

Additional file 10: Supplementary Information

**Structural variant analysis of a cancer reference cell line sample using
multiple sequencing technologies**

Keyur Talsania^{1,2#}, Tsai-wei Shen^{1,2#}, Xiongfong Chen^{1,2#}, Erich Jaeger^{3#}, Zhipan Li^{4#}, Zhong Chen⁵, Wangqiu Chen⁵, Bao Tran⁶, Rebecca Kusko⁷, Limin Wang⁸, Andy Wing Chun Pang⁹, Zhaowei Yang¹⁰, Sulbha Choudhari^{1,2}, Michael Colgan¹¹, Li Tai Fang¹², Andrew Carroll¹³, Jyoti Shetty⁶, Yuliya Kriga⁶, Oksana German⁶, Tatyana Smirnova⁶, Tiantain Liu⁵, Jing Li¹⁰, Ben Kellman⁹, Karl Hong⁹, Alex Hastie⁹, Aparna Natarajan³, Ali Moshrefi³, Anastasiya Granat³, Tiffany Truong³, Robin Bombardi³, Veronnica Mankinen¹⁴, Daoud Meerzaman¹⁵, Christopher E. Mason¹⁶, Jack Collins^{1,2}, Eric Stahlberg², Chunlin Xiao¹⁷, Charles Wang^{5*}, Wenming Xiao^{11*}, Yongmei Zhao^{1,2*}

#Contributed equally

*Correspondence should be addressed to: Yongmei.Zhao@nih.gov, Wenming.Xiao@fda.hhs.gov, and oxwang@gmail.com

Supplementary Figures

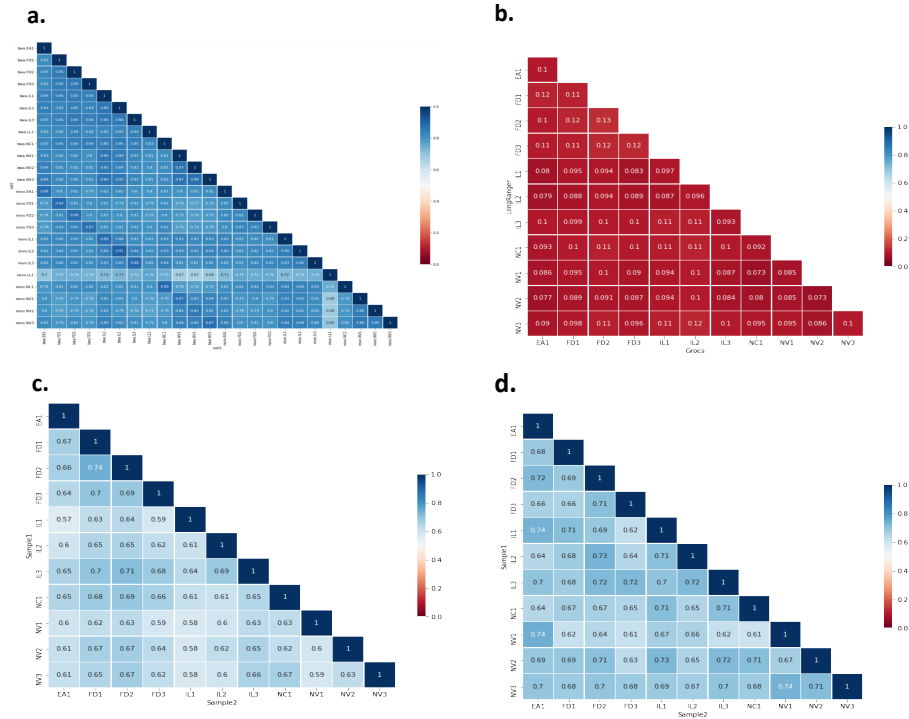


Fig. S1. Software Aligners and Caller Impact for SV Detection Sensitivity. (a) comparison of aligner effect. The matrix showed correlations between 24 pairs of replicates SV calling results which were generated by using BWA-mem or Novoalign for alignment and TNscope for SV calling. (b) comparison of variant caller effect between Longranger and Groc-SV. 11 pairs of 10x Chromium linked-read WGS replicates were compared. Observed low concordance calls between the two callers. (c) concordance among Longranger large SVs (>=30kbs) call results for 10x genomics data sets. The sequencing replicates were prepared by 5 different sequencing sites, same aligner (Lariat) and same caller (Long ranger) were used for data processing. (d) concordance among Groc-SVs results (large SVs includes DEL, INS, TRA, INV). The sequencing replicates were prepared by 5 sequencing sites, aligner (BWA-MEM) and caller (Groc-SVs) were used for mapping and variant calling.

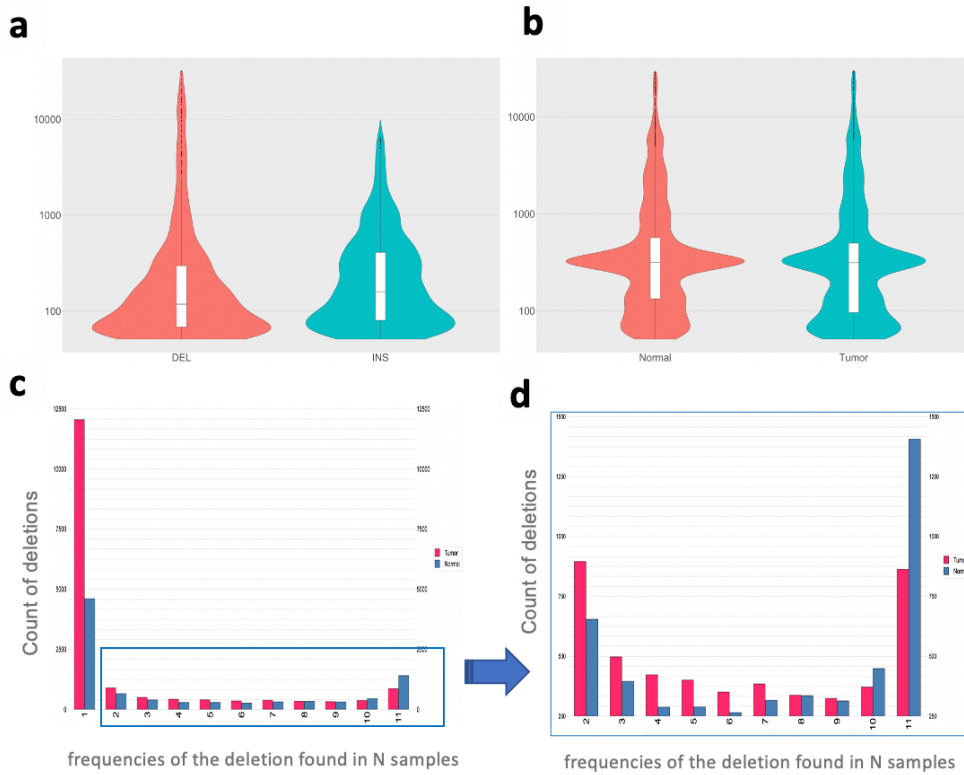


Fig. S2. 10x Genomics Short Deletion Reproducibility (a) PacBio somatic call of deletions and insertions (b) 10X Genomics – short deletions size between 50kb – 30kbs. The Y axis for both plots denote the deletion sizes. High number of short deletions less than 1kb observed for 10x Long Ranger calls.(c) All short deletions with size between 50kb - 30kb called by 10x Long Ranger software (d) Short deletions called by 10X Long Ranger in at least two samples. The X axis denotes the frequencies of the deletion found in 1 sample, 2 sample, N sample. Y axis denotes the count number of short deletions in a particular frequency. Sample specify private calls are the highest frequency calls observed for 10x Long Ranger short deletions.

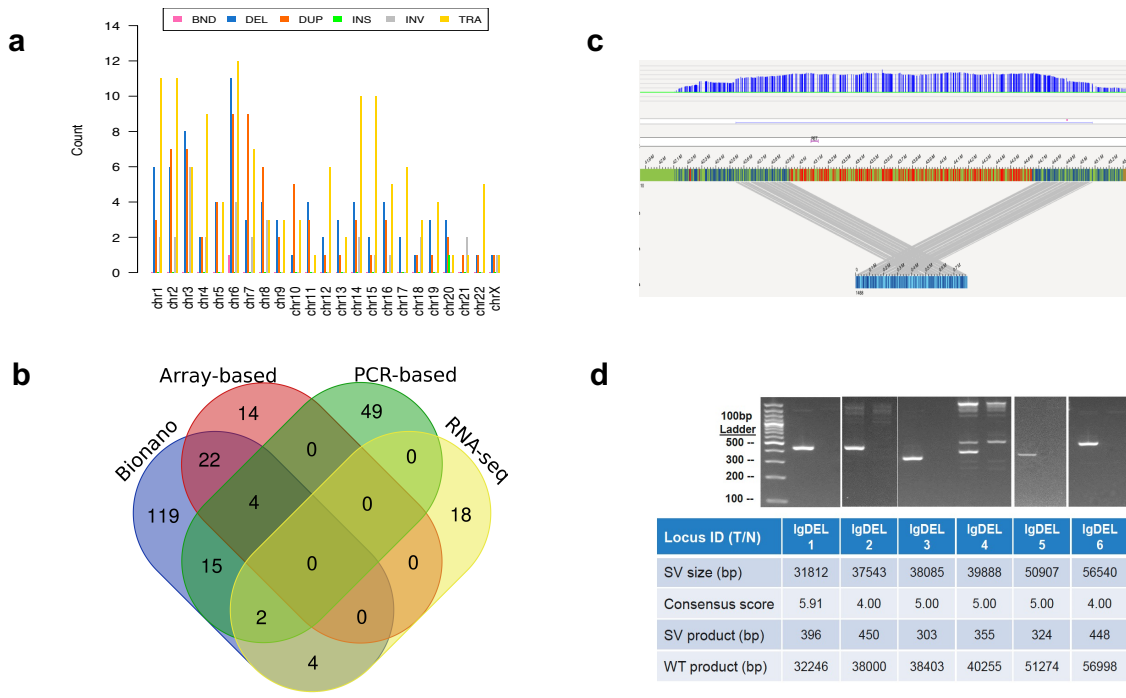


Fig. S3. Validation of the SV Consensus Call Set for HCC1395 Cancer Cell Line. (a) Total number of SVs were validated by different methods for each SV type (b) SVs validated by each platforms and multiple NGS platforms. (c) one 2.5 Mbp amplification impacting *RET* gene was detected by multiple NGS platform including PacBio/Illumina and 10x linked-read technologies at chr10: 42528,663 to chr10: 45,086,964 was confirmed by BioNano map on chromosome 10: 42,012,629-45,493,007, 3 - 6 copies; it is also confirmed by Affymetrix array at chr10:31813672 to chromosome 10: 45,092,199 with 3 copies gain. (d) PCR-based validation. The top showed the gel picture of the PCR product size bands and bottom showed the size of each SV and its corresponding PCR product size.

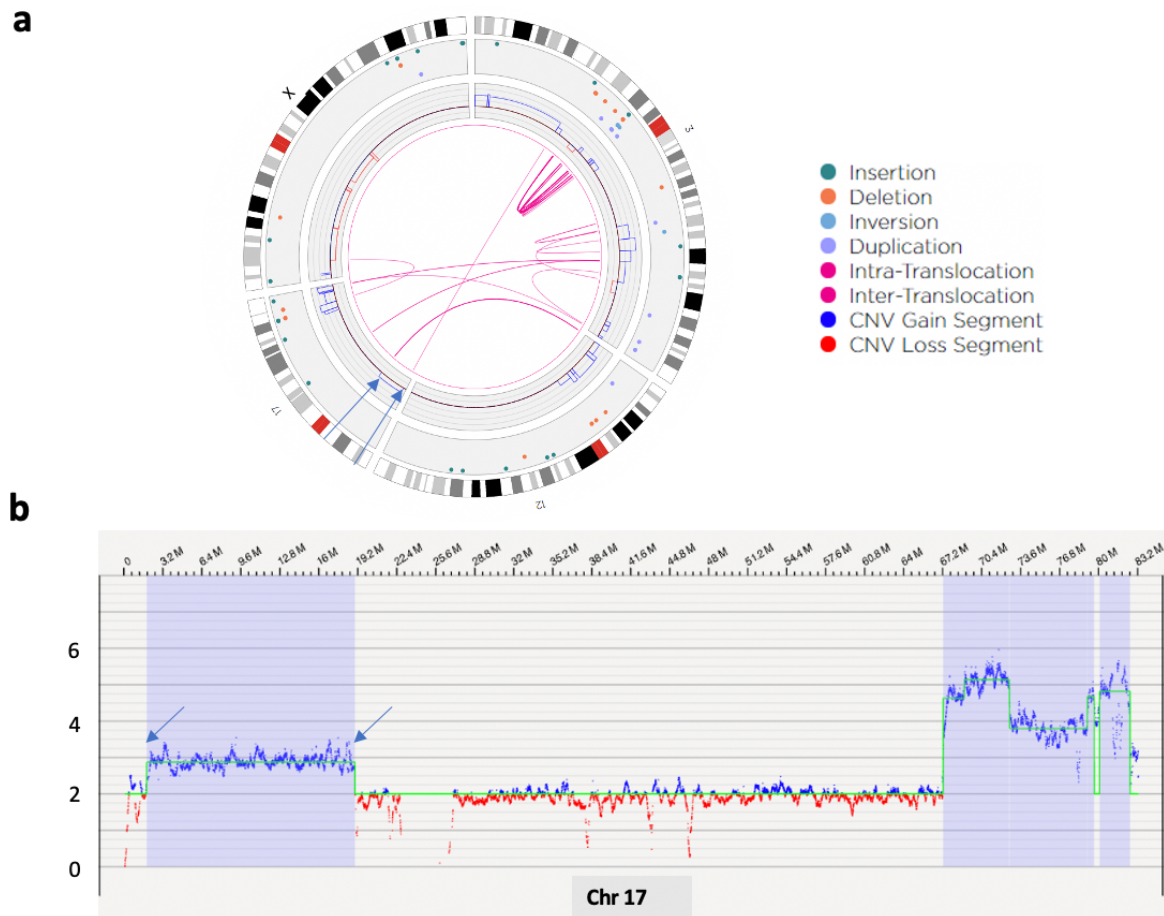


Fig. S4. Complex SV on Chromosome 17p of HCC1395 Cell Line. A 16.1 Mbp duplication on 17p (chr17:1,864,623-18,055,721) uniquely detected by Bionano genomics. This 3-copy duplication impacts multiple genes USP6, RABEP1, TP53, PER1, GAS7, MAP2K4, NCOR1 and FLCN. Furthermore, the duplication breakpoints coincide with t(3;17) and t(12;17). **(a)** The circos plot shows all the somatic variants on chr3, 12,17 and X. The breakpoints of the duplication of interest are indicated by the two blue arrows. **(b)** A view of the coverage profile of chromosome 17. The shaded regions are calls made by the CNV algorithm, and the one with two flanking blue arrows are the one 16.1 Mbp duplication.

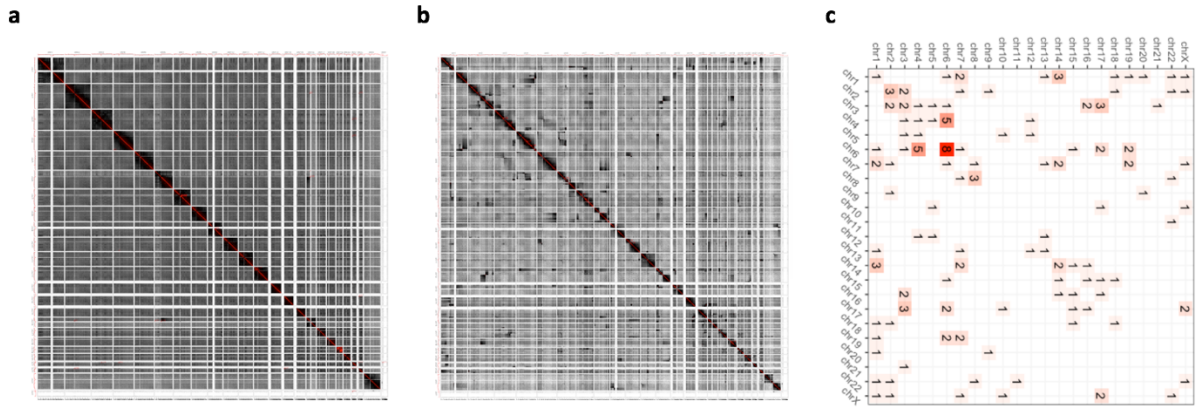


Fig. S5. Chromothripsis in the HCC1395 Tumor Cell line (a) Dovetail Selva generated Dot Plot for chromosomal view of HCC1395BL normal cell line. (b) Dovetail Selva generated Dot Plot for chromosomal view of HCC1395 tumor cell line. (c) Overview of the number of translocation events in HCC139 tumor cell line

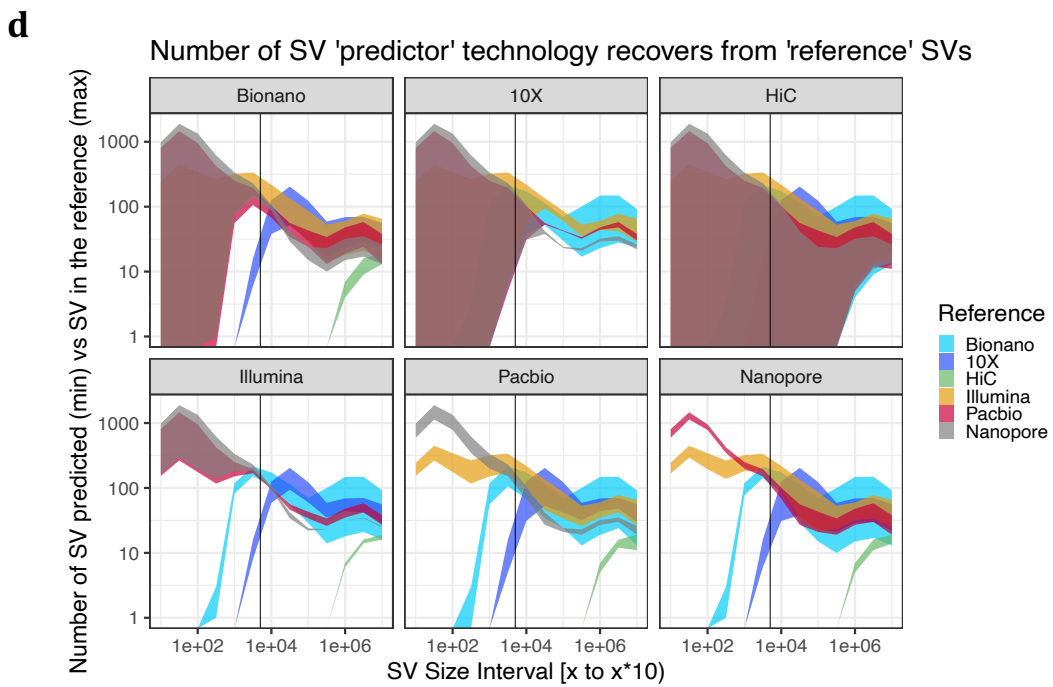
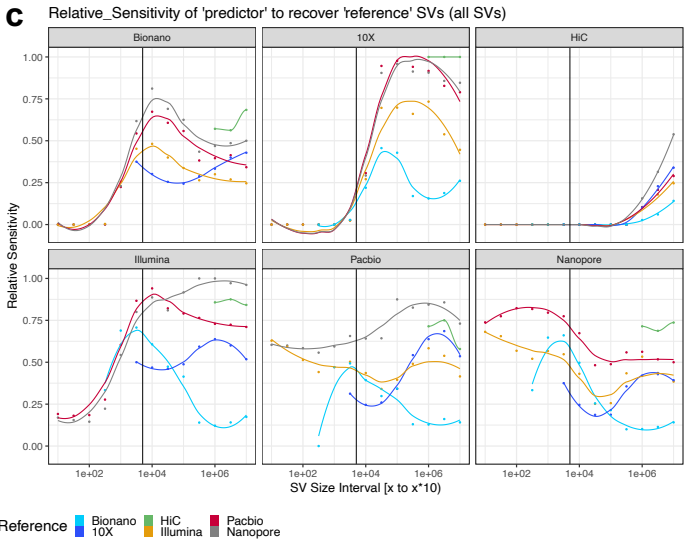
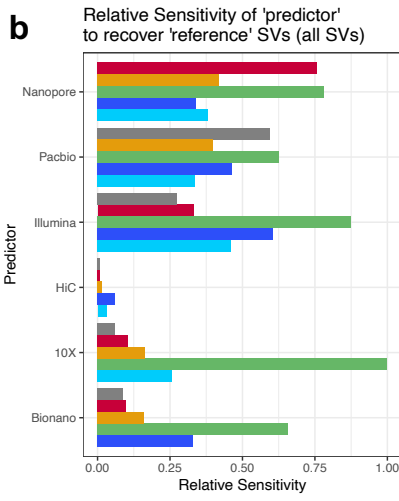
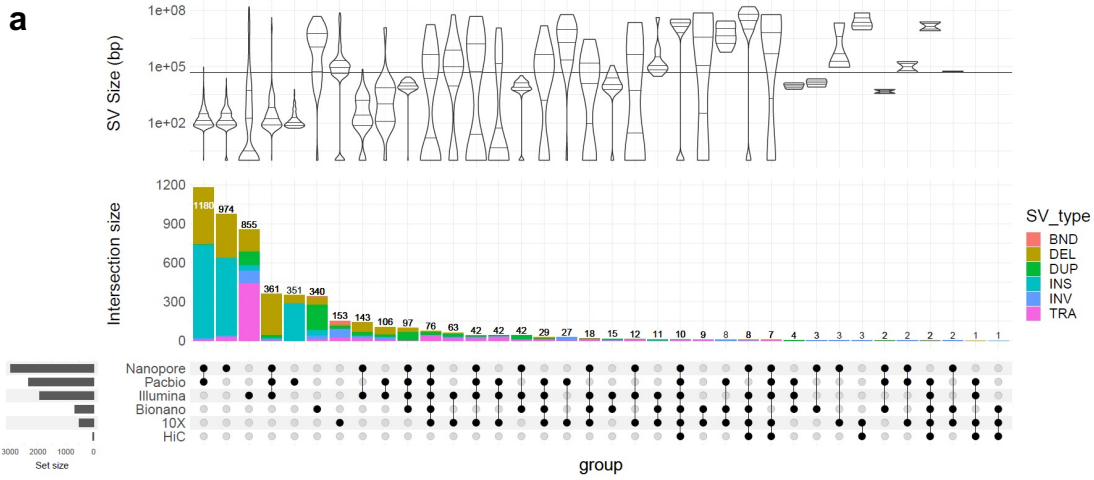


Fig. S6. Relative Sensitivity for SV Calls Across Platforms. (a) Comparison of SVs called by different platform including Nanopore, PacBio, Illumina short-read, Bionano, 10X Chromium linked-read, and HiC. Upset plot expanding on data similarities highlighted in **Figure 2a**. Intersection size is split into SV-type including multiple breakpoints SV events (BND), deletions (DEL), duplications (DUP), insertions (INS), inversions (INV), and inter-chromosomal translocations (TRA). Violin plots indicate the distribution of SV-sizes, in base-pairs (bp), within the vertically aligned with overlap-sets; horizontal lines within the violin plots indicate the 1st quantile, median, and 3rd quantile. A horizontal line is included in the violin plots to indicate 100kb size. (b) Relative sensitivity was calculated for all SV calls. (c) Relative sensitivity was calculated for all SV calls and stratified by size, between x and $x*10$ base pairs (x -axis) where x is the vertically aligned value on the x -axis. Relative sensitivity (x -axis in **b**, y -axis in **c**) is the proportion of all SV calls include SVs from initial call set defined in **Additional file 3: Table S3** and Bionano call set from **Additional file 15: Table S14** called by a reference platform (color) recoverable using another technology (y -axis in **b**, sub-panels in **c**) across different platforms. (b) Relative sensitivity of “predictor” to recover ‘reference’ SVs. Ribbon plot expansion on relative sensitivity to show the numerator (smallest shaded y -coordinate; number of reference SVs (color) found by the predictor platform (subpanel)) and the denominator (largest shaded y -coordinate; number of reference SVs (color)); the numerator at each size interval is the smallest vertically aligned y -coordinate within the shaded region while the denominator is the largest vertically aligned y -coordinate.

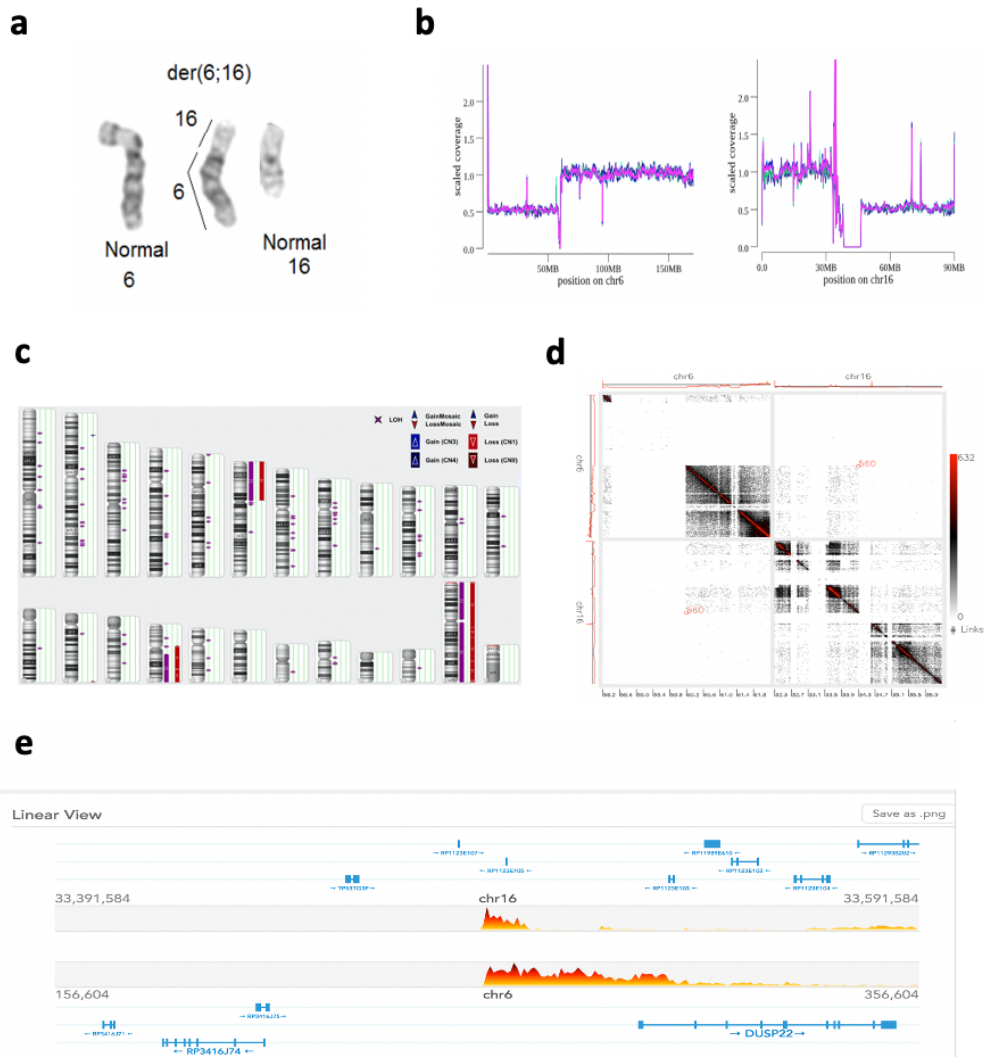


Fig. S7. Translocation and Complex Events in HCC1395 BL Cell Line. (a) Karyotyping showed fusion events of Chromosome 6 (chr6 p22) and Chromosome 16 (chr16 q21). (b) Global whole genome read-depth based analysis from all 3 WGS data sets (Illumina, PacBio and 10x Genomics) showed copy number change on Chr 6 p22 and Chr16 q21 (c) CytoScan® HD Assay confirmed HCC1395 translocation complex events occurred between Chr 6 p22 and Chr16 q21. (d) HiC detected a complex translocation event at chr6:60,243,000 and chr16:34,260,000. (e) 10x Genomics technology detected a 100kb translocation event between chr6:257,000 and chr16:32,280,000.

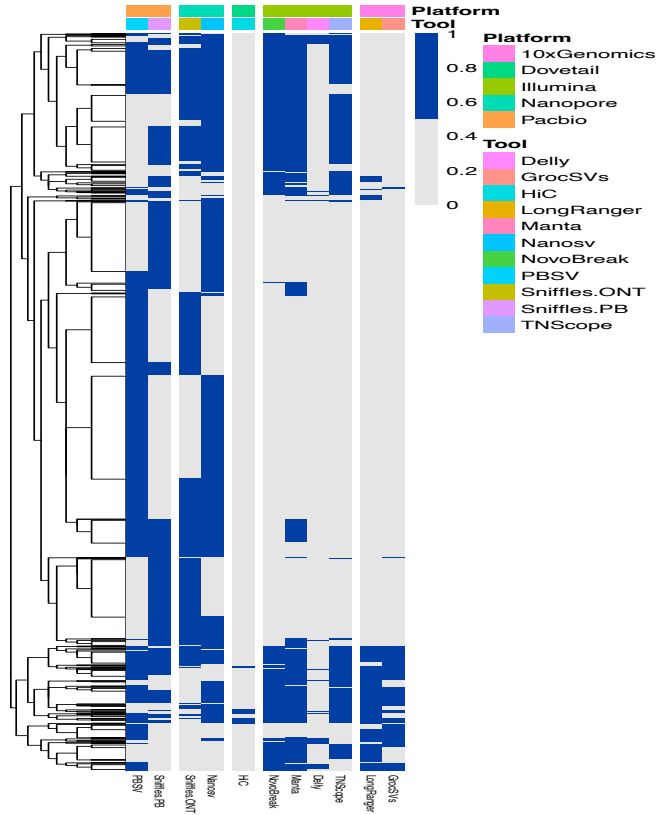
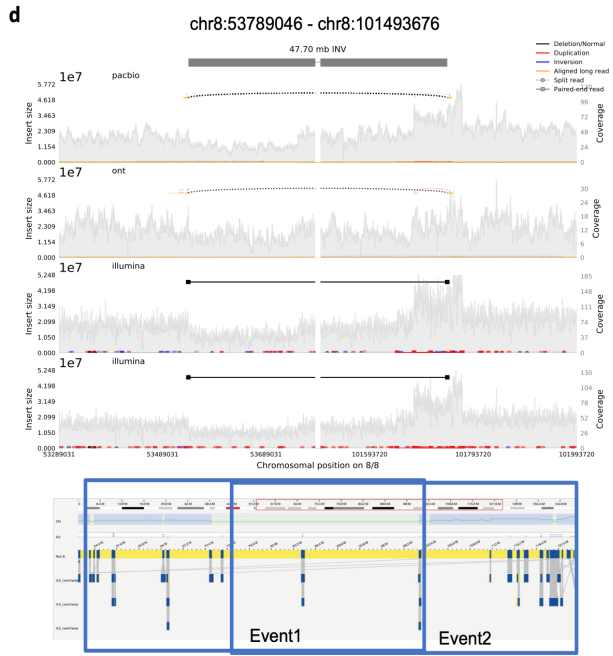
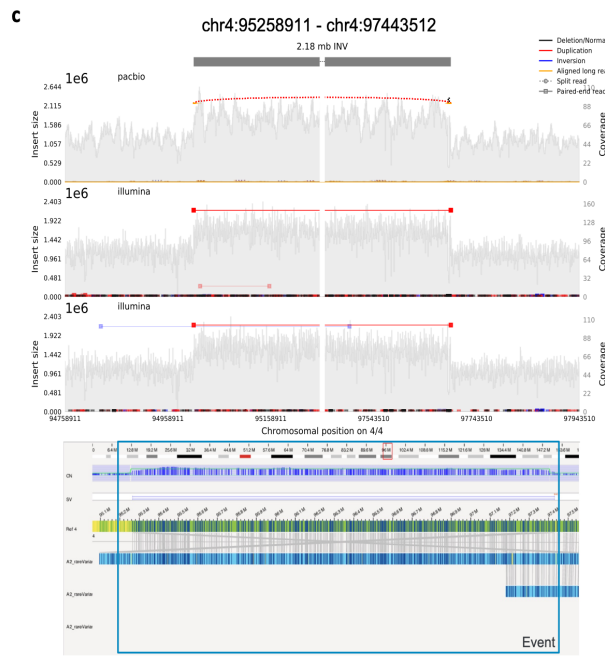
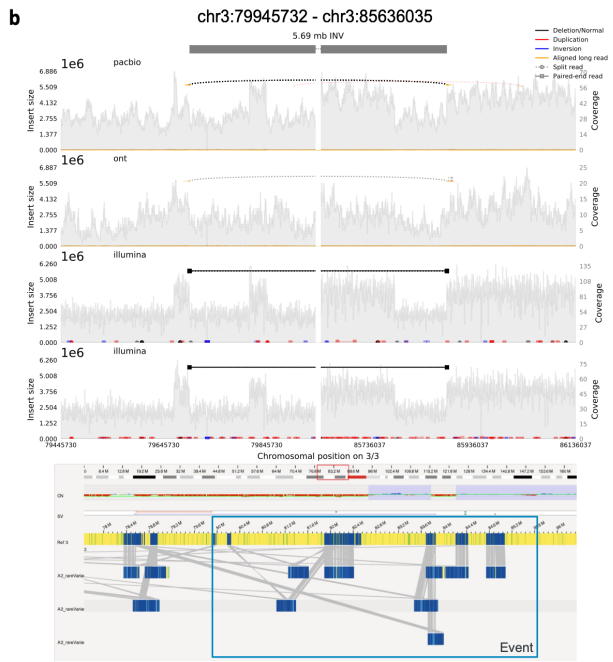
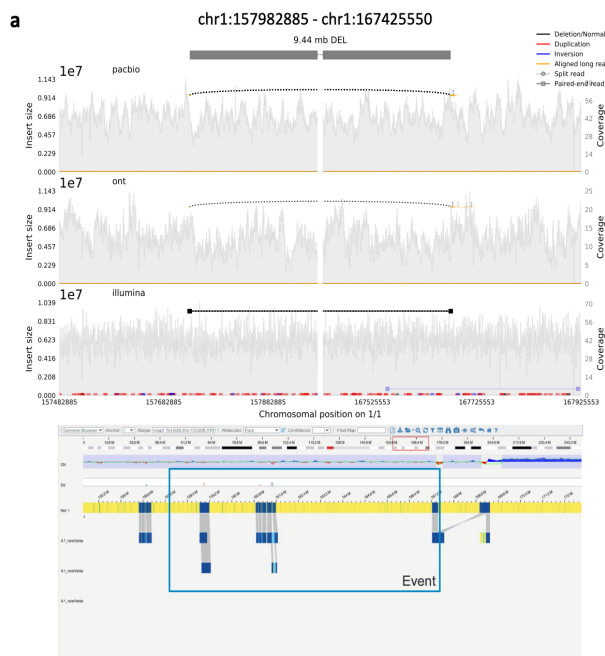


Fig. S8. Reproducibility of SVs in Consensus Call Set from 5 NGS Platforms and 11 Callers. Heatmap was generated based on SV location on genome and SV frequency which was calculated based on the method section of Calculating SV Calling Frequency and Consensus Call Set. Software tools displayed in the plots are: GrocSVs, Long Ranger, TNScope, Delly, Manta, Novobreak, HiC (Selva), Novasv, Sniffles.ONT for Nanopore data, Sniffles.PB for PacBio data, PBSV. The heatmap color denotes the SV frequencies detected by each tool for a technology, the dendrograms along the side of the heatmap show similarity and variability how the SVs are clustered.



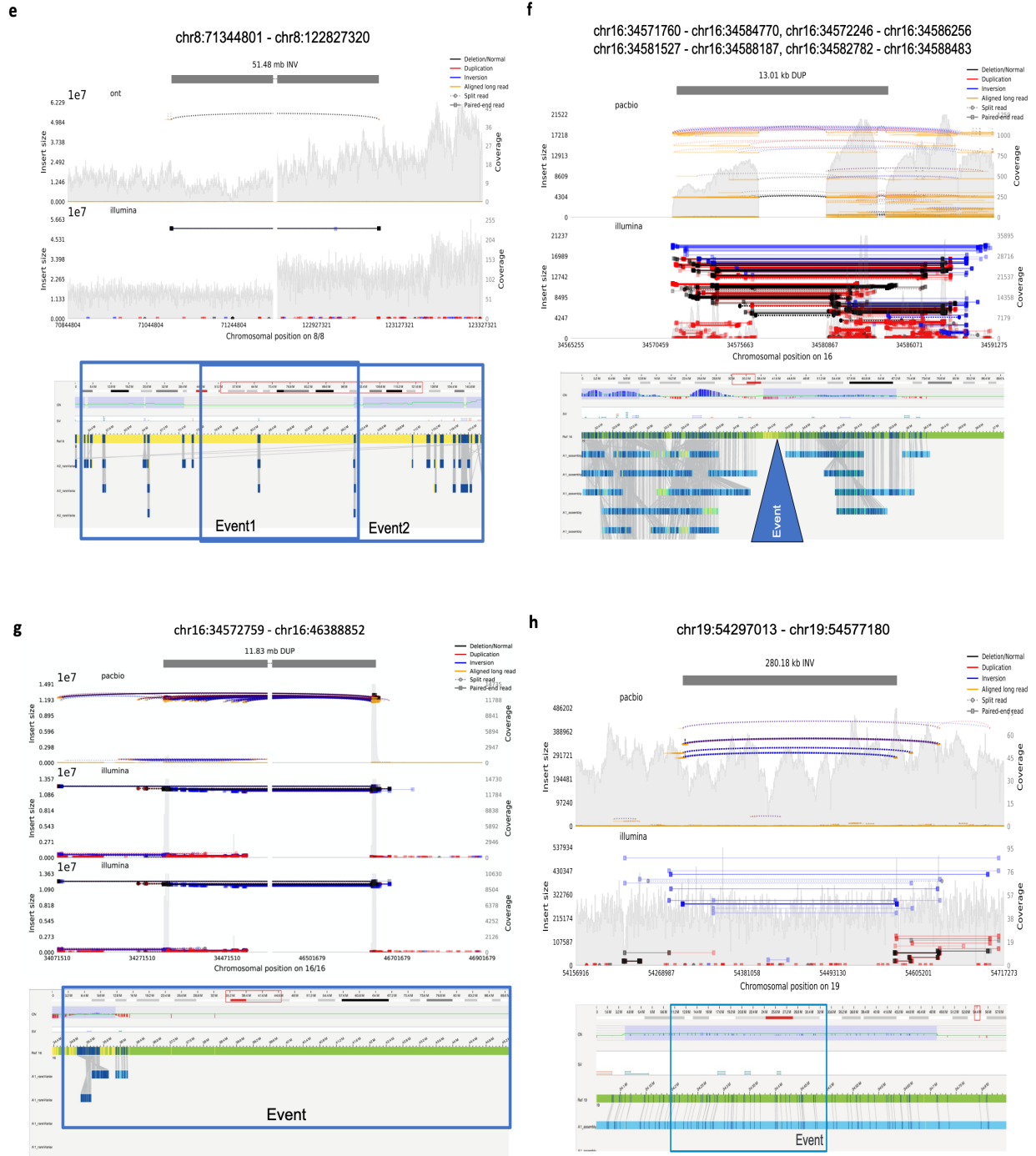


Fig. S9. Examples of Complex Break Point Events (BND) Called by Different Technology and Software Tools. Plots were generated by using SamPlot tool to show multiple break point events in regions displayed in each panel. The genome read coverage and structural variant events were profiled from PacBio, Nanopore (ONT), Illumina, and Bionano data. The summary of each BND event characterization was described in **Additional File 16: Table S15**.