BJ. Identifying structural variants using linked-read sequencing data. *Bioinformatics* 2017.
16. Bishara A, Liu Y, Weng Z, Kashef-Haghighi D, Newburger DE, West R, Sidow A, Batzoglou S. Read clouds uncover variation in complex regions of the human genome. *Genome Res* 2015, 25(10):1570-1580.
17. Campbell PJ, Stephens PJ, Pleasance ED, O'Meara S, Li H, Santarius T, Stebbings LA, Leroy C, Edkins S, Hardy C, Teague JW, Menzies A, Goodhead I, Turner DJ, Clee CM, Quail MA, Cox A, Brown C, Durbin R, Hurles ME, Edwards PA, Bignell GR, Stratton MR, Futreal PA. Identification of somatically acquired rearrangements in cancer using genome-wide massively parallel paired-end sequencing. *Nat Genet* 2008, 40(6):722-729.
18. Mitelman F, Johansson B, Mertens F. The impact of translocations and gene fusions on cancer causation. *Nat Rev Cancer* 2007, 7(4):233-245.
19. Stephens PJ, McBride DJ, Lin ML, Varela I, Pleasance ED, Simpson JT, Stebbings LA, Leroy C, Edkins S, Mudie LJ, Greenman CD, Jia M, Latimer C, Teague JW, Lau KW, Burton J, Quail MA, Swerdlow H, Churcher C, Natrajan R, Sieuwerts AM, Martens JW, Silver DP, Langerod A, Russnes HE, Foekens JA, Reis-Filho JS, van 't Veer L, Richardson AL, Borresen-Dale AL, Campbell PJ, Futreal PA, Stratton MR. Complex landscapes of somatic rearrangement in human breast cancer genomes. *Nature* 2009, 462(7276):1005-1010.
20. Lakich D, Kazazian HH, Jr., Antonarakis SE, Gitschier J. Inversions disrupting the factor VIII gene are a common cause of severe haemophilia A. *Nat Genet* 1993, 5(3):236-241.
21. De Brasi CD, Bowen DJ. Molecular characteristics of the intron 22 homologs of the coagulation factor VIII gene: an update. *J Thromb Haemost* 2008, 6(10):1822-1824.
22. Layer RM, Chiang C, Quinlan AR, Hall IM. LUMPY: a probabilistic framework for structural variant discovery. *Genome Biol* 2014, 15(6):R84.

23. Rozowsky J, Abyzov A, Wang J, Alves P, Raha D, Harmanci A, Leng J, Bjornson R, Kong Y, Kitabayashi N, Bhardwaj N, Rubin M, Snyder M, Gerstein M. AlleleSeq: analysis of allele-specific expression and binding in a network framework. *Mol Syst Biol* 2011, 7:522.
24. Zook JM, Chapman B, Wang J, Mittelman D, Hofmann O, Hide W, Salit M. Integrating human sequence data sets provides a resource of benchmark SNP and indel genotype calls. *Nat Biotechnol* 2014, 32(3):246-251.
25. Startek M, Szafranski P, Gambin T, Campbell IM, Hixson P, Shaw CA, Stankiewicz P, Gambin A. Genome-wide analyses of LINE-LINE-mediated nonallelic homologous recombination. *Nucleic Acids Res* 2015, 43(4):2188-2198.

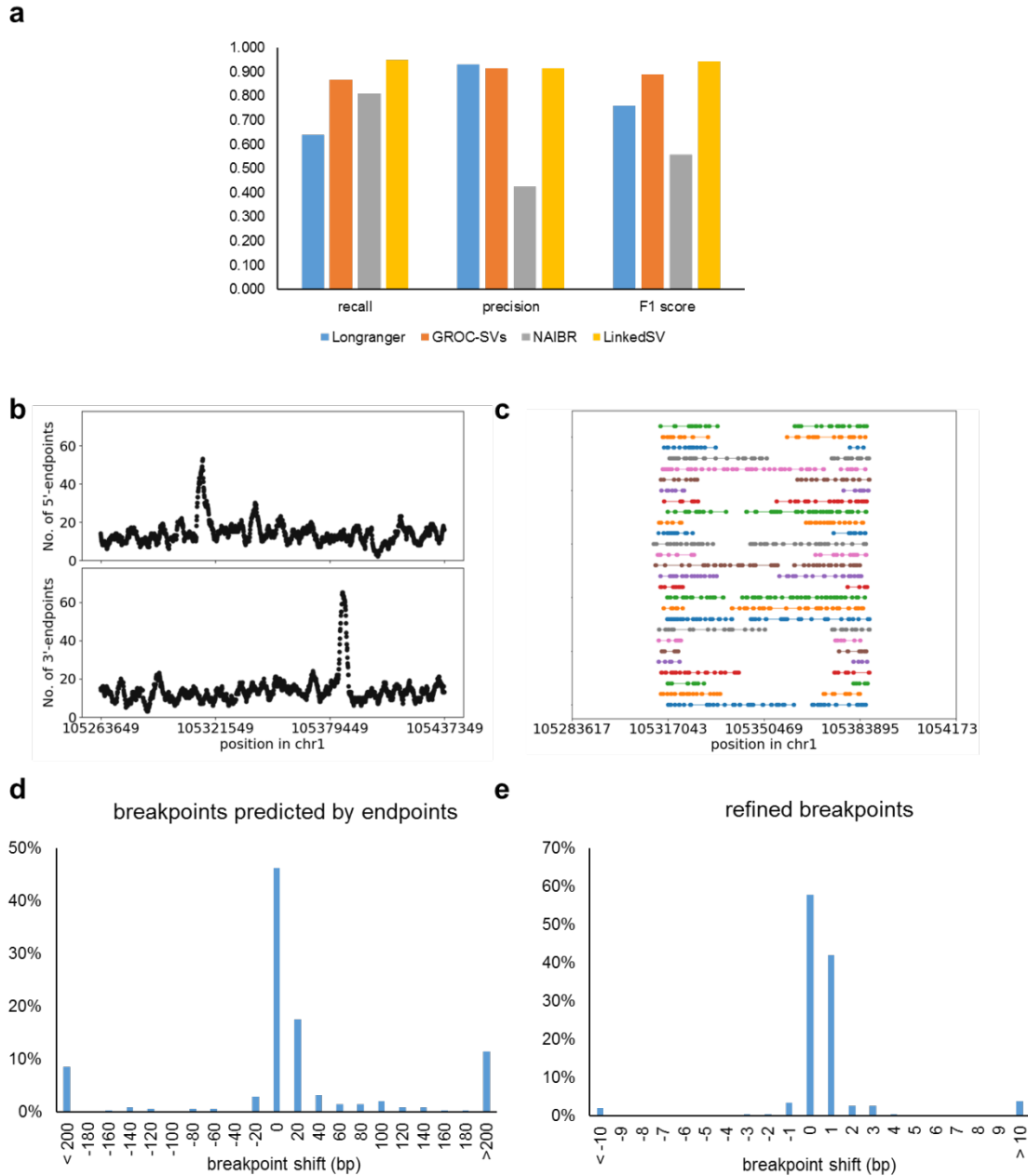
## Figures



**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

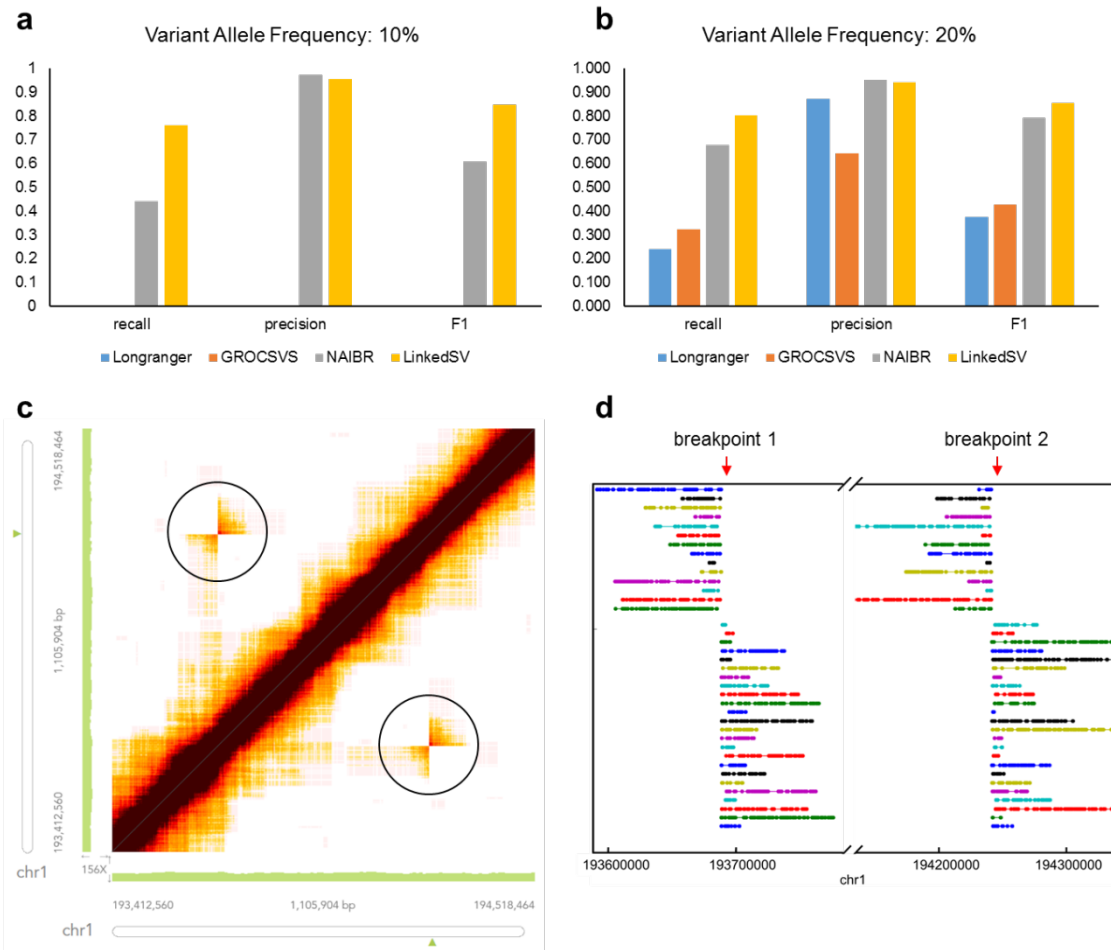
observed fragments and two sets of newly introduced fragment endpoints (solid dots). **b)** The patterns of enriched fragment endpoints indicate the SV types. **c)** Enriched fragment endpoints detected near two breakpoints of a deletion on NA12878 genome. 5'-endpoints and 3'-endpoints are plotted separately. 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 3'-endpoint and the y-value of the dot is the position of the second fragment's 5'-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**

**Performance of LinkedSV on the simulated WGS data set.** **a)** Recalls, precisions and F1 scores of four linked-read SV callers on the simulated WGS data set. **b)** Fragment endpoint signals of a small duplication that was missed by GROC-SVs. The peaks indicate the approximate breakpoint positions. **c)** Supporting fragments of the duplications. Horizontal lines represent linked reads with the same barcode; dots represent reads; colors indicate barcodes. **d)** Precision of breakpoints predicted by LinkedSV without checking short-read information. **e)** Precision of LinkedSV refined breakpoints using discordant read-pairs and split-reads.





**Figure 3**

**Performance of LinkedSV on the simulated WGS data with low variant allele frequencies. a, b** Recalls, precisions and F1 scores of four linked-read SV callers on the simulated WGS data set with VAF of 10% and 20%. **c** An inversion that was missed by Longranger, and NAIBR (VAF = 10%). The overlapped barcodes between the two inversion breakpoints can be clearly visualized (in the black circle) using the Loupe software developed by 10X Genomics. **d** Supporting fragments of the inversion detected by LinkedSV. Horizontal lines represent linked reads with the same barcode; dots represent reads; colors indicate barcodes. Predicted breakpoint positions are marked by red arrows.

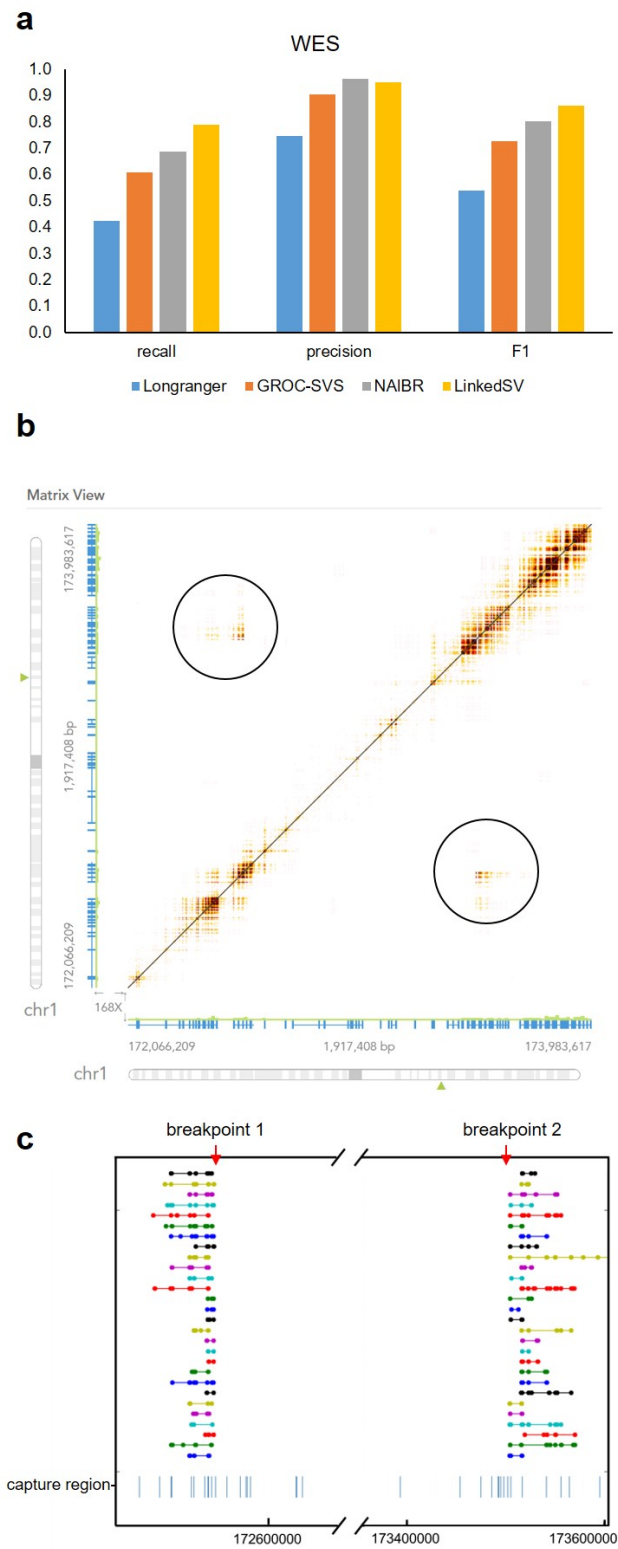
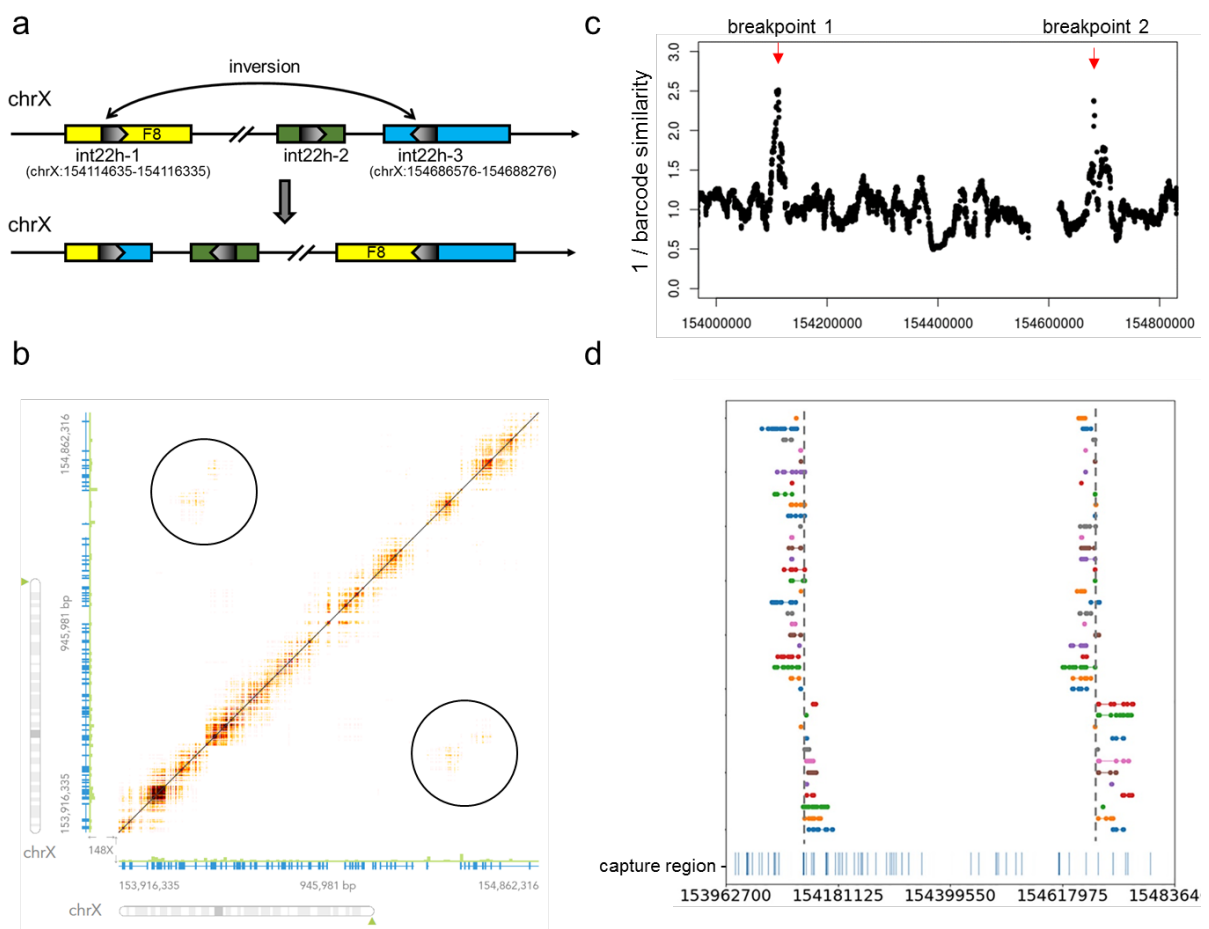


Figure 4

**Performance of LinkedSV on the simulated WES data set.** a) Recalls, precisions and F1 scores of four linked-read SV callers on the simulated WES data set. b) A deletion that was missed by NAIBR. The overlapped barcodes between the two breakpoints can be clearly visualized (in the black circle) using the Loupe software. c) Supporting fragments of the deletion detected by LinkedSV. Horizontal lines represent linked reads with the same barcode; dots represent reads; colors indicate barcodes. Predicted breakpoint positions are marked by red arrows.

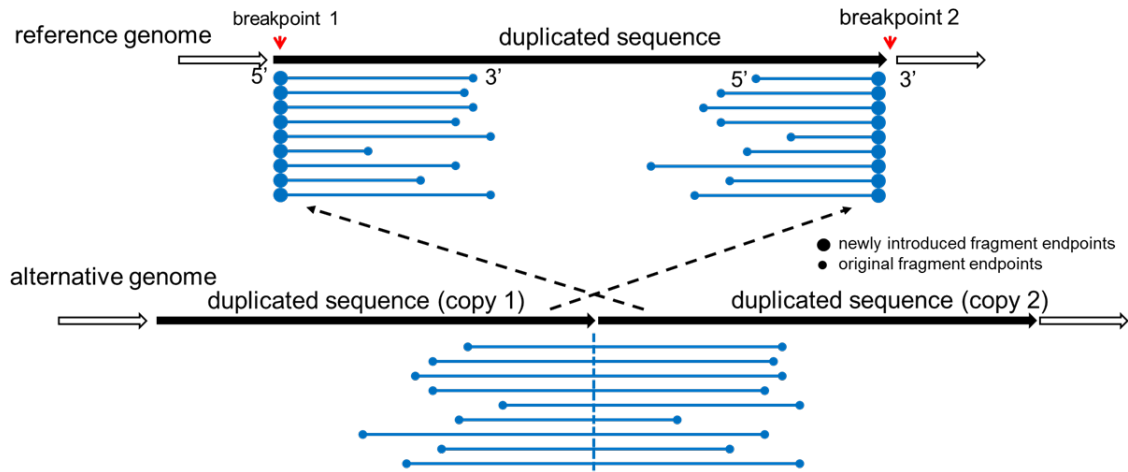


**Figure 5**

**Detection SVs from clinical exome sequencing data.** (a) Illustration of type I inversion of F8 gene. A portion of intron 22 is has three copies in chrX (int22h-1, int22h-2, int22h-3). The inversion is induced by the homologous recombination between two inverted copies int22h-1 and int22h-3. Int22h-1 is located in intron 22 of F8 gene and int22h-3 is located in the intergenic regions. (b) Barcode overlapping between

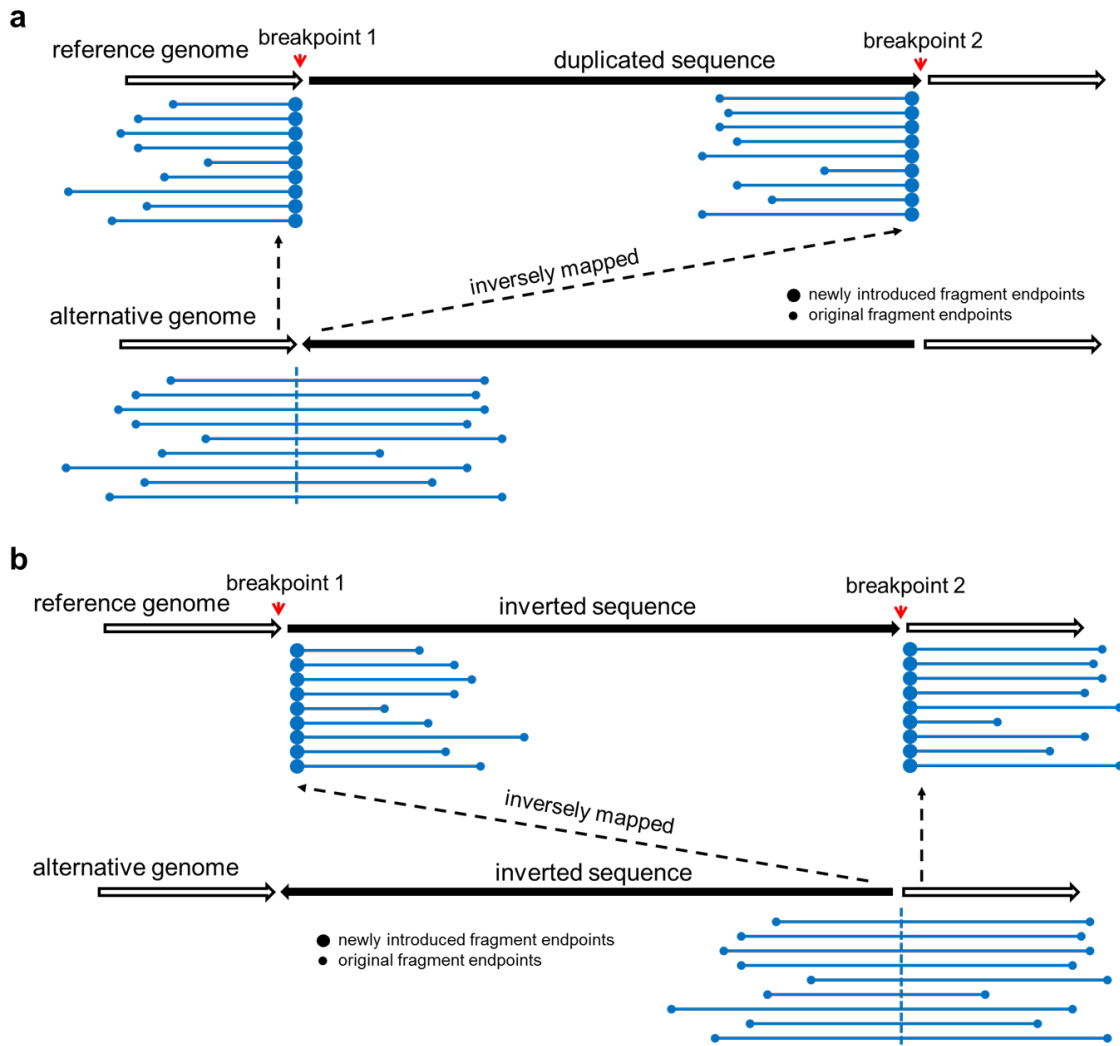
the two breakpoints can be visualized using the Loupe software (black circle). (c) Decreased barcode similarity at breakpoints detected by the twin window method of LinkedSV. Window size = 40 kb (d) Supporting fragments detected by LinkedSV. Horizontal lines represent linked reads with the same barcode; dots represent reads; colors indicate barcodes. Dashed vertical black lines represent breakpoints. Capture regions were shown as vertical bars in the bottom.

## Supplementary Figures



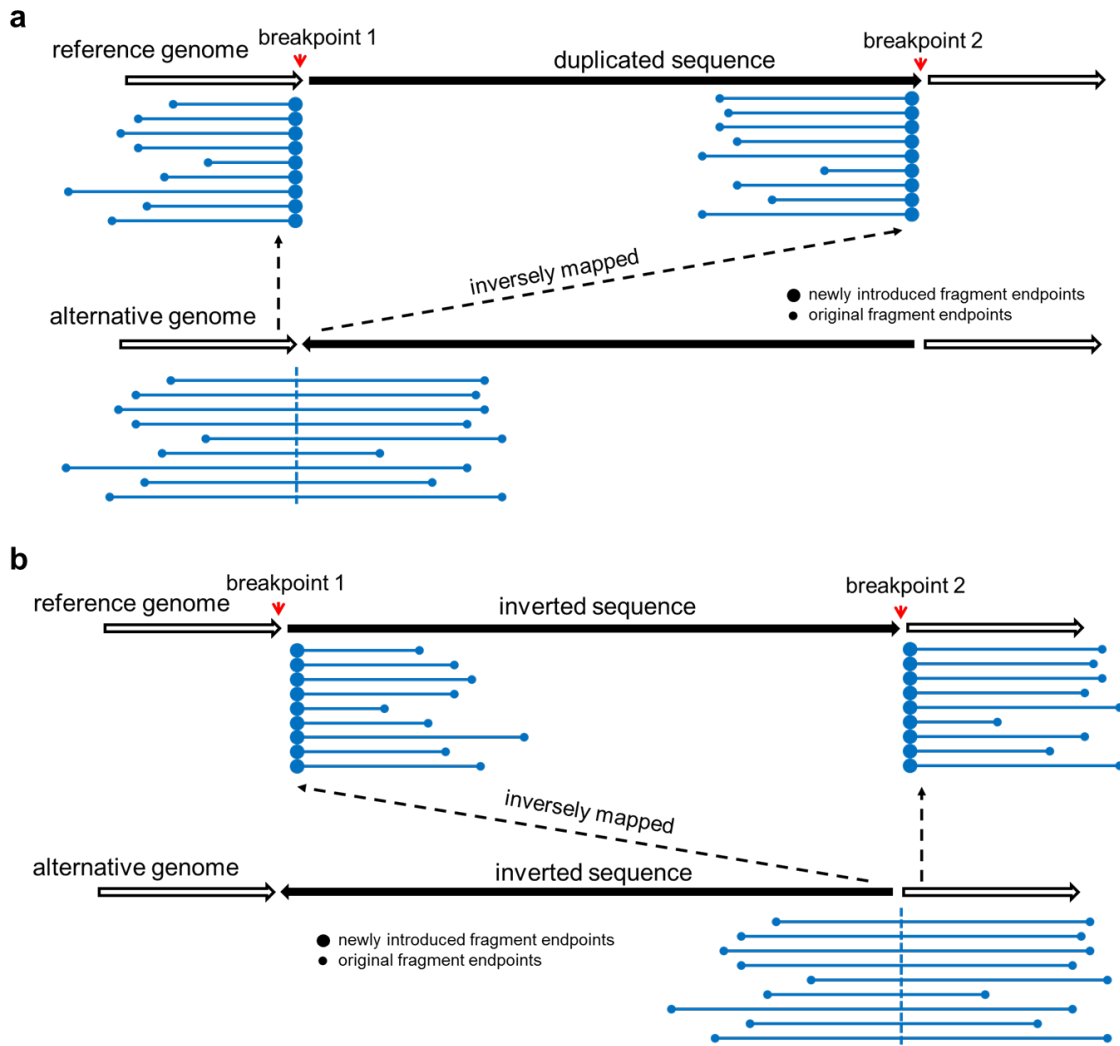
### Supplemental Figure 1

Pattern of enriched fragment endpoints for tandem duplication.



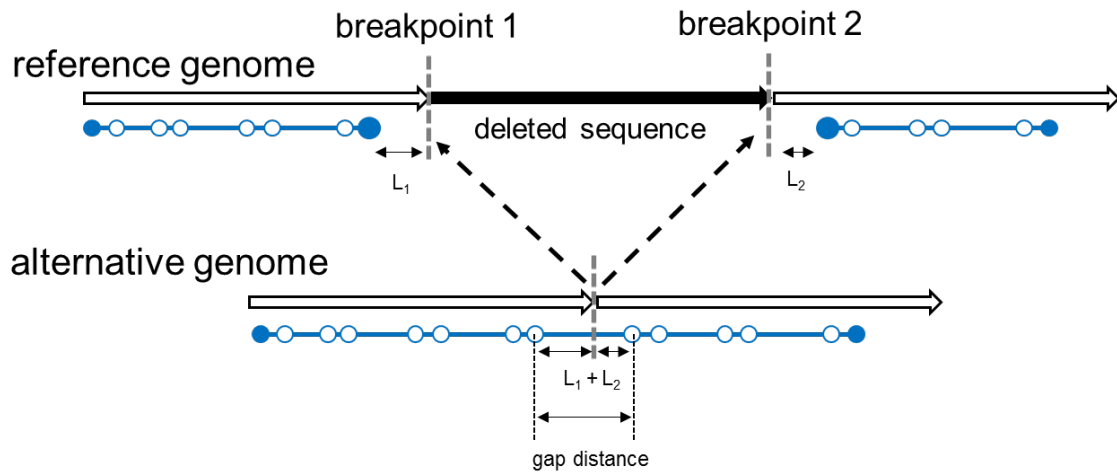
## Supplemental Figure 2

Pattern of enriched fragment endpoints for inversion.



### Supplemental Figure 3

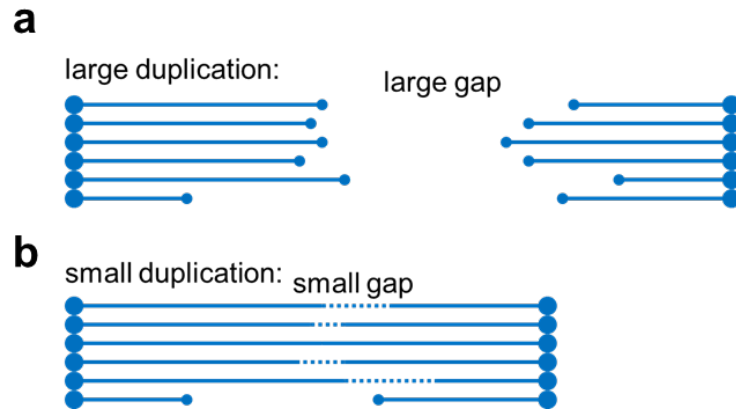
Pattern of enriched fragment endpoints for balanced translocation.



#### Supplemental Figure 4

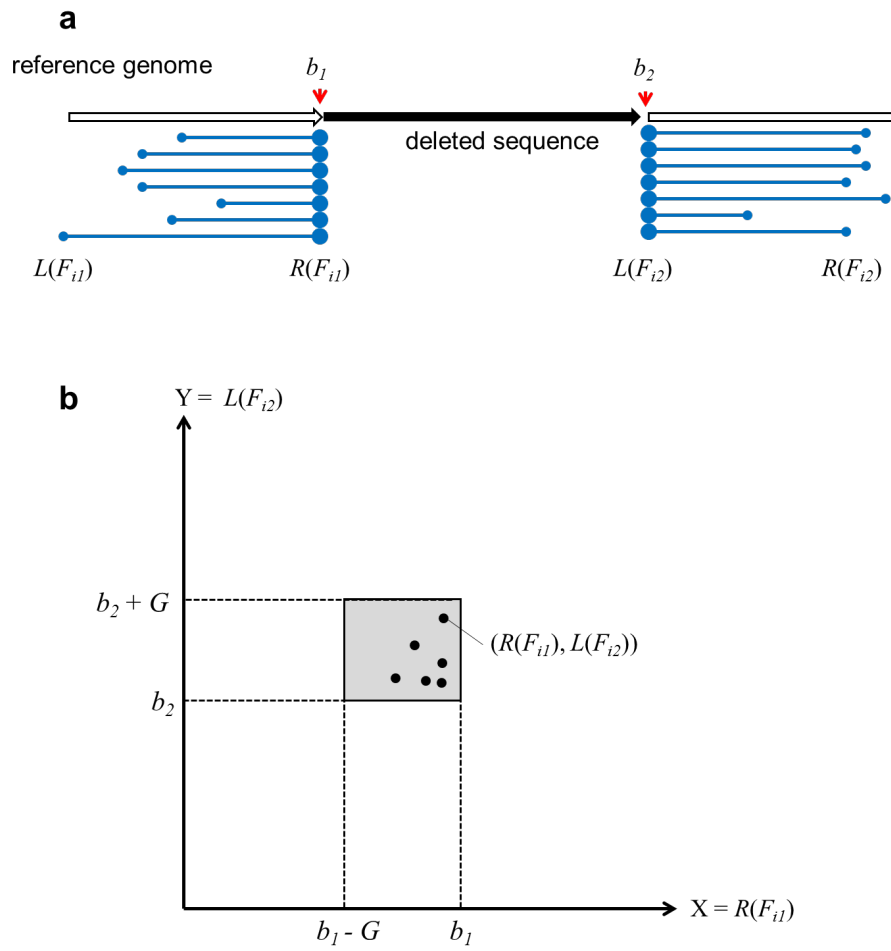
The distance between the fragment endpoint and the corresponding breakpoint is not more than the distance between two nearby reads (i.e. gap distance). This explains the enrichment of fragment endpoints near the breakpoints.





### Supplemental Figure 5

**a)** For large duplications, the reads of the alternative allele are separated by a large gap so that we can observe two sets of fragments with the same set of barcodes, which indicate a SV. **b)** If the duplication is not large enough, the reads will be probably clustered into one fragment.



### Supplemental Figure 6

Detection of type 1 evidence.

## Supplementary Tables

### Supplemental Table 1

SV calls detected by Delly on the F8 inversion sample (10 kb upstream/downstream of the inversion).

Chrom	Start	End	SV type
chrX	154002063	154002446	INV
chrX	154004508	154004709	INV
chrX	154005670	154005855	INV
chrX	154011709	154012062	DUP
chrX	154012155	154012297	INV
chrX	154014432	154014542	INV
chrX	154014535	154014759	INV
chrX	154017795	154017948	INV
chrX	154020109	154020407	INV
chrX	154045287	154045525	INV
chrX	154045633	154045781	INV
chrX	154057649	154057833	INV
chrX	154057785	154057923	INV
chrX	154083952	154084185	INV
chrX	154091092	154091440	INV
chrX	154115196	154115584	INV
chrX	154115667	154115957	INV
chrX	154124101	154124495	DUP
chrX	154129769	154129919	INV
chrX	154130086	154130222	INV
chrX	154132781	154133165	DUP
chrX	154146261	154146462	INV
chrX	154146565	154146775	INV
chrX	154152621	154152772	INV
chrX	154157646	154157941	DUP
chrX	154158133	154158450	DUP
chrX	154158358	154158630	INV
chrX	154159826	154160143	INV
chrX	154175838	154176055	INV
chrX	154176073	154176451	INV
chrX	154194502	154194721	INV
chrX	154208905	154209847	DUP
chrX	154209557	154209679	INV
chrX	154212475	154212843	INV
chrX	154212596	154212733	INV
chrX	154227431	154227667	INV
chrX	154272329	154272476	INV
chrX	154272335	154272592	DUP

chrX	154275546	154275664	INV
chrX	154275671	154275934	INV
chrX	154290114	154290323	INV
chrX	154344170	154344430	INV
chrX	154387813	154429972	INV
chrX	154387928	154429779	INV
chrX	154391126	154391291	INV
chrX	154426121	154426264	INV
chrX	154428351	154428890	DUP
chrX	154464614	154464716	INV
chrX	154508570	154508865	INV
chrX	154540689	154540872	INV
chrX	154609922	154610441	DUP
chrX	154610132	154610342	INV
chrX	154611498	154611783	INV
chrX	154611510	154688681	DEL
chrX	154612459	154612925	DUP
chrX	154613017	154613356	INV
chrX	154652935	154653319	DUP
chrX	154687101	154687535	DUP
chrX	154687312	154687480	INV
chrX	154688521	154688851	INV
chrX	154688714	154688969	INV
chrX	154689667	154689984	INV
chrX	154735308	154735475	INV
chrX	154755098	154755225	INV

---

## Supplemental Table 2

SV calls detected by Lumpy on the F8 inversion sample (all calls in chrX).

Chrom	Start	End	SV type
chrX	1413852	1414637	DEL
chrX	17520630	17520931	DEL
chrX	28749622	28749739	DUP
chrX	31189715	31189837	DUP
chrX	38145097	38145387	DUP
chrX	38145261	38145866	DUP
chrX	38145731	38145851	DEL
chrX	38145513	38145876	DEL
chrX	38145369	38145943	DUP
chrX	38145283	38145943	DUP
chrX	38145180	38145975	DUP
chrX	38145091	38145976	DUP
chrX	38145835	38145920	DEL
chrX	38145729	38145981	DUP
chrX	38145874	38145914	DEL
chrX	38145321	38145988	DUP
chrX	38145471	38145995	DUP
chrX	38145929	38146016	DUP
chrX	38145971	38146003	DEL
chrX	38145885	38146041	DUP
chrX	38145642	38146055	DUP
chrX	38146077	38146257	DUP
chrX	39932139	39932332	DUP
chrX	41788764	41788966	DUP
chrX	47308661	47308783	DUP
chrX	63411862	63412017	DUP
chrX	64063531	64063690	DUP
chrX	76890054	76890207	DUP
chrX	101912168	101912314	DUP
chrX	142794887	142795158	DUP
chrX	149745015	149745169	DUP

### Supplemental Table 3

SVs on chrX detected by Longranger on the F8 inversion sample.

<b>Chrom 1</b>	<b>Position 1</b>	<b>Chrom 2</b>	<b>Position 2</b>	<b>SV type</b>	<b>Distance to int22h-1</b>	<b>Distance to int22h-3</b>
chrX	3735199	chrX	3855112	Unknown	-150380286	-150832314
chrX	7810783	chrX	8151613	Unknown	-146304702	-146535813
chrX	7810783	chrX	8113116	Unknown	-146304702	-146574310
chrX	134855892	chrX	134985727	Unknown	-19259593	-19701699
chrX	152225317	chrX	152352113	Unknown	-1890168	-2335313
chrX	154091033	chrX	154660067	Unknown	-24452	-27359
chrX	154131339	chrX	154735755	Unknown	15854	48329

### Supplemental Table 4

SVs detected by GROC-SVs on the F8 inversion sample.

Chrom 1	Position 1	Chrom 2	Position 2	Orientation
chr1	12889258	chr1	12938883	+-
chr1	12919205	chr1	12853022	--
chr1	21754487	chr1	21794667	+-
chr1	148026036	chr1	144622101	++
chr1	223725671	chr1	223797843	+-
chr2	98162357	chr2	97860318	--
chr2	149687145	chr2	149790581	+-
chr2	234053740	chr2	234002965	-+
chr3	129809651	chr3	129762840	-+
chr4	9452594	chr4	9485059	+-
chr4	69681327	chr4	69893173	--
chr5	155137378	chr5	155188725	+-
chr5	180430554	chr5	180375038	-+
chr6	29909733	chr6	29843849	--
chr6	29913575	chr6	29844437	++
chr6	160956444	chr6	160877743	--
chr7	100550750	chr7	100609409	--
chr7	100555813	chr7	100610572	++
chr11	1162747	chr11	1212758	+-
chr11	5809264	chr11	5777102	-+
chr12	7239875	chr12	7189849	-+
chr12	11545335	chr12	11503243	-+
chr12	18018173	chr12	17922871	--
chr12	109372790	chr12	109423610	+-
chr12	132926655	chr16	86452753	+-
chr12	133041064	chr2	231869678	+-
chr13	114325993	chr13	114425990	+-
chr14	24436900	chr14	24474868	+-
chr15	20740860	chr15	23406226	+-
chr15	22743596	chr15	23572115	+-
chr15	23407964	chr15	20739466	+-
chr15	23573198	chr15	22742499	+-
chr15	28597075	chr15	28806596	++
chr15	28804860	chr15	28595402	--
chr15	83003956	chr15	82934463	--
chr15	83014625	chr15	82936703	++
chr15	84960511	chr15	84859997	++
chr16	14988609	chr16	15031370	--
chr16	70009791	chr16	74426101	-+

chr16	86453332	chr12	132926283	+-
chr17	36297237	chr17	36337601	+-
chr20	1600172	chr20	1559471	-+
chr20	56771729	chr12	2828954	--
chr20	56771963	chr12	2830313	++
chr22	18666166	chr22	18737108	--
chr22	18737933	chr22	18686213	++

---



### Supplemental Table 5

SVs on chrX detected by NAIBR on the F8 inversion sample.

<b>Chrom 1</b>	<b>Position 1</b>	<b>Chrom 2</b>	<b>Position 2</b>	<b>Orientation</b>	<b>Distance to int22h-1</b>	<b>Distance to int22h-3</b>
chrX	1399792	chrX	1400312	+-	-152715693	-153287114
chrX	26179787	chrX	26212951	++	-127935698	-128474475
chrX	49162006	chrX	49180528	++	-104953479	-105506898
chrX	49208629	chrX	49209218	+-	-104906856	-105478208
chrX	49218203	chrX	49218751	+-	-104897282	-105468675
chrX	52830446	chrX	52830826	+-	-101285039	-101856600
chrX	57147470	chrX	57162963	++	-96968015	-97524463
chrX	129651865	chrX	129652250	+-	-24463620	-25035176
chrX	140140043	chrX	140140427	+-	-13975442	-14546999
chrX	154387911	chrX	154430070	--	272426	-257356
chrX	155245001	chrX	155245507	+-	1129516	558081
chrX	155251108	chrX	155252593	+-	1135623	565167

## Supplementary Notes

### Supplementary Note 1 Explanation of the model for detecting type 2 evidence.

Barcode similarity between two nearby regions very high because the reads originate from almost the same set of HMW DNA molecules. However, the barcode similarity between the left side and right side of the breakpoint are dramatically reduced. We call this evidence as type 2 evidence. To detect type 2 evidence, LinkedSV uses two adjacent sliding windows (window 1 and window 2) to scan the genome and calculate the barcode similarity between the window 1 and window 2.

In WES data sets, the numbers of reads in the sliding windows vary a lot due to capture bias and the length of capture regions. To detect type 2 evidence from both WGS and WES data sets, our model considers the variation of sequencing depth and capture regions. The barcode similarity is calculated as:

$$S = \frac{x}{m_1^a m_2^b} n e^{-\alpha d} \quad (1),$$

where:

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

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

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

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

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

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

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

$n$ ,  $\alpha$ ,  $a$  and  $b$  are estimated from the data using regression.

Suppose there are  $n$  different HMW DNA molecules span both window 1 and window 2, each of which has a different barcode and generates a number of read pairs in the library. The read pairs in the library may or may not be sequenced. We assume the  $n$  HMW DNA molecules have the same rate to generate read pairs in the library, so the  $n$  HMW DNA molecules have the same chance to be sequenced (have at least 1 read). Let  $m_1$  be the number of HMW DNA molecules sequenced in window 1,  $m_2$  be the number of HMW DNA molecules sequenced in window 2.  $m_1$  and  $m_2$  can be different due to the bias of target enrichment and the total length of target regions in each window. Let  $X$  be the number of HMW DNA molecules sequenced in both window 1 and window 2.  $X$  follows the hypergeometric distribution:

$$P(X = x | m_1, m_2, n) = \frac{C_{m_1}^x C_{n-m_1}^{m_2-x}}{C_n^{m_2}} \quad (2).$$

The expectation of  $X$  is:

$$E(X) = \frac{m_1 m_2}{n} \quad (3).$$

However, the length of sliding windows may be as long as 40 kb and not all the  $n$  HMW DNA molecules are long enough to span both windows. In addition, the capture regions in window 1 and window 2 may be close to each other or far away from each other. Therefore, we need to adjust  $n$  to be approximately  $ne^{-\alpha d}$ .  $d$  is calculated using the following equation:

$$d = w_2 - w_1 \quad (4),$$

where  $w_1$  is the mean mapping position of all reads in window 1 and  $w_2$  is the mean mapping position of all reads in window 2. The larger  $d$ , the smaller number of HMW DNA molecules can span a region of length  $d$ . We choose exponential distribution because the length of HMW DNA molecules follows exponential distribution and thus the number of HMW DNA molecules longer than  $d$  also follows exponential distribution.

$m_1$  also need to be adjusted because not all HMW DNA molecules being sequenced in window 1 span both windows. We adjust  $m_1$  to be approximately  $m_1^a$  and similarly adjust  $m_2$  to be approximately  $m_2^b$ .

After the adjustment, the expectation of  $X$  is:

$$E(X) = \frac{m_1^a m_2^b}{ne^{-\alpha d}} \quad (5).$$

We define barcode similarity as:

$$S = \frac{x}{E(X)} = \frac{x}{m_1^a m_2^b} ne^{-\alpha d} \quad (6),$$

where  $x$  is the number of shared barcodes between window 1 and 2,  $E(X)$  is the expected number of shared barcodes between window 1 and 2.

Take the log of both sides equation (6), we have:

$$\log(S) = \log(x) - a \log(m_1) - b \log(m_2) + \log(n) - \alpha d \quad (7).$$

Assuming most regions in the genome do not have breakpoints, we can replace  $S$  with 1 and estimate  $a$ ,  $b$ ,  $n$ ,  $\alpha$  from the data using linear regression.