

# SQUID: Transcriptomic Structural Variation Detection from RNA-seq

Cong Ma<sup>1</sup>, Mingfu Shao<sup>1</sup>, and Carl Kingsford<sup>\*1</sup>

<sup>1</sup>Computational Biology Department, School of Computer Science, Carnegie Mellon University,  
5000 Forbes Ave., Pittsburgh, PA

July 10, 2017

## Abstract

Transcripts are frequently modified by structural variations, which leads to either a fused transcript of two genes (known as fusion gene) or an insertion of intergenic sequence into a transcript. These modifications are termed transcriptomic structural variants (TSV), and they can lead to drastic change of a downstream translation product. Detecting TSVs, especially in cancer tumor sequencing where they are known to frequently occur, is an important and challenging computational problem. This problem is made even more challenging in that often only RNA-seq measurements are available from the sample. We introduce SQUID, a novel algorithm and its implementation, to accurately and comprehensively predict both fusion-gene and non-fusion-gene TSVs from RNA-seq alignments. SQUID takes the unique approach of attempting to reconstruct an underlying genome sequence that best explains the observed RNA-seq reads. By unifying both concordant alignments and discordant read alignments into one model, SQUID achieves high sensitivity with many fewer false positives than other approaches. We detect TSVs on TCGA tumor samples using SQUID, and observe that breast cancer samples are more likely to contain a large number of TSVs than several other cancer types. We further find that non-fusion-gene TSVs are more likely to be intra-chromosomal than fusion-gene TSVs while the breakpoint separation distance tends to be larger than that of fusion-gene TSVs in intra-chromosomal case. We also identify several novel TSVs involving tumor suppressor genes, which may lead to loss-of-function of corresponding genes and play a role in tumorigenesis.

---

\*To whom correspondence should be addressed: [carlk@cs.cmu.edu](mailto:carlk@cs.cmu.edu)

## 25 **1 Introduction**

26 Large-scale transcriptome sequence changes are known to be associated with cancer [21, 34]. Those changes  
27 are usually a consequence of genomic structural variation (SV). By pulling different genomic regions to-  
28 gether or separating one region into pieces, structural variants can potentially cause severe alteration to  
29 transcribed or translated products. Transcriptome changes induced by genomic SVs, called transcriptomic  
30 structural variants (TSVs), can have a particularly large impact on disease genesis and progression. In some  
31 cases, TSVs bring regions from one gene next to regions of another, causing exons from both genes to be  
32 transcribed into a single transcript (known as a fusion gene). Domains of the corresponding RNA or proteins  
33 can be fused, inducing new functions or causing loss of function, or the transcription or translation levels can  
34 be altered, leading to disease states. For example, BCR-ABL1 is a well-known fusion oncogene for chronic  
35 myeloid leukemia [8], and the TMPRSS2-ERG fusion product leads to over-expression of ERG and helps  
36 triggers prostate cancer [35]. These fusion events are used as biomarkers for early diagnosis or treatment  
37 targets [36]. In other cases, TSVs can affect genes by causing a previously non-transcribed region to be  
38 incorporated into a gene, causing disruption to the function of the altered gene. There are fewer studies on  
39 these TSVs between transcribed and non-transcribed regions, but their ability to alter downstream RNA and  
40 protein structure is likely to lead to similar results as fusion gene TSVs, and contribute to tumor genesis and  
41 progression.

42 Genomic SVs are typically detected from whole-genome sequencing (WGS) data by identifying reads and  
43 read pairs that are incompatible with a reference genome [e.g., 5, 14, 17, 27, 28]. However, WGS data are  
44 not completely suitable to infer TSVs since they neither inform which region is transcribed nor reveal how  
45 transcribed sequence will change if SVs alter a splicing site or the stop codon. In addition, WGS data is more  
46 scarce and more expensive to obtain than RNA-seq [31] measurements, which sequence transcribed regions  
47 directly. RNA-seq is relatively inexpensive, high-throughput, and widely available in many existing and  
48 growing data repositories. For example, The Cancer Genome Atlas (TCGA, <https://cancergenome.nih.gov>)  
49 contains RNA-seq measurements from thousands of tumor sample across various cancer types, but 80% of  
50 tumor samples in TCGA have RNA-seq data but no WGS data (Supp. Figure S1). While methods exist to  
51 detect fusion genes from RNA-seq measurements [e.g., 7, 15, 20, 26, 41], fusion genes are only a subset  
52 of TSVs, and existing fusion gene detection methods rely heavily on current gene annotations and are

53 generally not able or at least not optimized to predict non-fusion-gene TSV events. This motivates the need  
54 for a method to detect all types of TSVs directly from RNA-seq data.

55 We present SQUID, the first computational tool that comprehensively and accurately predicts TSVs from  
56 RNA-seq data. SQUID divides the reference genome into segments and builds a genome segment graph  
57 from both concordant and discordant RNA-seq read alignments. In this way, it can detect both fusion-gene  
58 events and TSVs incorporating previously non-transcribed regions into transcripts. Using an efficient, novel  
59 integer linear program (ILP), SQUID rearranges the segments of the reference genome so that as many  
60 read alignments as possible are concordant with the rearranged sequence. TSVs are represented by pairs  
61 of breakpoints realized by the rearrangement. Discordant reads that cannot be made concordant through  
62 the optimal rearrangement given by the ILP are discarded as false positive discordant reads, likely due to  
63 misalignments. By building a consistent model of the entire rearranged genome and maximizing the number  
64 of overall concordant read alignments, SQUID drastically reduces the number of spurious TSVs reported  
65 compared with other methods.

66 SQUID features high accuracy. SQUID is usually  $> 20\%$  more accurate than applying WGS-based SV  
67 detection methods to RNA-seq data directly. It is similarly more accurate than a pipeline that uses de novo  
68 transcript assembly and transcript-to-genome alignment to detect TSVs. We also show that SQUID is able  
69 to detect more TSVs involving non-transcribed regions than any existing fusion gene detection method.

70 We use SQUID to detect TSVs within 401 TCGA tumor samples of four cancer types (99–101 samples each  
71 of breast invasive carcinoma [22], bladder urothelial carcinoma [23], lung adenocarcinoma [24], and prostate  
72 adenocarcinoma [25]). SQUID's predictions suggest that breast invasive carcinoma has more fusion-gene  
73 TSVs and more non-fusion-gene TSVs than other cancer types. We also characterize the differences between  
74 fusion-gene TSVs and non-fusion-gene TSVs. Non-fusion-gene TSVs, for example, are more likely to be  
75 intra-chromosomal events, and within those intra-chromosomal events the two breakpoints of non-fusion-  
76 gene TSVs tend to be farther apart from each other than of fusion-gene TSVs. We show that breakpoints  
77 can occur in multiple samples, and among those that do repeatedly occur, their breakpoint partners are  
78 also often conserved. Finally, we identify several novel non-fusion-gene TSVs that affect known tumor  
79 suppressor genes, which may result in loss-of-function of corresponding proteins and play a role in tumor  
80 genesis.

## 81 2 Methods

### 82 2.1 The computational problem: rearrangement of genome segments

83 We formulate the TSV detection problem as the optimization problem of rearranging genome segments to  
84 maximize the number of observed reads that are consistent (termed *concordant*) with the rearranged genome.  
85 This approach requires defining the genome segments that can be independently rearranged. It also requires  
86 defining what reads are consistent with a particular arrangement of the segments. We will encode both of  
87 these (segments and read consistency) within a *Genome Segment Graph* (GSG). See Figure 1 as an example.

88 **Definition 1** (Segment). A segment is a pair  $s = (s_h, s_t)$ , where  $s$  represents a continuous sequence in  
89 reference genome and  $s_h$  represents its head and  $s_t$  represents its tail in reference genome coordinates. In  
90 practice, segments will be derived from the read locations (Section 2.4).

91 **Definition 2** (Genome Segment Graph (GSG)). A genome segment graph  $G = (V, E, w)$  is an undirected  
92 weighted graph, where  $V$  contains both endpoints of each segment in a set of segments  $S$ , i.e.,  $V = \{s_h : s \in S\} \cup \{s_t : s \in S\}$ . Thus, each vertex in the GSG represents a location in the genome. An edge  
93  $(u, v) \in E$  indicates that there is evidence that the location  $u$  is in fact adjacent to location  $v$ . Weight  
94 function,  $w : E \rightarrow \mathbb{R}^+$ , represents the reliability of an edge. Generally speaking, the weight is the number  
95 of read alignments supporting, but we allow a multiplier to calculate edge weight which will be discussed  
96 below. In practice,  $E$  and  $w$  will be derived from split-aligned and paired-end reads (Section 2.5).

98 Defining vertices by endpoints of segments is required to avoid ambiguity. Only knowing that segment  $i$  is  
99 connected with segment  $j$  is not enough to recover the sequence, since different relative positions of  $i$  and  
100  $j$  spell out different sequences. Instead, for example, an edge  $(i_t, j_h)$  indicates that the tail of segment  $i$  is  
101 connected head of segment  $j$ , and this specifies a unique desired local sequence with only another possibility  
102 of the reverse complement (i.e. it could be that the true sequence is  $i \cdot j$  or  $rev(j) \cdot rev(i)$ ; here  $\cdot$  indicates  
103 concatenation and  $rev(i)$  is the reverse complement of segment  $i$ ).

104 The GSG is similar to the breakpoint graph [2] but with critical differences. A breakpoint graph has edges  
105 representing both connections in reference genome and in target genome. While edges in the GSG only  
106 represents the target genome, and they can be either concordant or discordant. In addition, the GSG does

107 not require that the degree of every vertex is two, and thus alternative splicing and erroneous edges can exist  
108 in the GSG.

109 Our goal is to reorder and reorient the segments in  $S$  so that as many edges in  $G$  are compatible with the  
110 rearranged genome as possible.

111 **Definition 3** (Permutation). *A permutation  $\pi$  on a set of segments  $S$  projects a segment in  $S$  to a set of*  
112 *integers from 1 to  $|S|$  (the size of  $S$ ) representing the indices of the segments in an ordering of  $S$ . In other*  
113 *words, each permutation  $\pi$  defines a new order of segments in  $S$ .*

114 **Definition 4** (Orientation Function). *An orientation function  $f$  maps both ends of segments to 0 or 1:*

$$f : \{s_h : s \in S\} \cup \{s_t : s \in S\} \longrightarrow \{0, 1\}$$

115 *subject to  $f(s_h) + f(s_t) = 1$  for all  $s = (s_h, s_t) \in S$ . An orientation function specifies the orientations of*  
116 *all segments in  $S$ . Specifically,  $f(s_h) = 1$  means  $s_h$  goes first and  $s_t$  next, corresponding to forward strand*  
117 *of segment, and  $f(s_t) = 1$  corresponds to the reverse strand of the segment.*

118 With a permutation  $\pi$  and an orientation function  $f$ , the exact and unique sequence of genome is determined.  
119 The reference genome also corresponds to a permutation and an orientation function, where the permutation  
120 is the identity permutation, and the orientation function maps all  $s_h$  to 1 and all  $s_t$  to 0.

121 **Definition 5** (Edge Compatibility). *Given a set of segments  $S$ , a genome segment graph  $G = (V, E, w)$ , a*  
122 *permutation  $\pi$  on  $S$ , and an orientation function  $f$ , an edge  $e = (u_i, v_j) \in E$ , where  $u_i \in \{u_h, u_t\}$  and*  
123  *$v_j \in \{v_h, v_t\}$ , is compatible with permutation  $\pi$  and orientation  $f$  if and only if*

$$1 - f(v_j) = \mathbf{I}[\pi(v) < \pi(u)] = f(u_i) \quad (1)$$

124 *where  $\mathbf{I}[x]$  is the indicator function that is 1 if  $x$  is true and 0 otherwise. We write  $e \sim (\pi, f)$  if  $e$  is*  
125 *compatible with  $\pi$  and  $f$ .*

126 The above two edge compatibility equations (1) require that, in order for an edge to be compatible with  
127 the rearranged and reoriented sequence determined by  $\pi$  and  $f$ , the edge needs to connect the right side  
128 of the segment in front to the left side of segment following it. As we will see in Section 2.5, edges of

129 GSG are derived from reads alignments. An edge being compatible with  $\pi$  and  $f$  is essentially equivalent to  
130 the statement that the corresponding read alignments are concordant (Section 2.3) with respect to the target  
131 genome determined by  $\pi$  and  $f$ . When  $(\pi, f)$  is clear, we refer to edges that are compatible as concordant  
132 edges, and edges that are incompatible as discordant edges.

133 With the above definitions, we formulate an optimization problem as follows:

134 **Problem 1. Input:** A set of segments  $S$  and a GSG  $G = (V, E, w)$ .

135 **Output:** Permutation  $\pi$  on  $S$  and orientation function  $f$  that maximizes:

$$\max_{\pi, f} \sum_{e \in E} w(e) \cdot \mathbf{I}[e \sim (\pi, f)] \quad (2)$$

136 This objective function tries to find a rearrangement of genome segments  $(\pi, f)$ , such that when aligning  
137 reads to the rearranged sequence, as many reads as possible will be aligned concordantly. This objec-  
138 tive function includes both concordant alignments and discordant alignments and sets them in competition,  
139 which will be effective in reducing false positives when tumor transcripts out-number normal transcripts.  
140 There is the possibility that some rearranged tumor transcripts are out-numbered by normal counterparts. In  
141 order to be able to detect TSV in this case, we weight discordant read alignments more than concordant read  
142 alignments. Specifically, for each discordant edge  $e$ , we multiply the weight  $w(e)$  by a constant  $\alpha$ , which  
143 represents our estimate of the ratio of normal transcripts over tumor counterparts.

144 The final TSVs are modeled as pairs of breakpoints. Denote the permutation and orientation corresponding  
145 to an optimally rearranged genome as  $(\pi^*, f^*)$  and those that correspond to reference genome as  $(\pi_0, f_0)$ .  
146 An edge  $e$  can be predicted as a TSV if  $e \sim (\pi^*, f^*)$  and  $e \not\sim (\pi_0, f_0)$ .

## 147 2.2 Integer linear programming formulation

148 We use integer linear programming (ILP) to compute an optimal solution  $(\pi^*, f^*)$  of Problem 1. To do this,  
149 we introduce the following boolean variables:

- 150 •  $x_e$ :  $x_e = 1$  if edge  $e \sim (\pi^*, f^*)$ , and  $x_e = 0$  if not.
- 151 •  $z_{uv}$ :  $z_{uv} = 1$  if segment  $u$  is before  $v$  in the permutation  $\pi^*$ , and 0 otherwise.

- 152 •  $y_u$ :  $y_u = 1$  if  $f^*(u_h) = 1$  for segment  $u$ .

153 With this representation, the objective function can be rewritten as

$$\max_{x_e, y_u, z_{uv}} w(e) \cdot x_e \quad (3)$$

154 We add constraints to the ILP derived from edge compatibility equations (1). Without loss of generality,  
 155 we first suppose segment  $u$  is in front of  $v$  in the reference genome, and edge  $e$  connects  $u_t$  and  $v_h$  (which  
 156 is a tail-head connection). Plugging in  $u_t$ , the first equation in (1) is equivalent to  $1 - \mathbf{1}[\pi(u) > \pi(v)] =$   
 157  $1 - f(u_t)$ , and can be rewritten as  $\mathbf{1}[\pi(u) < \pi(v)] = f(u_h) = y_u$ . Note that  $\mathbf{1}[\pi(u) < \pi(v)]$  has the  
 158 same meaning as  $z_{uv}$ ; it leads to the constraint  $z_{uv} = y_u$ . Similarly, the second equation in (1) indicates  
 159  $z_{uv} = y_v$ . Therefore,  $x_e$  can only reach 1 when  $y_u = y_v = z_{uv}$ . This is equivalent to the inequalities (4)  
 160 below. Analogously, we can write constraints for other three types of edge connections: tail-tail connec-  
 161 tions impose inequalities (5); head-head connections impose inequalities (6); head-tail connections impose  
 162 inequalities (7):

$$\begin{array}{ll} x_e \leq y_u - y_v + 1 & x_e \leq y_u - (1 - y_v) + 1 \\ x_e \leq y_v - y_u + 1 & x_e \leq (1 - y_v) - y_u + 1 \\ x_e \leq y_u - z_{uv} + 1 & x_e \leq y_u - z_{uv} + 1 \\ x_e \leq z_{uv} - y_u + 1 & x_e \leq z_{uv} - y_u + 1 \end{array} \quad \begin{array}{l} (4) \\ (5) \end{array}$$

$$\begin{array}{ll} x_e \leq (1 - y_u) - y_v + 1 & x_e \leq (1 - y_u) - (1 - y_v) + 1 \\ x_e \leq y_v - (1 - y_u) + 1 & x_e \leq (1 - y_v) - (1 - y_u) + 1 \\ x_e \leq (1 - y_u) - z_{uv} + 1 & x_e \leq (1 - y_u) - z_{uv} + 1 \\ x_e \leq z_{uv} - (1 - y_u) + 1 & x_e \leq z_{uv} - (1 - y_u) + 1 \end{array} \quad \begin{array}{l} (6) \\ (7) \end{array}$$

163 We also add constraints to enforce that  $z_{uv}$  forms a valid topological ordering. For each pair of nodes  $u$  and  
 164  $v$ , one must be in front of other, that is  $z_{uv} + z_{vu} = 1$ . In addition, for each triple of nodes,  $u$ ,  $v$  and  $w$ , they  
 165 cannot be all in front of another; one must be at the beginning of these three and one must be at the end.

166 Therefore we add  $1 \leq z_{uv} + z_{vw} + z_{wu} \leq 2$ .

167 Solving an ILP in theory takes exponential time, but in practice, solving the above ILP to rearrange genome  
168 segments is very efficient. The key is that we can solve for each connected component separately. Because  
169 the objective maximizes the sum of compatible edge weight, the best rearrangement of one connected com-  
170 ponent is independent from the rearrangement of another because by definition there are no edges between  
171 connected components.

### 172 **2.3 Concordant and discordant alignments**

173 Discordant alignments are alignments of reads that contradict library preparation in sequencing. Concordant  
174 alignments are alignments of reads that agree with the library preparation. Take Illumina sequencing as an  
175 example. In order for a paired-end read alignment to be concordant, one end should be aligned to the forward  
176 strand and the other to the reverse strand, and the forward strand aligning position should be in front of the  
177 reverse strand aligning position (Figure 2a). Concordant alignment traditionally used in WGS also requires  
178 that a read cannot be split and aligned to different locations. But these requirements are invalid in RNA-seq  
179 alignments because alignments of reads can be separated by an intron with unknown length.

180 We define concordance criteria separately for split-alignment and paired-end alignment. If one end of a  
181 paired-end read is split into several parts and each part is aligned to a location, the end has split-alignments.  
182 Denote the vector of the split alignments of an end to be  $R = [A_1, A_2, \dots, A_r]$  ( $r$  depends on the number  
183 of splits). Each alignment  $R[i] = A_i$  is comprised of 4 components: chromosome (Chr), alignment starting  
184 position (Spos), alignment ending position (Epos) and orientation (Ori, with value either + or -). We  
185 require that the alignments  $A_i$  are sorted by their position in read. A split-aligned end  $R = [A_1, A_2, \dots, A_r]$   
186 is concordant if all the following conditions hold:

$$\begin{aligned} A_i.Chr &= A_j.Chr && \forall i, \forall j \\ A_i.Ori &= A_j.Ori && \forall i, \forall j \\ A_i.Spos &< A_j.Spos && \text{if } A_i.Ori = + \text{ for all } i < j \\ A_i.Spos &> A_j.Spos && \text{if } A_i.Ori = - \text{ for all } i < j \end{aligned} \tag{8}$$

187 Note that if the end is not split, but continuous aligned, the alignment automatically satisfy equation (8).



188 Denote the alignments of  $R$ 's mate as  $M = [B_1, B_2, \dots, B_m]$ . An alignment of the paired-end read is  
189 concordant if the following conditions all hold:

$$\begin{aligned} A_i.Chr &= B_j.Chr \\ A_i.Ori &\neq B_j.Ori \\ A_1.Spos < B_m.Spos &\text{ if } A_1.Ori = + \\ A_m.Spos > B_1.Spos &\text{ if } A_1.Ori = - \end{aligned} \tag{9}$$

190 We only require the left-most split of the forward read  $R$  be in front of the left-most split of the reverse read  
191  $M$  since the two ends in a read pair may overlap. In order for a paired-end read to be concordant, each  
192 end should satisfy split-read alignment concordance (8), and the pair should satisfy paired-end alignment  
193 concordance (9).

## 194 **2.4 Splitting the genome into segments $S$**

195 We use a set of breakpoints to partition the genome. The set of breakpoints contains two types of positions:  
196 (1) the start position and end position of each interval of overlapping discordant alignments, (2) an arbitrary  
197 position in each 0-coverage region.

198 Ideally, both ends of a discordant read should be located in separate segments, otherwise, the discordant  
199 read contained in a single segment will always be discordant no matter how the segments are rearranged.  
200 Assuming discordant read alignments of each TSV pile up around the breakpoints and do not overlap with  
201 discordant alignments of other TSVs, we set a breakpoint on the start and end positions of each contiguous  
202 interval of overlapping discordant alignments.

203 For each segment that contains discordant read alignments, it may also contain concordant alignments that  
204 connect the segment to its adjacent segments. To avoid having all segments in GSG connected to their  
205 adjacent segments and thus creating one big connected component, we pick the starting point of each 0-  
206 coverage region as a breakpoint. By adding those breakpoint, different genes will be in separate connected  
207 components unless some discordant reads support their connection. Overall, the size of each connected  
208 component is not very large: the number of nodes generated by each gene is approximately the number of  
209 exons located in them and these gene subgraphs are connected only when there is a potential TSV between

210 them.

## 211 **2.5 Defining edges in the genome segment graph**

212 In a GSG, an edge is added between two vertices when there are reads supporting the connection. For each  
213 read spanning different segments, we build an edge such that when traversing the segments along the edge,  
214 the read is concordant with the new sequence (equations (8) and (9)). Examples of deriving an edge from a  
215 read alignment are given in Figure 2. In this way, concordance of an alignment and compatibility of an edge  
216 with respect to a genome sequence is equivalent.

217 The weight of a concordant edge is the number of read alignments supporting the connection, while the  
218 weight of a discordant edge is the number of alignments supporting multiplied by discordant edge weight  
219 coefficient  $\alpha$ . Edges with very low read support are likely to be a result of alignment error, therefore we filter  
220 out edges with weight lower than a threshold  $\theta$ . Segments with too many connections to other regions are  
221 likely to have low mappability, so we also filter out segments connecting to more than  $\gamma$  other segments. The  
222 parameters  $\alpha$ ,  $\theta$ , and  $\gamma$  are the most important user-defined parameters to SQUID (Supplementary Table 1  
223 and Supplementary Figure S2).

## 224 **2.6 Identifying TSV breakpoint locations**

225 Edges that are discordant in the reference genome indicate potential rearrangements in transcripts. Among  
226 those edges, some are compatible with the permutation and orientation from ILP. These edges are taken to be  
227 the predicted TSVs. For each edge that is discordant initially but compatible with the optimal rearrangement  
228 found by the ILP, we examine the discordant read alignments to determine the exact breakpoint located  
229 within related segments. Specifically, for each end of a discordant alignment, if there are 2 other read  
230 alignments that start or end in the same position and support the same edge, then the end of the discordant  
231 alignment is predicted to be the exact TSV breakpoint. Otherwise, the boundary of the corresponding  
232 segment will be output as the exact TSV breakpoint.

## 233 **2.7 Simulation methodology**

234 Simulations with randomly added structural variations and simulated RNA-seq reads were used to evalu-  
235 ate SQUID's performance in situations with a known correct answer. RSVsim [3] was used to simulate  
236 SV on the human genome (Ensembl 87 or hg38) [40]. We use the 5 longest chromosomes for simulation  
237 (chromosome 1 to chromosome 5). RSVsim introduces 5 different types of SVs: deletion, inversion, inser-  
238 tion, duplication, and inter-chromosomal translocation. To vary the complexity of the resulting inference  
239 problem, we simulated genomes with 200 SVs of each type, 500 SVs of each type, and 800 SVs of each  
240 type. We generated 4 replicates for each level of SV complexity (200, 500, 800). For inter-chromosomal  
241 translocations, we only simulate 2 events because only 5 chromosomes were used.

242 In the simulated genome with SVs, the original gene annotations are not applicable, and we cannot simulate  
243 gene expression from the rearranged genome. Therefore, for testing purposes, we interchange the role  
244 of the reference (hg38) and rearranged genome, and use the new genome as the reference genome for  
245 alignment, and hg38 with the original annotated gene positions as the target genome for sequencing. Flux  
246 Simulator [12] was used to simulate RNA-seq reads from the hg38 genome using the Ensembl annotation  
247 version 87 [1].

248 After simulating SVs on genome, we need to transform SVs into a set of TSVs, because not all SVs affect  
249 transcriptome, and thus not all SVs can be detected by RNA-seq. To derive the list of TSVs, we compare  
250 the positions of simulated SVs with the gene annotation. If a gene is affected by an SV, some adjacent  
251 nucleotides in the corresponding transcript may be located far part in the RSVsim-generated genome. The  
252 adjacent nucleotides can be consecutive nucleotides inside an exon if the breakpoint breaks the exon, or the  
253 end points of two adjacent exons if the breakpoint hits the intron. So for each SV that hits a gene, we find  
254 the pair of nucleotides that are adjacent in transcript and separated by the breakpoints, and converted them  
255 into coordinate of the RSVsim-generated genome, thus deriving the TSV.

256 Since there are no existing methods for annotation-free TSV detection, we compare SQUID to the pipeline  
257 of de novo transcriptome assembly and transcript-to-genome alignment. We also use the same set of simu-  
258 lations to test whether existing WGS-based SV detection methods can be directly applied to RNA-seq data.  
259 For the de novo transcriptome assembly and transcript-to-genome alignment pipeline, we use all combi-  
260 nations of the existing software Trinity [11], Trans-ABYSS [29], GMAP [37] and MUMmer3 [16]. For

261 WGS-based SV detection methods, we test LUMPY [17] and DELLY2 [28]. We test both STAR [9] and  
262 SpeedSeq [6] (which is based on BWA-MEM [18]) to align RNA-seq reads to the genome. LUMPY is only  
263 compatible with SpeedSeq output, so we do not test it with STAR alignments.

## 264 **3 Results**

### 265 **3.1 SQUID is accurate on simulation data**

266 Overall, SQUID's predictions of TSVs are far more precise than other approaches at similar sensitivity  
267 on simulated data (Section 2.7). SQUID achieves 60% to 80% percent precision and about 50% percent  
268 sensitivity on simulation data (Figure 3). SQUID's precision is  $\approx 40\%$  higher than all combinations of  
269 de novo transcriptome assembly and transcript-to-genome alignment pipeline, and the precision of WGS-  
270 based SV detection methods on RNA-seq data is even lower. The sensitivity of SQUID is similar to de novo  
271 assembly with MUMmer3, but a little lower than DELLY2 and LUMPY with SpeedSeq aligner. The overall  
272 sensitivity is not as high as precision, which is probably because there are not enough supporting reads  
273 aligned correctly to some TSV breakpoints. The fact that assembly and WGS-based SV detection methods  
274 achieve similar sensitivity corroborates the hypothesis that it is the data limiting the achievable sensitivity.

275 The low specificity of the pipeline- and WGS-based methods shows neither of these types of approaches  
276 are suitable for TSV detection from RNA-seq data. WGS-based SV detection methods are able to detect  
277 TSV signals, but not able to filter out false positives. Assembly-based approaches require solving the tran-  
278 scriptome assembly problem which is a harder and more time-consuming problem, and thus errors are more  
279 easily introduced. Further, the performance of assembly pipelines depends heavily on the choice of software  
280 — for example, MUMmer3 is better at discordantly aligning transcripts than GMAP.

281 SQUID is likely effective due to its unified model of both concordant reads and discordant reads. Coverage  
282 in RNA-seq alignment is proportional to the expression level of the transcript, and using one read count  
283 threshold for TSV evidence is not appropriate. Instead, the ILP in SQUID sets concordant and discordant  
284 alignments into competition and selects the winner as the most reliable TSVs.

## 285 **3.2 SQUID is able to detect non-fusion-gene TSV on two previously-studied cell lines**

286 Fusion gene events are a strict subset of TSVs where the two breakpoints are each be within a gene region  
287 and the fused sequence corresponds to the sense strand of both genes. Fusion genes thus exclude TSV events  
288 where a gene region is fused with a intergenic region or an anti-sense strand of another gene. Nevertheless,  
289 fusion genes have been implicated (likely because of available methods to detect them) in playing a role in  
290 cancer.

291 To probe SQUID's ability to detect this subclass of TSVs, we use two cell lines, HCC1954 and HCC1395, for  
292 which previous studies have experimentally validated predicted SVs and fusion gene events. Specifically,  
293 we compile results from Bignell et al. [4], Galante et al. [10], Stephens et al. [33], Zhao et al. [42] and  
294 Robinson et al. [30] for HCC1954, and results from Stephens et al. [33] and Zhang et al. [41] for HCC1395.  
295 After removing short deletions and overlapping structural variations among different studies, we have 326  
296 validated structural variations for HCC1954 cell line, in which 245 of them have at least one breakpoint  
297 outside gene region, and the rest 81 have both breakpoints within gene region; we have 256 validated  
298 true structural variations for HCC1395 cell line, in which 94 have at least one breakpoint outside gene  
299 region, while the rest 162 have both breakpoints within gene. For a predicted structural variation to be  
300 true positive, both predicted breakpoints should be within a window of 30kb of true breakpoints and the  
301 predicted orientation should agree with true orientation. We use a relatively large window since the true  
302 breakpoints can be located within an intron or other non-transcribed region, while the observed breakpoint  
303 from RNA-seq reads will be at a nearby coding or expressed region.

304 We use publicly available RNA-seq data from the NIH Sequencing Read Archive (SRA; accessions:  
305 SRR2532344 and SRR925710 for HCC1954, SRR2532336 for HCC1395). Because the data are from  
306 an pool of experiments, the sample from which RNA-seq was collected may be different from those used  
307 for experimental validation. We align reads to the reference genome using STAR.

308 When restricted to fusion gene events, SQUID achieves similar precision and sensitivity compared to fusion  
309 gene detection tools (Figure 4A). SQUID has the highest accuracy in HCC1954 cell line, with very similar  
310 sensitivity as all fusion gene detection tools. For HCC1395, SQUID is in the middle of fusion gene detection  
311 methods, while INTEGRATE and JAFFA are the best performers on this sample.

312 For non-fusion-gene TSVs, it is even harder to predict them accurately, since current annotations cannot be  
313 used to limit the search space for potential read alignments or TSV events. Only SQUID and deFuse are  
314 able to detect non-fusion-gene events. Between these two methods, SQUID is able to predict more known  
315 non-fusion-gene TSVs correctly (Figure 4B). By considering both fusion-gene and non-fusion-gene TSVs  
316 in SQUID predictions, the number of correct predictions greatly increases compared to considering fusion-  
317 gene TSVs only, since a considerable proportion of validated TSVs are non-fusion-gene TSVs. At the same  
318 time, precision does not decrease very much by considering both fusion-gene and non-fusion-gene TSVs.

### 319 **3.3 Charactering TSVs on four types of TCGA cancer samples**

320 To compare the distributions and characteristics of TSVs among cancer types and between TSV types, we  
321 arbitrarily selected 99 to 101 tumor samples from TCGA for each of four cancer types: breast invasive  
322 carcinoma (BRCA), bladder urothelial carcinoma (BLCA), lung adenocarcinoma (LUAD), and prostate  
323 adenocarcinoma (PRAD).

324 To estimate the accuracy of SQUID's prediction on selected TCGA samples, we use WGS data of the  
325 same patients to validate TSV junctions. There are in total 72 WGS experiments available for the 400  
326 samples (20 BLCA, 10 BRCA, 31 LUAD, 11 PRAD). For each TSV prediction, we extract a 25Kb sequence  
327 around both breakpoints and concatenate them according to the predicted TSV orientation. We then map  
328 the WGS reads against these junction sequences using SpeedSeq. If a paired-end WGS read can only be  
329 mapped concordantly to a junction sequence but not the reference genome, that paired-end read is marked  
330 as supporting the TSV. If at least 3 WGS reads support a TSV, the TSV is considered as validated. Using  
331 this approach, SQUID's overall validation rate is 88.21%, and this indicates that SQUID is quite accurate  
332 and reliable on TCGA data.

333 We find that most samples have  $\approx 15 - 20$  TSVs including  $\approx 3 - 5$  non-fusion-gene TSVs among all four  
334 cancer types (Figure 5A,B). BRCA has more samples with a larger number of TSVs: there are 37 BRCA  
335 samples with more than 20 TSVs, while for other cancer types there are at most 26 samples with  $> 20$  TSVs.  
336 The same trend is observed when restricted to non-fusion-gene TSVs, where there are 29 BRCA samples  
337 with more than 8 non-fusion-gene TSVs, while any of the other cancer types has at most 11 samples with  
338  $> 8$  non-fusion-gene TSVs. This observation agrees with Yang et al. [39], where they observe BRCA has  
339 more somatic SVs than PRAD.

340 Inter-chromosomal TSVs are more prevalent than intra-chromosomal TSVs for all cancer types (Figure 5C),  
341 although this difference is much more pronounced in bladder and prostate cancer. Non-fusion-gene TSVs  
342 are more likely to have intra-chromosomal events than fusion gene TSVs (Figure 5D), and in fact in  
343 bladder, breast, and lung cancer, we detect more intra-chromosomal non-fusion-gene TSVs than inter-  
344 chromosomal non-fusion-gene TSVs. Prostate cancer is an exception in that for non-fusion-gene TSVs,  
345 inter-chromosomal events are observed more often than intra-chromosomal events. Nevertheless, it also  
346 holds true that non-fusion-gene TSVs are more likely to be intra-chromosomal than fusion-gene, because  
347 the percentage of intra-chromosomal TSVs within non-fusion-gene TSVs is higher than that within all TSVs.

348 For a large proportion of breakpoints occurring multiple times within a cancer type, their partner in the TSV  
349 is likely to be fixed and to reoccur every time that breakpoint is used. To quantify this, for each breakpoint  
350 that occurred  $\geq 3$  times, we compute the entropy of its partner promiscuity. Specifically, we derive a  
351 discrete, empirical probability distribution of partners for each breakpoint and compute the entropy of this  
352 distribution. This measure thus represents the uncertainty of the partner given one breakpoint, with higher  
353 entropy corresponding to a less conserved partnering pattern. In Figure 5E, we see that there there is a high  
354 peak near 0 for all cancer types, which indicates that for a large proportion of recurring breakpoints, we are  
355 certain about its rejoined partner once we know the breakpoint. However, there are promiscuous breakpoints  
356 with entropy larger than 0.5.

357 Finally, we consider the span of distance between breakpoints of intra-chromosomal TSVs. We find that  
358 generally the two breakpoints are most likely to be separated by between  $10^5 - 10^7$  nt. The separation  
359 distance for non-fusion-gene breakpoints, tends to be on the higher end of this range ( $\approx 10^7$  nt). The full  
360 distributions of breakpoint separation for intra-chromosomal TSVs are given in Figures 5F and 5G. There  
361 are some differences among cancer types in these distributions. In BLCA, BRCA, and LUAD, breakpoints  
362 of intra-chromosomal TSVs are more likely to be separated by around  $10^5$  nt than  $10^7$  nt, while when only  
363 looking at non-fusion-gene events, number of TSVs with distance  $10^7$  nt is greater than or equal to the  
364 number of events with distance  $10^5$ . Thus for these three cancer types, non-fusion-gene TSVs occur more  
365 at longer distances while fusion-gene TSVs more at shorter distances. PRAD behaves differently, where for  
366 both overall TSV events and non-fusion-gene events, the distance is most likely to be around  $10^7$  nt rather  
367 than  $10^5$  nt.

### 368 **3.4 Tumor suppressor genes can undergo TSV and generate altered transcript**

369 Tumor suppressor genes (TSG) protect cells from becoming cancer cells. Usually their functions involve  
370 inhibiting cell cycle, facilitating apoptosis, and so on [32]. Mutations in TSGs may lead to loss of function of  
371 the corresponding proteins and benefit tumor growth. For example, homozygous loss-of-function mutation  
372 in p53 is found in about half of cancer samples across various cancer types [13]. TSVs are likely to cause  
373 loss of function of TSGs as well. Indeed, we observe several TSGs that are affected by TSVs, both of the  
374 fusion-gene type and the non-fusion-gene type.

375 The ZFH3 gene encodes a transcription factor that transactivates cyclin-dependent kinase inhibitor 1A  
376 (aka p21CIP1), a cell cycle inhibitor [19]. We find that in one BLCA and one BRCA sample, there are  
377 TSVs affecting ZFH3. These two TSVs events are different from each other in terms of the breakpoint  
378 partner outside of ZFH3. In the BLCA tumor sample, a intergenic region is inserted after the third exon of  
379 ZFH3 (Figure 6A). The fused transcript stops at the inserted region, causing the ZFH3 transcript to lose  
380 the rest of exons. In the BRCA tumor sample, a region of the anti-sense strand of gene MYLK3 is inserted  
381 after the third exon of ZFH3 gene (Figure 6B). Because codons and splicing sites are not preserved on the  
382 anti-sense strand, the transcribed insertion region does not correspond to known exons of MYLK3 gene, but  
383 covers the range of first exon of MYLK3 and extend to the first intron and 5' intergenic region. Transcription  
384 stops within inserted region, and causes the ZFH3 transcript to lose exons after exon 3, which resembles  
385 the fusion with intergenic region in BLCA sample.

386 Another example is given by the ASXL1 gene, which is essential for activating INK4B to inhibit tumor-  
387 genesis [38]. We observe two distinct TSVs related to ASXL1 from BLCA and BRCA samples. The first  
388 TSV merges the first 11 exons and half of exon 12 of ASXL1 with a intergenic region on chromosome 4  
389 (Figure 6C). Transcription stops at the inserted intergenic region, leaving the rest of exon 12 not transcribed.  
390 The breakpoint within the ASXL1 is before the 3' UTR, so the downstream protein sequence from exon 12  
391 will be affected. The other TSV involving ASXL1 is a typical fusion-gene TSV where the first three exons  
392 of ASXL1 are fused with the last three exons from the PDRG1 gene (Figure 6D). Protein domains after  
393 ASXL1 exon 4 and before PDGR1 exon 2 are lost in the fused transcript.

394 These examples are novel predicted TSV events that are not typically detectable via traditional fusion-gene  
395 detection methods using RNA-seq data. They suggest that non-fusion-gene events can also be involved in



396 tumorigenesis by causing disruption of tumor suppressor genes.

## 397 **4 Discussion**

398 We developed SQUID, the first algorithm for accurate and comprehensive TSV detection, spanning both  
399 traditional fusion-gene detection and the much broader class of general TSVs. SQUID exhibits far higher  
400 precision at similar sensitivities compared with WGS-based SV detection methods and pipeline of de novo  
401 transcriptome assembly and transcript-to-genome alignment. In addition, it has the ability to detect non-  
402 fusion-gene TSVs. These features are derived from its unique approach to predicting TSVs, whereby it  
403 constructs a consistent model of the underlying rearranged genome that explains as much of the data as pos-  
404 sible. In particular, it simultaneously considers both concordant and discordant reads, and by rearranging  
405 genome segments to maximize the number of concordant reads, SQUID generates a set of compatible TSVs  
406 that are most reliable in terms of the numbers of reads supporting them. Instead of a universal read support  
407 threshold, the objective function in SQUID naturally balances reads supporting and not supporting a candi-  
408 date TSV. This design is efficient in filtering out sequencing and alignment noise in RNA-seq, especially in  
409 the annotation-free context for predicting non-fusion-gene TSV events.

410 We use SQUID to analyze TCGA RNA-seq data of tumor samples. We identify BRCA to have more TSVs  
411 within a typical sample than the other cancer types studied. We observe that non-fusion-gene TSVs are  
412 more likely to have intra-chromosomal TSVs but the intra-chromosomal breakpoint distance tends to be  
413 larger than fusion-gene TSVs. This is likely due to the different sequence composition features in gene vs.  
414 non-gene regions. PRAD also stands out because the percentage of inter-chromosomal TSVs is the largest,  
415 and it is the least likely to have breakpoint distances less than  $10^5$ . Overall, these findings continue to  
416 suggest that different cancer types have different preferred patterns of TSVs, although the question remains  
417 whether these differences will hold up as more samples are analyzed and whether the different patterns are  
418 causal, correlated, or mostly non-functional randomness.

419 We also use SQUID to observe both non-fusion-gene and fusion-gene TSVs involving known tumor sup-  
420 pressor genes ZFH3 and ASXL1. In these cases, transcription usually stops within the inserted region  
421 of the non-fusion-gene TSVs, which causes TSG transcript to lose some of its exons, reasonably leading  
422 to downstream loss of function. These non-fusion-gene TSVs related to TSG may provide an alternative

423 reason or contributor to tumor genesis.

424 Other important uses and implications for general TSVs have yet to be explored and represent possible  
425 directions for future work. TSVs will impact accuracy of transcriptome assembly and expression quantifi-  
426 cation, and methodological advancements are needed to correct those downstream analyses for the effect  
427 of TSVs. For example, current reference-based transcriptome assemblers are not able to assemble from  
428 different chromosomes to handle the case of inter-chromosomal TSVs. In addition, TSV-affected transcripts  
429 cannot be quantified if they are not present in the transcript database. Incorporating TSVs into transcriptome  
430 assembly and expression quantification can potentially improve their accuracy. SQUID's ability to provide a  
431 new genome sequence that is as consistent as possible with the observed reads will facilitate its use as a pre-  
432 processing step for transcriptome assembly and expression quantification, though optimizing this pipeline  
433 remains a task for future work.

434 Several natural directions exist for extending SQUID. First, SQUID is not able to predict small deletions,  
435 instead, it treats the small deletions the same as intron-splitting events. This is to some extent a limitation of  
436 using RNA-seq data: introns and deletions are difficult to distinguish, as both result in concordant split reads  
437 or stretched mate pairs. The use of gene annotations can somewhat address this problem. Second, when  
438 the RNA-seq reads are derived from a highly heterogeneous sample, SQUID is likely not able to predict  
439 all TSVs occurring in the same region if they are conflicting since it seeks a single, consistent genome  
440 model. Instead, SQUID will only pick the dominating one that is compatible with other predicted TSVs.  
441 One approach to handle this would be to iteratively re-run SQUID, removing reads that are explained at each  
442 step. Again, this represents an attractive avenue for future work.

443 SQUID is open source and available at <http://www.github.com/Kingsford-Group/squid>.

444 **Acknowledgements.** We thank Jacob West-Roberts for useful discussions. This research is funded in part  
445 by the Gordon and Betty Moore Foundation's Data-Driven Discovery Initiative through Grant GBMF4554  
446 to C.K., by the US National Science Foundation (CCF-1256087, CCF-1319998) and by the US National  
447 Institutes of Health (R21HG006913, R01HG007104). C.K. received support as an Alfred P. Sloan Research  
448 Fellow. This project is funded, in part, under a grant (#4100070287) with the Pennsylvania Department  
449 of Health. The Department specifically disclaims responsibility for any analyses, interpretations or conclu-  
450 sions.

## 451 **References**

- 452 [1] Bronwen L Aken, Sarah Ayling, Daniel Barrell, Laura Clarke, Valery Curwen, Susan Fairley, Julio  
453 Fernandez Banet, Konstantinos Billis, Carlos García Girón, Thibaut Hourlier, et al. The Ensembl gene  
454 annotation system. *Database*, 2016, 2016.
- 455 [2] Vineet Bafna and Pavel A Pevzner. Genome rearrangements and sorting by reversals. *SIAM Journal*  
456 *on Computing*, 25(2):272–289, 1996.
- 457 [3] Christoph Bartenhagen and Martin Dugas. RSVSim: an R/Bioconductor package for the simulation of  
458 structural variations. *Bioinformatics*, 29(13):1679–1681, 2013.
- 459 [4] Graham R Bignell, Thomas Santarius, Jessica CM Pole, Adam P Butler, Janet Perry, Erin Pleasance,  
460 Chris Greenman, Andrew Menzies, Sheila Taylor, Sarah Edkins, et al. Architectures of somatic ge-  
461 nomic rearrangement in human cancer amplicons at sequence-level resolution. *Genome Research*, 17  
462 (9):1296–1303, 2007.
- 463 [5] Ken Chen, John W Wallis, Michael D McLellan, David E Larson, Joelle M Kalicki, Craig S Pohl,  
464 Sean D McGrath, Michael C Wendl, Qunyuan Zhang, Devin P Locke, et al. BreakDancer: an algorithm  
465 for high-resolution mapping of genomic structural variation. *Nature Methods*, 6(9):677–681, 2009.
- 466 [6] Colby Chiang, Ryan M Layer, Gregory G Faust, Michael R Lindberg, David B Rose, Erik P Garrison,  
467 Gabor T Marth, Aaron R Quinlan, and Ira M Hall. SpeedSeq: ultra-fast personal genome analysis and  
468 interpretation. *Nature Methods*, 2015.
- 469 [7] Nadia M Davidson, Ian J Majewski, and Alicia Oshlack. JAFFA: High sensitivity transcriptome-  
470 focused fusion gene detection. *Genome Medicine*, 7(1):43, 2015.
- 471 [8] Michael WN Deininger, John M Goldman, and Junia V Melo. The molecular biology of chronic  
472 myeloid leukemia. *Blood*, 96(10):3343–3356, 2000.
- 473 [9] Alexander Dobin, Carrie A Davis, Felix Schlesinger, Jorg Drenkow, Chris Zaleski, Sonali Jha, Philippe  
474 Batut, Mark Chaisson, and Thomas R Gingeras. STAR: ultrafast universal RNA-seq aligner. *Bioinfor-*  
475 *matics*, 29(1):15–21, 2013.
- 476 [10] Pedro AF Galante, Raphael B Parmigiani, Qi Zhao, Otávia L Caballero, Jorge E De Souza, Fábio CP

- 477 Navarro, Alexandra L Gerber, Marisa F Nicolás, Anna Christina M Salim, Ana Paula M Silva, et al.  
478 Distinct patterns of somatic alterations in a lymphoblastoid and a tumor genome derived from the same  
479 individual. *Nucleic Acids Research*, 39(14):6056–6068, 2011.
- 480 [11] Manfred G Grabherr, Brian J Haas, Moran Yassour, Joshua Z Levin, Dawn A Thompson, Ido Amit,  
481 Xian Adiconis, Lin Fan, Raktima Raychowdhury, Qiandong Zeng, et al. Trinity: reconstructing a  
482 full-length transcriptome without a genome from RNA-Seq data. *Nature Biotechnology*, 29(7):644,  
483 2011.
- 484 [12] Thasso Griebel, Benedikt Zacher, Paolo Ribeca, Emanuele Raineri, Vincent Lacroix, Roderic Guigó,  
485 and Michael Sammeth. Modelling and simulating generic RNA-Seq experiments with the flux simu-  
486 lator. *Nucleic Acids Research*, 40(20):10073–10083, 2012.
- 487 [13] Monica Hollstein, David Sidransky, Bert Vogelstein, and Curtis C Harris. p53 mutations in human  
488 cancers. *Science*, 253(5015):49–54, 1991.
- 489 [14] Fereydoun Hormozdiari, Iman Hajirasouliha, Phuong Dao, Faraz Hach, Deniz Yorukoglu, Can Alkan,  
490 Evan E Eichler, and S Cenk Sahinalp. Next-generation VariationHunter: combinatorial algorithms for  
491 transposon insertion discovery. *Bioinformatics*, 26(12):i350–i357, 2010.
- 492 [15] Matthew K Iyer, Arul M Chinnaiyan, and Christopher A Maher. ChimeraScan: a tool for identifying  
493 chimeric transcription in sequencing data. *Bioinformatics*, 27(20):2903–2904, 2011.
- 494 [16] Stefan Kurtz, Adam Phillippy, Arthur L Delcher, Michael Smoot, Martin Shumway, Corina Antonescu,  
495 and Steven L Salzberg. Versatile and open software for comparing large genomes. *Genome Biology*, 5  
496 (2):R12, 2004.
- 497 [17] Ryan M Layer, Colby Chiang, Aaron R Quinlan, and Ira M Hall. LUMPY: a probabilistic framework  
498 for structural variant discovery. *Genome Biology*, 15(6):1, 2014.
- 499 [18] Heng Li and Richard Durbin. Fast and accurate short read alignment with Burrows–Wheeler transform.  
500 *Bioinformatics*, 25(14):1754–1760, 2009.
- 501 [19] Donna Maglott, Jim Ostell, Kim D Pruitt, and Tatiana Tatusova. Entrez Gene: gene-centered informa-  
502 tion at NCBI. *Nucleic Acids Research*, 39(suppl 1):D52–D57, 2011.

- 503 [20] Andrew McPherson, Fereydoun Hormozdiari, Abdalnasser Zayed, Ryan Giuliany, Gavin Ha, Mark GF  
504 Sun, Malachi Griffith, Alireza Heravi Moussavi, Janine Senz, Nataliya Melnyk, et al. deFuse: an  
505 algorithm for gene fusion discovery in tumor RNA-Seq data. *PLoS Comput. Biol.*, 7(5):e1001138,  
506 2011.
- 507 [21] Fredrik Mertens, Bertil Johansson, Thoas Fioretos, and Felix Mitelman. The emerging complexity of  
508 gene fusions in cancer. *Nature Reviews Cancer*, 15(6):371–381, 2015.
- 509 [22] Cancer Genome Atlas Network et al. Comprehensive molecular portraits of human breast tumors.  
510 *Nature*, 490(7418):61, 2012.
- 511 [23] Cancer Genome Atlas Research Network et al. Comprehensive molecular characterization of urothelial  
512 bladder carcinoma. *Nature*, 507(7492):315–322, 2014.
- 513 [24] Cancer Genome Atlas Research Network et al. Comprehensive molecular profiling of lung adenocar-  
514 cinoma. *Nature*, 511(7511):543–550, 2014.
- 515 [25] Cancer Genome Atlas Research Network et al. The molecular taxonomy of primary prostate cancer.  
516 *Cell*, 163(4):1011–1025, 2015.
- 517 [26] Daniel Nicorici, Mihaela Satalan, Henrik Edgren, Sara Kangaspeska, Astrid Murumagi, Olli Kallion-  
518 iemi, Sami Virtanen, and Olavi Kilkku. FusionCatcher - a tool for finding somatic fusion genes in  
519 paired-end RNA-sequencing data. *bioRxiv*, 2014. doi: 10.1101/011650.
- 520 [27] Aaron R Quinlan, Royden A Clark, Svetlana Sokolova, Mitchell L Leibowitz, Yujun Zhang, Matthew E  
521 Hurles, Joshua C Mell, and Ira M Hall. Genome-wide mapping and assembly of structural variant  
522 breakpoints in the mouse genome. *Genome Research*, 20(5):623–635, 2010.
- 523 [28] Tobias Rausch, Thomas Zichner, Andreas Schlattl, Adrian M Stütz, Vladimir Benes, and Jan O Korbel.  
524 DELLY: structural variant discovery by integrated paired-end and split-read analysis. *Bioinformatics*,  
525 28(18):i333–i339, 2012.
- 526 [29] Gordon Robertson, Jacqueline Schein, Readman Chiu, Richard Corbett, Matthew Field, Shaun D Jack-  
527 man, Karen Mungall, Sam Lee, Hisanaga Mark Okada, Jenny Q Qian, et al. De novo assembly and  
528 analysis of RNA-seq data. *Nature Methods*, 7(11):909–912, 2010.

- 529 [30] Dan R Robinson, Shanker Kalyana-Sundaram, Yi-Mi Wu, Sunita Shankar, Xuhong Cao, Bushra Ateeq,  
530 Irfan A Asangani, Matthew Iyer, Christopher A Maher, Catherine S Grasso, et al. Functionally recur-  
531 rent rearrangements of the MAST kinase and Notch gene families in breast cancer. *Nature Medicine*,  
532 17(12):1646–1651, 2011.
- 533 [31] Andrea Sboner, Xinmeng Jasmine Mu, Dov Greenbaum, Raymond K Auerbach, and Mark B Gerstein.  
534 The real cost of sequencing: higher than you think! *Genome Biology*, 12(8):125, 2011.
- 535 [32] Charles J Sherr. Principles of tumor suppression. *Cell*, 116(2):235–246, 2004.
- 536 [33] Philip J Stephens, David J McBride, Meng-Lay Lin, Ignacio Varela, Erin D Pleasance, Jared T Simp-  
537 son, Lucy A Stebbings, Catherine Leroy, Sarah Edkins, Laura J Mudie, et al. Complex landscapes of  
538 somatic rearrangement in human breast cancer genomes. *Nature*, 462(7276):1005–1010, 2009.
- 539 [34] A Sveen, S Kilpinen, A Ruusulehto, RA Lothe, and RI Skotheim. Aberrant RNA splicing in cancer;  
540 expression changes and driver mutations of splicing factor genes. *Oncogene*, 2015.
- 541 [35] Scott A Tomlins, Daniel R Rhodes, Sven Perner, Saravana M Dhanasekaran, Rohit Mehra, Xiao-Wei  
542 Sun, Sooryanarayana Varambally, Xuhong Cao, Joelle Tchinda, Rainer Kuefer, et al. Recurrent fusion  
543 of TMPRSS2 and ETS transcription factor genes in prostate cancer. *Science*, 310(5748):644–648,  
544 2005.
- 545 [36] Jianghua Wang, Yi Cai, Wendong Yu, Chengxi Ren, David M Spencer, and Michael Ittmann.  
546 Pleiotropic biological activities of alternatively spliced TMPRSS2/ERG fusion gene transcripts. *Cancer*  
547 *Research*, 68(20):8516–8524, 2008.
- 548 [37] Thomas D Wu and Colin K Watanabe. GMAP: a genomic mapping and alignment program for mRNA  
549 and EST sequences. *Bioinformatics*, 21(9):1859–1875, 2005.
- 550 [38] Xudong Wu, Ida Holst Bekker-Jensen, Jesper Christensen, Kasper Dindler Rasmussen, Simone Sidoli,  
551 Yan Qi, Yu Kong, Xi Wang, Yajuan Cui, Zhijian Xiao, et al. Tumor suppressor ASXL1 is essential  
552 for the activation of INK4B expression in response to oncogene activity and anti-proliferative signals.  
553 *Cell Research*, 25(11):1205–1218, 2015.
- 554 [39] Lixing Yang, Lovelace J Luquette, Nils Gehlenborg, Ruibin Xi, Psalm S Haseley, Chih-Heng Hsieh,

- 555 Chengsheng Zhang, Xiaojia Ren, Alexei Protopopov, Lynda Chin, et al. Diverse mechanisms of so-  
556 matic structural variations in human cancer genomes. *Cell*, 153(4):919–929, 2013.
- 557 [40] Andrew Yates, Wasiu Akanni, M Ridwan Amode, Daniel Barrell, Konstantinos Billis, Denise  
558 Carvalho-Silva, Carla Cummins, Peter Clapham, Stephen Fitzgerald, Laurent Gil, et al. Ensembl  
559 2016. *Nucleic Acids Research*, page gkv1157, 2015.
- 560 [41] Jin Zhang, Nicole M White, Heather K Schmidt, Robert S Fulton, Chad Tomlinson, Wesley C War-  
561 ren, Richard K Wilson, and Christopher A Maher. INTEGRATE: gene fusion discovery using whole  
562 genome and transcriptome data. *Genome Research*, 26(1):108–118, 2016.
- 563 [42] Qi Zhao, Otavia L Caballero, Samuel Levy, Brian J Stevenson, Christian Iseli, Sandro J De Souza,  
564 Pedro A Galante, Dana Busam, Margaret A Leversha, Kalyani Chadalavada, et al. Transcriptome-  
565 guided characterization of genomic rearrangements in a breast cancer cell line. *Proceedings of the*  
566 *National Academy of Sciences, USA*, 106(6):1886–1891, 2009.

567 **Main Figures**

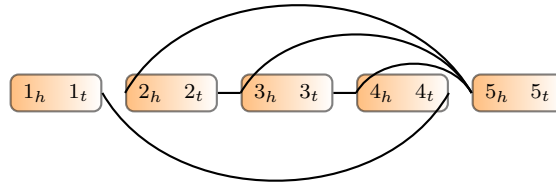


Figure 1: Example of genome segment graph. Boxes are genome segments, each of which has two ends subscripted by  $h$  and  $t$ . The color gradient indicates the orientation from head to tail. Edges connect ends of genome segments.

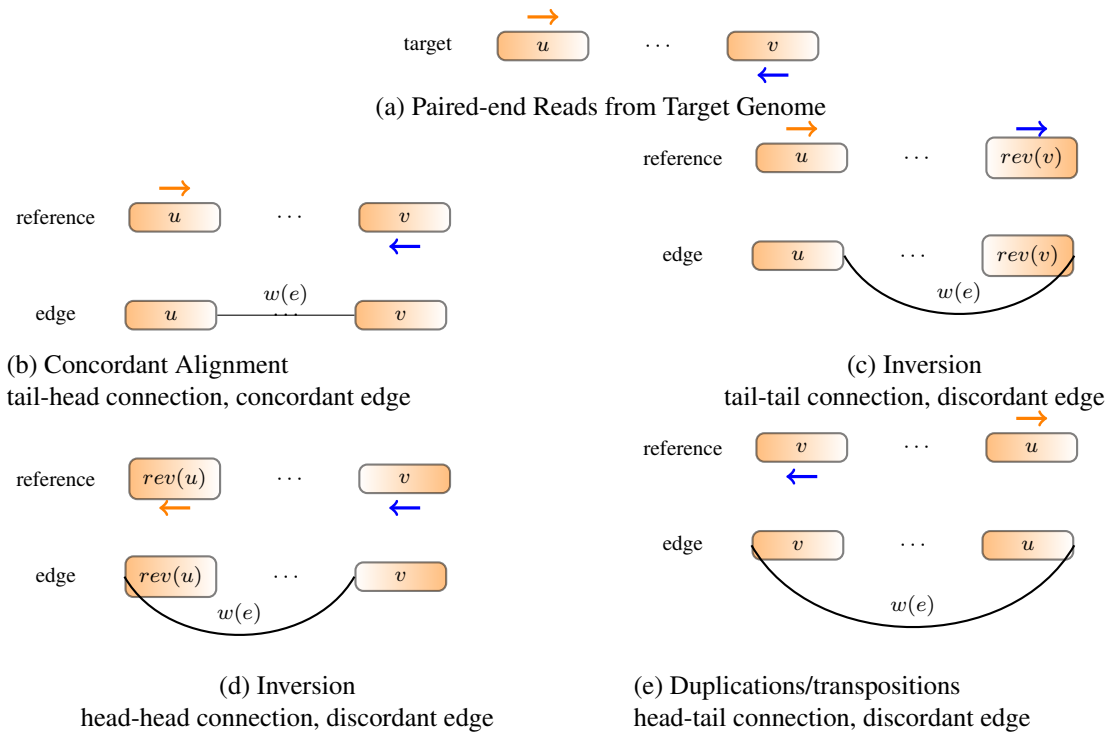


Figure 2: Constructing edges from alignment. (a) Read positions and orientations generated from the target genome. (b) If the reference genome does not have rearrangements, the read should be concordantly aligned to reference genome. An edge is added to connect the right end of  $u$  to the left end of  $v$ . Traversing the two segments along the edge reads out  $u \cdot v$ , which is the same as reference. (c) Both ends of the read align to forward strand. An edge is added to connect the right end of  $u$  to the right end of  $rev(v)$ . Traversing the segments along the edge reads out sequence  $u \cdot rev(rev(v)) = u \cdot v$ , which recovers the target sequence and the read can be concordantly aligned to. (d) If both ends align to the reverse strand, an edge is added to connect the left end of front segment to the left end of back segment. (e) If two ends of a read point out of each other, an edge is added to connect the left end of front segment to the right end of back segment.



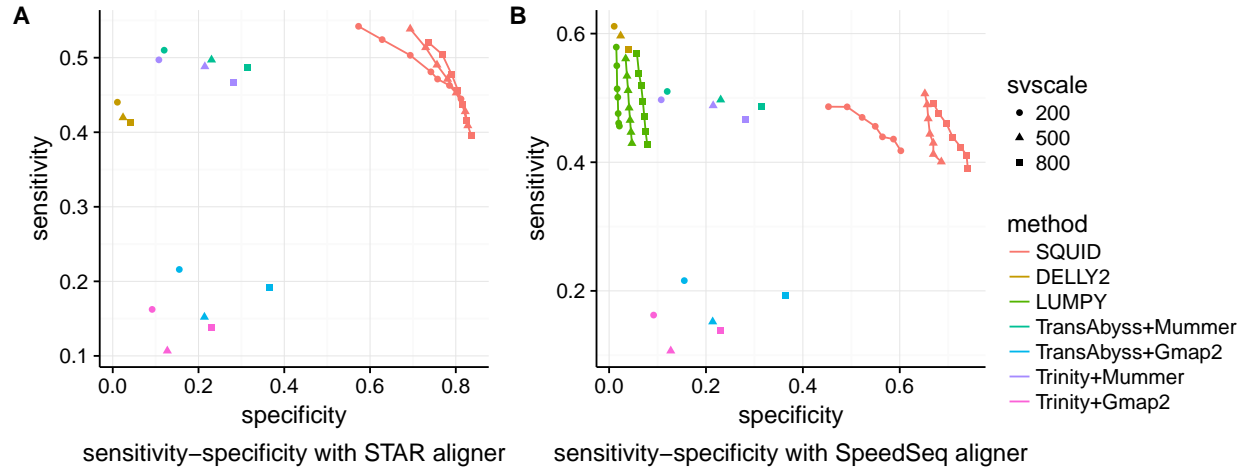


Figure 3: Performance of SQUID and other methods on simulation data. Different number of SVs (200, 500, 800 SVs) are simulated in each dataset. Each simulated read is aligned with both (A) STAR and (B) SpeedSeq aligner. If the method allows for user-defined minimum read support for prediction, we vary the threshold from 3 to 9, and plot a curve on sensitivity-specificity curve (SQUID and LUMPY), otherwise it is shown as a single point

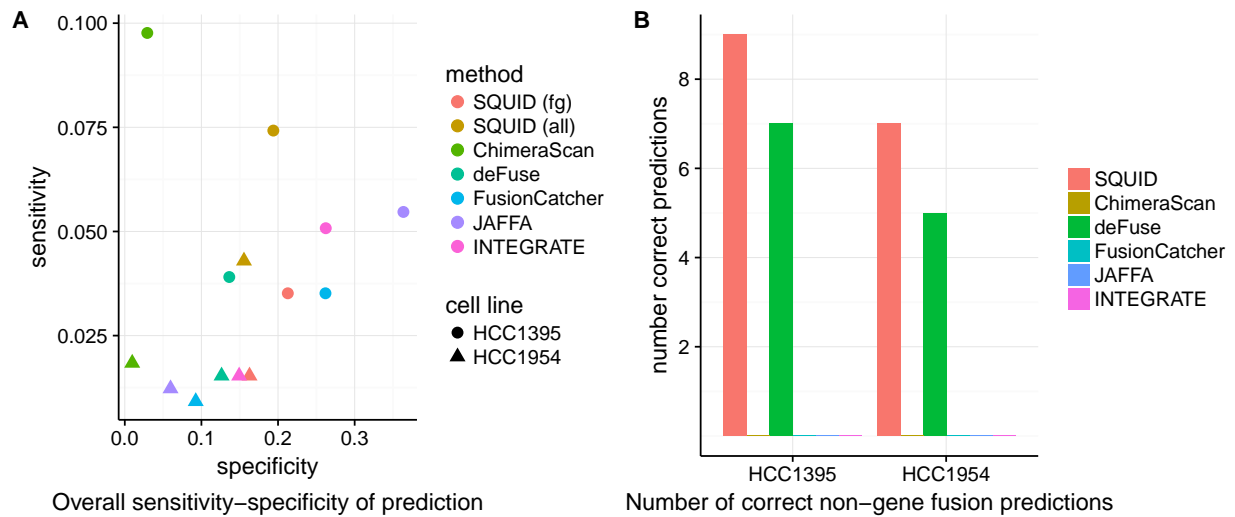


Figure 4: Performance of SQUID and fusion gene detection methods on breast cancer cell lines HCC1954 and HCC1395. Predictions are evaluated by previously validated SVs and fusions. (A) Sensitivity-specificity of different methods on both cell lines. SQUID (fg) represents the sensitivity and specificity when restricting SQUID prediction result to be fusion-gene TSVs. SQUID (all) is the performance of SQUID when considering all predictions. (B) Number of correct non-fusion-gene TSV predictions that correspond to previously validated SVs.

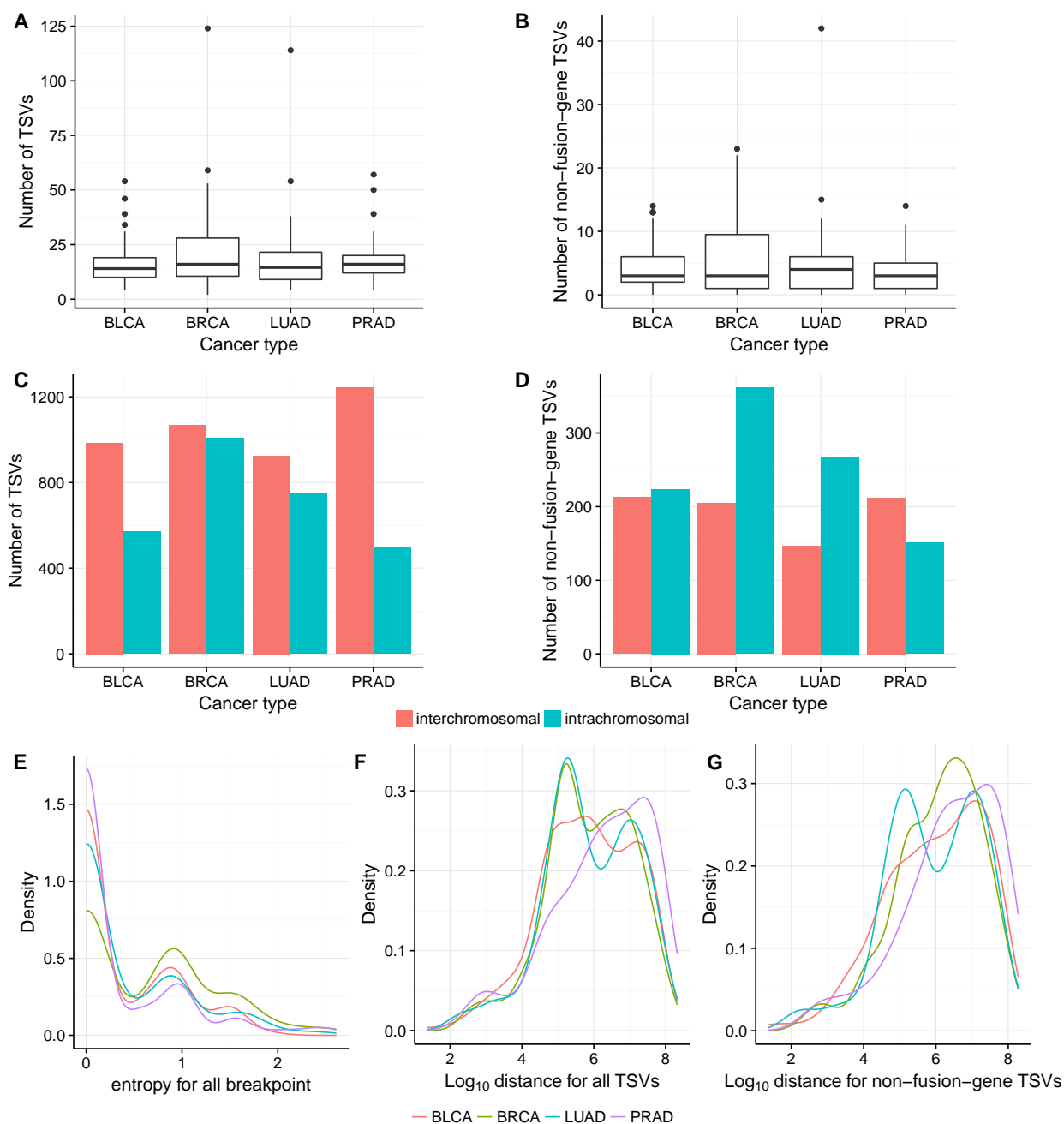


Figure 5: (A,B) Number of TSVs and non-fusion-gene TSVs in each sample in different cancer types. BRCA has slightly more samples with larger number of (non-fusion-gene) TSVs, thus showing a longer tail on y axis. (C,D) Number of inter-chromosomal and intra-chromosomal TSVs within all TSVs and within non-fusion-gene TSVs. Non-fusion-gene TSVs contain more intra-chromosomal events than fusion-gene TSVs. (E) For breakpoints occurring more than 3 times in the same cancer type, the distribution of the entropy of its TSV partner. The lower the entropy, the more likely the breakpoint has a fixed partner in TSV. The peak near 0 indicates a large portion of breakpoints are likely to be rejoined with the same partner in TSV. However, there are still some breakpoints that have multiple rejoined partners. (F,G) Distance between the pair of breakpoints in a TSV for intra-chromosomal TSV. Overall, breakpoint distance of intra-chromosomal TSVs is likely to have magnitude of  $10^5$  or  $10^7$ ; but non-fusion-gene TSVs contribute more to peak  $10^7$  than fusion-gene TSVs.

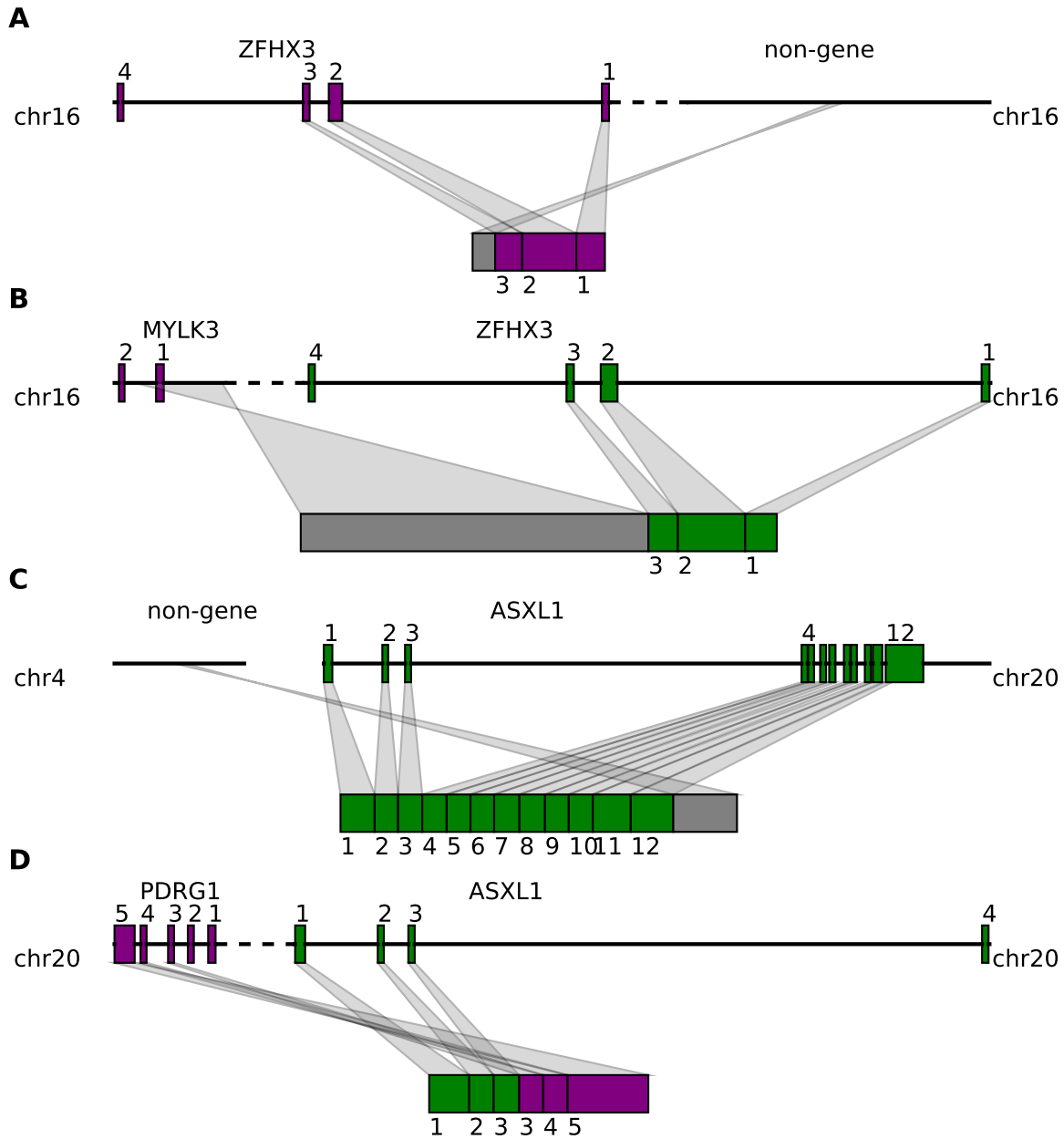


Figure 6: Tumor suppressor genes are affected by both fusion-gene and non-fusion-gene TSVs and generate transcripts with various features. A. ZFH X3 is fused with a intergenic region after exon 3. Transcript stops at the inserted region, and losing the rest of exons. B. ZFH X3 is fused with a part of MYLK3 anti-sense strand after exon 3. Codon and splicing signals are not preserved on anti-sense strand, thus MYLK3 anti-sense insertion acts the same as intergenic region insertion, and cause transcription stop before reaching the rest of ZFH X3 exons. C. ASXL1 is fused with intergenic region in the middle of exon 12. Resulting transcript contains a truncated ASXL1 exon 12 and intergenic sequence. D. First 3 exons of ASXL1 gene is joined with last 3 exons of PDRG1, resulting in a fused transcript containing 6 complete exons from both ASXL1 and PDRG1.

## 568 **Supplementary Text**

### 569 **Using de novo assembly and transcript to genome alignment to predict TSV**

570 For the pipeline of de novo transcriptome assembly and transcript-to-genome alignment, the direct output is  
571 a series of alignment pieces for each assembled transcript. To derive TSV from the pieces of alignment of  
572 each transcript, we still need to use the split-read alignment concordance criteria (8) and the edge-building  
573 approach. In the case of no TSV, equation (8) still holds, since a transcript is generated from one strand of  
574 one chromosome, without rearrangements but only deletion of introns. Any violation of (8) is treated as a  
575 TSV. Here TSVs are still able to be represented by edges in GSG, where segments are the intervals of each  
576 piece of alignment, and edges are added in the same principle that traversing segments along the edges will  
577 result in a concordant alignment of the assembled transcript. The positions of both breakpoints in a TSV are  
578 exactly the two positions linked by the discordant edge, and the orientations corresponds to the connection  
579 type of the edge.

### 580 **Processing TCGA RNA-seq data**

581 In order to be cautious about TSV prediction, we reprocess RNA-seq alignment data in the following way.  
582 We use STAR aligner to align TCGA RNA-seq reads to Ensemble genome 87 with corresponding gene  
583 annotation. STAR aligner is set with the option of outputting chimeric alignments with hanging length  
584 15bp. The chimeric alignments generated by STAR are further filtered out if the paired-end reads can be  
585 aligned concordantly by SpeedSeq aligner.

586 SQUID is applied to concordant alignment generated by STAR and filtered chimeric alignment. The dis-  
587 cordant edge weight coefficient  $\alpha$  is set to be 1, that is, we require tumor transcripts to dominate normal  
588 transcripts in order to predict corresponding TSVs. Only when reads supporting one TSV compose more  
589 than 50% of reads at the junction, can the TSV be treated as a candidate.

590 A large number of fusions between immunoglobulin genes are predicted by SQUID. However, there is  
591 possibility that B cells are in the mixture of sequencing and have very high expression of immunoglobulin  
592 genes (Ig). We cannot tell whether Ig rearrangements are generated by tumor cells or B cells. Therefore, we  
593 exclude Ig TSVs during post-processing, and exclude them from the descriptive statistics. Note that SQUID

594 does not exclude Ig TSVs internally, because Ig expression and VDJ recombination have been observed to  
595 exist in tumor cells, and revealing the role of Ig in tumor can deepen our understanding of cancer. When  
596 normal cells are removed from tumor samples, using SQUID to predict Ig TSV will help the study of Ig and  
597 tumor.

## 598 **SQUID parameters**

Table 1: Value of SQUID's parameters used in experiments

Symbol	Description	Value
$\gamma$	segment degree threshold	4
$\theta$	edge weight threshold	5
$\alpha$	discordant edge weight coefficient	8 (simulation and HCC cell line), 1 (TCGA)

599 **Supplementary Figures**

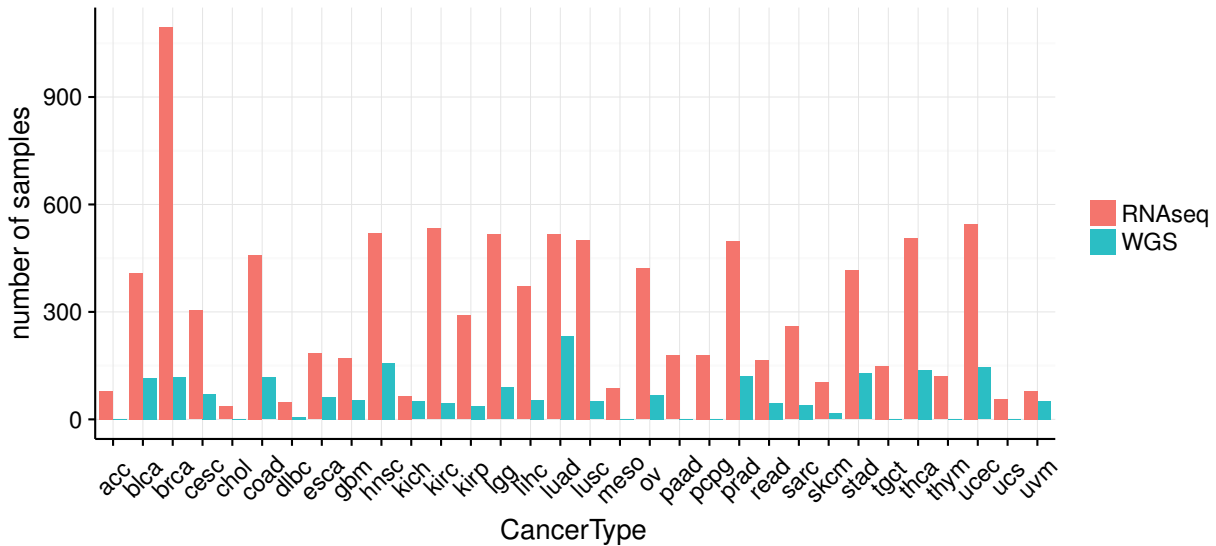


Figure S1: Number of samples with RNA-seq or WGS data in TCGA

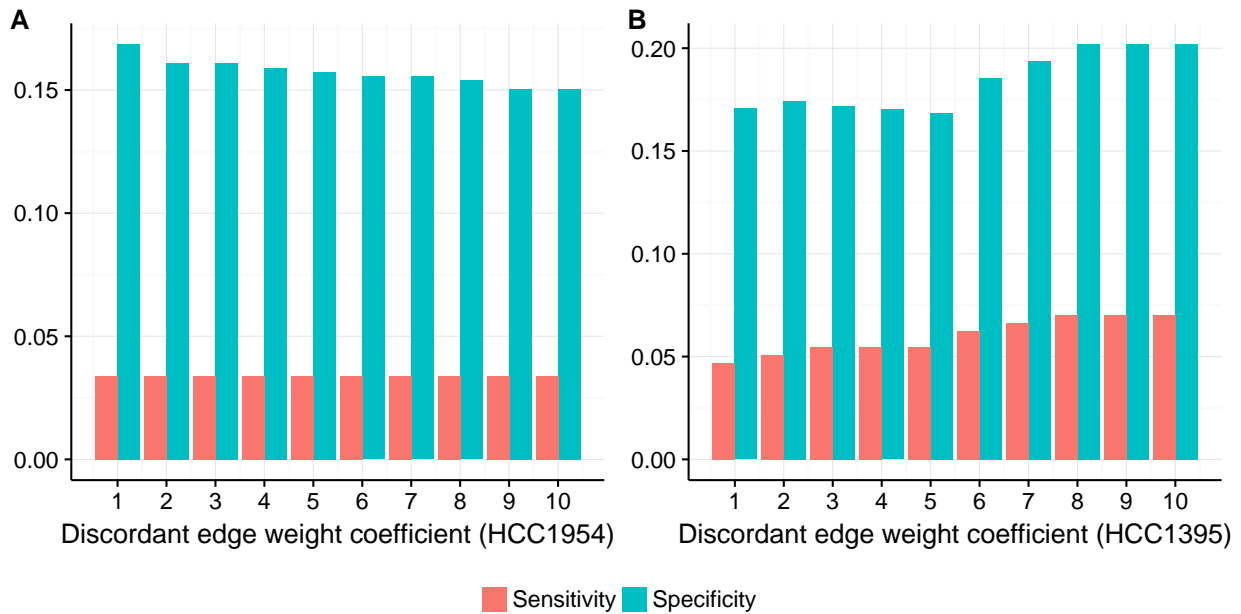


Figure S2: Specificity and sensitivity of SQUID against different value of discordant edge weight coefficient. (A) HCC1954 cell line. Sensitivity does not change when increasing discordant edge weight coefficient, indicating rearranged tumor transcripts out-number their normal counterparts. Specificity decreases slightly because SQUID predicts more as discordant edge weight coefficient increases. (B) HCC1395 cell line. Sensitivity and specificity reach the highest at discordant edge weight coefficient 8 and remain unchanged at 9 and 10. Some normal transcripts out-number the rearranged tumor transcripts, increasing this parameter allows SQUID to capture these TSVs.