

 Open access • Posted Content • DOI:10.1101/2021.03.01.433465

Mako: a graph-based pattern growth approach to detect complex structural variants

— [Source link](#) 

Jiadong Lin, Xiaofei Yang, Walter A. Kusters, Tun Xu ...+12 more authors

Institutions: Xi'an Jiaotong University, Leiden University, University of Maryland, Baltimore, University of Washington ...+1 more institutions

Published on: 02 Mar 2021 - [bioRxiv](#) (Cold Spring Harbor Laboratory)

Related papers:

- [Mako: A Graph-based Pattern Growth Approach to Detect Complex Structural Variants.](#)
- [A variant selection framework for genome graphs.](#)
- [MpBsmi: A new algorithm for the recognition of continuous biological sequence pattern based on index structure.](#)
- [Rapid and enhanced remote homology detection by cascading hidden Markov model searches in sequence space.](#)
- [A space and time-efficient index for the compacted colored de Bruijn graph](#)

Share this paper:    

View more about this paper here: <https://typeset.io/papers/mako-a-graph-based-pattern-growth-approach-to-detect-complex-2woxlpomy>

1 **Mako: a graph-based pattern growth approach to detect** 2 **complex structural variants**

3 Jiadong Lin^{1,2,3,4,#,a}, Xiaofei Yang^{2,5,#,b}, Walter Kusters^{4,c}, Tun Xu^{1,d}, Yanyan Jia^{1,e},
4 Songbo Wang^{1,f}, Qihui Zhu^{6,g}, Mallory Ryan^{6,h}, Li Guo^{2,8,i}, Chengsheng Zhang^{6,7,j}, The
5 Human Genome Structural Variation Consortium⁸, Charles Lee^{6,7,k}, Scott E. Devine^{9,l},
6 Evan E. Eichler^{10,11,m}, Kai Ye^{1,2,3,12,*,n}

7 ¹School of Automation Science and Engineering, Faculty of Electronic and Information
8 Engineering, Xi'an Jiaotong University, Xi'an, 710049 China.

9 ²MOE Key Lab for Intelligent Networks & Networks Security, Faculty of Electronic
10 and Information Engineering, Xi'an Jiaotong University, Xi'an, 710049 China.

11 ³Genome Institute, the First Affiliated Hospital of Xi'an Jiaotong University, Xi'an,
12 710061 China.

13 ⁴Leiden Institute of Advanced Computer Science, Faculty of Science, Leiden
14 University, Leiden, Netherland.

15 ⁵School of Computer Science and Technology, Faculty of Electronic and Information
16 Engineering, Xi'an Jiaotong University, Xi'an, 710049 China.

17 ⁶The Jackson Laboratory for Genomic Medicine, Farmington, CT 06032.

18 ⁷Precision Medicine Center, the First Affiliated Hospital of Xi'an Jiaotong University,
19 Xi'an, 710061 China.

20 ⁸Consortium authors are enumerated at the end of this document.

21 ⁹Institute for Genome Sciences, University of Maryland School of Medicine, Baltimore,
22 MD 21201, USA.

23 ¹⁰Department of Genome Sciences, University of Washington School of Medicine,
24 Seattle, WA 98119-5065, USA.

25 ¹¹Howard Hughes Medical Institute, University of Washington, Seattle, WA 98195,
26 USA.

27 ¹²The School of Life Science and Technology, Xi'an Jiaotong University, Xi'an,
28 710049 China.

29 [#]These authors contributed equally to this work.

30 ^{*}To whom correspondence should be addressed. E-mail: kaiye@xjtu.edu.cn (Ye K)

31
32

33 **Running title:** *Jiadong Lin et al / Mako: a graph-based pattern growth approach to*
34 *detect complex structural variants*

35

36

37 Total word counts (from “Introduction” to “Conclusions” or “Materials and methods”):

38 6,827

39 Total references: 59

40 Total figures: 6

41 Total tables: 3

42 Total supplementary figures: 23

43 Total supplementary tables: 9

44 Total supplementary files: 4

45 Total characters in title: 71

46 Total characters in running title: 86

47 Total keywords: 5

48 Total words in abstract: 209

49 **Abstract**

50 Complex structural variants (CSVs) are genomic alterations that have more than two
51 breakpoints and are considered as simultaneous occurrence of simple structural variants.
52 However, detecting the compounded mutational signals of CSVs is challenging through
53 a commonly used model-match strategy. As a result, there has been limited progress for
54 CSV discovery compared with simple structural variants. We systematically analyzed
55 the multi-breakpoint connection feature of CSVs, and proposed Mako, utilizing a
56 bottom-up guided model-free strategy, to detect CSVs from paired-end short-read
57 sequencing. Specifically, we implemented a graph-based pattern growth approach,
58 where the graph depicts potential breakpoint connections and pattern growth enables
59 CSV detection without predefined models. Comprehensive evaluations on both
60 simulated and real datasets revealed that Mako outperformed other algorithms. Notably,
61 validation rates of CSV on real data based on experimental and computational
62 validations as well as manual inspections are around 70%, where the medians of
63 experimental and computational breakpoint shift are 13bp and 26bp, respectively.
64 Moreover, Mako CSV subgraph effectively characterized the breakpoint connections
65 of a CSV event and uncovered a total of 15 CSV types, including two novel types of
66 adjacent segments swap and tandem dispersed duplication. Further analysis of these
67 CSVs also revealed impact of sequence homology in the formation of CSVs. Mako is
68 publicly available at <https://github.com/jiadong324/Mako>.

69 **KEYWORDS:** Next-generation sequencing; Complex structural variants; Pattern
70 growth; Graph mining; Formation mechanism

71

72

73 **Introduction**

74 Computational methods based on next-generation-sequencing (NGS) have provided an
75 increasingly comprehensive discovery and catalog of simple structure variants (SVs)
76 that usually have two breakpoints, such as deletions and inversions [1-7]. In general,
77 these approaches follow a model-match strategy, where a specific SV model and its
78 corresponding mutational signal model is proposed. Afterwards, the mutational signal
79 model is used to match observed signals for the detection (**Figure 1A**). This model-
80 match strategy has been proved effective for detecting simple SVs, providing us with
81 prominent opportunities to study and understand genome evaluation and disease
82 progression [8-11]. However, recent research has revealed that some rearrangements
83 have multiple, compounded mutational signals and usually cannot fit into the simple
84 SV models [8, 12-16] (**Figure 1B**). For example, in 2015, Sudmant et al. systematically
85 categorized 5 types of complex structural variants (CSVs) and found that a remarkable
86 80% of 229 inversion sites were complex events [8]. Collins et al. used long-insert size
87 whole genome sequencing (liWGS) on autism spectrum disease (ASD) and
88 successfully resolved 16 classes of 9,666 CSVs from 686 patients [17]. In 2019, Lee et
89 al. revealed that 74% of known fusion oncogenes of lung adenocarcinomas were caused
90 by complex genomic rearrangements, including *EML4-ALK* and *CD74-ROSI* [16].
91 Though less frequently reported compared with simple SVs, these multiple breakpoint
92 rearrangements were considered as punctuated events, leading to severe genome
93 alterations at once [10, 18-21]. This dramatic change of genome provided distinctive
94 evidence to study formation mechanisms of rearrangement and to understand cancer
95 genome evolution [13, 14, 17, 19, 21-25].

96 However, due to lack of effective CSV detection algorithms, most CSV related
97 studies screen these events from the “sea” of simple SVs through computational
98 expensive contig assembly and realignment, incomplete breakpoints clustering or even
99 targeted manual inspection [8, 12, 16]. In fact, many CSVs have been already neglected
100 or misclassified in this “sea” because of the incompatibility between complicated
101 mutational signals and existing SV models. Although the importance and challenge for
102 CSV detection have been recognized, only a few dedicated algorithms were proposed
103 for CSVs discovery, and they followed two major approaches guided by the model-
104 match strategy. TARDIS and SVelter utilizes the top-down approach, where they
105 attempt to model all the mutational signals of a CSV event instead of modeling specific

106 parts of signals. In particular, TARDIS [26] proposed sophisticated abnormal alignment
107 models to depict the mutational signals reflected by dispersed duplication and inverted
108 duplication. The pre-defined models were then used to fit observed signals from
109 alignments for the detection of the two specific CSV types. Indeed, this was
110 complicated and greatly limited by the diversity types of CSV. To solve this, SVelter
111 [27] replaced the modeling process for specific CSVs with a randomly created virtual
112 rearrangement. And CSVs were detected by minimizing the difference between the
113 virtual rearrangement and the observed signals. Whereas GRIDSS [28] represents the
114 assembly-based approach, which detected CSVs through extra breakpoints discovered
115 from contig-assembly and realignment. Though assembly-based approach is sensitive
116 for breakpoint detection, it lacks certain regulations to constrain or classify these
117 breakpoints and leave them as independent events. As a result, these model-match
118 guided approaches would substantially break-up or misinterpret the CSVs because of
119 partially matched signals (**Figure 1B**). Moreover, graph is another approach that has
120 been widely used for simple [2, 29] and complex [19, 30] SV detection. Notably, ARC-
121 SV [30] uses clustered discordant read-pairs to construct an adjacency graph and adopts
122 a maximum likelihood model to detection complex SVs, showing great potential of
123 using graph to detect complex SVs. Accordingly, there is an urgent demand of a new
124 strategy, enabling CSV detection without predefined models as well as maintaining the
125 completeness of a CSV event.

126 In this study, we proposed a bottom-up guided model-free strategy, implemented as
127 Mako, to effectively discover CSVs all at once based on short-read sequencing.
128 Specifically, Mako uses a graph to build connections of mutational signals derived from
129 abnormal alignment, providing the potential breakpoint connections of CSVs.
130 Meanwhile, Mako replaces model fitting with the detection of maximal subgraphs
131 through a pattern growth approach. Pattern growth is a bottom-up approach, which
132 captures the natural features of data without sophisticated model generation, allowing
133 CSV detection without predefined models. We benchmarked Mako against five widely
134 used tools on a series of simulated and real data. The results show that Mako is an
135 effective and efficient algorithm for CSV discovery, which will provide more
136 opportunities to study genome evolution and disease progression from large cohorts.
137 Remarkably, the analysis of subgraphs detected by Mako highlights the unique strength
138 of Mako, where Mako was able to effectively characterize the CSV breakpoint
139 connections, confirming the completeness of a CSV event. Moreover, we

140 systematically analyzed the CSVs detected by Mako on three healthy samples,
141 revealing a novel role of sequence homology in CSV formation.

142 **Results**

143 In this section, we give an overview of the Mako algorithm, with full details available
144 in the Methods section. For performance comparison, we propose all-breakpoint match
145 and unique-interval match to evaluate Mako against five published methods on both
146 simulated and real data. The detailed explanation of the evaluation measurements, CSV
147 simulation and real data CSV benchmarks are described in the Methods section.
148 Additionally, we describe our observations of Mako's CSV discovery from HG00514,
149 HG00733 and NA19240. These samples are sequenced by Human Structural Variants
150 Consortium (HGSVC) and publicly available.

151 **Overview of the Mako algorithm**

152 Given the fact that a CSV is a single event with multiple breakpoint connections, we
153 observe that either false positive breakpoints or breakpoints from other events will not
154 have connection with the breakpoints in the current CSV because of weak or non-exist
155 connections. Thus, we formulate the detection of CSVs as maximal subgraph pattern
156 detection in a signal graph. To detect the CSV subgraphs, Mako comprises two major
157 steps (**Figure 2**). Firstly, it collects abnormal aligned reads clusters as nodes and uses
158 two types of edges to build the so-called signal graph. To build the high-quality graph,
159 we filtered discordance alignments based on procedure described in BreakDancer [4]
160 (**Methods**). The resulting signal graph is formally defined as follows:

161 $G = (V, E)$ with $V = \{v_1, v_2, \dots, v_n\}$ and $E = \{E_{pe}, E_{ae}\}$, where each node $v \in V$ is
162 represented as $v = (type, pos, weight)$, and each edge in E_{pe} and E_{ae} is represented
163 as either $e_{pe} = (v_i, v_j, rp)$ or $e_{ae} = (v_i, v_j, dist)$, with $v_i, v_j \in V$. In particular, Mako uses
164 weight and the ratio between weight and coverage at pos to filter nodes, which are
165 created separately by clustering discordant read-pairs, clipped reads and split reads
166 (**Methods**). For the edge set, E_{pe} contains the paired edges that represent connections
167 between two signals on the genome derived from paired-reads or split-alignment, while
168 E_{ae} consists of adjacency edges that indicate distances between signals along the
169 genome. Afterwards, Mako applies a pattern growth search strategy to efficiently
170 discover these subgraphs as potential CSVs at whole genome scale. Meanwhile, the
171 attributes of the subgraph are used to measure the complexity and to define the types of

172 CSVs. Specifically, the CSVs types are given by the edge connection types of the
173 corresponding subgraphs (**Figure 2**).

174 **Mako effectively characterizes multiple breakpoints of CSV**

175 The most important feature for a CSV is the presence of multiple breakpoints in a single
176 event. Thus, we first examined the performance of breakpoint detection for Mako,
177 Lumpy, Manta, SVelter, TARDIS and GRIDSS. The results were evaluated according
178 to the all-breakpoint match criteria on both reported and randomized CSV type
179 simulations (**Methods**). For convenience, we used the terms reported CSV and
180 randomized CSV throughout this study. Overall, for the heterozygous (HET) (**Figure**
181 **3A**) and homozygous (HOM) (**Figure 3B**) simulation, Mako was comparable to
182 GRIDSS and they outperformed other algorithms. For example, GRIDSS, Mako and
183 Lumpy detected 50%, 51% and 46% for reported HET CSV breakpoints, while they
184 reported 53%, 54% and 44% for randomized ones. Because the graph encoded both
185 multiple breakpoints and their substantial connections for each CSV, Mako achieved
186 better performance on randomized events, which included more subcomponents than
187 the reported ones. Indeed, by comparing reported and randomized simulation, the
188 breakpoint detection sensitivity (**Figure 3A** and **Figure 3B**) of Mako increased, while
189 that of other algorithms dropped except for GRIDSS. Although the assembly-based
190 method, GRIDSS, is as effective as Mako for breakpoint detection, it lacks a proper
191 procedure to resolve the connections among breakpoints.

192 **Mako precisely discovers CSV unique-interval**

193 CSV is considered as a single event consisted of connected breakpoints and we have
194 demonstrated that Mako was able to detect CSV breakpoints effectively. However, the
195 breakpoint detection evaluation only assesses the discovery of basic components for a
196 CSV and lacks examination for CSV completeness. We then investigated whether
197 Mako could precisely capture the entire CSV interval even with missing breakpoints.
198 In general, according to the unique-interval match (**Methods**) criteria, Mako
199 consistently outperformed other algorithms for both reported and randomly created
200 CSVs, while SVelter and GRIDSS ranked second and third, respectively. For the
201 reported CSVs at 30X coverage (**Figure 3C** and **Figure 3D**), the recall of Mako was
202 94% and 92%, which was significantly higher than SVelter (49% and 57%) for both
203 reported HET and HOM CSVs, respectively. Due to the model guided top-down
204 approach, SVelter was able to discover some complete CSV events. However, the
205 virtual rearrangement generation may not fully explore all possibilities. Remarkably,

206 we noted that Mako's superior sensitivity was most significant for randomized
207 simulation (**Figure 3E** and **Figure 3F**), which was consistent with our previous
208 observation (**Figure 3A**, **Figure 3B**). In particular, at 30X coverage, Mako's recall
209 (88%) was much higher than SVelter (29%) for the HET CSVs (**Figure 3E**). This was
210 due to the complementary nature of the graph edges (adjacent and paired), from which
211 the subgraph can be expanded alternatively through one or the other, enabling the
212 complete CSVs discovery even with missing breakpoints.

213 **Performance on real data**

214 Since Mako outperformed other methods on simulated data, we further compared Mako
215 with SVelter, GRIDSS and TARDIS on whole genome sequencing data of NA19240
216 and SKBR3. Firstly, we obtained 6,060, 7,733, 6,426 and 15,358 calls for NA19240,
217 and 2,962, 2,468, 3,077 and 4,010 for SKBR3 predicted by Mako, SVelter, GRIDSS
218 and TARDIS, respectively (**Methods**, **Supplementary Figure S1-S2**). By comparing
219 their predictions, we found Mako and GRIDSS showed similar performance (**Figure**
220 **4A** and **4B**) which was consistent with our observation in simulated data (**Figure 3**).
221 Furthermore, we examined the discovery completeness of 59 (NA19240) and 21
222 (SKBR3) benchmark CSVs (**Table 1**, **Supplementary File 1**, **Supplementary Table**
223 **S1**). Because Manta and Lumpy contributed to the CSV benchmarks, they were
224 excluded from the comparison. The results showed that Mako performed the best for
225 the two benchmarks with different CXS thresholds, while TARDIS ranked second
226 (**Figure 4C**). Given that inverted duplication and dispersed duplication dominated the
227 benchmark set and that TARDIS has designed specific models for these two types,
228 TARDIS detected more events of these two duplication types than others did (**Table 1**).
229 SVelter only detected a few benchmark CSVs for SKBR3, because the procedure of
230 randomly created rearrangement was not optimized, leading to either incorrect events
231 or inaccurate breakpoints. Based on the above observation, we concluded that either
232 randomized model (SVelter) or specific model (TARDIS) was far from comprehensive
233 to cover the large diversity of real CSV types.

234 **CSV subgraph illustrates breakpoints connections**

235 CSVs from autosomes were selected from Mako's callset with more than one edge
236 connection type observed in the subgraph, leading to 403, 609, and 556 events for
237 HG00514, HG00733 and NA19240, respectively (**Figure 5A**, **Supplementary Table**
238 **S2**). We systemically evaluated all CSV events in HG00733 via experimental and
239 computational validation as well as manual inspection. For experimental validation, we

240 successfully designed primers for 107 CSVs (**Supplementary Table S3**), where 15 out
241 of 21 (71%, **Table 2**) successfully amplified were validated by Sanger sequencing
242 (**Supplementary Table S4**). The computational validation (**Supplementary Figure S3**)
243 showed up to 87% accuracy, indicating a combination of methods and external data is
244 necessary for comprehensive CSV validation (**Table 3, Methods**). Further analysis
245 showed that the medians of breakpoint shift were 13bp and 26bp comparing to
246 breakpoints given by experimental and computational evaluation (**Supplementary**
247 **Figure S4**). We observed that approximately 54% of CSVs were found in either STR
248 or VNTR regions, contributing to 75% of all events inside the repetitive regions (**Figure**
249 **5A**). For the connection types, more than half of the events contains DUP and INS
250 edges in the graph, indicating duplication involved sequence insertion. Moreover,
251 around 40% of the events contain DEL edges (**Figure 5A**), showing two distant
252 segment connections derived from either duplication or inversion events. We further
253 examined whether the CSV subgraph depict the connections for each CSV via
254 discordant read-pairs. Interestingly, we observed two representative events with four
255 breakpoints at chr6:128,961,308-128,962,212 (**Figure 5B**) and chr5:151,511,018-
256 151,516,780 (**Figure 5C**) from NA19240 and SKBR3, respectively. Both events were
257 correctly detected by Mako, but missed by SVelter and reported more than once by
258 GRIDSS and TARDIS (**Supplementary Table S5**). In particular, the CSV at
259 chr6:128,961,308-128,962,212 that consists of two deletions and an inverted spacer
260 was reported twice and five times by GRIDSS and TARDIS. The event at chromosome
261 5 that consists of a deletion and dispersed duplication was reported four and three times
262 by GRDISS and TARDIS. These redundant predictions complicate and mislead
263 downstream functional annotations. On the contrary, Mako was able to completely
264 detect the above two CSV events, and also capable of revealing the breakpoint
265 connections of CSVs encoded in the subgraphs. The above observations suggested that
266 the subgraphs detected through pattern growth are interpretable, from which we can
267 characterize the breakpoint connections for a given CSV event.

268 **Contribution of homology sequence in CSV formation**

269 Ongoing studies have revealed that genome alterations are mainly caused by the
270 inaccurate DNA repair and the 2-33bp long microhomology sequence at breakpoint
271 junctions plays an important role in CSV formation [18, 31-34]. To further characterize
272 CSVs' internal structure and examine the impact of homology sequence on CSV
273 formation, we manually reconstructed (**Methods**) 1,052 high-confident CSV calls

274 given by Mako (252/403 from HG00514, 440/609 from HG00733 and 360/556 from
275 NA19240) via PacBio HiFi reads (**Figure 6A, Supplementary Table S6,**
276 **Supplementary Figure S5, Supplementary File 2**). The percentage of successfully
277 reconstructed events was similar to the orthogonal validation rate, showing CSVs
278 detected by Mako were accurate and the validation method was effective. The high-
279 confident CSV callset contains 816 insDup events with both insertion and duplication
280 edge connections. Further investigation revealed that these events contains irregular
281 repeat sequence expansion, making them different from simple insertion or duplications
282 (**Supplementary Figure S6**). Besides, we found two novel types, which were named
283 adjacent segments swap and tandem dispersed duplication (**Figure 6B, Supplementary**
284 **Figure S7-S8**). We inferred that homology sequence mediated inaccuracy replication
285 was the major cause for these two types. Furthermore, we observed that 134 CSVs
286 contains either inverted or dispersed duplications (**Supplementary Table S6**). These
287 duplications involved CSVs were mainly caused by Microhomology Mediated Break-
288 Induced Replication (MMBIR) according to previous studies[18, 32, 35]. It was known
289 that different homology patterns cause distinct CSV types (**Figure 6C and Figure 6D**).
290 Surprisingly, one particular pattern of homology sequence yielded multiple CSV types
291 (**Figure 6E**). In particular situations of the three different homology patterns, DNA
292 double strand break (DSB) occurred after replication of fragment *C*. According to the
293 MMBIR mechanism and template switch [23, 32-34], pattern I (**Figure 6C**) and pattern
294 II (**Figure 6D**) can only have one output but pattern III (**Figure 6E**) produces three
295 different outcomes. The results provided additional evidence for understanding the
296 impact of sequence contents on DNA DSB repair, leading to better understanding of
297 diversity variants produced by CRISPR [36, 37].

298 **Discussion**

299 Currently, short read sequencing is significantly reduced in cost and has been applied
300 to clinical diagnostics and large cohort studies [16, 38, 39]. However, CSVs from short
301 read data are not fully explored due to the methodology limitations. Though long read
302 sequencing technologies bring us promising opportunities to characterize CSVs [13, 14,
303 40], their application is currently limited to small-scale projects and the methods for
304 CSV discovery are also underdeveloped. As far as we know, NGMLR combined with
305 Sniffles is the only pipeline that utilizes the model-match strategy to discover two
306 specific forms of CSVs, namely deletion-inversion and inverted duplication. Therefore,

307 there is a strong demand in the genomic community to develop effective and efficient
308 algorithms to detect CSV using short read data. It should be noted that CSV breakpoints
309 might come from either single haplotype or different haplotypes, where two simple SVs
310 from different haplotypes lead to false positives (**Supplementary Figure S9**). This may
311 substantially increase the false discoveries because Mako currently is not able to
312 determine the exact haplotype of each breakpoints. However, Mako can be extended to
313 differentiate such false positives by adding additional features to the graph, e.g. phased
314 reads. Given that short read sequencing is not able to span all breakpoints of a CSV,
315 Mako could only infer the CSV types based on the edge connections from the subgraph,
316 while it is difficult to characterize the exact components of CSVs. Therefore, our next
317 work will integrate both short and long reads to the signal graph for CSV discovery and
318 characterization.

319 To sum up, we developed Mako, utilizing the graph-based pattern growth approach,
320 to discover CSVs. Meanwhile, the intensive experimental and computational
321 validations as well as manual inspections showed around 70% accuracy and 20bp
322 median breakpoint shift. Besides the improvement of CSV detection performance, the
323 optimized pattern growth algorithm on sequentially constrained subgraph detection is
324 not restricted to CSV detection and can be generalized to other graph problems with
325 similar constraints. Most importantly, to the best of our knowledge, Mako is the first
326 algorithm that utilizes the bottom-up guided model-free strategy for SV discovery,
327 avoiding the complicated model and match procedures. Given the fact that CSVs are
328 largely unexplored, Mako presents opportunities to broaden our knowledge of genome
329 evolution and disease progression.

330

331

332 **Materials and methods**

333 **Materials**

334 The short read aligned BAM files for NA19240, HG00514 and HG00733 were obtained
335 from the HGSVC [9] (**Supplementary Note**). The PacBio HiFi reads were provided
336 by HGSVC and we aligned these reads with pbmm2
337 (<https://github.com/PacificBiosciences/pbmm2>) and NGLMR [40] under default
338 settings (**Supplementary Note**). The haploid assembly of HG00733 were obtained
339 from HGSVC and aligned with pbmm2 (**Supplementary Note**). Both short reads and
340 long reads were aligned to the human reference genome GRCh38. The coverage was
341 approximately 70X and 30X for short and long reads, respectively (**Supplementary**
342 **Note**). The simple SV callset for NA19240 is publicly available from HGSVC, and was
343 contributed by Manta [7], Lumpy [3], Pindel [1] and etc. Alignment files and SV callset
344 for the SK-BR-3 cell line were obtained from a recent publication [13] (**Supplementary**
345 **Note**). The SK-BR-3 callset (**Supplementary Note**) was merged by SURVIVOR from
346 contributions by Manta [7], Lumpy [3], Delly [2] and PopIns [41], and contains 627
347 inversions (INV), 2,776 deletions (DEL), 483 duplications (DUP) and 1,160
348 translocations (TRA).

349 **Building signal graph**

350 To create the signal graph G , Mako collects mutational signals satisfying one of the
351 following criteria from the alignment file to create the signal nodes set V of G : 1)
352 clipped portion with minimum 10% size fraction of the overall read length; 2) split
353 reads with high mapping quality; 3) discordant read-pairs. Notice that a discordant
354 alignment will create two nodes correspondingly. Meanwhile, each node is represented
355 by a cluster of mutational signals and is given three attributes `type`, `pos` and `weight`.
356 Mako uses two types of signal clusters. One of the clusters is single-nucleotide
357 resolution cluster created by clipped reads or split reads, namely Mako clusters these
358 reads at the same location to create node. Another cluster is formed by discordant read-
359 pairs, where the clustering distance is set as estimated average insert size minus two
360 times read length. To avoid using randomly occurred discordant alignment clusters, we
361 followed the procedure introduced by Chen [4]. Specifically, it assumed one type of
362 discordant alignment at the genomic location is uniformly distributed under the null
363 hypothesis of no variant. For locations that have more than one type of discordant

364 alignment, the number of such alignments at particular location forms a mixture Poisson
365 distribution with each mixture component representing one of the discordant types.
366 Thus, we summarize the statistics of clustering of a particular type i as the probability
367 of having more than observed number of discordant alignments in a given region:

$$368 \quad P(n_i \geq k_i)$$

369 where n_i denotes the Poisson random variable with mean equal to λ_i , and k_i is the
370 number of observed type i discordant alignment. The estimation of λ_i can be
371 calculated based on the uniform assumption:

$$372 \quad \lambda_i = \frac{sN_i}{G}$$

373 where s represents the cumulative size of the regions that discordant alignments
374 anchored, N_i the total number of type i alignment in the BAM and G the length of
375 reference genome.

376 It should be noted that some discordant read-pairs may contain two types of signals, e.g.
377 abnormal insert size and incorrect mapping orientation, which are clustered separately
378 to create nodes. Moreover, split reads created nodes not only provide precise location
379 but also complement edges for discordant read-pairs. Therefore, Mako's performance
380 will not be dramatically affected by the skewed insert size distribution because skewed
381 distribution only affects estimation of abnormal insert size. The attribute `weight` and
382 `pos` indicate the number of abnormal reads and approximate position on the genome,
383 respectively; and `type` denotes the type of abnormal alignment, such as *MS*, indicating
384 the node consists of reads clipped at the right part. Importantly, we consider nodes with
385 the same `type` as identical nodes. For the edge set $E = \{E_{pe}, E_{ae}\}$ of signal graph G ,
386 the paired edges from E_{pe} are derived from read-pairs or split-reads between two
387 signal nodes, where `rp` indicates the number of paired reads involved. Adjacency
388 edges from E_{ae} measure the distance `dist` between two adjacent signals. However,
389 adjacent edges are virtual links compared with the paired edges derived from
390 alignments, thus the pattern growth through adjacent edges is constrained by `dist` to
391 avoid pointless pattern expansion. It should be noted that both types of edges might co-
392 exist between two nodes. To achieve efficient subgraph detection and avoid
393 overlapping subgraphs, we use a linearized database to store the graph and this graph
394 can be built efficiently in linear time by reading the input file once.

395 **Detecting CSVs with pattern growth**

396 Pattern growth is an efficient heuristic approach for frequent pattern discovery in strings
397 and graphs [42], which has been widely used in many areas [43-48], such as INDEL
398 detection in DNA sequences [1, 24]. Compared with statistical methods, pattern growth
399 discloses the intrinsic features of the data without sophisticated model generation.
400 Meanwhile, the output of the pattern growth approach is usually interpretable, which is
401 very important for specific applications [49].

402 In the CSV detection, the subgraph pattern starts at a single node and grows by adding
403 more nodes until it cannot find a proper node (**Algorithm I, Supplementary Figure**
404 **S10**). In addition, to avoid overlapping subgraphs, we only allow the subgraph to grow
405 according to the increasing order of pos value for each node. Meanwhile,
406 backtracking is only allowed for nodes involved in the current subgraph. For example
407 (**Figure 2**), Mako detects the maximal subgraph by visiting nodes *A*, *C*, *B*, and *D*,
408 respectively. Since the edge distances between *A* and *B* as well as *D* and *E* is larger than
409 the distance (minDist) threshold, Mako grows the subgraph through *C* and backtrack
410 node *B* to expand the subgraph, whereas edge between *D* and *E* is constrained.

411 Given the fact that the signal graph contains millions of nodes at whole genome scale,
412 we use a strategy similar to “seed-and-extension” that has been utilized by sequence
413 alignment algorithms [50, 51] to accelerate the subgraph detection process. Meanwhile,
414 we only keep the index of each node in the database to save memory for subgraph
415 detection (**Supplementary Figure S11**). Moreover, as we assigned attributes to each
416 node, the discovered subgraphs not only differ in edge connections but also in the type
417 of signal nodes in the subgraph. Therefore, we propose an algorithm that starts at
418 multiple signal nodes of the same type and extends locally for efficient subgraph
419 detection (**Algorithm II**). It should be noted that sequence alignment usually results in
420 one best alignment [50, 51], whereas our algorithm is also encouraged to discover
421 multiple maximal subgraphs that share the same edge connections but different node
422 attributes. To avoid missing subgraphs or incomplete detection, minFreq = 1 is a
423 default parameter for subgraph detection, but this could also be time consuming and
424 affected by graph noise. Thus, Mako allows users to set larger minFreq to avoid
425 random subgraphs and detects the connected components of subgraphs to ensure
426 complete detection. In particular, a larger minFreq value allows multiple identical
427 subgraphs to be discovered, and edges between these subgraphs are kept and used to
428 build connections between subgraphs. These edges can be reliably marked, because the

429 frequency of the current subgraph becomes smaller than the `minFreq` value by adding
430 those edges. Then, a local maximal subgraph represented by a connected component
431 can be discovered from the subgraph connection graph. A significant feature of
432 discovering CSVs from a graph is that it provides the connections between multiple
433 breakpoints of a CSV, so that the attributes of the discovered subgraph can be directly
434 used as a measure for CSV. Namely, if the subgraph contains more non-identical nodes
435 and E_{pe} edges, this subgraph is more likely to indicate a complex event. Therefore,
436 Mako defines the boundary of CSVs using the leftmost and rightmost `pos` value of the
437 nodes involved, and utilizes the number of identical node types multiplied by the
438 number of E_{pe} edges as a complexity score `CXS` (default=2). For example (**Figure 2**),
439 the discovered CSV subgraph has a `CXS` score of 8, because of four identical nodes
440 and two paired edges.

Algorithm 1: Detect maximal subgraphs

Input: Signal graph $G = (V, E)$, **parameters** `minFreq`, `minDist`

Output: A set of CSV subgraphs $O = \{g_1, g_2, \dots, g_n\}$, with $freq(g_i) \geq minFreq$

1: **procedure** `findMaximalSubgraph(G, minFreq, minDist)`

2: Initialize `freq_types` equals to `type` frequency of node in V ;

3: Build index-projection $G|_{\emptyset}$ of G ;

4: **for** α **in** `freq_types` **do**:

5: Build index-projection $G|_{\alpha}$;

6: $g_i = \alpha$;

7: **if** $freq(g_i) > minFreq$ **then**

8: `multiLocPatternGrowth(O, g_i, G|_{\alpha}, minFreq, minDist)`;

9: **end if**

10: **end for**

11: **end procedure**

441

Algorithm II: Multi-location subgraph growth

```
1: procedure multiLocPatternGrowth( $O, g, G|_g, minFreq, minDist$ )
2: Initialize adj_list with adjacent node direct after  $g$  through  $E$ ;
3: for node in adj_list do:
4:   if nodeInRange( $g, node$ ) then
5:      $g' = g + node$ ;
6:      $O.append(g')$ ;
7:     multiLocPatternGrowth( $O, g', G|_{g'}, minFreq, minDist$ );
8:   end if
9: end for
10: end procedure
11: procedure nodeInRange( $g, v$ )
12:   Set the nodes in  $g$  with respect increasing order of pos value:
    $v_0, v_1, \dots, v_n$ ;
13:   Set  $v' = v_n$ ;
14:   if  $freq(v) > minFreq$  then
15:     if  $dist(v', v) < minDist$  then
16:       return True
17:     else:
18:       for  $i = n$  to 0 do
19:         if  $\exists e_{pe}$  between  $v$  and  $v_i$  then
20:           return True
21:         end if
22:       end for
23:     return False
24: end procedure
```

442

443

Design of simulation studies

444

To create CSVs, we follow the simulation strategy introduced by the Sniffles[40]. In general, simple SVs generated by VISOR[52] are randomly selected and combined to make CSVs (Supplementary Figure S12). In this study, we first create deletion, inversion, inverted tandem duplication, tandem duplication and translocation copy-paste with 5000bp average size and 500bp standard deviation (Supplementary Note).

445

446

447

448

449

450

451

452

453

454

455

456

We only consider focal translocations, where the distance between source sequence and insert position is smaller than 100Kbp. These events are created using reference genome GRCh38 and collected as basic operations for further random combination usage. For example, suppose segments on the reference genome are ABCDE and the following criteria are considered for CSV simulation:

- 1) The deletion (C) associated with inversion (D') ABD'DE can be generated by first creating a deletion event and adding the inversion to a flanking region of the deletion.

457 2) The dispersed duplication and inverted duplication are produced through
458 translocation copy-paste, and the orientations at the paste position distinguish these
459 two types of duplication. For example, if we copy-paste segment B and insert it
460 after D, a dispersed duplication ABCDBE will be created.

461 3) Additionally, to create translocation copy-paste involved CSVs, we only
462 manipulate segments adjacent to the insert position of the source segment. For
463 instance, a deletion can be associated with the dispersed duplication ABCDBE by
464 removing D or E, leading to ABCBE or ABCDB.

465 To produce homozygous or heterozygous CSVs, we use the purity parameter
466 introduced by VISOR to control the ratio of reads sequenced from variation genome
467 and reference genome. After the variation genome is created, VISOR used wgsim
468 (<https://github.com/lh3/wgsim>) to simulate paired-end reads and applied BWA-MEM
469 [51] to align the simulated reads to the reference genome (**Supplementary Note**).
470 Overall, VISOR has efficient functions for creating basic operations, building variation
471 genome with simulated CSVs, simulating reads and alignment. We add the random
472 selection and combination step as part of VISOR.

473 We first evaluate whether Mako is able to capture reported CSV types published by
474 previous studies [8, 17], such as deletion flanked by inversion, inverted duplication,
475 dispersed duplication and etc. This was termed as reported CSV. For the reported CSV,
476 we only randomly select and combine deletion, inversion, inverted tandem duplication
477 and tandem duplication, but leave translocation copy-paste unchanged
478 (**Supplementary Note**). In total, we simulated 300 reported CSV types on chromosome
479 1. The reported CSVs usually have four to six breakpoints, which are still feasible to be
480 detected by model-based methods. However, we emphasize that limited knowledge of
481 CSV variety and the complex mutational signals produced by breakpoint connections
482 are the major challenges for CSV discovery. From this perspective, we made another
483 set of randomly simulated CSV types on autosomes, termed as randomized CSV, where
484 we created 4,500 CSVs with 4~10 breakpoints through random combinations of at least
485 two basic operations including translocation copy-paste (**Supplementary Note**).

486 **Creating CSV benchmark from real data**

487 It has been recognized that the most significant feature of CSVs is simultaneous
488 appearance of multiple breakpoints[8, 12, 27, 53, 54]. However, the development of
489 robust tools for screening complex events is a difficult and unsolved problem because
490 there are currently no well-defined rules for constraining the expected breakpoint

491 patterns[12]. In order to study CSVs, researchers follow four major steps[12, 20] to
492 resolve CSVs from an enormous number of simple SVs: 1) breakpoint clustering; 2)
493 clustered breakpoints enrichment test; 3) contig assembly and realignment; 4) manually
494 inspection from visualization. Fortunately, PacBio reads provide us with the direct
495 evidence to validate and categorize CSVs, which can be used to screen each simple SV
496 site for CSVs. But to avoid the intensive manually investigation of each simple SV, we
497 first cluster simple SVs and only inspect clusters with at least two SVs. In particular,
498 we treat each SV as an interval and apply the hierarchical clustering to find interval
499 clusters. The distance measure for clustering is defined as follows:

$$\min(|Iter1.start - Iter2.start|, |Iter1.end - Iter2.end|, |Iter1.center - Iter2.center|) / 1000$$
$$Iter1.center = (Iter1.start + Iter1.end) / 2$$
$$Iter2.center = (Iter2.start + Iter2.end) / 2$$

501 where *Iter* is an SV breakpoint interval, and *Iter.start*, *Iter.end* and *Iter.center*
502 indicate the start, end and center of the interval, respectively. We then use the average
503 method to calculate distance between intervals in two clusters *u* and *v*, which is
504 assigned by:

$$d(u, v) = \sum_{i,j} \frac{d(u[i], v[j])}{(|u| \times |v|)}$$

506 To select a proper threshold for merging clusters from the hierarchical clustering results,
507 we use the threshold from a set of values that could produce the most clusters for each
508 chromosome independently (**Supplementary Table S7, Supplementary Note,**
509 **Supplementary Figure S13-S16**).

510 We further utilize the sequence dot-plot to resolve CSVs based on PacBio long reads.
511 Sequence dot-plot is a classic way to investigate genome rearrangement between
512 species or chromosomes[55]. It applies a k-mer match approach between sequences and
513 keeps matches in a similarity matrix. Thus, we can define the breakpoints and type of a
514 CSV by visualizing the similarity matrix. We use the publicly available interactive
515 sequence dot-plot tool Gepard[56] for this process. Since CSVs are rare and might
516 appear at the minor allele, we create a dot-plot for each long read that spans the
517 corresponding SV cluster. Afterwards, we manually inspect all these dot-plots to
518 identify CSVs, and their breakpoints can be easily obtained from Gepard's interactive
519 user interface (**Supplementary Figure S17**).

520 **Parameter selection for Mako and other methods**

521 Mako run with `minAf = 0.2`, `minFreq = 1`, `minWeight = 10` for real data

522 (NA19240, HG00514, HG00733) and all simulated data (**Supplementary Note**). The
523 minFreq was set to 1 to detect rare events. The minDist is set four times the
524 estimated library fragment average size. And these values are all default settings for
525 Mako. For the cancer cell line (SKBR3), considering the coverage and highly
526 rearranged nature compared with the normal genome, we reduce the cutoff from 0.2 to
527 0.1 and 10 to 5 for minAf and minWeight, respectively, so that the graph could
528 involve more nodes. Signal nodes satisfying either the minAf or minWeight
529 threshold will be included to create the graph. The other selected tools are run under
530 default settings for both simulated and real data (**Supplementary Note**). We use the
531 latest version of TARDIS [26] and the SVelter callset for NA19240 is provided by
532 HGSVC [9] (**Supplementary Note**). For the CSV detection evaluation, all predictions
533 larger than 50p are involved and additional filtering has been done according to the
534 recommended procedures [26-28]. In particular, GRIDSS's callset is filtered by a filter
535 field in VCF header such as ASSEMBLY_TOO_FEW_READ and SVs with
536 coordinates like [57] and [p2, p1] are kept only once. The prediction of SVelter is
537 filtered by a validation score of -1 (**Supplementary Note**).

538 **Performance evaluation**

539 Typically, a correct discovery is defined as a best match between benchmark and
540 predictions, and thus the closest event to the benchmark CSVs with similar size is
541 considered as true positive [58]. However, performance comparison of CSVs is less
542 straightforward than that of simple SVs because of multiple breakpoints involved [27].
543 To address the demand of detecting CSVs as a single event and avoiding redundant
544 predictions [12], the performance is evaluated from two aspects. For example, a CSV
545 with inversion flanked by two deletions is evaluated as three components. Correct
546 prediction of all breakpoints for the three components is considered as all-breakpoint
547 match. Meanwhile, if only one prediction is close to the leftmost and rightmost
548 breakpoints of the CSV with similar size, this prediction is treated as unique-interval
549 match. In the evaluation, the closeness bpDist and size similarity sim between
550 prediction and benchmark are 500bp and 0.7. For example, assume a benchmark
551 [b.start, b.end, b.size], and a prediction [p.start, p.end, p.size]; then a correct
552 prediction will satisfy the following equations:

$$553 \min(|b.start - p.start|, |b.end - p.end|) \leq bpDist$$
$$b.size \times sim \leq p.size \leq b.size \times (2 - sim)$$

554 For simulated data, true positive (TP) is defined as the nearest prediction with similar
555 size to the benchmark, while predictions not in the benchmark are treated as false
556 positives (FP). False negatives (FN) are events in the benchmark set that are not
557 matched by predictions (**Supplementary Note**). Then, the usual measurements can be
558 calculated as follows:

$$recall = TP / (TP + FN)$$

559 $precision = TP / (TP + FP)$

$$F1 = (recall \times precision) / (recall + precision)$$

560 Since it is usually hard to measure the false positives of each tool for real data, we only
561 consider the number of correct discoveries. To fully characterize Mako's performance,
562 we further evaluate it on NA19240 based on PacBio reads by using sensitivity and
563 specificity (**Supplementary Note**). Additionally, because the breakpoints are not as
564 precise as that in the simulation, we relax the size similarity threshold `sim` to 0.5 for
565 real data sensitivity evaluation. To examine Mako's CSV breakpoint offset, we first
566 manually labeled the breakpoints of each CSV from HG00514, HG00733 and
567 NA19240 based on PacBio reads create sequence Doplots (**Supplementary Note**).
568 Secondly, we compare manually labeled breakpoints to Mako reported ones to calculate
569 the offset.

570 **Orthogonal validation of Mako detected CSVs**

571 To evaluate detected CSVs, we used experimental and computational validation as well
572 as manual inspections of HG00733. The raw CSV calls from HG00733 was obtained
573 by selecting events with more than one link types observed in the subgraph, resulting
574 in 609 CSVs. To design primers, Primer3 (<https://github.com/primer3-org/primer3>)
575 was used in conjunction with internal software to design and select PCR primers, where
576 the optimal primer size was set to 23bp. In particular, we extend Mako detected
577 breakpoints by 500bp to select primers with average GC contents close to 50% and a
578 predicted melting temperature 60 °C. Primers were then selected within the extended
579 distance but 200bp outside of the boundaries of the breakpoints defined by Mako
580 (**Supplementary Figure S18**). If duplication and inversion like edges were found in
581 the subgraph, primers were also designed on the reverse complementary strand. All
582 primer pairs were tested for their uniqueness across the human genome using In Silico
583 PCR from UCSC Genome Browser. BLAT (<https://users.soe.ucsc.edu/~kent/>) search
584 was also performed at the same time to make sure all primer candidates have only one
585 hit in the human genome. If the above procedure does not result in a valid primer pair,

586 the size of the regions for which primers are designed was increased from 500bp to
587 650bp and all process were repeated to search for primers (primers are in
588 **Supplementary Table S3**). PCR amplifications were performed in a volume of 25 ul
589 concentration of reagents, consisting of 1) 1x of 10x *Ex Taq* Buffer (Mg²⁺ Plus); 2) 0.4
590 mM of dNTP mix, 0.4 uM for each primer; 3) 0.75 units of Ex Taq DNA polymerase
591 (TakaRa, Japan) and 4) 30 ng of DNA. The amplification cycle was performed in
592 Mastercycler® nexus gradient (Eppendorf, Germany), including 1) 5 minutes' pre-
593 denaturation at 94°C; 2) 35 cycles of denaturation at 94°C for 45 seconds, annealing 45
594 seconds according to different T_M value of each primer and elongation at 72°C for 90
595 seconds; 3) followed by 10 minutes' extension at 72°C. The amplification products were
596 separated by electrophoresis in 1.5% agarose gels with CellProTMDNA-Red
597 (InCellGene LLC, USA) and bands were visualized under the UV light. Then, we
598 selected products with the expected product size and bright electrophoretic bands
599 (**Supplementary Figure S19**, all results in **Supplementary File 3**), which were further
600 purified and cloned into the expression vector pEASY-T1 (Transgene, China). The
601 positive clones containing the targeted fragments were send to TsingKe Biological
602 Technology Company for Sanger sequencing. The Sanger sequencing data were aligned
603 against the reference allele of the CSV site and visualized with Gepard for breakpoint
604 inspection.

605 We used HiFi reads from HGVC to manually reconstruct each CSVs. Similar to the
606 procedure of creating the benchmark CSV for NA19240 and SKBR3, SAMtools was
607 used to get the HiFi reads spanning the breakpoints. Afterwards, Gepard was applied to
608 create the sequence dotplot between each read and the reference genome. We than go
609 through all the sequence dotpot to validate CSVs detected by Mako (**Supplementary**
610 **Figure S17, Supplementary Note, Supplementary File 2**). The validation rate
611 measured whether Mako detected subgraphs contained different types of breakpoint
612 connection edges. For dotplots with 'messy' regions, they could produce duplication
613 and insertion like breakpoint connections based on short-read sequencing. Therefore, it
614 was difficult or even impossible for short reads to distinguish between distinct complex
615 events and those detected at repeat regions. To characterize these events based on long-
616 read sequencing, we introduced a three steps workflow as follows:

617 Step 1. Identifying event breakpoints inside the 'messy' regions in the dotplots.

618 Those outside the 'messy' regions were considered as distinct complex events.

619 Step 2. We defined 3 dotplot patterns (**Supplementary Figure S20**) to classify
620 ‘messy’ events to CSVs, where the x-axis and y-axis are REF and ALT sequence,
621 respectively. Region 1, 2 and 3 indicates regions where extra segments could be
622 found. Especially, region 2 in each case indicates the ‘messy’ region caused by
623 repeats.

624 ○ Case A: Blue segments indicate an insertion event with single
625 breakpoint on the reference. A CSV should contain at least one
626 duplicated segments (purple) in region 1, 2 or 3. Example events include
627 chr1:206,924,211-206,924,525 (**Supplementary File 2, page 89**) and
628 etc.

629 ○ Case B: Blue segments indicate repeat expansion on the ALT sequence.
630 A CSV should contain extra segments in region 1, 2 or 3. Example
631 events include chr1:1,382,295-1,382,470 (**Supplementary File 2, page**
632 **145**) and etc.

633 ○ Case C: Blue segments overlap on the REF, but have a gap on the ALT
634 sequence. This type of events could be interpreted as insertion with
635 duplications, which is considered as complex event. We also observed
636 some CSV contained segments (purple) in region 1, 2 or 3. Example
637 events include chr3:50,311,835-50,312,092 (**Supplementary File 2,**
638 **page 226**) and etc.

639 Step 3. We further investigate events that failed the examination in Step 2
640 according to 2 dotplot patterns (**Supplementary Figure S21**).

641 ○ Case A: A simple insertion event, where the breakpoint locates
642 inside tandem repeats (region 2) and other segments cannot be found.

643 ○ Case B: Regular repeat expansion (purple segments) in ALT
644 sequence.

645 For computational validation, we obtained ONT reads of HG00733 from HGSVC and
646 applied VaPoR [59], an independent structural variants validation method, to validate
647 these CSVs (**Supplementary Note**). VaPoR is able to validate calls based predicted
648 region and types with a confidence score. VaPoR labeled NAs and 0 to some of the
649 inconclusive events due to highly repetitive sequence and unclear recurrent pattern that
650 can be observed (**Supplementary Figure S22**). We termed the above procedure as
651 ONT validation. Besides, we obtained HiFi assemblies from HGSVC and applied a K-
652 mer based breakpoint examination and calculate the breakpoint shifts. Specifically,

653 CSV spanning H1 or H2 contig sequence (ALT) and reference (REF) sequence were
654 extracted from alignment and GRCh38, respectively. We first identified the matched
655 segments between ATL and REF through K-mer (k=32bp) realignment as well as sorted
656 these segments according to their position on reference. Afterwards, we marked the
657 unmatched or gap regions, from which, we calculated the breakpoints and size
658 similarity. A CSV was considered valid if both left and right breakpoint difference are
659 smaller than 500bp. This constrain was used by Truvari
660 (<https://github.com/spiralgenetics/truvari/>), a standard benchmarking tool used by
661 Genome In A Bottle (GIAB). The implementation of K-mer validation is available at
662 Mako GitHub site. Breakpoint comparison of experimental and K-mer validation were
663 listed in **Supplementary Table S8**, which was used to calculate the breakpoint
664 resolution. Because VaPoR is able to report Valid, NA and 0 events but not to report
665 the breakpoint based on ONT (**Supplementary Table S9**), we did not include VaPoR's
666 results in the breakpoint shift analysis.

667 **Code availability**

668 Mako is implemented in Java 1.8, and it is available at
669 <https://github.com/jiadong324/Mako>. It is free for non-commercial use by academic,
670 government, and non-profit/not-for-profit institutions. A commercial version of the
671 software is available and licensed through Xi'an Jiao-tong University. All scripts used
672 in this study are also included in the Github repository, and a detailed description of
673 using these scripts and other tools is provided in **Supplementary Note**.

674 **Data availability**

675 All materials or datasets used in this study are publicly available and their links are
676 listed in **Supplementary Note**.

677 **Authors' contributions**

678 In particular, KY conceived and designed the study; JL, XY and WK developed the
679 graph-based pattern growth algorithm for SV breakpoint discovery; JL and TX created
680 the CSV benchmarks for real data and manually reconstructed CSVs. YJ, CZ, QH and
681 MR performed the wet lab experimental validation; SW performed the computational
682 validation; JL, XY, CZ, LG, WK and KY wrote the manuscript; EE, SD, CL provides
683 the ONT reads and HiFi reads. HGSVC produced the HiFi assembly and all authors
684 contributed to the critical revision of the manuscript and approved the final version.

685 **Competing interests**

686 The authors have declared no competing interests.

687 **Acknowledgments**

688 This study was supported by the National Key R&D Program of China (grand NO.
689 2018YFC0910400, 2017YFC0907500), the National Science and Technology Major
690 Project of China (grand NO. 2018ZX10302205), the National Science Foundation of
691 China (grand NO. 31671372, 61702406 and 31701739) and the “World-Class
692 Universities and the Characteristic Development Guidance Funds for the Central
693 Universities”. Supported by Shanghai Municipal Science and Technology Major
694 Project (Grant No.2017SHZDZX01).

695 **ORCID**

696 ^aORCID: 0000-0002-8116-5901 (Lin J)

697 ^bORCID: 0000-0002-5118-7755 (Yang X)

698 ^cORCID: 0000-0001-8860-0390 (Kosters W)

699 ^dORCID: 0000-0003-3194-1834 (Xu T)

700 ^eORCID: 0000-0002-4966-0574 (Jia Y)

701 ^fORCID: 0000-0003-4482-8128 (Wang S)

702 ^gORCID: 0000-0003-2401-8443 (Zhu Q)

703 ^hORCID: 0000-0001-5428-0018 (Ryan M)

704 ⁱORCID: 0000-0001-6100-3481 (Guo L)

705 ^jORCID: 0000-0002-5238-083X (Zhang C)

706 ^kORCID: 0000-0001-7317-6662 (Lee C)

707 ^lORCID: 0000-0001-7629-8331 (Devine S)

708 ^mORCID: 0000-0002-8246-4014 (Eichler E)

709 ⁿORCID: 0000-0002-2851-6741 (Ye K)

710

711 **References**

- 712 [1] Ye K, Schulz MH, Long Q, Apweiler R, Ning Z. Pindel: a pattern growth approach
713 to detect break points of large deletions and medium sized insertions from paired-end
714 short reads. *Bioinformatics* 2009;25:2865-71.
- 715 [2] Rausch T, Zichner T, Schlattl A, Stutz AM, Benes V, Korbel JO. DELLY: structural
716 variant discovery by integrated paired-end and split-read analysis. *Bioinformatics*
717 2012;28:i333-i9.
- 718 [3] Layer RM, Chiang C, Quinlan AR, Hall IM. LUMPY: a probabilistic framework
719 for structural variant discovery. *Genome Biol* 2014;15:R84.
- 720 [4] Chen K, Wallis JW, McLellan MD, Larson DE, Kalicki JM, Pohl CS, et al.
721 BreakDancer: an algorithm for high-resolution mapping of genomic structural variation.
722 *Nat Methods* 2009;6:677-81.
- 723 [5] Cameron DL, Di Stefano L, Papenfuss AT. Comprehensive evaluation and
724 characterisation of short read general-purpose structural variant calling software. *Nat*
725 *Commun* 2019;10:3240.
- 726 [6] Kosugi S, Momozawa Y, Liu X, Terao C, Kubo M, Kamatani Y. Comprehensive
727 evaluation of structural variation detection algorithms for whole genome sequencing.
728 *Genome Biol* 2019;20:117.
- 729 [7] Chen X, Schulz-Trieglaff O, Shaw R, Barnes B, Schlesinger F, Kallberg M, et al.
730 Manta: rapid detection of structural variants and indels for germline and cancer
731 sequencing applications. *Bioinformatics* 2016;32:1220-2.
- 732 [8] Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, et al.
733 An integrated map of structural variation in 2,504 human genomes. *Nature*
734 2015;526:75-81.
- 735 [9] Chaisson MJP, Sanders AD, Zhao X, Malhotra A, Porubsky D, Rausch T, et al.
736 Multi-platform discovery of haplotype-resolved structural variation in human genomes.
737 *Nat Commun* 2019;10:1784.
- 738 [10] Gao R, Davis A, McDonald TO, Sei E, Shi X, Wang Y, et al. Punctuated copy
739 number evolution and clonal stasis in triple-negative breast cancer. *Nat Genet*
740 2016;48:1119-30.
- 741 [11] Yates LR, Knappskog S, Wedge D, Farmery JHR, Gonzalez S, Martincorena I, et
742 al. Genomic Evolution of Breast Cancer Metastasis and Relapse. *Cancer Cell*
743 2017;32:169-84 e7.

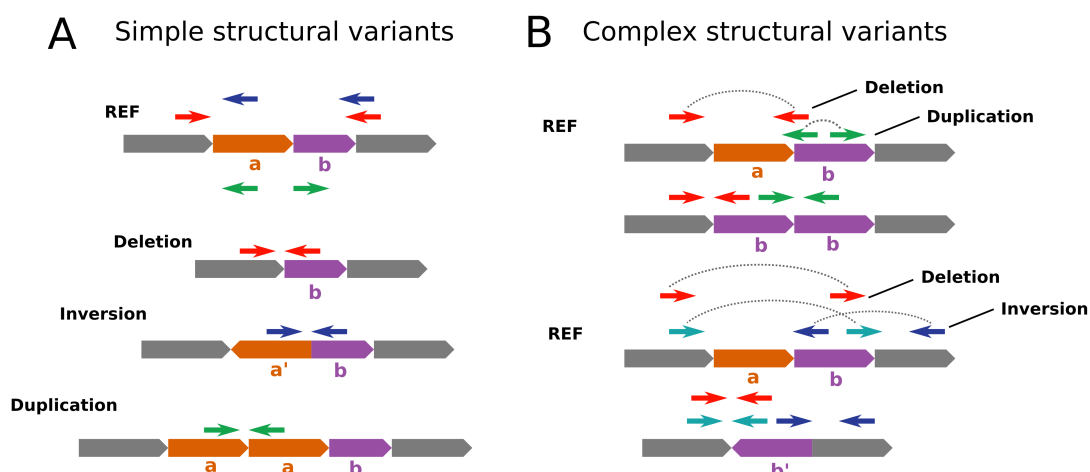
- 744 [12] Quinlan AR, Hall IM. Characterizing complex structural variation in germline and
745 somatic genomes. *Trends Genet* 2012;28:43-53.
- 746 [13] Nattestad M, Goodwin S, Ng K, Baslan T, Sedlazeck FJ, Rescheneder P, et al.
747 Complex rearrangements and oncogene amplifications revealed by long-read DNA and
748 RNA sequencing of a breast cancer cell line. *Genome Res* 2018;28:1126-35.
- 749 [14] Sanchis-Juan A, Stephens J, French CE, Gleadall N, Megy K, Penkett C, et al.
750 Complex structural variants in Mendelian disorders: identification and breakpoint
751 resolution using short- and long-read genome sequencing. *Genome Med* 2018;10:95.
- 752 [15] Greer SU, Nadauld LD, Lau BT, Chen J, Wood-Bouwens C, Ford JM, et al. Linked
753 read sequencing resolves complex genomic rearrangements in gastric cancer metastases.
754 *Genome Med* 2017;9:57.
- 755 [16] Lee JJ, Park S, Park H, Kim S, Lee J, Lee J, et al. Tracing Oncogene
756 Rearrangements in the Mutational History of Lung Adenocarcinoma. *Cell*
757 2019;177:1842-57 e21.
- 758 [17] Collins RL, Brand H, Redin CE, Hanscom C, Antolik C, Stone MR, et al. Defining
759 the diverse spectrum of inversions, complex structural variation, and chromothripsis in
760 the morbid human genome. *Genome Biol* 2017;18:36.
- 761 [18] Carvalho CM, Lupski JR. Mechanisms underlying structural variant formation in
762 genomic disorders. *Nat Rev Genet* 2016;17:224-38.
- 763 [19] Baca SC, Prandi D, Lawrence MS, Mosquera JM, Romanel A, Drier Y, et al.
764 Punctuated evolution of prostate cancer genomes. *Cell* 2013;153:666-77.
- 765 [20] Korbel JO, Campbell PJ. Criteria for inference of chromothripsis in cancer
766 genomes. *Cell* 2013;152:1226-36.
- 767 [21] Sanders AD, Meiers S, Ghareghani M, Porubsky D, Jeong H, van Vliet M, et al.
768 Single-cell analysis of structural variations and complex rearrangements with tri-
769 channel processing. *Nat Biotechnol* 2019.
- 770 [22] Carvalho CMB, Lupski JR. Mechanisms underlying structural variant formation
771 in genomic disorders. *Nature Reviews Genetics* 2016;17:224-38.
- 772 [23] Malhotra A, Lindberg M, Faust GG, Leibowitz ML, Clark RA, Layer RM, et al.
773 Breakpoint profiling of 64 cancer genomes reveals numerous complex rearrangements
774 spawned by homology-independent mechanisms. *Genome Res* 2013;23:762-76.
- 775 [24] Ye K, Wang J, Jayasinghe R, Lameijer EW, McMichael JF, Ning J, et al.
776 Systematic discovery of complex insertions and deletions in human cancers. *Nat Med*
777 2016;22:97-104.

- 778 [25] Zhang CZ, Leibowitz ML, Pellman D. Chromothripsis and beyond: rapid genome
779 evolution from complex chromosomal rearrangements. *Genes Dev* 2013;27:2513-30.
- 780 [26] Soylev A, Le TM, Amini H, Alkan C, Hormozdiari F. Discovery of tandem and
781 interspersed segmental duplications using high-throughput sequencing. *Bioinformatics*
782 2019;35:3923-30.
- 783 [27] Zhao X, Emery SB, Myers B, Kidd JM, Mills RE. Resolving complex structural
784 genomic rearrangements using a randomized approach. *Genome Biol* 2016;17:126.
- 785 [28] Cameron DL, Schroder J, Penington JS, Do H, Molania R, Dobrovic A, et al.
786 GRIDSS: sensitive and specific genomic rearrangement detection using positional de
787 Bruijn graph assembly. *Genome Res* 2017;27:2050-60.
- 788 [29] Marschall T, Costa IG, Canzar S, Bauer M, Klau GW, Schliep A, et al. CLEVER:
789 clique-enumerating variant finder. *Bioinformatics* 2012;28:2875-82.
- 790 [30] Arthur JG, Chen X, Zhou B, Urban AE, Wong WH. Detection of complex
791 structural variation from paired-end sequencing data. *bioRxiv* 2017:200170.
- 792 [31] Ottaviani D, LeCain M, Sheer D. The role of microhomology in genomic structural
793 variation. *Trends Genet* 2014;30:85-94.
- 794 [32] Kramara J, Osia B, Malkova A. Break-Induced Replication: The Where, The Why,
795 and The How. *Trends Genet* 2018;34:518-31.
- 796 [33] Hartlerode AJ, Willis NA, Rajendran A, Manis JP, Scully R. Complex Breakpoints
797 and Template Switching Associated with Non-canonical Termination of Homologous
798 Recombination in Mammalian Cells. *PLoS Genet* 2016;12:e1006410.
- 799 [34] Zhou W, Zhang F, Chen X, Shen Y, Lupski JR, Jin L. Increased genome instability
800 in human DNA segments with self-chains: homology-induced structural variations via
801 replicative mechanisms. *Hum Mol Genet* 2013;22:2642-51.
- 802 [35] Yang L, Luquette LJ, Gehlenborg N, Xi R, Haseley PS, Hsieh CH, et al. Diverse
803 mechanisms of somatic structural variations in human cancer genomes. *Cell*
804 2013;153:919-29.
- 805 [36] Chen W, McKenna A, Schreiber J, Haeussler M, Yin Y, Agarwal V, et al.
806 Massively parallel profiling and predictive modeling of the outcomes of CRISPR/Cas9-
807 mediated double-strand break repair. *Nucleic Acids Res* 2019;47:7989-8003.
- 808 [37] Allen F, Crepaldi L, Alsinet C, Strong AJ, Kleshchevnikov V, De Angeli P, et al.
809 Predicting the mutations generated by repair of Cas9-induced double-strand breaks. *Nat*
810 *Biotechnol* 2018.

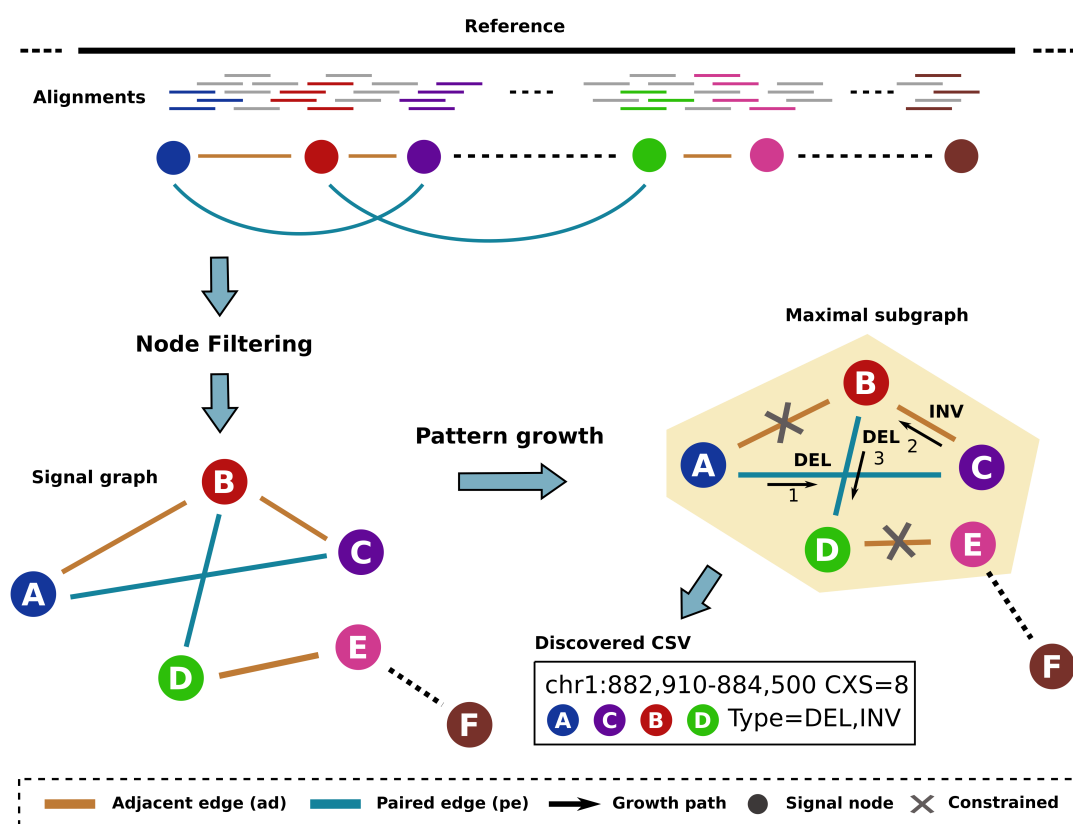
- 811 [38] Quigley DA, Dang HX, Zhao SG, Lloyd P, Aggarwal R, Alumkal JJ, et al.
812 Genomic Hallmarks and Structural Variation in Metastatic Prostate Cancer. *Cell*
813 2018;175:889.
- 814 [39] Fraser M, Sabelnykova VY, Yamaguchi TN, Heisler LE, Livingstone J, Huang V,
815 et al. Genomic hallmarks of localized, non-indolent prostate cancer. *Nature*
816 2017;541:359-64.
- 817 [40] Sedlazeck FJ, Rescheneder P, Smolka M, Fang H, Nattestad M, von Haeseler A,
818 et al. Accurate detection of complex structural variations using single-molecule
819 sequencing. *Nat Methods* 2018;15:461-8.
- 820 [41] Kehr B, Melsted P, Halldorsson BV. PopIns: population-scale detection of novel
821 sequence insertions. *Bioinformatics* 2016;32:961-7.
- 822 [42] Han J. *Data Mining: Concepts and Techniques*. Morgan Kaufmann Publishers Inc.,
823 2005.
- 824 [43] Liao VCC, Chen MS. DFSP: a Depth-First SPelling algorithm for sequential
825 pattern mining of biological sequences. *Knowledge and Information Systems*
826 2014;38:623-39.
- 827 [44] Tsai HP, Yang DN, Chen MS. Mining Group Movement Patterns for Tracking
828 Moving Objects Efficiently. *Ieee Transactions on Knowledge and Data Engineering*
829 2011;23:266-81.
- 830 [45] Huang Y, Zhang LQ, Zhang PS. A framework for mining sequential patterns from
831 spatio-temporal event data sets. *Ieee Transactions on Knowledge and Data Engineering*
832 2008;20:433-48.
- 833 [46] Ye K, Kusters WA, Ijzerman AP. An efficient, versatile and scalable pattern
834 growth approach to mine frequent patterns in unaligned protein sequences.
835 *Bioinformatics* 2007;23:687-93.
- 836 [47] Pei J, Han J, Wang W. Constraint-based sequential pattern mining: the pattern-
837 growth methods. *Journal of Intelligent Information Systems* 2007;28:133-60.
- 838 [48] Pei J, Han JW, Mortazavi-Asl B, Wang JY, Pinto H, Chen QM, et al. Mining
839 sequential patterns by pattern-growth: The PrefixSpan approach. *Ieee Transactions on*
840 *Knowledge and Data Engineering* 2004;16:1424-40.
- 841 [49] Hu H, Yan X, Huang Y, Han J, Zhou XJ. Mining coherent dense subgraphs across
842 massive biological networks for functional discovery. *Bioinformatics* 2005;21 Suppl
843 1:i213-21.

- 844 [50] Li H, Homer N. A survey of sequence alignment algorithms for next-generation
845 sequencing. *Brief Bioinform* 2010;11:473-83.
- 846 [51] Li H, Durbin R. Fast and accurate short read alignment with Burrows-Wheeler
847 transform. *Bioinformatics* 2009;25:1754-60.
- 848 [52] Bolognini D, Sanders A, Korbel JO, Magi A, Benes V, Rausch T. VISOR: a
849 versatile haplotype-aware structural variant simulator for short and long read
850 sequencing. *Bioinformatics* 2019.
- 851 [53] McPherson A, Wu C, Wyatt AW, Shah S, Collins C, Sahinalp SC. nFuse:
852 discovery of complex genomic rearrangements in cancer using high-throughput
853 sequencing. *Genome Res* 2012;22:2250-61.
- 854 [54] Dzamba M, Ramani AK, Buczkowicz P, Jiang Y, Yu M, Hawkins C, et al.
855 Identification of complex genomic rearrangements in cancers using CouGaR. *Genome*
856 *Res* 2017;27:107-17.
- 857 [55] Delcher AL, Phillippy A, Carlton J, Salzberg SL. Fast algorithms for large-scale
858 genome alignment and comparison. *Nucleic Acids Res* 2002;30:2478-83.
- 859 [56] Krumsiek J, Arnold R, Rattei T. Gepard: a rapid and sensitive tool for creating
860 dotplots on genome scale. *Bioinformatics* 2007;23:1026-8.
- 861 [57] Carvalho CM, Ramocki MB, Pehlivan D, Franco LM, Gonzaga-Jauregui C, Fang
862 P, et al. Inverted genomic segments and complex triplication rearrangements are
863 mediated by inverted repeats in the human genome. *Nat Genet* 2011;43:1074-81.
- 864 [58] Audano PA, Sulovari A, Graves-Lindsay TA, Cantsilieris S, Sorensen M, Welch
865 AE, et al. Characterizing the Major Structural Variant Alleles of the Human Genome.
866 *Cell* 2019;176:663-75 e19.
- 867 [59] Zhao X, Weber AM, Mills RE. A recurrence-based approach for validating
868 structural variation using long-read sequencing technology. *Gigascience* 2017;6:1-9.
- 869
- 870
- 871
- 872

873 **Figure legends**

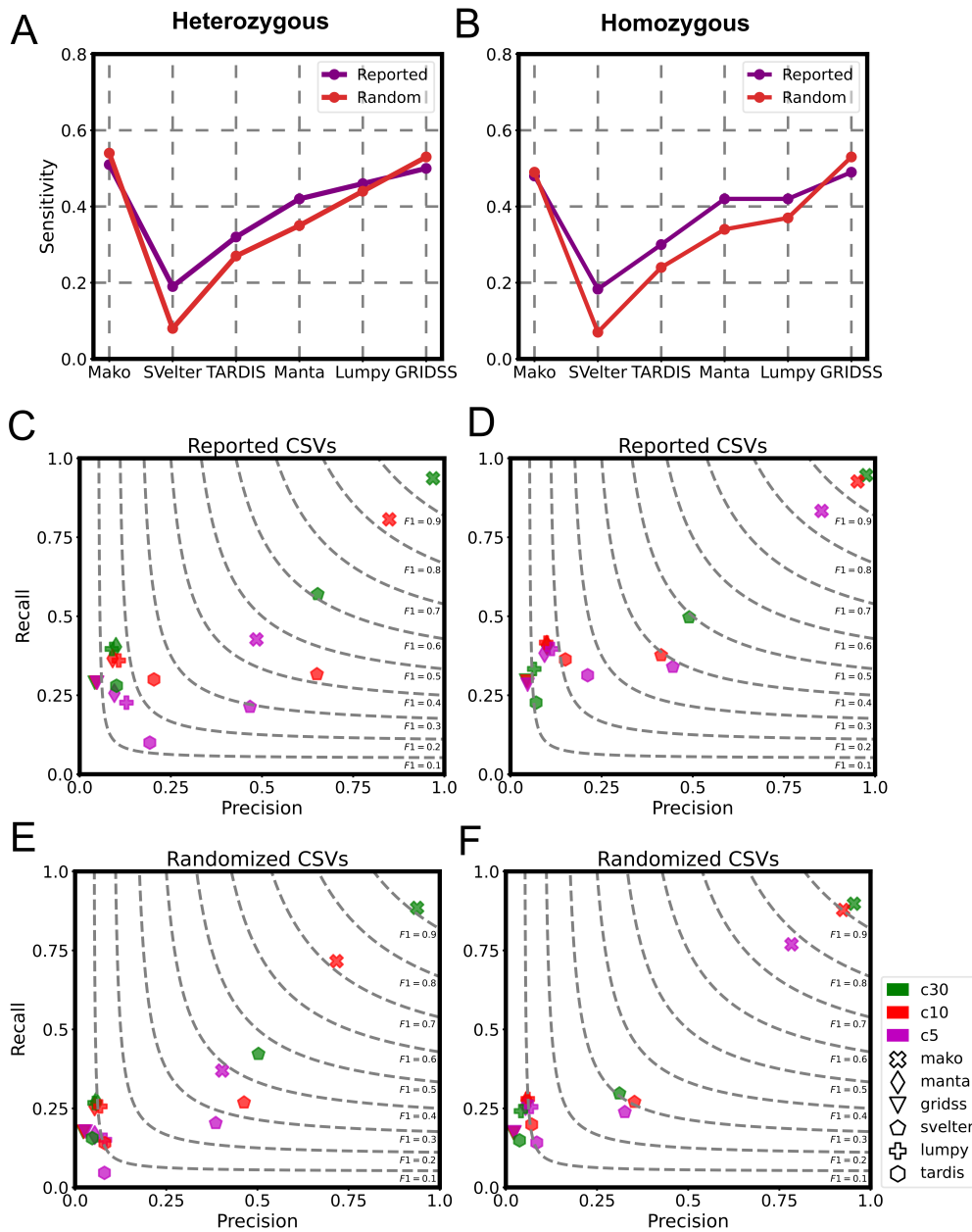


875 **Figure 1. Explanation of simple and complex structure variants alignment**
 876 **models derived from abnormal read pairs**
 877 (A) Three common simple SV and their corresponding abnormal read pair alignment
 878 on the reference genome, representing by red, blue and green arrows. (B) The alignment
 879 signature of two CSVs, each of them involves two types of signature that can be
 880 matched by simple SV alignment model.



882 **Figure 2. Overview of Mako for identifying CSVs from NGS data.**

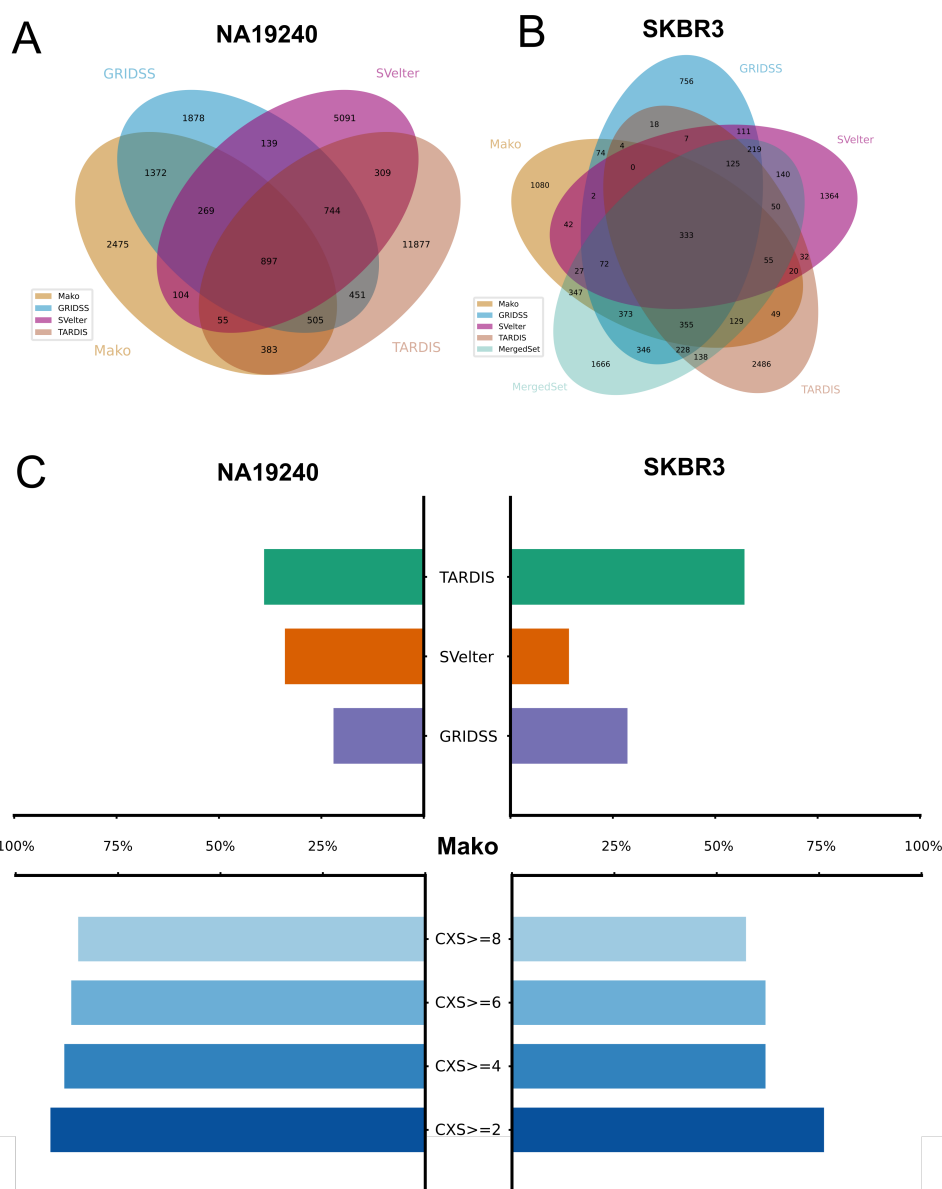
883 Mako first builds a signal graph by collecting abnormal aligned reads as nodes and their
 884 edge connections are provided by paired-end alignment and split alignment. Afterwards,
 885 Mako utilizes the pattern growth approach to find a maximal subgraph as potential CSV
 886 site. In the example output, the maximal subgraph contains A, B, C, D, whereas F is not
 887 able to be appended because of none existing edge (dashed line). The CSV is derived from
 888 this subgraph with estimate breakpoints and CXS score, where the discovered CSV
 889 subgraph contains four different nodes, one E_{ae} edge and two E_{pe} edges of type DEL
 890 and INV, thus $CXS = 8$.



891

892 **Figure 3. Performance comparison on simulated CSVs with different match**
 893 **criteria.**

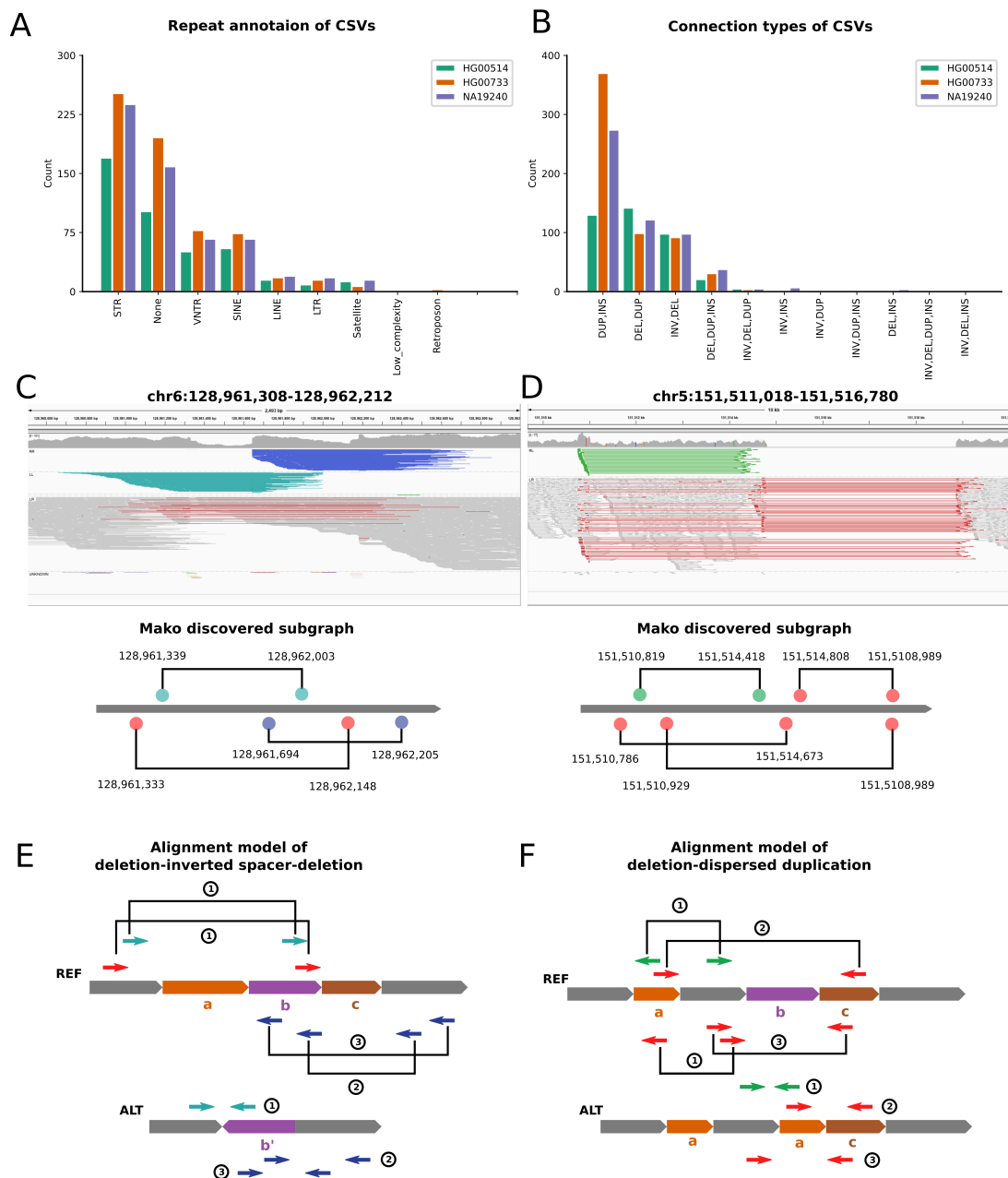
894 All-breakpoint match (**A** and **B**) and unique-interval match (**C-F**) evaluation of selected
 895 tools on simulated CSVs. (**A**) The sensitivity of detecting heterozygous CSVs
 896 breakpoints. (**B**) The sensitivity of detecting homozygous CSVs breakpoints. The red
 897 and purple curve indicate randomized and reported CSVs, respectively. (**C**) Evaluation
 898 of reported heterozygous CSV simulation. (**D**) Evaluation of reported homozygous
 899 CSV simulation. (**E**) Evaluation of randomized heterozygous CSV simulation. (**F**)
 900 Evaluation of randomized homozygous CSV simulation. From (**C**) to (**F**), the
 901 performance is evaluated by recall (y-axis), precision (x-axis) and F1-score (dotted
 902 lines). The right top corner of the plot indicates better performance. The c5-c30
 903 indicates coverage, e.g. c5 indicates 5X coverage.



905 **Figure 4. Overview of performance on NA19240 and SKBR3 for Mako, GRIDSS,**

906 **SVelter and TARDIS.**

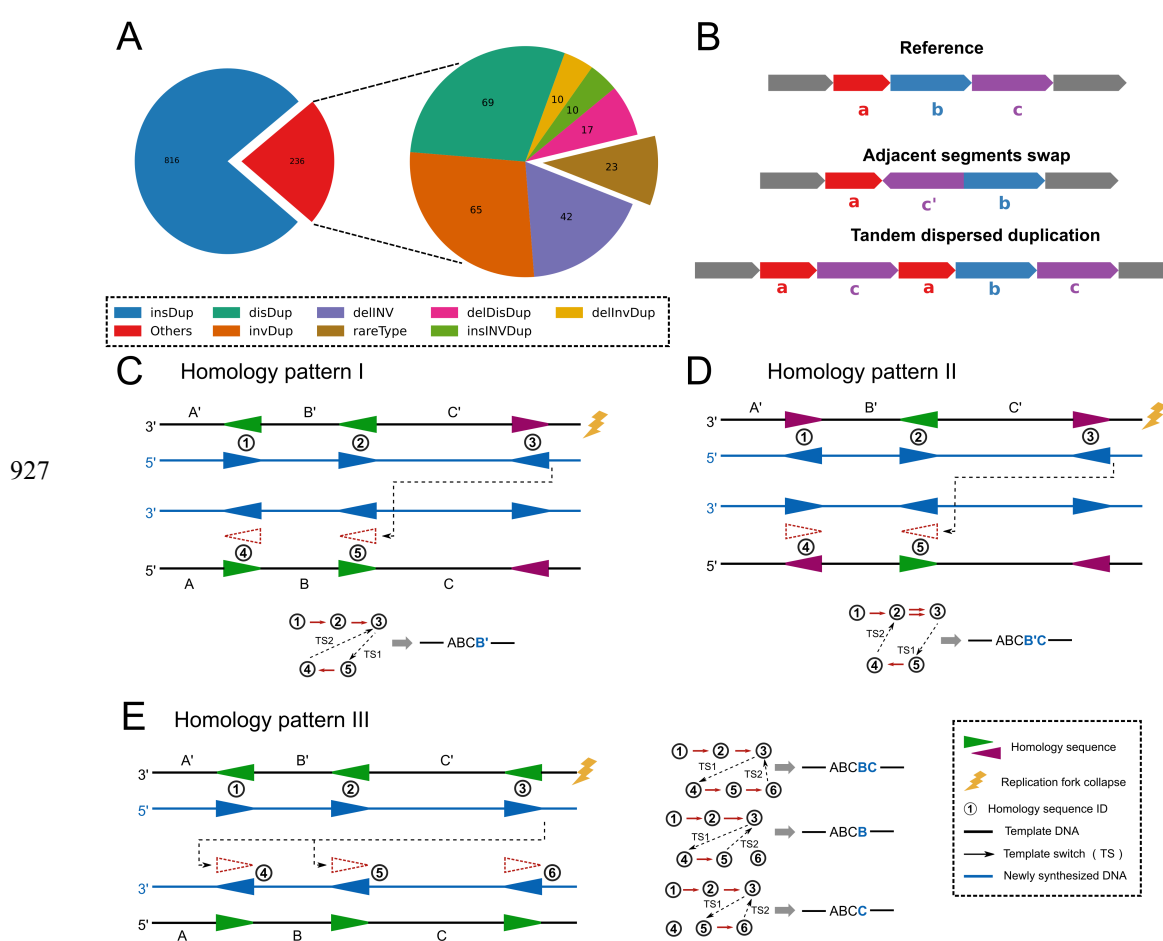
907 (A) and (B) are the Venn diagram of 50% reciprocal overlap between callsets for both
 908 NA19240 and SKBR3. They are created by a publicly available tool Intervene with –
 909 bedtools-options enabled. (B) The MergedSet is the callset provided by the publication.
 910 (C) The percentage of completely and uniquely discovered CSVs from the NA19240
 911 and SKBR3. The results of Mako (bottom panel) are shown according to different CXS
 912 threshold.



914 **Figure 5. Repeat annotation and types of CSVs with two representative examples**
 915 **identified by Mako.**

916 (A) is repeat annotation and (B) is detected connection types of CSVs, respectively. The

917 top panel of (C) and (D) are IGV view of the two events and the alignments are grouped
 918 by pair orientation. The dark blue shows reverse-reverse alignments, light blue is the
 919 forward-forward alignments, green is the reverse-forward alignments and the red
 920 indicates alignment of large insert size. The bottom panel of (C) and (D) are sub-graph
 921 structure discovered by Mako. The colored circles and solid lines are nodes and edges
 922 in the sub-graph. (E) The alignment model of deletions with inverted spacer. (F) The
 923 alignment model of deletion associated with dispersed duplication. In (E) and (F), short
 924 arrows are paired-end reads that span breakpoint junctions, and their alignment are
 925 shown on the reference genome with corresponding ID in circle. Noted that a single ID
 926 may have more than one corresponding abnormal alignment types on the reference.



928 **Figure 6. Overview of Mako's CSV discoveries from three healthy samples and**
 929 **proposed CSV formation mechanisms.**

930 (A) Summary of discovered CSV types, these types are reconstructed by HiFi PacBio
 931 reads, where a type with less than 10 events was summarized as rareType. (B) Diagrams
 932 of two novel and rare CSV types discovered by Mako. In particular, Mako finds three
 933 events of adjacent segments swap and only one tandem dispersed duplication. (C-E)

934 Replication diagram explains the impact of homology pattern for MMBIR produced
935 CSVs. In these diagrams, sequence *ABC* has been replicated before the replication fork
936 collapse (flash symbol). The single strand DNA at the DNA double strand break (DSB)
937 starts searching for homology sequence (purple and green triangle) to repair. The above
938 procedure is explicitly explained as a replication graph, from which, nodes are
939 homology sequences and edges keep track of the template switch (dotted arrow lines)
940 as well as the normal replication at different strand (red lines). If there are two red lines
941 between two nodes, the sequence between these two nodes will be replicate twice as
942 shown in **(D)**.
943

944 **Tables**

945 **Table 1. Summary of benchmark CSVs. The CSV type abbreviations and their**
 946 **corresponding descriptions are also listed.**

Benchmark summaries			
Type	NA19240		Description
	NA19240	SKBR3	
Disdup	15	12	Dispersed duplication
Invdup	18	-	Inverted duplication
DelInv	7	5	Deletion associated with inversion
DelDisdup	5	1	Deletion associated with dispersed duplication
DelInvdup	1	-	Deletion associated with inverted duplication
DisdupInvdup	2	2	Dispersed duplication with inverted duplication
InsInv	1	-	Insertion associated with inversion
Tantrans	1	-	Adjacent segments swap
DelSpaDel	8	1	Two deletions with inverted or non-inverted spacer
TanDisdup	1	-	Tandem dispersed duplications

947

948 **Table 2. Summary of experimentally validated CSVs.**

Chromosome	Start	End	Mako Type
Chr1	81,194,398	81,195,874	DEL, INV
Chr2	119,659,504	119,661,322	DUP, INS
Chr3	146,667,093	146,667,284	DEL, DUP
Chr5	141,480,327	141,483,116	DEL, DUP
Chr7	1,940,931	1,941,009	DUP, INS
Chr9	29,591,409	29,593,057	DEL, INV
Chr10	14,568,488	14,568,677	DUP, INS
Chr12	71,315,482	71,316,928	DEL, INV
Chr12	77,989,900	77,994,324	DEL, INV
Chr13	74,340,759	74,342,810	DEL, DUP
Chr16	78,004,459	78,007,456	DEL, DUP
Chr17	34,854,438	34,855,851	DEL, INV
Chr17	48,538,270	48,540,171	DEL, DUP
Chr18	72,044,575	72,045,937	DEL, DUP
Chr21	26,001,844	26,001,844	DEL, INV

949

950

951

952

953

954

955 **Table 3. Summary of experimental and computational validation as well as manual**
956 **inspection for CSVs.**

Validation Strategy	Total	Valid	Invalid	Inconclusive
Experimental (PCR succeeded)	21	15 (71%)	6 (29%)	-
ONT reads		256 (42%)	-	353 (58%)
Computational	609	414 (68%)	191 (32%)	-
HiFi contig		414 (68%)	191 (32%)	-
ONT reads or HiFi contig		544 (87%)	76 (13%)	-
Manual	609	440 (72%)	169 (28%)	-
HiFi reads		440 (72%)	169 (28%)	-

957

958 **Supplementary material**

959 **Supplementary Note** contains supplementary information for MATERIALS and
960 METHODS.

961 **Supplementary Figures** contains the supplementary figures for this study.

962 **Supplementary Table S1** provides the benchmark CSVs, SV clustering summary and
963 examples used to illustrate Mako CSV subgraph.

964 **Supplementary Table S2** provides Mako detected CSVs for HG00733, HG00514 and
965 NA19240.

966 **Supplementary Table S3** provides events with successfully designed primers.

967 **Supplementary Table S4** provides the summary of experimental and computational
968 validation as well as manual inspections of HG00733.

969 **Supplementary Table S5** provides the details of breakpoints for the two examples in
970 Figure 5C to 5F.

971 **Supplementary Table S6** provides the results of manual inspections of HG00733,
972 HG00514 and NA19240 based on PacBio HiFi reads.

973 **Supplementary Table S7** provides parameters used for creating the CSV benchmarks
974 for NA19240 and SKBR3.

975 **Supplementary Table S8** provides experimental and computational evaluated
976 breakpoints, which was used for breakpoint shift analysis.

977 **Supplementary Table S9** provides the details of VaPoR results of HG00733.

978 **Supplementary File 1** provides the IGV view and PacBio reads dotplot of each
979 benchmark CSVs.

980 **Supplementary File 2** provides the PacBio HiFi reads dotplots for manual inspections
981 of HG00733.

982 **Supplementary File 3** provides the PCR results and visualization of CSV breakpoint
983 validated through Sanger sequencing.

984

985 **Authors from HGSVC**

986 Mark B. Gerstein¹, Ashley D. Sanders², Micheal C.Zody³, Michael E. Talkowski⁴, Ryan
987 E. Mills⁵, Jan O. Korbel², Tobias Marschall⁶, Peter Ebert⁶, Peter A. Audano⁷, Bernardo
988 Rodriguez-Martin⁸, David Porubsky⁷, Marc Jan Bonder^{8,9}, Arvis Sulovari⁷, Jana Ebler⁶,
989 Weichen Zhou⁵, Rebecca Serra Mari⁶, Feyza Yilmaz¹⁰, Xuefang Zhao⁴, PingHsun
990 Hsieh⁷, Joyce Lee¹¹, Sushant Kumar¹, Tobias Rausch⁸, Yu Chen¹², Zechen Chong¹²,

991 Jingwen Ren¹³, Martin Santamarina¹⁴, Wolfram Jops⁸, Hufsa Sshraf⁶, Katherine
992 M.Munson⁷, Mark J.P. Chaisson¹³, Junjie Chen¹⁵, Xinghua Shi¹⁵, Harrison Bran¹⁶,
993 Aaron M.Wenger¹⁷, William T.Harvey⁷, Patrick Hansenfeld⁸, Allison Regier¹⁸, Haley
994 Abel¹⁷, Ira Hall¹⁷, Paul Flicek¹⁸, Alex R. Hastie¹¹, Susan Fairely¹⁸

995 ¹Program in Computational Biology and Bioinformatics, Yale University, BASS
996 432&437, 266 Whitney Avenue, New Haven, CT 06520, USA. ²European Molecular
997 Biology Laboratory (EMBL), Genome Biology Unit, Meyerhofstr. 1, 69117 Heidelberg,
998 Germany. ³New York Genome Center, New York, NY 10013, USA. ⁴Center for
999 Genomic Medicine, Massachusetts General Hospital, Department of Neurology,
1000 Harvard Medical School, Boston, MA 02114, USA. ⁵Department of Computational
1001 Medicine & Bioinformatics, University of Michigan, 500 S. State Street, Ann Arbor,
1002 MI 48109, USA. ⁶Heinrich Heine University, Medical Faculty, Institute for Medical
1003 Biometry and Bioinformatics, Moorenstr. 20, 40225 Düsseldorf, Germany. ⁷Department
1004 of Genome Sciences, University of Washington School of Medicine, 3720 15th Ave
1005 NE, Seattle, WA 98195-5065, USA. ⁸European Molecular Biology Laboratory (EMBL),
1006 Genome Biology Unit, Meyerhofstr. 1, 69117 Heidelberg, Germany. ⁹Division of
1007 Computational Genomics and Systems Genetics, German Cancer Research Center
1008 (DKFZ), 69120 Heidelberg, Germany. ¹⁰The Jackson Laboratory for Genomic
1009 Medicine, 10 Discovery Dr, Farmington, CT 06030, USA. ¹¹Bionano Genomics, San
1010 Diego, CA 92121, USA. ¹²Department of Genetics and Informatics Institute, School of
1011 Medicine, University of Alabama at Birmingham, Birmingham, AL 35294, USA.
1012 ¹³Molecular and Computational Biology, University of Southern California, Los
1013 Angeles, CA 90089, USA. ¹⁴Genomes and Disease, Centre for Research in Molecular
1014 Medicine and Chronic Diseases (CIMUS), Universidade de Santiago de Compostela,
1015 Santiago de Compostela, Spain. ¹⁵Department of Computer & Information Sciences,
1016 Temple University, Philadelphia, PA 19122, USA. ¹⁶Center for Genomic Medicine,
1017 Massachusetts General Hospital, Department of Neurology, Harvard Medical School,
1018 Boston, MA 02114, USA. ¹⁷Pacific Biosystems of California, Inc., Menlo Park, CA
1019 94025, USA. ¹⁷Washington University, St. Louis, MO 63108, USA. ¹⁸European
1020 Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome
1021 Campus, Hinxton, Cambridge, CB10 1SD, United Kingdom.

1022

1023