

1 Detecting Transcriptomic Structural Variants in Heterogeneous 2 Contexts via the Multiple Compatible Arrangements Problem

3 Yutong Qiu^{*1}, Cong Ma^{*1}, Han Xie¹, and Carl Kingsford^{†1}

4 ¹Computational Biology Department, School of Computer Science, Carnegie Mellon University,
5 5000 Forbes Ave., Pittsburgh, PA

6 June 24, 2019

7 Abstract

8 Transcriptomic structural variants (TSVs) — structural variants that affect expressed regions — are
9 common, especially in cancer. Detecting TSVs is a challenging computational problem. Sample het-
10 erogeneity (including differences between alleles in diploid organisms) is a critical confounding factor
11 when identifying TSVs. To improve TSV detection in heterogeneous RNA-seq samples, we introduce
12 the MULTIPLE COMPATIBLE ARRANGEMENT PROBLEM (MCAP), which seeks k genome rearrange-
13 ments to maximize the number of reads that are concordant with at least one rearrangement. This directly
14 models the situation of a heterogeneous or diploid sample. We prove that MCAP is NP-hard and provide
15 a $\frac{1}{4}$ -approximation algorithm for $k = 1$ and a $\frac{3}{4}$ -approximation algorithm for the diploid case ($k = 2$)
16 assuming an oracle for $k = 1$. Combining these, we obtain a $\frac{3}{16}$ -approximation algorithm for MCAP
17 when $k = 2$ (without an oracle). We also present an integer linear programming formulation for general
18 k . We completely characterize the graph structures that require $k > 1$ to satisfy all edges and show
19 such structures are prevalent in cancer samples. We evaluate our algorithms on 381 TCGA samples and
20 2 cancer cell lines and show improved performance compared to the state-of-the-art TSV-calling tool,
21 SQUID.

22 *Keywords:* transcriptomic structural variation, integer linear programming, sample heterogeneity

*These authors contributed equally to this work

†To whom correspondence should be addressed: carlk@cs.cmu.edu

23 **1 Introduction**

24 Transcriptomic structural variations (TSVs) are genomic structural variants (SVs) that disturb the transcrip-
25 tome. TSVs may cause the joining of parts from different genes, which are fusion-gene events. Fusion genes
26 are known for their association with various types of cancer. For example, the joint protein products of *BCR-*
27 *ABL1* genes are prevalently found in leukemia [4]. In addition to fusion genes, the joining of intergenic and
28 genic regions, called non-fusion-gene events, are also related to cancer [22].

29 TSV events are best studied with RNA-seq data. Although SVs are more often studied with whole genome
30 sequencing (WGS) [2, 5, 9, 12, 18, 20], the models built on WGS data lack the flexibility to describe
31 alternative splicing and differences in expression levels of transcripts affected by TSVs. In addition, RNA-
32 seq data is far more common [14] than WGS data, for example, in The Cancer Genome Atlas (TCGA,
33 <https://cancergenome.nih.gov>).

34 Many methods have been proposed that identify fusion genes with RNA-seq data. Generally, these tools
35 identify candidates of TSV events through investigation into read alignments that are inconsistent with the
36 reference genome (e.g. [3, 10, 11, 15, 16, 21]). A series of filtering or scoring functions are applied on each
37 TSV candidate to eliminate the errors in alignment or data preparation. The performance of filters often
38 relies heavily on a large set of method parameters and requires prior annotation [13]. Furthermore, most of
39 the fusion-gene detection methods limit the scope to the joining of protein-coding regions and ignore the
40 joining of intergenic regions that could also affect the transcriptome. An approach that correctly models
41 both fusion-gene and non-fusion-gene events without a large number of ad hoc assumptions is desired.

42 An intuitive TSV model is the one that describes directly the rearrangement of the genome. For example,
43 when an inversion happens, two double-strand breaks (DSB) are introduced to the genome and the segment
44 between the DSBs is flipped. After a series of TSVs are applied to a genome, a rearranged genome is
45 produced. In order to identify the TSVs, we can attempt to infer the rearranged genome from the original
46 genome and keep track of the rearrangements of genome segments. Since a model of the complete genome is
47 produced, both fusion-gene and non-fusion-gene events can be detected. A recently published TSV detection
48 tool, SQUID [14], models TSV events in this way by determining a single rearrangement of a reference
49 genome that can explain the maximum number of observed sequencing reads. SQUID finds one arrangement

50 of genome segments such that the total number of consistent read alignments is maximized. The originally
51 discordant edges that are made consistent in the rearranged genome are output as predicted TSVs and the
52 other discordant edges are regarded as sequencing or alignment errors.

53 Despite the generally good performance of SQUID, it relies on the assumption that the sample is homo-
54 geneous, i.e. the original genome contains only one allele that can be represented by a single rearranged
55 string. This assumption is unrealistic in diploid (or high ploidy) organisms. When TSV events occur within
56 the same regions on different alleles, the set of inconsistent read alignments may appear conflicting with
57 each other if they are coerced to be on the same allele. Under the homogeneous assumption, conflicting
58 TSV candidates are regarded as errors. Therefore, this assumption leads to discarding the conflicting TSV
59 candidates that would be compatible on separate alleles and therefore limits the discovery of true TSVs.
60 Conflicting SV candidates are addressed in a few SV detection tools such as VariationHunter-CR [9]. How-
61 ever, VariationHunter-CR assumes a diploid genome, and its model is built for WGS data that lacks ability
62 to handle RNA-seq data.

63 We present an improved model of TSV events in heterogeneous contexts. We address the limitation of the
64 homogeneous assumption by extending the assumption to k alleles. We introduce the MULTIPLE COMPAT-
65 IBLE ARRANGEMENT PROBLEM (MCAP), which seeks, for a given k , an optimal set of k arrangements of
66 segments from GSG such that number of consistent read alignments is maximized, where each arrangement
67 describes the permutation of all segments and orientation of each segment. The originally discordant edges
68 that are concordant in any of the k arrangements are predicted as TSVs, and those edges are regarded as
69 errors otherwise. We show that MCAP is NP-hard. To address NP-hardness, we propose a $\frac{1}{4}$ -approximation
70 algorithm for the $k = 1$ case and a $\frac{3}{4}$ -approximation solution to the $k = 2$ case using an oracle for $k = 1$.
71 Combining these, we obtain a $\frac{3}{16}$ -approximation algorithm for MCAP when $k = 2$ (without an oracle). We
72 also present an integer linear programming (ILP) formulation that gives an optimal solution for general k .

73 We completely characterize the patterns of reads that result in conflicting TSV candidates under a single-
74 allele assumption. We show that these patterns are prevalent in both cancer cell lines and TCGA samples,
75 thereby further motivating the importance of SV detection approaches that directly model heterogeneity.

76 We apply our algorithms to 381 TCGA samples from 4 cancer types and show that many more TSVs can
77 be identified under a diploid assumption compared to a haploid assumption. We also evaluate an exact ILP

78 formulation under a diploid assumption (D-SQUID) on previously annotated cancer cell lines HCC1395
79 and HCC1954, identifying several previously known and novel TSVs. We also show that, in most of the
80 TCGA samples, the performance of the approximation algorithm is very close to optimal and the worst case
81 of $\frac{3}{16}$ -approximation is rare.

82 2 The Genome Segment Graph (GSG)

83 A Genome Segment Graph, similar to a splice graph [8], encodes relationships between genomic segments
84 and a set of reads. We say a read alignment is *consistent* with the reference genome when the orientation
85 of its two ends are consistent with the reference genome (e.g. 5'-to-3' for the forward read and the reverse
86 for the mate read in Illumina sequencing). A *segmentation* S of the genome is a partition of the genome
87 into disjoint intervals according to consistent and inconsistent paired-end alignments with respect to the
88 reference genome.

89 **Definition 1** (Genome Segment Graph). A genome segment graph is a weighted, undirected graph $G =$
90 (V, E, \mathbf{w}) derived from a segmentation S of the genome and a collection of reads. The vertex set, $V =$
91 $\{s_h \in S\} \cup \{s_t \in S\}$, includes a vertex for both endpoints, head (h) and tail (t), for each segment $s \in S$.
92 The head of a segment is the end that is closer to the 5' end of the original genome. The tail is the end that
93 is close to the 3' end. Pairs of reads that span more than one segment are represented by edges. There are
94 four types of connections: head-head, head-tail, tail-head and tail-tail. Each edge $e = (u_i, v_j) \in E$, where
95 $i, j \in \{h, t\}$, is undirected and connects endpoints of two segments. The weight ($w_e \in \mathbf{w}$) is the number of
96 read alignments that support edge e .

97 We also define the weight of a subset $E' \subseteq E$ of edges $w(E') = \sum_{e \in E'} w_e$. (More details on the GSG
98 provided in Ma et al. [14].)

99 **Definition 2** (Permutation and Orientation functions). A *permutation function* is a function where $\pi(u) = i$,
100 where i is the index of segment $u \in S$ in an ordering of a set S of segments. We also define orientation
101 function $f(u) = 1$ if segment u should remain the original orientation, or 0 if it should be inverted.

102 If $\pi(u) < \pi(v)$, we say that segment u is closer to the 5' end of the rearranged genome than segment v . We
103 call a pair of permutation and orientation functions (π, f) an *arrangement*.

104 **Definition 3** (Concordant and discordant edges). *An edge e is concordant if it connects the right end of a*
105 *segment s_i and the left end of a segment s_j with $\pi(s_i) < \pi(s_j)$. Otherwise, e is discordant.*

106 A discordant edge represents a set of inconsistent read alignments. In other words, each discordant edge is
107 a candidate TSV. Examples of tail-tail and head-head connections are shown in Figure 1a.

108 Segments connected by discordant edges can be arranged to make the edge concordant by either flipping the
109 segments or changing the ordering of the segments. For example, a head-head edge $e' = (u_h, v_h)$ can be
110 made concordant by flipping u into $-u$, or by flipping v into $-v$ and moving $-v$ to the front of u (here, the
111 negation of a segment denotes a segment flipped relative to the reference). Biologically, a flip represents an
112 inversion and a change of ordering represents an insertion or translocation.

113 **3 The Multiple Compatible Arrangements Problem (MCAP)**

114 **3.1 Problem Statement**

Given an input GSG $G = (V, E, w)$ and a constant k , the MULTIPLE COMPATIBLE ARRANGEMENTS
PROBLEM seeks a set of arrangements $A = \{(\pi_i, f_i)\}_{i=1}^k$, to optimize:

$$\max_A \sum_{e \in E} w(e) \cdot \mathbf{1}[e \sim A], \quad (1)$$

115 where $\mathbf{1}[e \sim A]$ is 1 if edge e is concordant in at least one $(\pi_i, f_i) \in A$, and 0 otherwise.

116 This objective function aims to find an optimal set of k arrangements of segments where the total number of
117 edges made concordant is maximized in the rearranged alleles. In the context of TSV calling, the objective
118 aims to find an optimal set of TSVs such that the resulting k rearranged alleles given by these TSVs have
119 the maximum number of consistent read alignments. In other words, MCAP separates the conflicting edges
120 onto k alleles as shown in an example in Figure 1.

121 When $k = 1$, the problem reduces to finding a single rearranged genome to maximize the number of
122 concordant reads, which is the problem that SQUID [14] solves. We refer to the special case when $k = 1$ as
123 SINGLE COMPATIBLE ARRANGEMENT PROBLEM (SCAP).

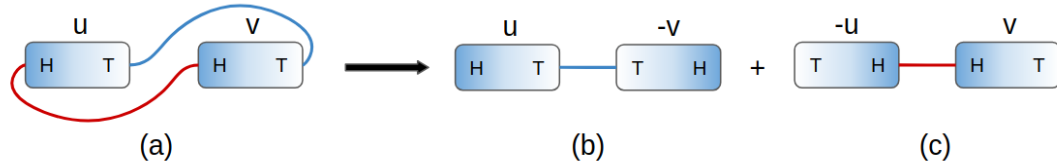


Figure 1: MCAP resolves conflicts. (a) Two conflicting edges connecting two segments u and v . If the sample is known to be homogeneous ($k = 1$), then the conflict is due to errors. If $k = 2$, MCAP seeks to separate two edges into two compatible arrangements as in (b) and (c). (b) In the first rearrangement, segment v is flipped, which makes the blue edge concordant. (c) In the second rearrangement, u is flipped to make the red edge concordant.

124 3.2 NP-hardness of SCAP and MCAP

125 **Theorem 1.** *SCAP is NP-hard.*

126 *Proof Sketch.* We prove the NP-hardness of SCAP by reducing from MAX-2-SAT. While 2-SAT can be
 127 solved in polynomial time, MAX-2-SAT, which asks for the maximum number of clauses that can be satis-
 128 fied, is NP-hard. For boolean variables and clauses in any MAX-2-SAT instance, we create gadget segments
 129 in the GSG so that the satisfaction of each clause is determined by the edge concordance and the boolean
 130 assignment is determined by segment inversion. The gadgets force the optimal sum of concordant edge
 131 weights to directly represent the number of satisfied clauses. Correspondingly, the optimal orientations of
 132 segments represent the assignment of boolean variables. See Appendix A for a complete proof. \square

133 **Corollary 1.** *MCAP is NP-hard.*

134 *Proof.* SCAP is a special case of MCAP with $k = 1$, so the NP-hardness of MCAP is immediate. \square

135 4 A $\frac{1}{4}$ -approximation Algorithm for SCAP

136 We provide a greedy algorithm for SCAP that achieves at least $\frac{1}{4}$ approximation ratio and takes $O(|V||E|)$
 137 time. The main idea of the greedy algorithm is to place each segment into the current order one by one
 138 by choosing the current “best” position. The current “best” position is determined by the concordant edge
 139 weights between the segment to be placed and the segments already in the current order.

140 **Theorem 2.** *Algorithm 1 approximates SCAP with at least $\frac{1}{4}$ approximation ratio.*

Data: Segment set S , genome segment graph $G = (V, E, w)$

Result: An arrangement of the segments and the sum of concordant edge weights

```

1 order = [];
2 orientation = [];
3 for i in 1 : |S| do
4     si = the ith segment in S;
      // choose from 4 possible order and orientation options
5     options = [(si in the beginning of order in forward strand), (si in the beginning of order in reverse
      strand), (si in the end of order in forward strand), (si in the end of order in reverse strand)];
6     for j in 1 : 4 do
7         weights[j] = w({e ∈ E : e connects si with sk and concordant in options[j], k < i});
8     end
      // update the current order and orientation
9     opt = argmax1 ≤ i ≤ 4, i ∈ ℕ weights[i];
10    order = update segment order as given by options[opt];
11    orientation = update segment orientation as given by options[opt];
12 end

```

Algorithm 1: Greedy algorithm for SCAP

141 *Proof.* Denote $E' \subset E$ as the concordant edges in the arrangement of Algorithm 1. Let OPT be the optimal
 142 value of SCAP. We are to prove $w(E') \geq \frac{1}{4}w(E) \geq \frac{1}{4}OPT$.

143 For iteration i in the for loop, the edges $E_i = \{e \in E : e \text{ connects } s^i \text{ with } s^j, i < j\}$ are considered when
 144 comparing the options. Each of the four options makes a subset of E_i concordant. These subsets are non-
 145 overlapping and their union is E_i . Specifically, the concordant edge subset is $\{e = (s_h^i, s_t^j)\}$ for the first
 146 option, $\{e = (s_h^i, s_h^j)\}$ for the second, $\{e = (s_t^i, s_h^j)\}$ for the third, and $\{e = (s_t^i, s_t^j)\}$ for the last.

By the selecting the option with the largest sum of concordant edge weights, the concordant edges E'_i in
 iteration i satisfies $w(E'_i) \geq \frac{1}{4}w(E_i)$. Therefore, the overall concordant edge weights of all iterations in the
 for loop satisfy

$$\sum_i w(E'_i) \geq \frac{1}{4} \sum_i w(E_i) = \frac{1}{4}w\left(\bigcup_i E_i\right).$$

147 Each edge $e \in E$ must appear in one and only one of E_i , and thus $\bigcup_i E_i = E$. This implies $\sum_i w(E'_i) \geq$
 148 $\frac{1}{4}w(E) \geq \frac{1}{4}OPT$. □

149 Algorithm 1 can be further improved in practice by considering more order and orientation options when in-
 150 serting a segment into current order. In the pseudo-code 1, only two possible insertion places are considered:
 151 the beginning and the end of the current order. However, a new segment can be inserted in between any pair

152 of adjacent segments in the current order. We provide an extended greedy algorithm to take into account the
 153 extra possible inserting positions (Algorithm 2). Algorithm 2 has a time complexity of $O(|V|^2|E|)$, but it
 154 may achieve a higher total concordant edge weight in practice.

<p>Data: Segment set S, genome segment graph $G = (V, E, w)$ Result: An arrangement of the segments and the sum of concordant edge weights</p> <pre> 1 order = []; 2 orientation = []; 3 for i in 1 : S do 4 s^i = the i^{th} segment in S; // choose from $i+1$ possible order and orientation options 5 options = [(s^i in the beginning of order in forward strand), (s^i in the beginning of order in reverse strand)]; 6 for j in 1 : $i - 1$ do 7 Append [(s^i right after order[j] in forward strand), (s^i right after order[j] in reverse strand)] to list of options ; 8 end 9 for j in 1 : $2i$ do 10 weights[j] = $w(\{e \in E : e \text{ connects } s^i \text{ with } s^k \text{ and concordant in options[j], } k < i\})$; 11 end // update the current order and orientation 12 opt = $\operatorname{argmax}_{1 \leq i \leq 2i, i \in \mathbb{N}} \text{weights}[i]$; 13 order = update segment order as given by options[opt] ; 14 orientation = update segment orientation as given by options[opt] ; 15 end </pre>
--

Algorithm 2: Extended greedy algorithm for SCAP

155 5 A $\frac{3}{4}$ -approximation of MCAP with $k = 2$ Using a SCAP Oracle

156 If an optimal SCAP solution can be computed, one way to approximate the MCAP's optimal solution is to
 157 solve a series of SCAP instances iteratively to obtain multiple arrangements. Here, we prove the iterative
 158 SCAP solution has an approximation ratio of $\frac{3}{4}$ for the special case of MCAP with $k = 2$.

Theorem 3. *Algorithm 3 is a $\frac{3}{4}$ -approximation of MCAP with $k = 2$. Denote the optimal objective sum of edge weights in MCAP with $k = 2$ as OPT , and the sum of edge weights in iterative SCAP as W , then*

$$W \geq \frac{3}{4}OPT$$

159 *Proof.* Denote MCAP with $k = 2$ as 2-MCAP. Let E_1^d and E_2^d be concordant edges in the optimal two

Data: A genome segment graph $G = (V, E, w)$

Result: a set of two arrangements, sum of weights of edges that are concordant in either arrangement

- 1 $a_1 =$ optimal SCAP arrangement on G ;
- 2 $E' = \{e \in E : e \text{ is discordant in } a_1\}$;
- 3 $G' = (V, E', w)$;
- 4 $a_2 =$ optimal SCAP arrangement on G' ;
- 5 $\tilde{E} = \{e \in E : e \sim A, A = \{a_1, a_2\}\}$;
- 6 $W = \sum_{e \in \tilde{E}} w(e)$;
- 7 **return** $(\{a_1, a_2\}, W)$;

Algorithm 3: $\frac{3}{4}$ -approximation for MCAP with $k = 2$

160 arrangements of 2-MCAP. It is always possible to make the concordant edges of the arrangements disjoint
 161 by removing the intersection from one of the concordant edge set, that is $E_1^d \cap E_2^d = \emptyset$. Let $E^d = E_1^d \cup E_2^d$.
 162 The optimal value is $w(E^d)$.

163 Denote the optimal set of concordant edges in the first round of Algorithm 3 as E_1^s . The optimal value of
 164 SCAP is $w(E_1^s)$. E_1^s can have overlap with the two concordant edge sets of the 2-MCAP optimal solution.
 165 Let the intersections be $I_1 = E_1^d \cap E_1^s$ and $I_2 = E_2^d \cap E_1^s$. Let the unique concordant edges be $D_1 = E_1^d - E_1^s$,
 166 $D_2 = E_2^d - E_1^s$ and $S = E_1^s - E_1^d - E_2^d$.

167 After separating the concordant edges in 2-MCAP into the intersections and unique sets, the optimal value
 168 of 2-MCAP can be written as $w(E^d) = w(I_1) + w(I_2) + w(D_1) + w(D_2)$, where the four subsets are
 169 disjoint. Therefore the smallest weight among the four subsets must be no greater than $\frac{1}{4}w(E^d)$. We prove
 170 the approximation ratio under the following two cases and discuss the weight of the second round of SCAP
 171 separately:

172 *Case (1): the weight of either D_1 or D_2 is smaller than $\frac{1}{4}w(E^d)$.* Because the two arrangements in 2-
 173 MCAP are interchangeable, we only prove for the case where $w(D_1) \leq \frac{1}{4}w(E^d)$. A valid arrangement
 174 of the second round of SCAP is the second arrangement in 2-MCAP, though it may not be optimal. The
 175 maximum concordant edge weights added by the second round of SCAP must be no smaller than $w(D_2)$.
 176 Combining the optimal values of two rounds of SCAP, the concordant edge weight is

$$W \geq w(E_1^s) + w(D_2) = w(S) + w(I_1) + w(I_2) + w(D_2) \geq w(E^d) - w(D_1) \geq \frac{3}{4}w(E^d). \quad (2)$$

177 *Case (2): both $w(D_1) \geq \frac{1}{4}w(E^d)$ and $w(D_2) \geq \frac{1}{4}w(E^d)$.* The subset with smallest sum of edge weights

178 is now either I_1 or I_2 . Without loss of generality, we assume I_1 has the smallest sum of edge weights and
 179 $w(I_1) \leq \frac{1}{4}w(E^d)$. Because the first round SCAP is optimal for the SCAP problem, its objective value
 180 should be no smaller than the concordant edge weights of either arrangement in 2-MCAP. Thus

$$w(E_1^s) \geq w(E_2^d) = w(D_2) + w(I_2). \quad (3)$$

181 A valid arrangement for the second round of SCAP can be either of the arrangements in 2-MCAP optimal
 182 solution. Picking the first arrangement of 2-MCAP as the possible (but not necessarily optimal) arrangement
 183 for the second round of SCAP, the concordant edge weights added by the second round of SCAP must be no
 184 smaller than $w(D_1)$. Therefore, the total sum of concordant edge weights of the optimal solutions of both
 185 rounds of SCAP is

$$W \geq w(E_1^s) + w(D_1) \geq w(D_2) + w(I_2) + w(D_1) = w(E^d) - w(I_1) \geq \frac{3}{4}w(E^d). \quad (4)$$

186

□

187 **Corollary 2.** *An approximation algorithm for MCAP with $k = 2$ can be created by using Algorithm 1 as*
 188 *the oracle for SCAP in Algorithm 3. This approximation algorithm runs in $O(|V||E|)$ time and achieves at*
 189 *least $\frac{3}{16}$ approximation ratio.*

190 The proof of the corollary is similar to the proof of iterative SCAP approximation ratio. By adding a
 191 multiplier of $\frac{1}{4}$ to the right of inequalities (3) and (4), the $\frac{3}{16}$ approximation ratio can be derived accordingly.

192 6 Integer Linear Programming Formulation for MCAP

193 MCAP, for general k , can be formulated as an integer linear programming (ILP) to obtain an optimal solu-
 194 tion. We rewrite the i^{th} permutation (π_i) , orientation (f_i) and decision $(\mathbf{1}[e \sim (\pi_i, f_i)])$ functions with three
 195 boolean variables y_e^i , z_e^i and x_e^i . For $i \in \{1, 2, \dots, k\}$ and $e \in E$, we have:

- 196 • $x_e^i = 1$ if edge $e \sim (\pi_i, f_i)$ and 0 otherwise.
- 197 • $y_u^i = 1$ if $f_i(u) = 1$ for segment u and 0 if $f_i(u) = 0$.

- 198 • $z_{uv}^i = 1$ if $\pi_i(u) < \pi_i(v)$, or segment u is in front of v in rearrangement i and 0 otherwise.

In order to account for the edges that are concordant in more than one arrangements in the summation in Equation 1, we define q_e such that $q_e = 1$ if edge e is concordant in one of the k arrangements and 0 if otherwise. The constraints for q_e are as follows:

$$q_e \leq \sum_i^k x_e^i \quad (5)$$

$$q_e \leq 1 \quad (6)$$

The objective function becomes

$$\max_{x_e^i, y_u^i, z_{uv}^i} \sum_{e \in E} w(e) \cdot q_e \quad (7)$$

199 We then add ordering and orientation constraints. If an edge is a tail-head connection, i.e. concordant to the
 200 reference genome, $x_e^i = 1$ if and only if $z_{uv}^i = y_u^i = y_v^i$. If an edge is a tail-tail connection, $x_e^i = 1$ if and
 201 only if $z_{uv}^i = -y_v^i = y_u^i$. If an edge is a head-tail connection, $x_e^i = 1$ if and only if $z_{uv}^i = -y_u^i = -y_v^i$. If
 202 an edge is a head-head connection, $x_e^i = 1$ if and only if $z_{uv}^i = -y_u^i = y_v^i$. The constraints for a tail-head
 203 connection are listed below in Equation 8, which enforce the assignment of boolean variables y_e^i , z_e^i and x_e^i .

$$\begin{aligned} x_e^i &\leq y_u^i - y_v^i + 1, \\ x_e^i &\leq y_v^i - y_u^i + 1, \\ x_e^i &\leq y_u^i - z_{uv}^i + 1, \\ x_e^i &\leq z_{uv}^i - y_u^i + 1, \end{aligned} \quad (8)$$

204 The constraints of other types of connections are similar and detailed in Ma et al. [14]. Additionally, con-
 205 straints are added so that all segments are put into a total order within each allele. For two segments u, v ,
 206 segment u will be either in front of or behind segment v , i.e. $z_{uv}^i + z_{vu}^i = 1$. For three segments u, v, w , if
 207 u is in front of v and v is in front of w , then u has to be in front of w : $1 \leq z_{uv}^i + z_{vw}^i + z_{wu}^i \leq 2$.

208 The total number of constraints as a function of k is $4k|E| + k \binom{|V|}{3} + 2|E| = O(k(|E| + V^3))$. When k
 209 increases, the number of constraints grows linearly. When $k = 1$, the ILP formulation reduces to the same

210 formulation as SQUID.

211 7 Characterizing the Conflict Structures That Imply Heterogeneity

212 In this section, we ignore edge weights and characterize the graph structures where homogeneous assump-
 213 tion cannot explain all edges. We add a set of segment edges, \hat{E} , to the GSG. Each $\hat{e} \in \hat{E}$ connects the two
 214 endpoints of each segment, i.e. $\hat{e} = \{s_h, s_t\}$ for $s \in S$. The representation of GSG becomes $G = (E, \hat{E}, V)$.

215 **Definition 4** (Conflict and Compatible Structures). *A conflict structure, $CS = (E', \hat{E}', V')$, is a subgraph
 216 of a GSG where there exists a set of edges that cannot be made concordant using any single arrangement.
 217 A compatible structure is a subgraph of a GSG where there exists a single arrangement such that all edges
 218 can be made concordant in it.*

219 **Definition 5** (Simple cycle in GSG). *A simple cycle, $C = (E', \hat{E}', \{v_0, \dots, v_{n-1}\})$, is a subgraph of a
 220 GSG, such that $E' \subseteq E, \hat{E}' \subseteq \hat{E}$ and $v_i \in V$, with $(v_i, v_{i+1 \bmod n}) \in E' \cup \hat{E}'$ and where $v_i \neq v_j$ when
 221 $i \neq j$ except $v_{n-1} = v_0$.*

222 **Definition 6** (Degree and special degree of a vertex in subgraphs of GSG). *Given a subgraph of GSG,
 223 $G' = (E', \hat{E}', V')$, $deg_{E'}(v)$ refers to the degree of vertex $v \in V'$ that counts only the edges $e \in E'$ that
 224 connect to v . $deg(v)$ refers to the number of edges $e \in E' \cup \hat{E}'$ that connect to v .*

225 **Theorem 4.** *Any acyclic subgraph of GSG is a compatible structure.*

226 **Theorem 5.** *A simple cycle $C = (E', \hat{E}', V')$ is a compatible structure if and only if there are exactly two
 227 vertices, v_j and v_i such that $deg_{E'}(v_i) = deg_{E'}(v_j) = 2$.*

228 The details of the proof of the above two theorems are in Appendix B.

229 **Corollary 3.** *A necessary condition for a subgraph (E', \hat{E}', V') to be a conflict structure is that it contains
 230 cycles. A sufficient condition for a subgraph (E', \hat{E}', V') to be a conflict structure is that it contains a
 231 simple cycle which is not a compatible structure. That is, there exists a simple cycle (E^*, \hat{E}^*, V^*) , such that
 232 $E^* \subset E', \hat{E}^* \subset \hat{E}', V^* \subset V'$ and $|\{v : deg_{E^*}(v) = 2\}| \neq 2$.*

233 The corollary is a direct derivation of theorem 4 and theorem 5 when considering general graph structures.

234 In practice, we determine if a discordant edge, $e = (u, v)$, is involved in a conflict structure by enumerating
235 all simple acyclic paths using a modified depth-first search implemented in Networkx [7, 19] between u and
236 v omitting edge e . We add e to each path and form a simple cycle. If the simple cycle satisfies Corollary
237 3, we stop path enumeration and label the e as discordant edge involved in conflict structure. If the running
238 time of path enumeration exceeds 0.5 seconds, we shuffle the order of DFS and repeat enumeration. If path
239 enumeration for e exceeds 1000 reruns, we label e as undecided.

240 **8 Experimental Results**

241 To produce an efficient, practical algorithm for TSV detection in diploid organisms, we use the following
242 approach, which we denote as D-SQUID: Run the ILP (Section 6) under the diploid assumption by setting
243 $k = 2$ on every connected component of GSG separately. If the ILP finishes or the running time of the ILP
244 exceeds one hour, output the current arrangements.

245 **8.1 D-SQUID Identifies More TSVs in TCGA Samples than SQUID**

246 We calculate the fraction of discordant edges involved in conflict structures (Figure 2a) in 381 TCGA sam-
247 ples from four types of cancers: bladder urothelial carcinoma (BLCA), breast invasive carcinoma (BRCA),
248 lung adenocarcinoma (LUAD) and prostate adenocarcinoma (PRAD). Among all samples, we found less
249 than 0.5% undecided edges out of all discordant edges. The distribution of fraction of discordant edges
250 within conflict structures are different among cancer types. The more discordant edges are involved in con-
251 flict structures, the more heterogeneous the sample is. Among four cancer types, PRAD samples exhibit
252 the highest extent of heterogeneity and BRCA samples exhibit the lowest. On average, more than 90% of
253 discordant edges are within conflict structures in all samples across four cancer types. This suggests that
254 TCGA samples are usually heterogeneous and may be partially explained by the fact that TCGA samples
255 are usually a mixture of tumor cells and normal cells [1].

256 We compare the number of TSVs found by D-SQUID and SQUID (Figure 2b). In all of our results, all of
257 the TSVs found by SQUID belong to a subset of TSVs found by D-SQUID. D-SQUID identifies many more
258 TSVs than SQUID on all four types of cancers.

259 A discordant edge is termed resolved if it is made concordant in one of the arrangements. Among all

260 discordant edges in all samples, D-SQUID is able to resolve most of them (Figure 2c), while SQUID is only
261 able to resolve fewer than 50% of them. The results demonstrate that D-SQUID is more capable of resolving
262 conflict structures in heterogeneous contexts, such as cancer samples, than SQUID.

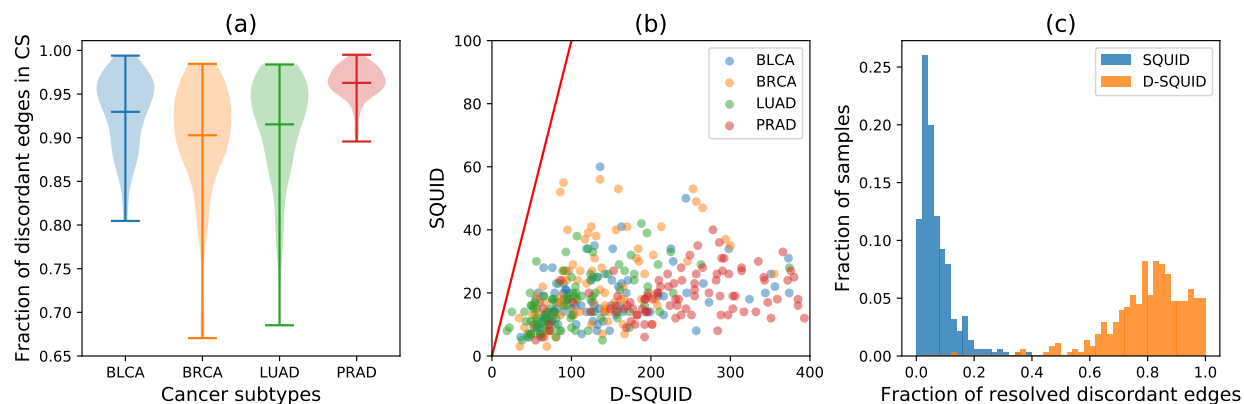


Figure 2: (a) The distribution of fractions of discordant edges that are involved in each identified conflict structure (CS) in four cancer subtypes. Minima, maxima and means of the distributions are marked by horizontal bars. (b) Number of TSVs identified by SQUID versus D-SQUID. (c) Histogram of fractions of resolved discordant edges by SQUID and D-SQUID.

263 8.2 D-SQUID Identifies More True TSV Events Than SQUID in Cancer Cell Lines

264 We compare the ability of D-SQUID and SQUID to detect fusion-gene and non-fusion-gene events on
265 previously studied breast cancer cell lines HCC1395 and HCC1954 [6]. The annotation of true SVs is taken
266 from Ma et al. [14]. In both cell lines, D-SQUID discovers more TSVs than SQUID. In HCC1954, D-
267 SQUID identifies the same number of known TSVs including fusions of gene (G) regions and intergenic
268 (IG) regions compared with SQUID. In HCC1395, D-SQUID identifies 2 more true TSV events that are
269 fusions of genic regions. We tally the fraction of discordant edges in conflict structures (Figure 3c) and
270 find similar fractions between HCC1395 and HCC1954, which indicates that the extent of heterogeneity in
271 two samples are similar. Compared to Figure 2a, the fraction in HCC samples is much lower than that in
272 TCGA samples. This matches the fact that two HCC samples contain the same cell type and are both cell
273 line samples, which are known to be less heterogeneous than TCGA samples.

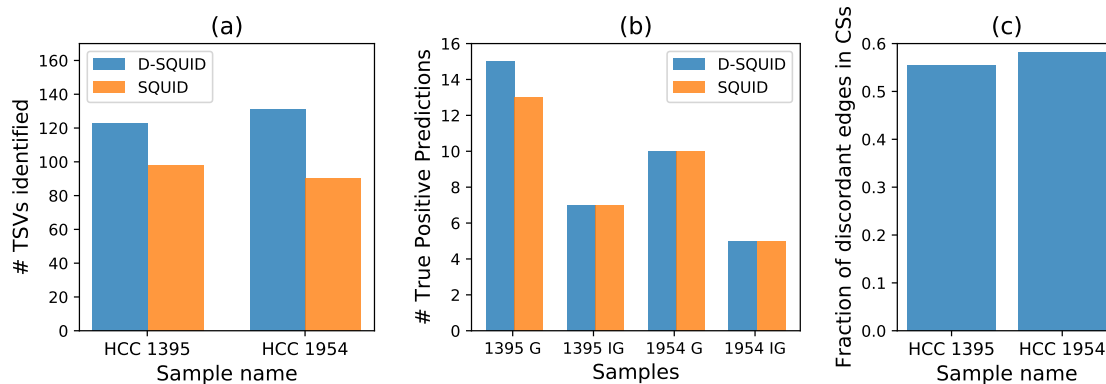


Figure 3: Performance of D-SQUID and SQUID on breast cancer cell lines with experimentally verified SV. (a) Total TSVs found. In both cell line samples, D-SQUID discovered more TSVs than SQUID. (b) Number of known fusion-gene and non-fusion-gene events recovered by D-SQUID and SQUID. G denotes TSVs that affect gene regions. IG denotes TSVs that affect intergenic regions. (c) Fraction of discordant edges in conflict structures.

274 8.3 Evaluation of approximation algorithms

275 We evaluate the approximation algorithms for diploid MCAP ($k = 2$) using two different subroutines de-
 276 scribed in Section 4. In this subsection, $A1$ refers to using Algorithm 1 with worst case runtime $O(|V||E|)$
 277 as subroutine and $A2$ refers to using Algorithm 2 with worst case runtime $O(|V|^2|E|)$ as subroutine. Both
 278 $A1$ and $A2$ solve SCAP by greedily inserting segments into the best position in the current ordering. While
 279 $A1$ only looks at the beginning and ending of the ordering, $A2$ looks at all the positions.

280 In order to compare the performance of approximations to the exact algorithm using ILP, we run D-SQUID,
 281 $A1$ and $A2$ on TCGA samples in Section 8.1. The algorithms are evaluated on runtime and total weight of
 282 concordant edges in the rearranged genomes. “Fold difference” on the axes of Figure 4 refers to the ratio of
 283 the axis values of D-SQUID over that of $A1$ or $A2$. Both $A1$ and $A2$ output results in a much shorter period
 284 of time than D-SQUID. $A2$ achieves better approximation than $A1$, demonstrated by closer-to-one ratio of
 285 total concordant edge weight, at a cost of longer run time.

286 The run time of D-SQUID ILP exceeds one hour on 4.5% of all connected components in all TCGA sam-
 287 ples. D-SQUID outputs sub-optimal arrangements in such cases. As a result, approximation algorithms,
 288 especially $A2$, appear to resolve more high-weight discordant edges than D-SQUID in some of the sam-
 289 ples in Figure 4, which is demonstrated by data points that fall below 1 on the y axes. $A1$ resolves more
 290 high-weight edges in 10 samples and $A2$ resolves more high-weight edges in 54 samples than D-SQUID.

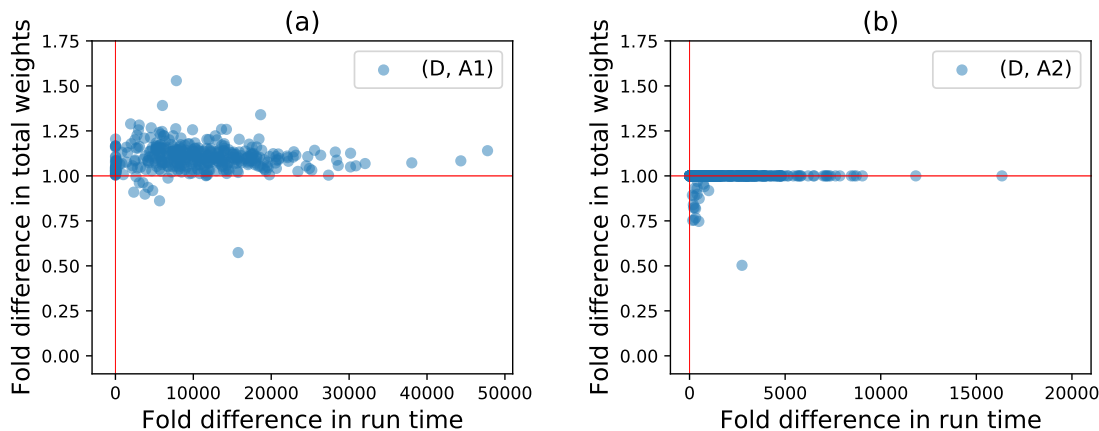


Figure 4: Fold differences (ILP/approx) in run time and total weights of concordant edges resolved by D-SQUID, A1 and A2 on TCGA samples. Horizontal and vertical red lines mark 1.0 on both axes. (a) shows fold differences between D-SQUID and A1. (b) shows fold differences between D-SQUID and A2.

291 9 Conclusion and Discussion

292 We present approaches to identify TSVs in heterogeneous samples via MULTIPLE COMPATIBLE AR-
293 RANGEMENT PROBLEM (MCAP). We characterize sample heterogeneity in terms of the fraction of dis-
294 cordant edges involved in conflict structures. In the majority of TCGA samples, the fractions of discordant
295 edges in conflict structures are high compared to HCC samples, which indicates that TCGA samples are
296 more heterogeneous than HCC samples. This matches the fact that bulk tumor samples often contain more
297 heterogeneous genomes than cancer cell lines, which suggests that fraction of conflicting discordant edges
298 is a valid measure of sample heterogeneity.

299 MCAP addresses this heterogeneity. In 381 TCGA samples, D-SQUID is able to resolve more conflicting
300 discordant edges than SQUID. In HCC cell lines, D-SQUID achieves better performance than SQUID. Since
301 D-SQUID solves MCAP by separating conflicting TSVs onto two alleles, D-SQUID's power to find TSVs
302 generally increases as the extent of heterogeneity increases.

303 We show that obtaining exact solutions to MCAP is NP-hard. We derive an integer linear programming
304 (ILP) formulation to solve MCAP exactly, of which the run time grows especially in more heterogeneous
305 samples. We provide a $\frac{3}{16}$ -approximation algorithm for MCAP when the number of arrangements is two
306 ($k = 2$), which runs in time $O(|V||E|)$. It approximates the exact solutions well in heterogeneous TCGA
307 samples.

308 Several open problems remain. MCAP relies on the number of arrangements (k) to make predictions. It
309 is not trivial to determine the optimal k for any sample. In addition, although MCAP is solved by separat-
310 ing TSVs onto different alleles, there are typically many equivalent phasings. Developing techniques for
311 handling these alternative phasings is an interesting direction for future work.

312 *Funding:* This work was supported in part by the Gordon and Betty Moore Foundations Data-Driven Discov-
313 ery Initiative [GBMF4554 to C.K.]; the US National Institutes of Health [R01GM122935]; and The Shurl
314 and Kay Curci Foundation. This project is funded, in part, by a grant (4100070287) from the Pennsylvania
315 Department of Health. The department specifically disclaims responsibility for any analyses, interpretations,
316 or conclusions.

317 *Acknowledgements:* The results shown here are in part based upon data generated by the TCGA Research
318 Network: <https://www.cancer.gov/tcga>. This work used the Extreme Science and Engineering Discovery
319 Environment (XSEDE), which is supported by National Science Foundation grant number ACI-1548562.
320 Specifically, it used the Bridges system, which is supported by NSF award number ACI-1445606, at the
321 Pittsburgh Supercomputing Center (PSC) [17]. C.K. is co-founder of Ocean Genomics, Inc.

322 **References**

- 323 [1] Dvir Aran, Marina Sirota, and Atul J Butte. Systematic pan-cancer analysis of tumour purity. *Nature*
324 *Communications*, 6:8971, 2015.
- 325 [2] Ken Chen, John W Wallis, Michael D McLellan, David E Larson, Joelle M Kalicki, Craig S Pohl,
326 Sean D McGrath, Michael C Wendl, Qunyuan Zhang, Devin P Locke, et al. BreakDancer: an algorithm
327 for high-resolution mapping of genomic structural variation. *Nature Methods*, 6(9):677, 2009.
- 328 [3] Nadia M Davidson, Ian J Majewski, and Alicia Oshlack. Jaffa: High sensitivity transcriptome-focused
329 fusion gene detection. *Genome Medicine*, 7(1):43, 2015.
- 330 [4] Michael WN Deininger, John M Goldman, and Junia V Melo. The molecular biology of chronic
331 myeloid leukemia. *Blood*, 96(10):3343–3356, 2000.
- 332 [5] Jesse R Dixon, Jie Xu, Vishnu Dileep, Ye Zhan, Fan Song, et al. Integrative detection and analysis of
333 structural variation in cancer genomes. *Nature Genetics*, 50(10):1388, 2018.

- 334 [6] Adi F Gazdar, Venkatesh Kurvari, Arvind Virmani, Lauren Gollahon, Masahiro Sakaguchi, et al. Char-
335 acterization of paired tumor and non-tumor cell lines established from patients with breast cancer.
336 *International Journal of Cancer*, 78(6):766–774, 1998.
- 337 [7] Aric Hagberg, Pieter Swart, and Daniel S Chult. Exploring network structure, dynamics, and function
338 using NetworkX. Technical report, Los Alamos National Lab.(LANL), Los Alamos, NM (United
339 States), 2008.
- 340 [8] Steffen Heber, Max Alekseyev, Sing-Hoi Sze, Haixu Tang, and Pavel A Pevzner. Splicing graphs and
341 EST assembly problem. *Bioinformatics*, 18(suppl_1):S181–S188, 2002.
- 342 [9] Fereydoun Hormozdiari, Iman Hajirasouliha, Phuong Dao, Faraz Hach, Deniz Yorukoglu, Can Alkan,
343 Evan E Eichler, and S Cenk Sahinalp. Next-generation variationhunter: combinatorial algorithms for
344 transposon insertion discovery. *Bioinformatics*, 26(12):i350–i357, 2010.
- 345 [10] Zhiqin Huang, David TW Jones, Yonghe Wu, Peter Lichter, and Marc Zapatka. confFuse: high-
346 confidence fusion gene detection across tumor entities. *Frontiers in Genetics*, 8:137, 2017.
- 347 [11] Wenlong Jia, Kunlong Qiu, Minghui He, Pengfei Song, Quan Zhou, et al. SOAPfuse: an algorithm for
348 identifying fusion transcripts from paired-end RNA-Seq data. *Genome Biology*, 14(2):R12, 2013.
- 349 [12] Ryan M Layer, Colby Chiang, Aaron R Quinlan, and Ira M Hall. LUMPY: a probabilistic framework
350 for structural variant discovery. *Genome Biology*, 15(6):R84, 2014.
- 351 [13] Silvia Liu, Wei-Hsiang Tsai, Ying Ding, Rui Chen, Zhou Fang, et al. Comprehensive evaluation of
352 fusion transcript detection algorithms and a meta-caller to combine top performing methods in paired-
353 end RNA-seq data. *Nucleic Acids Research*, 44(5):e47–e47, 2015.
- 354 [14] Cong Ma, Mingfu Shao, and Carl Kingsford. SQUID: transcriptomic structural variation detection
355 from RNA-seq. *Genome Biology*, 19(1):52, 2018.
- 356 [15] Andrew McPherson, Fereydoun Hormozdiari, Abdalnasser Zayed, Ryan Giuliany, Gavin Ha, et al.
357 deFuse: an algorithm for gene fusion discovery in tumor RNA-Seq data. *PLoS Computational Biology*,
358 7(5):e1001138, 2011.
- 359 [16] Daniel Nicorici, Mihaela Satalan, Henrik Edgren, Sara Kangaspeska, Astrid Murumagi, Olli Kallion-

- 360 iemi, Sami Virtanen, and Olavi Kilkku. FusionCatcher—a tool for finding somatic fusion genes in
361 paired-end RNA-sequencing data. *BioRxiv*, page 011650, 2014.
- 362 [17] Nicholas A Nystrom, Michael J Levine, Ralph Z Roskies, and J Scott. Bridges: a uniquely flexible
363 HPC resource for new communities and data analytics. In *Proceedings of the 2015 XSEDE Conference:
364 Scientific Advancements Enabled by Enhanced Cyberinfrastructure*, page 30. ACM, 2015.
- 365 [18] Tobias Rausch, Thomas Zichner, Andreas Schlattl, Adrian M Stütz, Vladimir Benes, and Jan O Korbel.
366 DELLY: structural variant discovery by integrated paired-end and split-read analysis. *Bioinformatics*,
367 28(18):i333–i339, 2012.
- 368 [19] Robert Sedgewick. *Algorithms in C, Part 5: Graph Algorithms, Third Edition*. Addison-Wesley
369 Professional, third edition, 2001.
- 370 [20] Fritz J Sedlazeck, Philipp Rescheneder, Moritz Smolka, Han Fang, Maria Nattestad, Arndt von Hae-
371 seler, and Michael C Schatz. Accurate detection of complex structural variations using single-molecule
372 sequencing. *Nature Methods*, 15(6):461–468, 2018.
- 373 [21] Wandaliz Torres-García, Siyuan Zheng, Andrey Sivachenko, Rahulsimham Vegesna, Qianghu Wang,
374 Rong Yao, Michael F Berger, John N Weinstein, Gad Getz, and Roel GW Verhaak. PRADA: pipeline
375 for RNA sequencing data analysis. *Bioinformatics*, 30(15):2224–2226, 2014.
- 376 [22] Xiaoke Wang, Renata Q Zamolyi, Hongying Zhang, Vera L Pannain, Fabiola Medeiros, Michele
377 Erickson-Johnson, Robert B Jenkins, and Andre M Oliveira. Fusion of HMGA1 to the LPP/TPRG1
378 intergenic region in a lipoma identified by mapping paraffin-embedded tissues. *Cancer Genetics and
379 Cytogenetics*, 196(1):64–67, 2010.

380 **A Proof of NP-hardness**

381 **Theorem 1.** *SCAP is NP-hard.*

382 *Proof.* To prove the NP-hardness, we reduce from MAX-2-SAT problem. It is necessary and sufficient to
383 show that for any MAX-2-SAT problem, a genome segment graph (GSG) can be constructed in polynomial
384 time, and the SCAP objective directly tells the objective of the MAX-2-SAT problem. For any MAX-2-
385 SAT instance, we are going to construct a GSG such that the satisfiability of a clause is indicated by the
386 concordance of an edge.

387 Given a MAX-2-SAT problem with n booleans $\{x_1, x_2, \dots, x_n\}$ and m clauses $\{c_1, c_2, \dots, c_m\}$, the key
388 gadget is the segments for boolean variables and clauses and the edges between them (Figure S1A). For
389 each boolean variable x_i , a segment X_i is constructed and termed as a boolean segment. For each clause
390 c_i , a segment C_i is constructed and termed as a clause segment. To ensure that the correspondence between
391 the edge concordance and the clause satisfiability as well as the correspondence between the orientation
392 of boolean segments and the assignment of boolean variables, we add edges between clause segments and
393 boolean segments in the following way. For clause c_i that involves boolean x_{i_1} , an edge is added between
394 the head of X_{i_1} and the head of C_i if clause c_i contains the negation of x_{i_1} , otherwise the edge is between
395 the tail of X_{i_1} and the head of C_i . When the literal is x_{i_1} , setting the orientation of segment X_{i_1} to be 1
396 indicates assigning True to variable x_{i_1} and leads to the concordance of the edge; when the literal is \bar{x}_{i_1} ,
397 setting the orientation of segment X_{i_1} to be 0 indicates assigning False to variable x_{i_1} and leads to the edge
398 concordance. The edge between clause c_i and the other involved boolean variable x_{i_2} is added in the same
399 principle. We call the edge between boolean segments and the clause segments as Type 1 edge. Type 1
400 edges have weight of 1.

401 Two extra edges between the two boolean segments involved in each clause are added. This is the
402 Type 2 edge with weight of 1. For each clause c_i that involves boolean x_{i_1} and x_{i_2} , two edges are added
403 between X_{i_1} and X_{i_2} as in Table S1. When both literals in c_i are True, there are two concordant Type 1
404 edges; when only one literal in c_i is True, one and only one of the two Type 2 edges is guaranteed to be
405 concordant, to compensate for the decrease of concordant Type 1 edges.

406 An extra $n + m + 1$ segments are added that we term blocking segments and denote as

407 $\{B_1, B_2, \dots, B_{n+m+1}\}$. Suppose w_1 and w_2 are large positive weights, and $w_2 \gg w_1 \gg 1$. Type 3 edges
 408 with edge weight w_2 are constructed between each adjacent pair of blocking segments, specifically between
 409 the tail of B_i and the head of B_{i+1} ($\forall i \in [1, n+m]$). Type 3 edges are used to enforce the order and orien-
 410 tation among blocking segments. Type 4 edges with weight w_1 are constructed between blocking segments
 411 and the other types of segments. Specifically, when $i \leq n$, an edge is added between the tail of segment B_i
 412 and both the head and the tail of X_i , as well as between the tail and the head of X_i and both the head of
 413 B_{i+1} . Similarly when $n < i \leq n+m$, two edges are added between the tail of B_i and C_{i-n} , and two other
 414 edges are added between the head and tail of C_{i-n} and B_{i+1} . Type 4 edges are used to enforce the relative
 415 order between blocking segments and the boolean and clause segments. But the orientation of the boolean
 416 and clause segments can be changed freely.

clause c_i	edge 1	edge 2
$x_{i_1} \vee x_{i_2}$	head of X_{i_1} to head of X_{i_2}	tail of X_{i_1} to tail of X_{i_2}
$\bar{x}_{i_1} \vee x_{i_2}$	tail of X_{i_1} to head of X_{i_2}	head of X_{i_1} to tail of X_{i_2}
$x_{i_1} \vee \bar{x}_{i_2}$	tail of X_{i_1} to head of X_{i_2}	head of X_{i_1} to tail of X_{i_2}
$\bar{x}_{i_1} \vee \bar{x}_{i_2}$	head of X_{i_1} to head of X_{i_2}	tail of X_{i_1} to tail of X_{i_2}

Table S1: Construction of Type 4 edges based on the clause.

417 We first prove that the order of the blocking segments in the optimal arrangement is $B_1 < B_2 < \dots <$
 418 B_{n+m+1} and the orientations of them are all in forward strand, where $<$ denotes the ordering between
 419 segments. Under the arrangement that uses the forward strand of all $\{B_i\}$ and have an order of $B_1 < B_2 <$
 420 $\dots < B_{n+m+1}$, the sum of concordant edge weights is at least $(n+m)w_2$. If the optimal arrangement
 421 contains any violations of the adjacencies between B_i and B_{i+1} , there will at least one Type 3 edge that does
 422 not connect blocking segments in a tail-to-head manner and become a discordant edge in the arrangement.
 423 Therefore, the optimal arrangement can at most have an objective value of $(n+m-1)w_2 + 4(n+m)w_1 +$
 424 $4m$. Since $w_2 \gg w_1 \gg 1$, the objective value is smaller than $(n+m)w_2$, and the arrangement is not
 425 optimal, which contradicts the assumption. Therefore assuming the whole chain of segments is not reverse
 426 complemented, the orientations of blocking segments are all in forward strand, and order is $B_1 < B_2 <$
 427 $\dots < B_{n+m+1}$ in the optimal arrangement.

We then prove that the Type 2 edges restrict the order of all segments but not the orientation of boolean and
 clause segments. The order between blocking segments and boolean segments must be $B_i < X_i < B_{i+1}$,
 the order between blocking and clause segments must be $B_i < C_{i-n} < B_{i+1}$, and all boolean segments

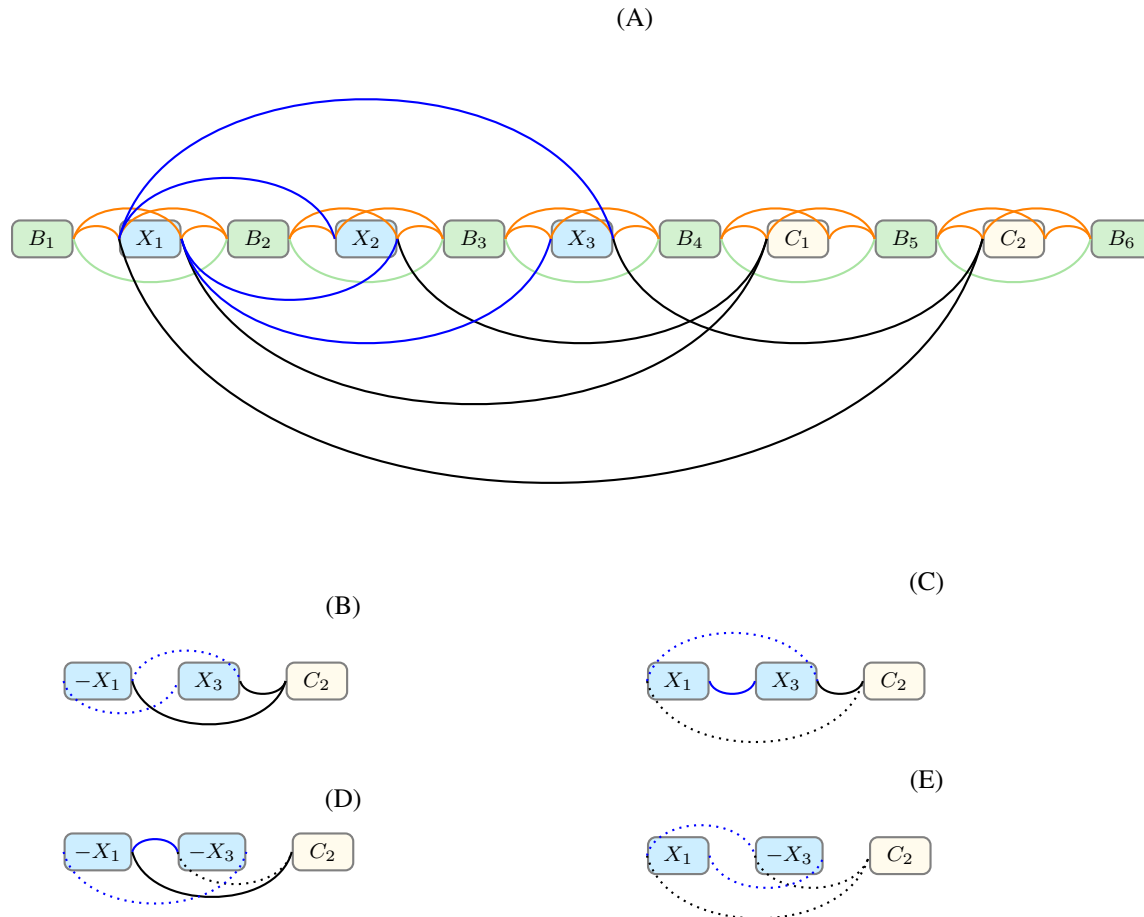


Figure S1: (A) Constructed GSG for boolean expression $(x_1 \vee x_2) \wedge (\bar{x}_1 \vee x_3)$. There is a segment for each boolean variable x_i (blue) and clause c_i (white), and 6 blocking segments (green) to separate between boolean segments and clause segments. Type 1 edges, black edges, are connecting between boolean segments and clause segments. Type 2 edges, blue edges, are connecting between a pair of boolean segments that appear in the same clause. Type 3 edges, green edges in the figure, are chaining the blocking segments. Type 4 edges, orange edges, are connecting between blocking and boolean / clause segments. (B-E) The subgraph corresponding clause $\bar{x}_1 \vee x_3$. $-X_1$ and $-X_3$ means the segment is inverted. Solid lines indicate the concordant edges in the arrangement, and dotted lines indicate the discordant edges. (B) The clause is satisfied with both literals satisfied. (C) The clause is satisfied with x_3 satisfied. (D) The clause is satisfied with \bar{x}_1 satisfied. (E) The clause is not satisfied.

must be before clause segments. When the order is $B_i < X_i < B_{i+1}$ among the three segments, and the orientations of B_i and B_{i+1} are both in forward strand, the concordant edge weights of Type two edge sum to $2w_1$ no matter whether X_i is in forward strand or inverted. The same weight can be achieved for order $B_i < C_{i-n} < B_{i+1}$. The arrangement with order $B_1 < X_1 < B_2 < \dots < B_n < X_n < B_{n+1} < C_1 < B_{n+2} < \dots < C_m < B_{n+m+1}$ and with all blocking segments in their forward strand will achieve a sum of concordant edge weight $(n+m)w_2 + 2(n+m)w_1$ at least. This concordant weight is summed over Type 3 and Type 4 edges. However, if the optimal arrangement violates any $B_i < X_i < B_{i+1}$ or $B_i < C_{i-n} < B_{i+1}$ order, the violated triplet can achieve at most w_1 of concordant edge weights, and thus the maximum sum of concordant edge weights is $(n+m)w_2 + 2(n+m-1)w_1 + w_1 + 4m$. Since $w_1 \gg 1$, the “optimal” arrangement objective is smaller than $(n+m)w_2 + 2(n+m)w_1$, which contradicts the optimality. Therefore, the order of all segments in the optimal arrangement must be

$$B_1 < X_1 < B_2 < \dots < B_n < X_n < B_{n+1} < C_1 < B_{n+2} < \dots < C_m < B_{n+m+1}.$$

428 Third, we prove that under the above segment order there are always two concordant edges of weight 1 when
 429 clause segment C_i has any concordant Type 1 edge. Suppose there is a clause c_i involving boolean variables
 430 x_{i_1} and x_{i_2} , segment C_i has one Type 1 edge between X_{i_1} and one Type 1 edge between X_{i_2} . When both
 431 Type 1 edges are concordant, both Type 2 edges between X_{i_1} and X_{i_2} are discordant (Figure S1B. When
 432 only one of the Type 1 edges is concordant, there is also one Type 2 edge between X_{i_1} and X_{i_2} that is
 433 concordant (Figure S1C,D). When neither of the Type 1 edges is concordant, both of the two Type 2 edges
 434 between X_{i_1} and X_{i_2} are discordant (Figure S1E). In this case, there is zero concordant edges of weight 1
 435 incident to C_i . Any arrangement solution of objective value W that satisfies the above segment order has
 436 $W - (n+m)w_2 - 2(n+m)w_1$ concordant edges of weight 1. Therefore, the arrangement solution will
 437 have $\frac{1}{2}(W - (n+m)w_2 - 2(n+m)w_1)$ clause segments with non-zero concordant Type 1 edges.

438 When multiple clauses involve the same pair of segment, multi-edges between X_{i_1} and X_{i_2} are constructed
 439 to make sure that two edges of weight 1 are contributed by any clause segment when it has non-zero con-
 440 cordant Type 1 edges.

441 Suppose the optimal number of satisfied clauses of the MAX-2-SAT instance is OPT_m and the optimal
 442 sum of concordant edge weights of the constructed SCAP instance is OPT_s , the following inequality holds:

443 $\frac{1}{2}(OPT_s - (n + m)w_2 - 2(n + m)w_1) \geq OPT_m$. Given the optimal solution of the MAX-2-SAT instance,
 444 a SCAP solution can be constructed by reversing segment X_i if x_i is assigned to False while keeping the
 445 order of $B_1 < X_1 < B_2 < \dots < B_n < X_n < B_{n+1} < C_1 < B_{n+2} < \dots < C_m < B_{n+m+1}$. By
 446 the construction of the Type 1 edges, a clause segment will have at least one concordant Type 1 edge if
 447 and only if it corresponds to a satisfied clause in the MAX-2-SAT solution. Denoting the objective value of
 448 the constructed solution of arrangement problem as W and applying the third proof, we have the following
 449 equality $OPT_m = \frac{1}{2}(W - (n + m)w_2 - 2(n + m)w_1)$. Since the optimal objective value of the arrangement
 450 problem is at least W ,

$$OPT_m = \frac{1}{2}(W - (n + m)w_2 - 2(n + m)w_1) \leq \frac{1}{2}(OPT_s - (n + m)w_2 - 2(n + m)w_1). \quad (S1)$$

451 Meanwhile $\frac{1}{2}(OPT_s - (n + m)w_2 - 2(n + m)w_1) \leq OPT_s$. Given the optimal solution of arrangement
 452 problem, there are $\frac{1}{2}(OPT_s - (n + m)w_2 - 2(n + m)w_1)$ clause segments with non-zero concordant
 453 Type 1 edges. Construct a MAX-2-SAT solution by assigning False to boolean variables if the corresponding
 454 boolean segment is reversed otherwise assigning True. The concordance of Type 1 edges guarantees that the
 455 corresponding literals in the MAX-2-SAT clauses are True. Thus the constructed MAX-2-SAT solution will
 456 have $\frac{1}{2}(OPT_s - (n + m)w_2 - 2(n + m)w_1)$ satisfied clauses, which is smaller than or equal to the optimal
 457 number of satisfied clauses. Therefore

$$\frac{1}{2}(OPT_s - (n + m)w_2 - 2(n + m)w_1) \leq OPT_m. \quad (S2)$$

458 Combining inequality (S1) and inequality (S2), the maximum number of satisfied clauses in MAX-2-SAT
 459 instance can be directly calculated from the optimal concordant edge weights in the arrangement problem,
 460 that is, $OPT_m = \frac{1}{2}(OPT_s - (n + m)w_2 - 2(n + m)w_1)$.

461 □

462 **B Proof of Characterization of Conflict Structures**

463 **Theorem 4.** *Any acyclic subgraph of GSG is a compatible structure.*

464 *Proof.* We show that any acyclic subgraph with N edges ($|E'| + |\hat{E}'| = N$), $G'_N = (E', \hat{E}', V')$, of GSG is
 465 a compatible structure by induction.

466 When $|E'| + |\hat{E}'| = 1$, G'_1 is a compatible structure because no other edge in G' is in conflict with the only
 467 edge $e \in E'$.

468 Assume the theorem hold for any acyclic subgraph that contains n edges. Let $G'_{n+1} = (E', \hat{E}', V')$ be an
 469 acyclic subgraph with $n + 1$ edges. Since G'_{n+1} is acyclic, there must be a leaf edge that is incident to a leaf
 470 node. Denote the leaf node as v_b and the leaf edge $e = (u_a, v_b) \in E' \cup \hat{E}'$ ($a, b \in \{h, t\}$). By removing edge
 471 e and leaf node v_b , the subgraph $G'_n = (E' - \{e\}, \hat{E}' - \{e\}, V' - \{v_b\})$ is also acyclic and contains n edges.
 472 According to the assumption, G'_n is a compatible structure and there is an arrangement of the segments in
 473 which all edges in $E' \cup \hat{e}' - \{e\}$ is concordant. Because no other edge in $E' \cup \hat{E}'$ except e connects to
 474 v_b , it is always possible to place segment v back to the arrangement such that e is concordant. Specifically,
 475 one of the four placing options will satisfy edge e : the beginning of the arrangement with orientation 1, the
 476 beginning with orientation 0, the end with orientation 1 and the end with orientation 0. Therefore, G'_{n+1} is
 477 a compatible structure.

478 By induction, acyclic subgraph G'_N of GSG with any $|E'|$ is a compatible structure. □

479 **Theorem 5.** A simple cycle $C = (E', \hat{E}', V')$ is a compatible structure if and only if there are exactly two
 480 vertices, v_j and v_i such that $deg_{E'}(v_i) = deg_{E'}(v_j) = 2$ and v_i and v_j belongs to different segments.

481 *Proof.* We prove sufficiency and necessity separately in Lemma 1 and Lemma 2. □

482 **Lemma 1.** If C is a compatible structure, there are exactly two vertices, v_i, v_j that belong to different
 483 segments, such that $deg_{E'}(v_i) = deg_{E'}(v_j) = 2$

484 *Proof.* We discuss compatibility in two cases:

485 *Case (1): All edges are concordant in C .* Sort the vertices by genomic locations in ascending order and
 486 label the first vertex v_1 and the last v_n , assuming $|V'| = n$. Similarly, sort the set of segments S' in C by
 487 the values of their permutation function π and label the first segment s^1 and the last s^m , assuming $|S'| = m$.
 488 Since concordant connections can only be right-left connections (e.g. Figure 1 b,c), $v_1 = s_r^1$ and $v_n = s_l^m$.
 489 Since C is a simple cycle, all vertices $v \in V'$ have $deg(v) = 2$. Because v_1 and v_n are the first and last
 490 vertices in this arrangement, the edges incident to v_1 or v_n must be in E' . It follows that the two edges

491 incident to v_1 connects to s_t^2 and s_t^m . Similarly, edges incident to v_n connects to s_r^1 and s_r^{n-1} . Therefore,
 492 we have $\deg_{E'}(v_1) = \deg_{E'}(v_n) = 2$. Any other vertex v_i ($1 < i < n$) is connected by one $e \in E'$ and one
 493 $\hat{e} \in \hat{E}'$ and thus has $\deg_{E'}(v_i) = 1$.

494 *Case (2): Some edges are discordant in C .* If discordant edges exist in cycle C , according to the definition
 495 of compatible structure, segments in C can be arranged such that all edges are concordant. This reduces to
 496 case (1). □

497 **Lemma 2.** *If there are exactly two vertices in V' that belong to different segments, v_i and v_j , such that*
 498 $\deg_{E'}(v_i) = \deg_{E'}(v_j) = 2$, *then C is a compatible structure.*

499 *Proof.* Let v_i and v_j be the one of the end points of segments s^i and s^j ($i \neq j$), respectively. We can arrange
 500 s^i and s^j such that $\pi(s^i) = \min_{s \in S'} \pi(s)$, $\pi(s^j) = \max_{s \in S'} \pi(s)$ and that $v_i = s_t^i$, $v_j = s_h^j$. Rename v_i to
 501 v_1 and v_j to v_n . Since C is a simple cycle, we can find two simple paths, P_1 and P_2 , between v_1 and v_n and
 502 there is no edge between P_1 and P_2 . Let P'_1 and P'_2 denote P_1 and P_2 that exclude v_1 and v_n and the edges
 503 incident to v_1 and v_n . Since P'_1 and P'_2 as acyclic subgraphs of GSG, according to Theorem 4, P'_1 and P'_2 are
 504 compatible structures and therefore segments in P'_1 and P'_2 can be arranged so that all edges are concordant.
 505 Denote the first and last vertices in the arranged P'_1 as v_2 and v_3 , and the first and last vertices in the arranged
 506 P'_2 as v_4 and v_5 . Because all the edges are concordant in P'_1 , v_2 and v_3 are the left and right ends of the first
 507 and last segments in P'_1 . Because only v_1 and v_n have $\deg_{E'} = 2$ in C , v_2 must be connected to v_1 or v_n
 508 and v_3 must be connected to v_n or v_1 . A similar argument applies to v_4 and v_5 . To ensure concordance of
 509 edges connected to v_1 and v_n , if v_n is connected to v_2 and v_1 is connected to v_3 , we flip all the segments in
 510 P'_1 . The similar operation is applied to v_4 , v_5 and P'_2 . Now we have a compatible structure. □