

1 Combining accurate tumour genome simulation 2 with crowd-sourcing to benchmark somatic 3 structural variant detection

4 Anna Y. Lee^{1,12}, Adam D. Ewing^{2,3,12}, Kyle Ellrott^{2,4,12}, Yin Hu⁵, Kathleen E. Houlahan¹, J.
5 Christopher Bare⁵, Shadrielle Melijah G. Espiritu¹, Vincent Huang¹, Kristen Dang⁵, Zechen
6 Chong^{6,7,8}, Cristian Caloian¹, Takafumi N. Yamaguchi¹, ICGC-TCGA DREAM Somatic Mutation
7 Calling Challenge Participants, Michael R. Kellen⁵, Ken Chen⁶, Thea C. Norman⁵, Stephen H.
8 Friend⁵, Justin Guinney⁵, Gustavo Stolovitzky⁹, David Haussler², Adam A. Margolin^{4,5*}, Joshua
9 M. Stuart^{2*}, Paul C. Boutros^{1,10,11*}

10

11 1 Ontario Institute for Cancer Research; Toronto, Ontario, Canada

12 2 Department of Biomolecular Engineering; University of California, Santa Cruz; Santa Cruz,
13 CA, USA

14 3 Mater Research Institute; University of Queensland; Woolloongabba, QLD, Australia

15 4 Computational Biology Program; Oregon Health & Science University; Portland, OR, USA

16 5 Sage Bionetworks; Seattle, WA, USA

17 6 Department of Bioinformatics and Computational Biology; University of Texas MD Anderson
18 Cancer Center; Houston, TX, USA

19 7 Department of Genetics; University of Alabama at Birmingham; Birmingham, AL, USA

20 8 Informatics Institute; University of Alabama at Birmingham; Birmingham, AL, USA

21 9 IBM Computational Biology Center; T.J.Watson Research Center; Yorktown Heights, NY,

22 USA

23 10 Department of Medical Biophysics; University of Toronto; Toronto, Ontario, Canada

24 11 Department of Pharmacology & Toxicology; University of Toronto; Toronto, Ontario, Canada

25 12 These authors contributed equally

26 * Correspondence: margolin@ohsu.edu; jstuart@ucsc.edu; paul.boutros@oicr.on.ca

27

28 **Abstract**

29 **Background**

30 The phenotypes of cancer cells are driven in part by somatic structural variants. Structural
31 variants can initiate tumors, enhance their aggressiveness and provide unique therapeutic
32 opportunities. Whole-genome sequencing of tumors can allow exhaustive identification of the
33 specific structural variants present in an individual cancer, facilitating both clinical diagnostics
34 and the discovery of novel mutagenic mechanisms. A plethora of somatic structural variant
35 detection algorithms have been created to enable these discoveries, however there are no
36 systematic benchmarks of them. Rigorous performance evaluation of somatic structural variant
37 detection methods has been challenged by the lack of gold-standards, extensive resource
38 requirements and difficulties arising from the need to share personal genomic information.

39 **Results**

40 To facilitate structural variant detection algorithm evaluations, we create a robust simulation
41 framework for somatic structural variants by extending the BAMSurgeon algorithm. We then
42 organize and enable a crowd-sourced benchmarking within the ICGC-TCGA DREAM Somatic
43 Mutation Calling Challenge (SMC-DNA). We report here the results of structural variant
44 benchmarking on three different tumors, comprising 204 submissions from 15 teams. In addition
45 to ranking methods, we identify characteristic error-profiles of individual algorithms and general

46 trends across them. Surprisingly, we find that ensembles of analysis pipelines do not always
47 outperform the best individual method, indicating a need for new ways to aggregate somatic
48 structural variant detection approaches.

49 **Conclusions**

50 The synthetic tumors and somatic structural variant detection leaderboards remain available as
51 a community benchmarking resource, and BAMSurgeon is available at
52 <https://github.com/adamewing/bamsurgeon>.

53 **Keywords**

54 somatic mutations, simulation, structural variants, benchmarking, cancer genomics, whole-
55 genome sequencing, crowd-sourcing

56

57 **Background**

58 Somatic structural variants (SVs) are mutations that arise in tumours involving rearrangements,
59 duplications or deletions of large segments of DNA. SVs are often defined to be events larger
60 than 100 bp in size, although with significant variability in this definition. Somatic SVs are critical
61 in driving and regulating tumour biology. They can initiate tumours [1,2] and because they are
62 unique to the cancer, can serve as highly-selective avenues for therapeutic intervention [3]. The
63 overall mutation load of somatic SVs serves as a proxy for genomic instability, and can robustly
64 predict tumour aggressiveness in multiple tumour types [4,5].

65 While somatic SVs that alter copy-number can be detected using microarray assays, the
66 resolution of such studies is limited, and many other important types of SVs cannot be detected.
67 As a result, high-throughput DNA sequencing is now a standard approach for detecting SVs in

68 cancer genomes. Although RNA-based assays are useful for detecting SVs that alter protein-
69 structure, DNA-based assays are required for most others. As a result, a broad range of
70 algorithms has been developed to detect SVs from short-read sequencing data using read
71 depth analysis, split read (*i.e.* a read that maps to multiple different parts of the reference
72 sequence) alignment, paired end mapping and de novo assembly techniques [6–9]. However,
73 the accuracy of existing methods is poorly described. There are no comprehensive benchmarks
74 of somatic SV detection approaches. Most comparison results are reported by the developers of
75 newly published methods. These developer-run benchmarks are potentially subject to several
76 types of selection biases. For example, the developers of one tool may be experts in
77 parameterizing and tuning it, but may lack the same skill in tuning methods developed by
78 others. Further, evaluating the accuracy of somatic SV detection is more challenging than
79 evaluating the accuracy of somatic single nucleotide variant (SNV) detection as validation data
80 is more difficult to generate for SVs. Even the metrics of measuring accuracy are not agreed
81 upon, with no community-accepted standards on how SV prediction accuracy should be
82 assessed, especially when predictions are close to, but not exactly at, the actual sequence
83 breakpoints. As a result, there are no robust estimates of the false positive and false negative
84 rates of somatic SV prediction tools on tumours of different characteristics.

85 To fill this gap, we created an open challenge-based assessment of somatic SV prediction tools
86 as part of the ICGC-TCGA DREAM Somatic Mutation Calling Challenge (the Challenge). The
87 lack of fully-characterized tumour genomes for building gold standard sets of SVs motivated our
88 simulation approach. Specifically, we first extended BAMSurgeon [10], a tool for adding
89 simulated mutations to existing reads, to generate somatic SVs. This approach is advantageous
90 because it permits flexibility with the added mutations while also capturing sequencing
91 technology biases through the use of existing reads. We created and distributed *in silico*
92 tumours (IS1-IS3), on which 204 submissions were made by 15 teams.

93 Results

94 Simulation of SVs with BAMSurgeon

95 In addition to point mutations [SNVs and short insertions or deletions (INDELs)], BAMSurgeon is
96 capable of creating simple SVs through read selection, local sequence assembly, manipulation
97 of assembled contigs, and simulation of sequence coverage over the altered contigs (Fig. 1a,
98 Additional file 1: Figure S1). This, combined with careful tracking of read depth, yields
99 approximations of SVs including insertions, deletions, duplication, and inversions into pre-
100 existing backgrounds of real sequence data. Here we present results based on simulations of
101 those SV types. Subsequent to the Challenge, BAMSurgeon was extended to support
102 translocations and more complex rearrangements. The BAMSurgeon manual, available online,
103 contains a full description of input formatting and available parameters. The input regions define
104 where local assembly will be attempted *via* Velvet [11]. For each region, the largest assembled
105 contig is selected and re-aligned to the reference genome using Exonerate [12]. The contig is
106 then trimmed to the length of its longest contiguous alignment and the alignment is used to
107 accurately track breakpoint locations within the contig in terms of reference coordinate space.
108 The location and identity of reads from the original BAM file in the assembled contig are tracked
109 *via* parsing of the AMOS [13] file output by Velvet [14], which also enables tracking of reads
110 included or excluded after contig trimming. If a suitable contig (*i.e.* sufficiently long, with a
111 sufficiently low number of discordant read pairs) is not available for a given input region, no
112 mutation is made for that region. For each segment where contig assembly succeeds, the contig
113 is rearranged according to the user specification (e.g. insertion, deletion, duplication, or
114 inversion of sequence). Then paired reads are simulated from the rearranged contig using
115 wgsim [15], with specific parameters controllable by the user. Because reads are simulated

116 using the rearranged contig, breakpoint-spanning reads and reads that will be discordant versus
117 the reference genome assembly will be created. The number of reads simulated (final coverage,
118 C_f) depends on the original coverage C_o , the difference in length between the original contig L_o
119 and the rearranged contig L_f , and a user-specified parameter controlling variant allele fraction
120 (VAF). Thus, $C_f = \text{VAF} * C_o * (L_f / L_o)$. Duplications and insertions result in larger contigs and require
121 new reads to be added to the final BAM, and deletions yielding a smaller contig require reads to
122 be removed from the final BAM. In the latter case, a list of reads to be deleted is maintained,
123 which correspond to reads covering the deleted region in the original BAM. BAMSurgeon
124 requires approximately 4GB of memory per thread if using the Burrows-Wheeler Aligner (BWA).
125 Its runtime varies depending on the number, variety and locations of the mutations, as well as
126 the depth of the original BAM. On average, runtime is about 2-3 minutes per SV per thread
127 followed by several hours to integrate all mutations into the output BAM, for a deeply sequenced
128 (e.g. 60x) genome. These are wallclock times, with the majority being spent in writing reads into
129 the BAM file.

130 **Validation of simulated somatic SVs**

131 To validate SVs simulated by BAMSurgeon, we performed a series of quality-control
132 experiments analogous to those performed to validate simulated SNVs [10]. Briefly, we used
133 BAMSurgeon to generate synthetic tumour-normal pairs, with the same set of target mutations,
134 that differ by the division of reads into tumour and normal sequence sets, aligner or cell line.
135 The target mutation set was designed to generate a synthetic tumour with a baseline level of
136 complexity and thus did not include insertions. We ran four SV callers using default parameters
137 on each pair: two widely used callers, CREST [16] and Delly [9], and two callers developed over
138 the course of the Challenge, Manta [17] and novoBreak [18]. We did not optimize parameters

139 for the callers since the goal of this validation was not to identify the best caller, but instead to
140 verify that the caller ranking is maintained across analogous datasets.

141 The definition of a SV suggests different scoring schemes for measuring the performance of a
142 caller. All SVs can be defined by at least one breakpoint; deletions, duplications and inversions
143 are SVs defined by a pair of breakpoints that in turn defines a genomic region. Thus, we
144 compared called SVs to gold-standard SVs based on i) region overlap or ii) breakpoint
145 closeness (Table 1, Additional file 1: Figure S2). The Challenge initially used a scoring scheme
146 based on region overlap (at least one or more bases in common; Additional file 1: Figure S2a).
147 Here we focus on the breakpoint closeness scheme since it is well suited for all types of SVs. A
148 called SV that is sufficiently similar to a known SV based on such criteria was considered a true
149 positive; otherwise, a false positive. We used such annotations to assess the performance of a
150 caller in terms of precision (fraction of calls that are true), recall (fraction of known SVs called)
151 and *F*-score (harmonic mean of precision and recall).

152 We performed several quality-control experiments. First, the caller ranking (by *F*-score) was
153 independent of the random division of reads: Manta > novoBreak > CREST > Delly (Additional
154 file 1: Figure S3a,b). Second, the same ranking was observed when alignments were generated
155 either using the BWA or NovoAlign with and without INDEL realignment (*i.e.* local realignment to
156 minimize mismatches across reads due to INDELS relative to the reference genome), indicating
157 that the ranking was independent of the aligner used (Fig. 1b, Additional file 1: Figure S3c).
158 Lastly, when the genomic background was varied by using HCC1143 BL or HCC1954 BL
159 sequence data, the caller ranking was largely independent of the cell line: Manta and novoBreak
160 retained first and second place, respectively, while CREST and Delly swapped places, although
161 their *F*-scores were very similar to each other when HCC1954 BL was used (Fig. 1c, Additional
162 file 1: Figure S3d). Overall, these results show that simulated SVs are robust to changes in the
163 read division, aligner and genomic background.

164 **Crowd-sourced benchmarking of somatic SV calling**

165 The SV component of the Challenge consisted of the same three synthetic tumour-normal data
166 sets used in the SNV component [10]. Briefly, the data sets were derived from existing cell line
167 sequence data (thus minimizing data access restrictions) and *in silico* tumours 1-3 (IS1-IS3)
168 were generated with increasing complexity (Fig. 1d). In terms of SVs, breakpoint locations were
169 randomly selected and the tumours had increasing mutation rates (371 vs. 2,886 somatic SVs in
170 IS1 and IS3, respectively). Moreover, IS1 contained deletions, duplications and inversions while
171 IS2 and IS3 additionally contained insertions. Like the SNV component, the SV component of
172 the Challenge was implemented using the Dialogue for Reverse Engineering Assessments and
173 Methods (DREAM) framework. Briefly, information about the Challenge was shared on its
174 website [19], participants registered online, downloaded a data set, applied their SV calling
175 pipelines to the data set and submitted the results in Variant Call Format (VCF) v4.1. IS1-IS3
176 were released sequentially, each data set had its own competition phase and participants could
177 make multiple submissions for each data set. Each tumour genome was divided into a training
178 set and a testing set by holding out a portion of the genome. During the competition phase,
179 leaderboards showed performance measures on the training set. After the competition closed,
180 leaderboards also showed performance measures on the whole genome (training + testing
181 sets).

182 The Challenge administration team prepopulated the leaderboards with two submissions and
183 the community provided 204 submissions from 15 teams (Table 2, Additional file 2). A list of all
184 submissions, and descriptions of pipelines used to generate them, can be found in Additional
185 files 3 and 4, respectively. The submissions were surprisingly discordant in format. For example,
186 between 5.5-11% of all submissions specified SV types that are not in the VCF controlled
187 vocabulary for types (Additional file 5). For this reason, and the ambiguity of specifying SV types

188 (i.e. the same SV can be specified with a specific type, or by specifying the breakpoints and
189 break-end adjacencies), type specifications were ignored when scoring submissions. Team
190 ranking varied with the stringency of the scoring (Additional file 1: Figure S2d-i). For simplicity,
191 we focused on scoring with $f = 100$ bp due to the balance between the median and variance of
192 the resulting F -scores (Additional file 1: Figure S4). While the top-performing teams achieved
193 near maximal precision on the simplest tumour, IS1, their recall remained less than 0.9 (Fig.
194 2a), and decreased further on the other tumours (Additional file 1: Figure S5a,b). On all three
195 tumours, F -scores on the training and testing sets were highly correlated (Spearman's rank
196 correlation coefficient (ρ) ≥ 0.98 ; Fig. 2b, Additional file 1: Figure S5c,d). However, the slightly
197 elevated F -scores in the training sets observed for IS1 and IS2 may reflect minor overfitting;
198 overfitting occurs when a statistical model is tuned to the training set, limiting generalizability.
199 Notably, the total number of somatic SV mutations in IS3 is $>4x$ that for IS1 and IS2 (Fig. 1d).
200 Conversely, the percentage of mutations used for training is greater for IS1 (93%) and IS2
201 (92%) vs. IS3 (89%). Sampling from the IS3 mutations, we simulated training and testing sets of
202 different sizes, and computed the differences between the F -scores on the training sets and the
203 F -scores on the testing sets. We found that that the differences tend to be greater when the
204 percentage of mutations used for training is greater (Additional file 1: Figure S5e). This suggests
205 that the F -score differences observed for IS1 and IS2 are at least in part an artefact of training
206 set size.

207 Pipeline optimization

208 The within-team variability in F -scores accounts for 21-43% of the total per-tumour variance in
209 F -scores. The large variability in submissions by certain teams highlights the impact of tuning
210 parameters during the Challenge (Fig. 3a, Additional file 1: Figure S6a,b). In comparing the
211 initial ("naive") and best ("optimized") submissions of each team, for each tumour, the maximum

212 *F*-score improvement was 0.75 (from 0.12 to 0.87 by Team 5 for IS1), and the median
213 improvements were 0.20, 0.01, and 0.07 for IS1, IS2 and IS3, respectively (Fig. 3b). At least
214 33% of teams improved their *F*-score by at least 0.05 and at least 25% of teams improved it by
215 more than 0.20, depending on the tumour. Despite these improvements by parameterization,
216 team rankings were only moderately changed: if a team's naive submission ranked in the top
217 three, their optimized submission remained in the top three 66% of the time (Fig. 3c).

218 Given the crowd-sourced nature of the Challenge, we explored “wisdom of the crowds” as an
219 approach to optimize performance [20,21]. Specifically, we aggregated SV calls into an
220 ensemble by first identifying sufficiently similar calls in the majority of the top *k* submissions.
221 Pairwise distances between calls from different submissions were computed (*i.e.* a breakpoint-
222 length distance that incorporates distances between breakpoints and differences in SV length,
223 Additional file 1: Figure S2c), and those calls with distances less than a selected threshold
224 (equal to *f*, for consistency) were considered to represent an equivalent called SV event. The
225 chromosome together with the median start and end positions of a set of similar calls would
226 then define a single ensemble SV prediction. We considered two variations of this approach: i) a
227 baseline approach with ensembles of the best submission from each team, and ii) a
228 conservative approach with ensembles of all submissions (where the top *k* may include multiple
229 submissions from the same team) and more stringent aggregation of called SVs (see Methods).
230 The baseline ensembles were found to have *F*-scores comparable to that of the best
231 submission (*e.g.* for IS1, the best ensemble and submission have *F*-scores of 0.92 and 0.91,
232 respectively; Fig. 3d, Additional file 1: Figure S7b). However, the ensembles had lower *F*-scores
233 than the best submission for IS2 (Additional file 1: Figure S7a). When *k* > 15, we found that the
234 conservative ensemble *F*-scores drop further below that of the best submission (Additional file
235 1: Figure S7c-e; *e.g.* for IS1, the best ensemble with *k* > 15 and the best submission have *F*-
236 scores of 0.83 and 0.91, respectively); these ensembles incorporate submissions from the top

237 three teams, at least. In contrast, the precision of all ensembles (range: 0.993-1.00) was similar
238 or slightly improved compared to that of the best submission. Thus, any changes in the
239 ensemble F -scores were mostly influenced by the changes in recall as k varied.

240 **Error characteristics**

241 We next exploited the large number of independent analyses to identify characteristics
242 associated with false negatives (FNs) and false positives (FPs). For example, error profiles
243 differed significantly between subclonal populations in IS3, with greater FN rates for mutations
244 present at lower VAFs (Additional file 1: Figure S8; one-sided Wilcoxon signed rank $P = 0.02$ for
245 $VAF = 0.2$ vs. 0.33 , $P = 0.04$ for $VAF = 0.33$ vs. 0.5 , $n = 7$). We also selected the best
246 submission from each team (by F -score) and focused on 14 variables associated with
247 breakpoint positions, representing factors like coverage and mapping quality (Additional file 6).
248 Several of these variables were associated with false-positive rates; in particular, tumour
249 coverage ($R > 0.24$), bridging reads count (the number of reads that bridge a putative
250 breakpoint, $R > 0.21$) and mapping quality ($R < -0.29$), have stronger associations with FPs for
251 both IS2 and IS3, compared to other variables (Additional file 1: Figure S9a, S10-S25). By
252 contrast, few were associated directly with false-negative rates ($0 \leq |R| \leq 0.15$; Additional file 1:
253 Figure S9b, S10-S25).

254 To evaluate whether these variables, and additional categorical variables, contribute
255 simultaneously to somatic SV prediction error, we generated two Random Forests (non-
256 parametric learning models that can trivially merge multiple data types) [22] for each team to
257 assess variable importance for FN and FP breakpoints separately. FN breakpoints are
258 associated with variables such as high bridging reads count and strand bias (Fig. 4a,c,e,g,i;
259 Additional file 1: Figure S26a). FP breakpoints are generally associated with variables such as
260 low mapping quality (Fig. 4b,d,f,h,j; Additional file 1: Figure S26b).

261 By executing specific SV callers, CREST (Fig. 4a,b), Delly (Fig. 4c,d) and Manta (Fig. 4e,f), with
262 the same parameters on all three tumours, we identified tumour-specific error profiles. For
263 example, the distance to the nearest germline INDEL tends to have stronger associations with
264 errors in IS2 and IS3 compared to IS1 (Fig. 4a-e). Team-specific error profiles are more
265 apparent with the FP breakpoint analysis. For example, Teams 8 and 10 have distinct FP
266 profiles for the same tumour, IS2 (Fig. 4h); FPs by Teams 8 and 10 are negatively and positively
267 associated with tumour coverage, respectively. Algorithmic approaches to SV calling from
268 sequencing data include i) read depth analysis, ii) paired end mapping, iii) split read alignment,
269 and iv) *de novo* assembly [23]. Some teams submitted sufficient algorithm details to determine
270 the general approaches used, as well as the choice of aligner (Fig. 4g-j). Based on the available
271 annotations, teams using the same aligner do not have error profiles that tightly cluster for all
272 three tumours, suggesting that the aligner is not as strong a driver of those profiles, compared
273 to the caller algorithm.

274

275 Discussion

276 Crowd-sourced benchmarking challenges are ideal for questions where significant diversity in
277 algorithmic approaches exists, particularly where individual methods are highly parameterized
278 or computationally intensive [24,25]. The detection of variants from high-throughput sequencing
279 data fits these criteria well: dozens of algorithms are in common use, many with complicated
280 sets of parameters to tune and most requiring tens to hundreds of CPU hours to execute. We
281 have quantified the critical importance of parameterization: it accounts for 21-43% of the
282 variability in performance across the 204 submissions evaluated. This is comparable to the 26%
283 of variability observed in somatic SNV detection benchmarking [10], and highlights the need for

284 algorithm developers to continue to optimize parameters, provide guidance for their tuning and
285 work toward automating their selection to make their software easier to use.

286 Scoring SV detection is complicated by the diversity of SVs. While some SV types may be well-
287 characterized by overlap-based scoring methods, others benefit more from breakpoint-based
288 scoring, and the choice of scoring metric and stringency parameters must be tuned to specific
289 biological questions of interest. For example, breakpoint identification is critical when
290 considering translocations -- especially those generating candidate fusion proteins -- while
291 overlap of the called and known regions is much more important for copy-number analyses.
292 Moreover, it may be useful to adapt scoring (*e.g.* by using a range of stringency parameter
293 values) to identify SVs in certain contexts (*e.g.* with breakpoints in repetitive regions) that are
294 still detectable by given tools, but with less precision. Taken together, SV diversity is an
295 important consideration for the development of standards for scoring SV detection.

296 The “wisdom of the crowds” is the idea that an ensemble of multiple algorithms can significantly
297 outperform any individual method. Several crowd-sourced benchmarking competitions from
298 diverse fields have shown great success in combining submissions from contestants to achieve
299 a high-performing meta-predictor including challenges for somatic SNV detection [10], gene
300 regulatory network inference [21] and mRNA-based prognostic signatures for breast cancer
301 [20]. By contrast, in somatic SV detection, we do not have clear evidence that an ensemble
302 improves on the best individual method consistently across different tumours. Specifically, the
303 majority vote approach works very well for somatic SNV detection, yet it appears to fail for
304 somatic SV detection. This may reflect the large diversity in the biases of each individual
305 algorithm (Fig. 4). Rather than focus on commonalities through a majority vote, it may be more
306 beneficial to leverage the strengths of individual algorithms. This might be achieved by using
307 machine learning to optimize the weighting of the algorithms for specific input patterns. For
308 example, an aggregating classifier could learn, if there is a sizable difference in coverage in the

309 tumour versus normal samples near given candidate breakpoints, the calling algorithms that use
310 read depth analysis should have more weight. The overall approach could involve the following
311 general steps: 1) apply all algorithms of interest to a given tumour-normal dataset and take the
312 union of all resulting call sets to define a list of candidate SVs; then for each candidate, 2)
313 compute sequence features (e.g. coverage) around the candidate breakpoints, and 3) provide
314 computed features and confidence scores from individual algorithms as input to an aggregating
315 classifier that will indicate whether or not the candidate is likely to be a true SV. In fact, a similar
316 approach is behind the SMC-DNA Meta-pipeline Challenge [26] for benchmarking pipelines that
317 aggregate calls from different SNV detection algorithms. In practise, analogous efforts for SV
318 detection would require additional considerations such as the identification of i) an optimal
319 method for merging similar yet different calls (due to imprecise breakpoint calling) when
320 compiling the list of candidate SVs, ii) the most informative sequence features for guiding the
321 relative weighting of individual algorithms (e.g. variables in Figure 4), and iii) an optimal scoring
322 method (as mentioned above). Thus there is a need for continued development of new, more
323 complex ways to integrate multiple somatic SV detection methods [27].

324 Given the paucity of gold-standard benchmarking data for somatic SVs, the creation of the
325 simulated datasets and the associated leaderboards constitutes a major contribution of this
326 Challenge. Ideally, a simulated dataset depicts realistic mutations through realistic sequence
327 reads. The synthetic tumours generated for the Challenge only represent straightforward SV
328 types (duplication, deletion, insertion, inversion) and cover relatively small regions. Subsequent
329 enhancements to BAMSurgeon have added support for additional SV types including
330 translocations and complex SV combinations, enabling simulations to more completely capture
331 the complexity of tumour genomes and by extension, challenge SV callers in different ways. For
332 each SV, simulated reads are generated (*via* wgsim) from a re-arranged contig, where the
333 original contig is constructed from real reads. Despite the basis on real reads, the simulated

334 reads do not necessarily reflect the non-uniform coverage that may arise during preparation of
335 real samples, for example [28]. There are other read simulators that learn biased-coverage
336 trends from real data and use them to generate reads (e.g. [29]) that could be used by
337 BAMSurgeon; however, it is an on-going challenge to simulate biases of real sequencing data
338 as sample preparation methods and sequencing technologies vary and/or advance. In fact, one
339 could sequence the same 'normal' sample twice to capture inter-sample variability, with one
340 replicate converted into a synthetic tumour sample using BAMSurgeon. Nevertheless, there are
341 distinct advantages to benchmarking on simulated datasets. It is dramatically easier to simulate
342 large numbers of tumours, or to create tumours with highly divergent mutational properties,
343 leading to well-supported estimates of per-tumour caller accuracy. This enables our strategy of
344 generating synthetic tumours of increasing complexity (e.g. with other SV types and/or
345 haplotype structure by using phased sequence data) whereby the impact of the complexity
346 introduced at each step can be assessed. With the three synthetic tumours described here, we
347 observed that caller ranking varied across tumours and we expect it to vary with a broad range
348 of tumour characteristics including coverage, normal contamination, complexity of the SVs, the
349 number of mutations adjacent to breakpoints and others, as they each present different
350 challenges. It is possible to identify strengths and weaknesses of an individual caller by
351 comparing its tumour-specific error profiles. Moreover, synthetic tumours can be designed to
352 test the limits of callers. These advantages highlight the usefulness of synthetic datasets for
353 benchmarking callers, and until synthetic datasets are completely realistic, they will serve as
354 important complements to real datasets.

355 While 15 teams participated in the actual competitive phase of the Challenge, 8 teams have
356 exploited the IS1-3 benchmarking resources since the competition, making 73 submissions to
357 benchmark their methods for pipeline evaluation and development. Evaluations based on the
358 first synthetic tumours, the simplest by design, provide lower-bounds on the error rates. As

359 subsequent updates to BAMSurgeon enable the generation of more complex and realistic
360 tumours, the corresponding error rates using these simulations will approach their upper-
361 bounds. We hope that journals and developers will begin to plan for benchmarking on these
362 standard datasets, including simulated ones, as a standard part of manuscripts reporting new
363 somatic SV detection algorithms.

364

365 Conclusions

366 Analysis of the error profiles of the Challenge submissions showed that somatic SV calling is a
367 distinctly harder problem than somatic SNV calling even given a relatively simple set of SVs,
368 with individual pipelines having complex and unique error profiles. Parameterization was a
369 critical factor in determining the performance of teams. Finally, we demonstrate that, unlike
370 almost every past DREAM Challenge, somatic SV prediction does *not* benefit from the “wisdom
371 of the crowds” -- simple voting of multiple prediction pipelines does not yield improved
372 predictions in this instance. The synthetic tumours and somatic SV detection leaderboards
373 remain available as a community benchmarking resource.

374

375 Methods

376 Simulation of SVs by BAMSurgeon

377 SV support in BAMSurgeon has evolved throughout the Challenge, largely as a result of
378 constructive feedback from participants. Our descriptions of BAMSurgeon's method for
379 simulating SVs is current as of commit (*i.e.*, version) b851573474 of the code available at [30].

380 As input, BAMSurgeon (addsv.py) requires an indexed reference genome, a pre-existing BAM
381 file (Additional file 1: Figure S1a), and a list of intervals (Additional file 1: Figure S1b) along with
382 the SV type and parameters (see manual [31]). The intervals should be wide enough that local
383 sequence assembly is successful in generating a contig that spans at least 2x the expected
384 library size in the input BAM file. Intervals for which a sufficiently long contig cannot be
385 generated are rejected, where the exact definition of 'sufficiently long' is an optional parameter.
386 Note that it may be less likely to obtain long contigs from genomic regions that are more difficult
387 to sequence, and by extension, less likely to simulate SVs in such regions. Intervals which
388 contain too many discordant read pairs (again, potentially indicating regions that are difficult to
389 sequence) are also rejected, subject to a parameter. Following local assembly, the contig is re-
390 arranged: the specific rearrangements for each supported SV type are illustrated in Fig. 1a (step
391 iii) and Additional file 1: Figure S1c,e,g. The assembled contig is then re-aligned to the target
392 interval in the reference genome (exonerate --bestn 1 -m ungapped) and is trimmed based on
393 the start and end coordinates of this alignment. Read pairs corresponding to trimmed contig
394 sequence are removed from further consideration.

395 Read coverage is generated over the rearranged contig using a read simulator (wgsim -e 0 -R 0
396 -r 0), to achieve the same average depth as the input BAM file, which has the effect of creating
397 split reads relative to the reference genome supporting SV detection. For a deletion, the number
398 of reads required to achieve (e.g.) 30x coverage is fewer than the number of reads required to
399 reach 30x coverage prior to the deletion, so reads must be removed from the original BAM (Fig.
400 1a, step iv). Inversely, for duplications and insertions additional reads need to be added to the
401 original BAM (Additional file 1: Figure S1d,h). Inversions generally do not affect coverage
402 (Additional file 1: Figure S1f). To ensure any reads removed actually correspond to the deleted
403 region of the contig, the locations of reads in the assembled contig are tracked. The number of
404 reads to be replaced, added, or deleted is scaled with the desired allele fraction. Finally, any

405 read pairs in the original BAM corresponding to reads altered in the simulated SV are replaced,
406 any read pairs marked for deletion are removed from the original BAM, and any additional read
407 pairs generated are added. It is recommended that the resulting altered BAM be post-processed
408 (with `postprocess.py`) to ensure compliance with the SAM format specification (see manual
409 [31]).

410 **Synthetic tumour generation**

411 Synthetic tumours were prepared by partitioning high-coverage BAMs from 'normal' cell lines
412 into two groups of reads, picking read pairs at random as described previously [10].
413 Alternatively, one could sequence the same 'normal' sample twice to capture inter-sample
414 variability, with one replicate converted into a synthetic tumour sample using BAMSurgeon. For
415 the three *in silico* challenges, non-overlapping regions were selected at random for SV addition,
416 resulting in 371 variants added for IS1, 655 for IS2, and 2,886 for IS3 (Fig. 1d). Variant input
417 files are available in Additional file 7. SVs were added using `addsv.py` with assembly
418 GRCh37/hg19 as the reference genome and default parameters except where noted. For IS3,
419 to simulate subclones a file specifying CNV fractions over SV regions was input via option `-c` to
420 specify the variant allele frequency (VAF) of the spiked-in variants at either 0.5, 0.33, or 0.2
421 (Additional file 7). The output BAMs were post-processed to account for any inconsistencies
422 introduced due to remapping and merging of reads supporting SVs using the script
423 `postprocess.py` included with BAMSurgeon. The BAMs were further adjusted with
424 `RealignerTargetCreator` and `IndelRealigner` from the Genome Analysis Toolkit (v.2.4.9). All
425 tumour-normal pairs generated via BAMSurgeon are verified for adherence to the SAM/BAM
426 format specification using the `ValidateSamFile` tool included in the Picard tool set [32]. Truth
427 VCF files, *i.e.* files specifying simulated mutations, for SVs were generated using the script
428 `etc/makevcf_sv.py` and merged with truth files for SNP and INDEL locations, where applicable.

429 SAMtools was used throughout to split, merge, sort, and index BAMs, and also index FASTA
430 files. Details on the specific BAMSurgeon commits used for generating each tumour, as well as
431 other tumour details are given at [33].

432 **Validation of BAMSurgeon**

433 To validate BAMSurgeon's ability to simulate somatic SVs, we compared the output of four
434 algorithms -- two widely used SV callers, CREST [16] and Delly [9], and two callers developed
435 over the course of the Challenge, Manta [17] and novoBreak [18] -- on the IS1 tumour-normal
436 data set, and analogous datasets generated with the same spike-in set of mutations, but with an
437 alternate aligner (NovoAlign v.3.00.05 [34]), cell line (HCC1954 BL) or read division. We did not
438 optimize parameters for the callers since the goal of this validation was not to identify the best
439 caller, but instead to verify that the caller ranking is maintained across analogous datasets.

440 Each tumour-normal pair was processed by CREST (v1.0) to extract soft clipping positions for
441 each chromosome separately, using default parameters. This data was then used by CREST to
442 call somatic SVs using the default protocol, and we converted the output into VCF v4.1 format.
443 Somatic SVs were called from each tumour-normal pair using Delly (v0.5.5) with default
444 parameters. Calling was performed on each chromosome separately for all supported SV types
445 except for translocations, and we converted the translocation output into VCFv4.1 format. Calls
446 with MAPQ < 20, PE < 5, or labeled as "LowQual" or "IMPRECISE" were filtered out. Somatic
447 SVs were called from each tumour-normal pair using Manta (v0.26.3) with the following
448 parameters: -m local -j 4 -g 10. Lastly, somatics SVs were called from each dataset using
449 novoBreak (v1.04) with a modification to ensure that sequence windows around breakpoints
450 never go beyond the start of the chromosome. All sets of SV calls were scored with $f = 100$ bp
451 and $j > 0$, callers were ranked based on F -score for each tumour-normal pair, and rankings were
452 compared across pairs (Fig. 1b,c and Additional file 1: Figure S3).

453 **Preprocessing VCF files**

454 We preprocess VCF files to parse out the SV-relevant details (e.g. the END coordinate in the
455 INFO value or from the length of the REF sequence; if the END coordinate cannot be
456 determined from those values, it is set to the POS coordinate), remove SVs that did not pass
457 filters (as indicated by the FILTER values) and ensure consistent formatting between files. To
458 ensure consistent formatting in accordance with the VCFv4.1 specification [35] we:

- 459 1. Add row entries to ensure that each MATEID specification has a corresponding pair of
460 entries, where only a single entry is provided
- 461 2. Re-assign IDs and MATEIDs to ensure unambiguous references to entries
- 462 3. Where possible, replace SVTYPE = BND entries with entries specifying SVTYPE =
463 {CNV, DEL, DUP, INS, INV} in accordance with REF, ALT and EVENT values

464 Testing set SVs are indicated in the truth VCF file with the addition of masked genomic regions
465 specified with CHROM, POS and END values indicating the chromosome, start and end
466 coordinates, and SVTYPE = MSK. Specifically, a SV where $\geq 50\%$ of the corresponding region
467 overlaps a masked region is allocated to the testing set; otherwise, it is in the training set.

468 **Structural variant scoring**

469 Our scoring approaches evaluate the accuracy of a set of called SVs and requires input VCF
470 files specifying: i) called SVs, and ii) true/known SVs. Generally, a called SV that is sufficiently
471 similar to a known SV based on specific criteria (Table 1) is considered a true positive (TP);
472 otherwise, a false positive (FP). Also, a known SV that is sufficiently similar to a called SV is
473 considered a TP; otherwise, a false negative (FN). Our scoring supports two ways of quantifying
474 similarity:

475 A. **Region overlap.** The Jaccard coefficient (j) is computed from the lengths (in bp) of
476 intersection and union regions (Additional file 1: Figure S2a).

477 B. **Breakpoint closeness.** The distance (Δ , in bp) between called and known breakpoints
478 is computed (Additional file 1: Figure S2b). If $\Delta \leq f$ (where f is a flank threshold
479 parameter), a relative closeness is computed, $c' = 1 - \Delta/f$. The overall closeness (c) is
480 defined as the geometric mean of the c' values for the start and end breakpoints. If only
481 one of the start and end breakpoints has $\Delta \leq f$, the called and known SVs are annotated
482 as partially matching.

483 Unless otherwise specified, we scored with $f = 100$ bp. If there is an ambiguous matching of
484 called SVs to known SVs by sufficient similarity, the similarity values (j/c) are used to identify an
485 optimal one-to-one matching. First, we restrict the matching to the best match(es) for each
486 called and known SV. If a SV has multiple best matches by similarity, we attempt to break the
487 tie by favouring SVs with the same SVTYPE, and/or test/training set membership. If the best
488 matching is still ambiguous, we then use corresponding similarity values together with the
489 Hungarian algorithm to obtain a one-to-one matching (with the clue v0.3-48 R package [36]).
490 Finally, SVs are annotated based on this matching. SVs that have sufficient similarity but are not
491 in the final matching are annotated as partially matching. Mated breakpoints are initially
492 annotated separately. If one is annotated as partially matching or as a TP, and the other is a FP,
493 the FP annotation is replaced by a partial match annotation. Subsequently, each set of mated
494 breakpoints is treated as a single SV.

495 These annotations are used to assess the performance of a SV caller in terms of precision =
496 $nTP/(nTP + nFP)$, recall = $nTP/(nTP + nFN)$ and F -score (specifically, F_1 -score) = $2 \times \text{precision} \times$
497 $\text{recall}/(\text{precision} + \text{recall})$, where nTP , nFP and nFN represent the numbers of TPs, FPs and
498 FNs, respectively. Partial matches are not counted in these computations. Unless otherwise

499 specified, the precision, recall and F -score values presented here were computed on the testing
500 and training sets combined. The best submission of a given team is defined as the team's
501 submission with the greatest F -score computed against all known SVs.

502 **Execution of challenge-based benchmarking**

503 The SV component of the Challenge was executed concurrently with the SNV component, and
504 the procedure has been described previously [10]. It was implemented using the Dialogue for
505 Reverse Engineering Assessments and Methods (DREAM) framework. Briefly, information
506 about the Challenge was shared on its website [19], participants registered online, downloaded
507 a data set, applied their SV calling pipelines to the data set and submitted the results in
508 VCFv4.1 format. IS1-IS3 were released sequentially, each data set had its own competition
509 phase and participants could make multiple submissions for each data set. Each tumour
510 genome was divided into a training set and a testing set. During the competition phase,
511 leaderboards showed performance measures on the training set. After the competition closed,
512 leaderboards also showed performance measures on the whole genome (training + testing
513 sets), thus benchmarking the SV calling pipelines. The SV leaderboards for IS1 and IS2 were
514 pre-populated with results from BreakDancer (v1.1.2_2013_03_08 [7]) run with default
515 parameters; a reference point submission indicated labeled as "Standard" in figures and tables.
516 Due to our exploration of multiple SV scoring methods in this manuscript, the leaderboard
517 results are not completely consistent with the results presented here, but all raw and
518 leaderboard data are available.

519 **Overfitting artefact analysis**

520 Due to the order of magnitude greater number of SVs spiked into IS3, we simulated training and
521 testing sets of different sizes by sampling from the IS3 training set. Specifically, we assessed
522 mutation totals of 100 to 1000 (by increments of 100), and training sets that were 80-95% (by

523 increments of 1%) of the total, by sampling each {mutation-total, training-set%} combination 100
524 times. For each sample, we computed $F_{train} - F_{test}$ for each IS3 submission where F_{train} and F_{test}
525 are F -scores computed on the simulated training and testing sets, respectively. We then
526 computed the median difference across samples to obtain a summary value for each
527 submission, and finally show the median across submissions in Additional file 1: Figure S5e.
528 ($F_{train} - F_{test}$) > 0 suggests overfitting but such values are an artefact of testing set size since no
529 fitting/training was done in this analysis.

530 **Team variation**

531 For each tumour-normal pair, we computed the percentage of variation in F -score, across all
532 submissions, that is accounted for by within-team variation. Specifically, we computed the
533 within-team sum of squares as a percentage of the total sum of squares.

534 **Definition of ensembles**

535 We aggregated SV calls from k submissions into an ensemble set with the following general
536 approach:

- 537 1. **BND filter.** Calls defined with SVTYPE = BND were excluded for simplicity.
- 538 2. **Compute call distances.** Pairwise distances (d , in bp) between remaining predictions
539 were computed (*i.e.* a breakpoint-length distance that incorporates distances between
540 breakpoints and differences in predicted SV length, Additional file 1: Figure S2c).
541 Distances were only computed between predictions from different submissions.
- 542 3. **Generate sets of similar calls.** A distance less than a selected threshold (100 for
543 consistency with f , see **Structural variant scoring**) indicated sufficiently similar calls.
544 We assessed two variations:

- 545 a. **Baseline.** We defined a graph such that vertices represented calls and edges
546 connected sufficiently similar calls. We identified the connected components to
547 define the sets of similar calls. Sets with median intra-set distances $> f$ were
548 refined. Specifically, the call with the greatest median distance to the other set
549 members was iteratively removed until the median intra-set distance dropped
550 below f , or the set became empty.
- 551 b. **Conservative.** We used the added constraint that called SVs overlap by ≥ 1 bp
552 to be treated as sufficiently similar. Sets of similar calls were constructed by
553 iterating over the sufficient similarity pairs from least to most distant. If a pair did
554 not contain a call in an existing call set, the pair was used to define a new call
555 set. Otherwise, one call was already in a set, and the other was a candidate for
556 addition to the same set via guilt-by-association. If the candidate came from a
557 submission that was not already covered by the set, and its median distance to
558 the existing set members $\leq f$, it was added to the set. Any unprocessed pairs
559 within or between the prediction sets at that stage were excluded from
560 consideration.
- 561 4. **Majority vote filter.** Sets with calls from $\leq k/2$ submissions were excluded; each
562 remaining set covered the majority of submissions.
- 563 5. **Aggregate sets to define ensemble calls.** The chromosome together with the median
564 start and end positions of each set of calls defined a single ensemble SV call.

565 An additional distinction between the baseline and conservative approaches is that the baseline
566 approach only involved the best submission from each team whereas all submissions were
567 used with the conservative approach. To investigate different ensembles of N submissions for
568 the same tumour-normal pair, we first ordered the submissions by overall F -score, computed
569 after excluding calls with SVTYPE = BND. We then generated an ensemble call set with the top
570 k submissions, for $k = 2..N$. The performance of ensembles was compared to that of the
571 individual submissions, after excluding calls with SVTYPE = BND (e.g. Fig. 3d).

572 **Error characterization**

573 To characterize the errors made by a team, we assessed the team's best submission for a given
574 tumour-normal pair. We also assessed errors made by CREST, Delly and Manta when run, with
575 the same protocols described in the **Validation of BAMSurgeon** section, on all three tumour-
576 normal pairs. Characterizing FNs and FPs involved comparisons to TPs and true negatives
577 (TNs), respectively. Moreover, we characterized errors at the level of breakpoints.

578 **Sampling true negatives.** Given a set of submissions for the same tumour-normal pair, we
579 identified the maximum number of FPs from a single submission, m . We then sampled $\geq m$ TNs
580 for each submission, by sampling regions from the reference genome that satisfied these
581 criteria:

- 582 1. length sampled from a log-normal distribution with mean and standard deviation equal to
583 that of the logged lengths of the known SVs
- 584 2. start position is not in known gap and repeat regions
- 585 3. region does not overlap with any known SVs
- 586 4. region does not overlap with any SVs called in the submission

587 Some sampled regions qualified as TNs for multiple submissions. For IS2, we excluded Team
588 14's submission because it had a very large number (17,806) of FPs, and thus was
589 computationally problematic for the subsequent Random Forest analysis.

590 **Breakpoint annotations based on scoring.** A single breakpoint may be associated with
591 multiple (called/known) SVs, and therefore may be associated with multiple annotations
592 depending on the scoring approach used, *i.e.* > 1 of {TP, FN, FP}. To remove ambiguity, we
593 choose a single annotation for each breakpoint by prioritizing as follows: TP > FN > FP. This
594 prioritization favours good performance (*i.e.* TP has highest priority) and then recall (*i.e.* FN >
595 FP) since it appears to be a greater challenge than precision for SV calling (Fig. 2a, Additional
596 file 1: Figure S5a,b). TN breakpoints should be unambiguous due to the way in which they were
597 sampled (see above).

598 **Genomic variables.** For each breakpoint position, we computed 16 genomic factors, 12 of
599 which were previously described [10]. The additional genomic variables were computed as
600 follows:

601 A. **Bridging reads count.** We used samtools v0.1.19 to identify reads in the tumour BAM
602 mapped to a genomic region containing the window defined by the breakpoint position
603 +/- 1 bp. The bridging read count was defined as the number of identified reads. Note
604 that a bridging read does not necessarily have a secondary mapping for part of the read,
605 as one might expect for a split read.

606 B. **Distance to nearest germline INDEL.** Germline calls were obtained as previously
607 described [10] and INDELS were parsed out. The distance of a breakpoint to the nearest
608 germline INDEL was computed using BEDTools closest v2.18.2.

609 C. **Nucleotide complexity.** The sequence for the window defined by the breakpoint
610 position +/- 50 bp was extracted from the reference fasta file. The nucleotide complexity

611 was defined as the entropy of the sequence: $-\sum p_x \log_2(p_x)$ over $x \in \{A, G, C, T\}$ where p_x
612 is the proportion of the sequence with x (case-insensitive).

613 D. **Strand bias.** We used samtools v0.1.19 to identify reads in the tumour BAM mapped to
614 a genomic region containing the breakpoint position. The strand bias was defined as the
615 proportion of these reads mapped to the + strand.

616 **Univariate analysis.** To assess the relationship between each non-categorical variable and
617 prediction error rates, we computed the Pearson correlation coefficient between the variable
618 values and the proportion of teams with a FN/FP at the breakpoints, as well as the
619 corresponding P value. Reference and alternative allele counts, base quality, tumour and
620 normal coverages, bridging reads counts and distances to the germline SNP and INDEL were
621 logged (base 10) prior to computing correlations (zeros were replaced with -1 instead of
622 logged). For the categorical variables, trinucleotide and genomic location, the P value measured
623 the significance of the variable in a fitted binomial model predicting the FN/FP rate at a
624 breakpoint. A binomial model was fitted because it is a relatively simple model (and thus less
625 prone to overfitting) to test the relationship between a categorical variable and a proportion
626 variable (*i.e.* an error rate).

627 **Multivariate analysis.** Random Forests were generated as previously described [10] with a few
628 alterations. Here, a total of 16 genomic variables (Fig. 4) were used to build: i) a classifier to
629 distinguish FN and TP breakpoints, and ii) a classifier to distinguish FP and TN breakpoints. A
630 FP classifier was not generated for Team 7 with respect to IS1 since the team produced only
631 one FP, and thus there was insufficient data to generate an accurate model. Conversely, a FP
632 classifier was not generated for Team 14 with respect to IS2 since the team produced a very
633 large number of FPs (17,806) that caused a failure to converge. Computation of the directional
634 effect of variables was also as previously described [10].

635 Non-parametric tests (*i.e.* Wilcoxon and Mann-Whitney tests) were used throughout to avoid
636 assumptions about the distributions of the tested populations; all tested populations had $n \geq 7$.
637 The BEDTools suite (v2.18.2 [37]) was used with the bedR R package (v0.5.3 [38]) throughout.
638 Plots were generated with the BPG (v5.3.9), lattice (v0.20-33) and latticeExtra (v0.6-26) R
639 packages and R (v3.2.1) was used throughout.

640 **Declarations**

641 **Availability of data and materials**

642 Sequences files are available at the Sequence Read Archive (SRA) under accession number
643 SRP042948. BAMSurgeon is available at Zenodo [39] and the code repository is available at
644 GitHub [30]. Submission (Synapse IDs syn12628575, syn12628576 and syn12628577 for IS1-
645 IS3, respectively) and known mutation (*i.e.* ground truth; Synapse IDs syn2354306, syn2399959
646 and syn2485207 for IS1-IS3, respectively) VCF files are available from the Challenge website
647 [19] following registration and subsequent login at Synapse.

648 **Acknowledgments**

649 The authors thank S.P. Shah, R.D. Morin and P.T. Spellman for helpful suggestions, L.E.
650 Heisler and B.F. Huang as well as all the members of the Boutros lab for insightful discussions
651 and technical support. The authors thank Google Inc. (in particular N. Deflaux) and Annai
652 Biosystems (in particular D. Maltbie and F. De La Vega) for their ongoing support of the ICGC-
653 TCGA DREAM Somatic Mutation Calling Challenge.

654 The ICGC-TCGA DREAM Somatic Mutation Calling Challenge Participants are: Bret D. Barnes,
655 Inanc Birol, Xiaoyu Chen, Readman Chiu, Anthony J. Cox, Li Ding, Markus H-Y. Fritz, Andrey
656 Grigoriev, Faraz Hach, Joseph K. Kawash, Jan O. Korbel, Semyon Kruglyak, Yang Liao,

657 Andrew McPherson, Ka M. Nip, Tobias Rausch, S. Cenk Sahinalp, Iman Sarrafi, Christopher T.
658 Saunders, Ole Schulz-Trieglaff, Richard Shaw, Wei Shi, Sean D. Smith, Lei Song, Difei Wang,
659 Kai Ye.

660 **Funding**

661 This study was conducted with the support of the Ontario Institute for Cancer Research to
662 P.C.B. through funding provided by the Government of Ontario. This work was supported by
663 Prostate Cancer Canada and is proudly funded by the Movember Foundation - Grant #RS2014-
664 01. This study was conducted with the support of Movember funds through Prostate Cancer
665 Canada and with the additional support of the Ontario Institute for Cancer Research, funded by
666 the Government of Ontario. This project was supported by Genome Canada through a Large-
667 Scale Applied Project contract to P.C.B., S.P. Shah and R.D. Morin. This work was supported
668 by the Discovery Frontiers: Advancing Big Data Science in Genomics Research program, which
669 is jointly funded by the Natural Sciences and Engineering Research Council (NSERC) of
670 Canada, the Canadian Institutes of Health Research (CIHR), Genome Canada, and the Canada
671 Foundation for Innovation (CFI). P.C.B. was supported by a Terry Fox Research Institute New
672 Investigator Award and a CIHR New Investigator Award. K.E.H. was supported by a CIHR
673 Computational Biology Undergraduate Summer Student Health Research Award. A.D.E was
674 supported by an Australian Research Council Discovery Early Career Researcher Award
675 DE150101117 and by the Mater Foundation. The following National Institutes of Health (NIH)
676 grants supported this work: R01-CA180778 (J.M.S.), and U24-CA143858 (J.M.S.). The funders
677 played no role in study design, data collection, data analysis, data interpretation or in writing of
678 this manuscript.

679 **Author's contributions**

680 A.A.M., J.M.S and P.C.B. initiated the project. A.D.E. created BAMSurgeon. A.D.E, K.E., Y.H.,
681 K.E.H., J.C.B., M.R.K., T.C.N., S.H.F., G.S., A.A.M., J.M.S. and P.C.B. created the ICGC-TCGA
682 DREAM Somatic Mutation Calling Challenge. A.Y.L., A.D.E., Y.H., K.E.H., S.M.G.E., V.H., K.D.,
683 Z.C., C.C., and T.N.Y. created datasets and analyzed sequencing data. A.Y.L., Y.H., K.E.H, and
684 P.C.B. were responsible for statistical modelling. Research was supervised by K.C., S.H.F.,
685 J.G., G.S., D.H., A.A.M., J.M.S. and P.C.B. The first draft of the manuscript was written by
686 A.Y.L. and P.C.B., extensively edited by A.D.E., K.E., A.A.M. and J.M.S. and approved by all
687 authors.

688 **Ethics approval and consent to participate**

689 Not applicable.

690 **Consent for publication**

691 Not applicable.

692 **Competing interests**

693 All authors declare that they have no competing interests.

694 **Additional files**

695 Additional file 1: Figures S1-S26. (PDF 3.5 MB)

696 Additional file 2: Table S1. Challenge participation. (XLS 5 KB)

697 Additional file 3: Table S2. All competition-phase submissions evaluated with $f = 100$ and $j > 0$.

698 (XLS 43 KB)

699 Additional file 4: Descriptions of pipelines used to generate submissions. (PDF 3.4 MB)

- 700 Additional file 5: Table S3. Invalid SV types. (XLS 8 KB)
- 701 Additional file 6: Table S4. Univariate error analysis. (XLS 14 KB)
- 702 Additional file 7: BAMSurgeon input files used to generate the three *in silico* tumour-normal
- 703 pairs (IS1-IS3). (TAR.GZ 122 KB)
- 704

705 References

- 706 1. Northcott PA, Lee C, Zichner T, Stütz AM, Erkek S, Kawauchi D, et al. Enhancer
707 hijacking activates GFI1 family oncogenes in medulloblastoma. *Nature*. 2014;511:428–
708 34.
- 709 2. Taub R, Kirsch I, Morton C, Lenoir G, Swan D, Tronick S, et al. Translocation of the c-
710 myc gene into the immunoglobulin heavy chain locus in human Burkitt lymphoma and
711 murine plasmacytoma cells. *Proc. Natl. Acad. Sci. U. S. A.* 1982;79:7837–41.
- 712 3. Huang M, Ye Y, Chen S, Chai J, Lu J, Zhou L, et al. Use of all-trans retinoic acid in the
713 treatment of acute promyelocytic leukemia. *Blood*. 1988;72:567–72.
- 714 4. Lalonde E, Ishkanian AS, Sykes J, Fraser M, Ross-Adams H, Erho N, et al. Tumour
715 genomic and microenvironmental heterogeneity for integrated prediction of 5-year
716 biochemical recurrence of prostate cancer: a retrospective cohort study. *Lancet. Oncol.*
717 2014;15:1521–32.
- 718 5. Vollan HKM, Rueda OM, Chin S-F, Curtis C, Turashvili G, Shah S, et al. A tumor DNA
719 complex aberration index is an independent predictor of survival in breast and ovarian
720 cancer. *Mol. Oncol.* 2015;9:115–27.
- 721 6. Medvedev P, Stanciu M, Brudno M. Computational methods for discovering structural
722 variation with next-generation sequencing. *Nat. Methods*. 2009;6:S13–20.
- 723 7. Chen K, Wallis JW, McLellan MD, Larson DE, Kalicki JM, Pohl CS, et al. BreakDancer:
724 an algorithm for high-resolution mapping of genomic structural variation. *Nat. Methods*.
725 2009;6:677–81.
- 726 8. Hormozdiari F, Hajirasouliha I, Dao P, Hach F, Yorukoglu D, Alkan C, et al. Next-
727 generation VariationHunter: combinatorial algorithms for transposon insertion discovery.
728 *Bioinformatics*. 2010;26:i350-7.

- 729 9. Rausch T, Zichner T, Schlattl A, Stütz AM, Benes V, Korbel JO. DELLY: structural variant
730 discovery by integrated paired-end and split-read analysis. *Bioinformatics*. 2012;28:i333–
731 9.
- 732 10. Ewing AD, Houlihan KE, Hu Y, Ellrott K, Caloian C, Yamaguchi TN, et al. Combining
733 tumor genome simulation with crowdsourcing to benchmark somatic single-nucleotide-
734 variant detection. *Nat. Methods*. 2015;12:623–30.
- 735 11. Zerbino DR, Birney E. Velvet: algorithms for de novo short read assembly using de Bruijn
736 graphs. *Genome Res*. 2008;18:821–9.
- 737 12. Slater GSC, Birney E. Automated generation of heuristics for biological sequence
738 comparison. *BMC Bioinformatics*. 2005;6:31.
- 739 13. Pop M, Phillippy A, Delcher AL, Salzberg SL. Comparative genome assembly. *Brief.*
740 *Bioinform*. 2004;5:237–48.
- 741 14. Zerbino DR, McEwen GK, Margulies EH, Birney E. Pebble and rock band: heuristic
742 resolution of repeats and scaffolding in the velvet short-read de novo assembler. *PLoS*
743 *One*. 2009;4:e8407.
- 744 15. GitHub Code Repository: wgsim. <https://github.com/lh3/wgsim>. Accessed 22 November
745 2017.
- 746 16. Wang J, Mullighan CG, Easton J, Roberts S, Heatley SL, Ma J, et al. CREST maps
747 somatic structural variation in cancer genomes with base-pair resolution. *Nat. Methods*.
748 2011;8:652–4.
- 749 17. Chen X, Schulz-Trieglaff O, Shaw R, Barnes B, Schlesinger F, Källberg M, et al. Manta:
750 rapid detection of structural variants and indels for germline and cancer sequencing
751 applications. *Bioinformatics*. 2016;32:1220–2.
- 752 18. Chong Z, Ruan J, Gao M, Zhou W, Chen T, Fan X, et al. novoBreak: local assembly for
753 breakpoint detection in cancer genomes. *Nat. Methods*. 2017;14:65–7.

- 754 19. ICGC-TCGA DREAM Mutation Calling challenge.
755 <https://www.synapse.org/#!Synapse:syn312572/wiki/58893>. Accessed 22 November
756 2017.
- 757 20. Margolin AA, Bilal E, Huang E, Norman TC, Ottestad L, Mecham BH, et al. Systematic
758 analysis of challenge-driven improvements in molecular prognostic models for breast
759 cancer. *Sci. Transl. Med.* 2013;5:181re1.
- 760 21. Marbach D, Costello JC, Küffner R, Vega NM, Prill RJ, Camacho DM, et al. Wisdom of
761 crowds for robust gene network inference. *Nat. Methods.* 2012;9:796–804.
- 762 22. Strobl C, Boulesteix A-L, Zeileis A, Hothorn T. Bias in random forest variable importance
763 measures: illustrations, sources and a solution. *BMC Bioinformatics.* 2007;8:25.
- 764 23. Tattini L, D’Aurizio R, Magi A. Detection of Genomic Structural Variants from Next-
765 Generation Sequencing Data. *Front. Bioeng. Biotechnol.* 2015;3:92.
- 766 24. Boutros PC, Margolin AA, Stuart JM, Califano A, Stolovitzky G. Toward better
767 benchmarking: challenge-based methods assessment in cancer genomics. *Genome Biol.*
768 2014;15:462.
- 769 25. Meyer P, Alexopoulos LG, Bonk T, Califano A, Cho CR, de la Fuente A, et al. Verification
770 of systems biology research in the age of collaborative competition. *Nat. Biotechnol.*
771 2011;29:811–5.
- 772 26. ICGC-TCGA SMC-DNA Meta Challenge.
773 <https://www.synapse.org/#!Synapse:syn4588939/wiki/233672>. Accessed 29 June 2018.
- 774 27. Mohiyuddin M, Mu JC, Li J, Bani Asadi N, Gerstein MB, Abyzov A, et al. MetaSV: an
775 accurate and integrative structural-variant caller for next generation sequencing.
776 *Bioinformatics.* 2015;31:2741–4.
- 777 28. Aird D, Ross MG, Chen W-S, Danielsson M, Fennell T, Russ C, et al. Analyzing and
778 minimizing PCR amplification bias in Illumina sequencing libraries. *Genome Biol.*

- 779 2011;12:R18.
- 780 29. Frampton M, Houlston R. Generation of artificial FASTQ files to evaluate the performance
781 of next-generation sequencing pipelines. PLoS One. 2012;7:e49110.
- 782 30. GitHub Code Repository: BAMSurgeon. <https://github.com/adamewing/bamsurgeon>.
783 Accessed 22 November 2017.
- 784 31. BAMSurgeon Manual.
785 <https://github.com/adamewing/bamsurgeon/blob/master/doc/Manual.pdf>. Accessed 22
786 November 2017.
- 787 32. Picard Tools - By Broad Institute. <http://broadinstitute.github.io/picard/>. Accessed 22
788 November 2017.
- 789 33. ICGC-TCGA DREAM Mutation Calling challenge: Synthetic Tumours.
790 <https://www.synapse.org/#!Synapse:syn312572/wiki/62018>. Accessed 22 November
791 2017.
- 792 34. Novocraft. <http://www.novocraft.com/>. Accessed 22 November 2017.
- 793 35. The Variant Call Format (VCF) Version 4.1 Specification. [https://samtools.github.io/hts-
794 specs/VCFv4.1.pdf](https://samtools.github.io/hts-specs/VCFv4.1.pdf). Accessed 22 November 2017.
- 795 36. Kuhn HW. The Hungarian method for the assignment problem. Nav. Res. Logist. Q.
796 1955;2:83–97.
- 797 37. Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic
798 features. Bioinformatics. 2010;26:841–2.
- 799 38. Haider S, Waggott D, Lalonde E, Fung C, Liu F-F, Boutros PC. A bedr way of genomic
800 interval processing. Source Code Biol. Med. 2016;11:14.
- 801 39. BAMSurgeon v1.1. 2018. <https://doi.org/10.5281/zenodo.1288359>. Accessed 29 June
802 2018.

803

804

805 **Table 1 | Caller scoring schemes.**

Basis of comparison	Region Overlap (Additional file 1: Figure S2a)	Breakpoint Closeness (Additional file 1: Figure S2b)
Description	SVs match if there is sufficient overlap, determined with a Jaccard threshold parameter, between the genomic region associated with the called SV and that of the known SV	SVs match if the breakpoints of the called SV are sufficiently close to the those of the known SV, <i>i.e.</i> breakpoints are within f bp of one another where f is a flank parameter
Strengths	<ul style="list-style-type: none"> • identifies genomic regions affected by the known SVs 	<ul style="list-style-type: none"> • suited to all types of SVs • evaluates precision of breakpoint predictions, facilitating subsequent breakpoint validation
Weaknesses	<ul style="list-style-type: none"> • some SVs are not accurately defined by genomic regions, e.g. an insertion may be characterized by a single breakpoint • need criteria to define sufficient overlap 	<ul style="list-style-type: none"> • need criteria to define sufficient closeness

806

807 **Table 2 | Teams.**

Team ID	Method name	Institute
Standard	BreakDancer	Challenge Administrators
Team 1*	Delly ¹	European Molecular Biology Laboratory (EMBL)
Team 2	Manta	Illumina, Inc.
Team 3	Meerkat	Harvard Medical School
Team 4	novoBreak	MD Anderson Cancer Center
Team 5	deStruct ²	BC Cancer Agency Research Centre
Team 6	SWAN	Wharton School
Team 7	SmuFin	Barcelona Supercomputing Center
Team 8	not available	Peking University
Team 9*	deStruct ²	Simon Fraser University
Team 10	GROM	Rutgers University
Team 11	Pindel	McDonnell Genome Institute
Team 12	PAVFinder	Canada's Michael Smith Genome Sciences Centre
Team 13	Delly ¹	Georgetown University Medical Center, National Cancer Institute
Team 14	Subread	Walter and Eliza Hall Institute of Medical Research
Team 15	CX	Wharton School

808 Teams that made submissions for IS1, IS2 and/or IS3, the names of the SV detection methods
 809 they used and the institutes to which they belong.

810 ¹Delly was developed by Team 1 and also used by Team 13.

811 ²deStruct was developed by Team 9 and also used by Team 5.

812

813 Figure Legends

814 **Fig. 1 | BAMSurgeon simulates SVs in genome sequences.**

815 Method for adding SVs to existing BAM alignments. **a** Overview of SV (e.g. deletion) spike-in:
816 Starting with an original BAM (i), a region (ii) is selected where a deletion is desired. iii) Contigs
817 are assembled from reads in the selected region, and the contig is rearranged by deleting the
818 middle. The amount of contig deleted is a user-definable parameter. Read coverage is
819 generated over the contig using wgsim to match the number of reads per base in the original
820 BAM. Since the deletion contig is shorter than the original, fewer reads will be required to
821 achieve the equivalent coverage. iv) Generated read pairs include discordant pairs (*i.e.* paired
822 reads that do not align to the reference genome with the expected relative orientation and inner
823 distance) spanning the deletion and clipped reads (*i.e.* reads that are only partially aligned to the
824 reference). Reads mapping to the deleted region of the contig are not included in the final BAM.
825 **b,c** To test the robustness of BAMSurgeon with respect to changes in **(b)** aligner and **(c)** cell
826 line, we compared the ranks of CREST, Delly, Manta and novoBreak on two new tumour-normal
827 data sets: one with an alternative aligner, NovoAlign, and the other on an alternative cell line,
828 HCC1954 BL. Callers were scored with $f = 100$ bp (Additional file 1: Figure S2b); Manta retained
829 the top position, independent of aligner and cell line. **d** Summary of the three *in silico* (IS)
830 tumours described here. Abbreviations: DEL, deletion; DUP, duplication; INV, inversion; INS,
831 insertion.

832

833 **Fig. 2 | Overview of the SV Calling Challenge submissions.**

834 **a** Precision-recall plot of IS1 submissions. Each point represents a submission, each colour
835 represent a team and the best submission from each team (top F -score) is circled. The
836 “Standard” point corresponds to the reference point submission provided by Challenge
837 organizers. **b** The F -scores of submissions on the training and testing sets are highly correlated
838 for IS1 (Spearman’s $\rho = 0.98$), falling near the plotted $y = x$ line.

839

840 **Fig. 3 | Performance optimization by parameterization and ensembles.**

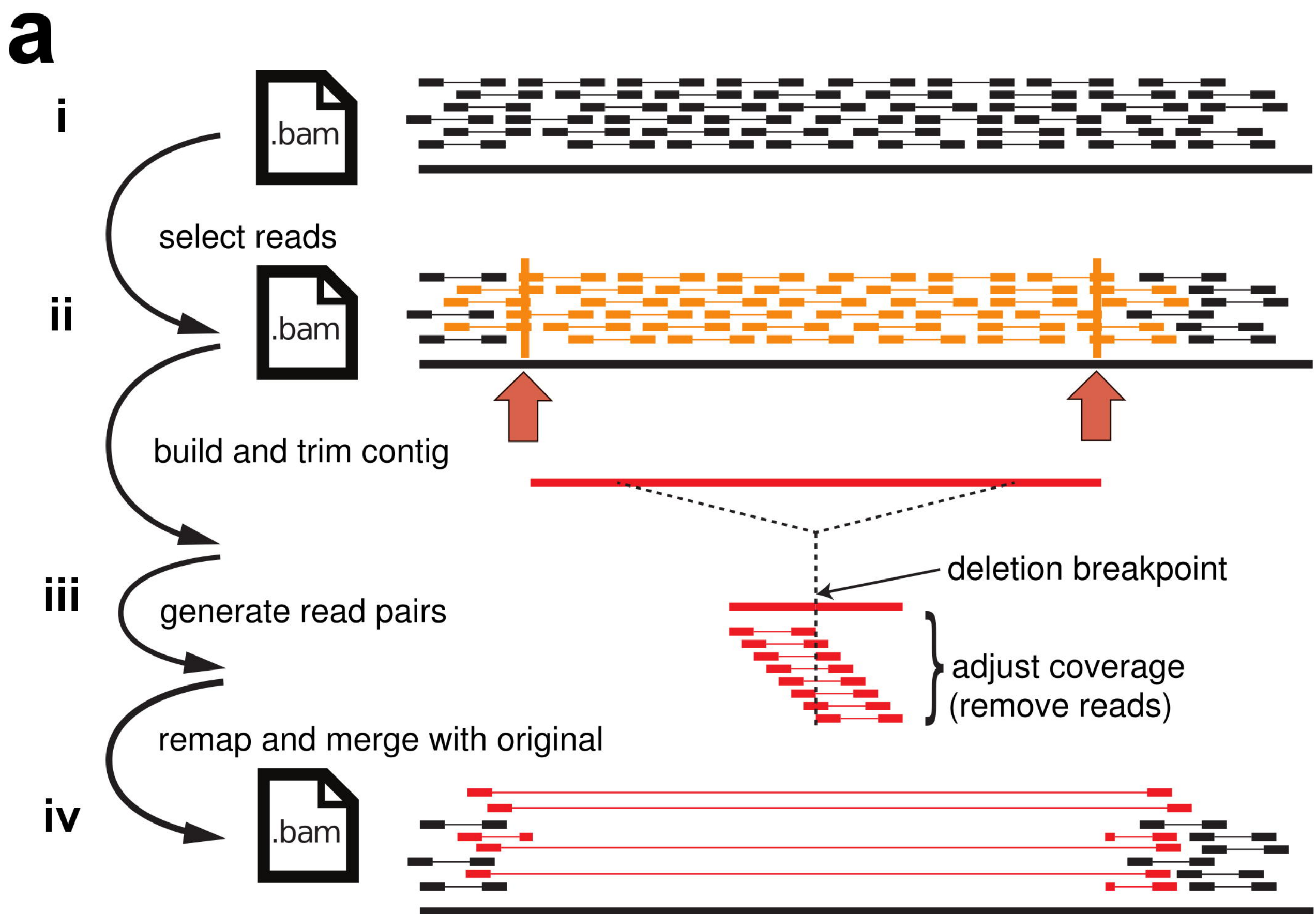
841 **a** Recall, precision and F -score of all IS1 submissions plotted by team, then submission order.
842 Teams were ranked by the F -score of their best submission, colour coding (top bar) as in Fig. 2.
843 The “Standard” lines correspond to the reference point submission provided by Challenge
844 organizers. **b** For each tumour, the improvement in F -score from the initial (“naive”) to the best
845 (“optimized”) submissions of each team. Darker shades of blue indicate greater improvement. **c**
846 For each tumour, team rankings based on their naive or optimized submissions. Larger dot
847 sizes indicate better ranks by F -score. **b,c** An “X” indicates that the team did not make a
848 submission for the specific tumour (or changed team name). **d** Recall, precision and F -score of
849 ensembles versus individual submissions for IS1. At the k th rank, the triangles indicate
850 performance of the ensemble of the top k submissions, and the circles indicate performance of
851 the k th ranked submission. The ensemble analysis focused on the best submission from each
852 team.

853

854 **Fig. 4 | Characteristics of prediction errors.**

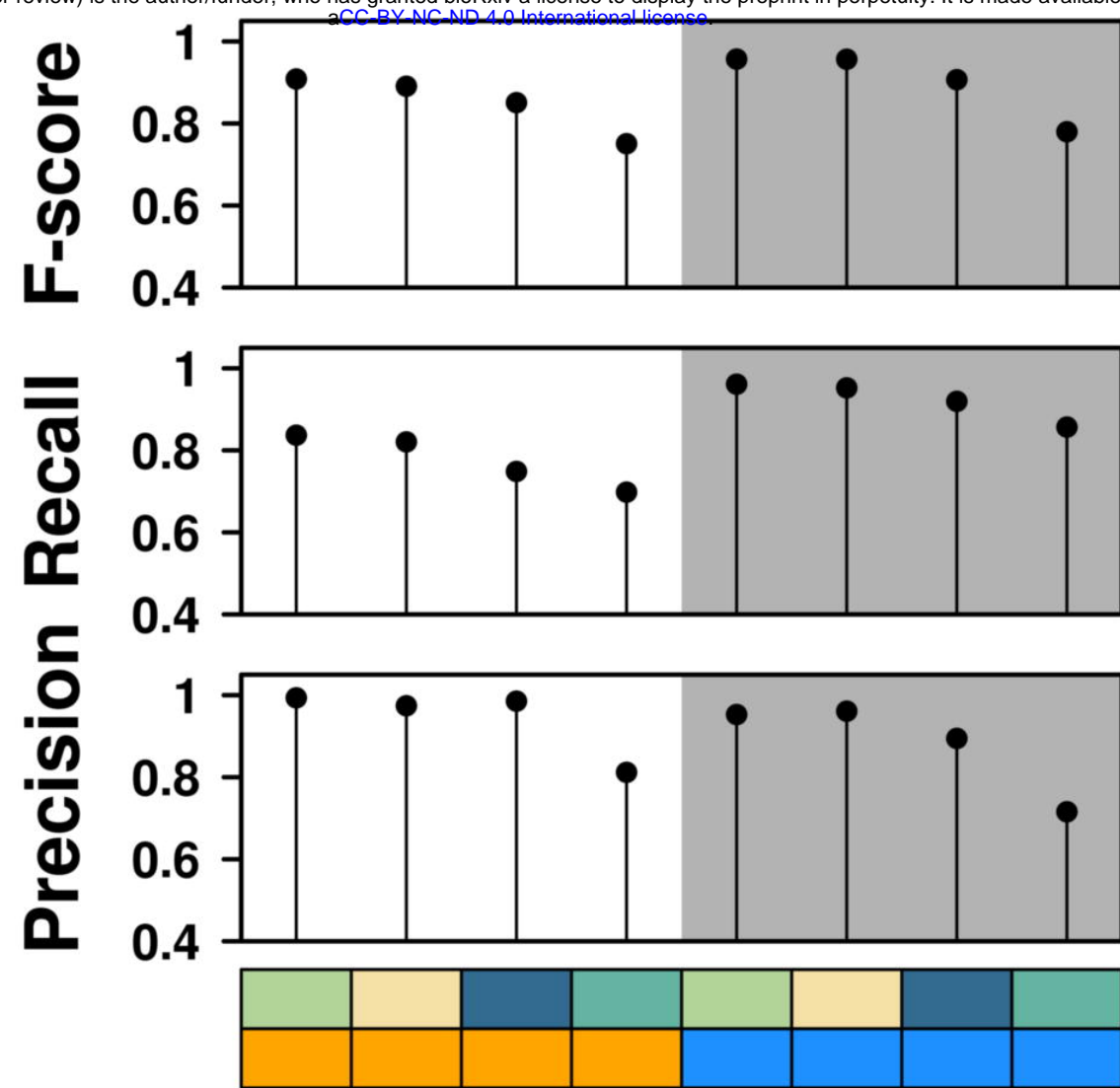
855 Random Forests assess the importance of 16 sequence-based variables for each caller’s FN

856 (a,c,e,g,i) and FP (b,d,f,h,j) breakpoints. Each panel shows variable importance on the left,
857 where each row represents the best performing set of predictions by the given team/caller (on
858 the given *in silico* tumour), and each column represents the indicated variable. Dot size reflects
859 variable importance, *i.e.* the mean change in accuracy caused by removing the variable from
860 the model (generated to predict erroneous breakpoints). Colour reflects the directional effect of
861 each variable (red and blue for greater and lower variable values, respectively, associated with
862 erroneous breakpoints; black for categorical variables or insignificant directional associations,
863 two-sided Mann-Whitney $P > 0.01$). Background shading indicates the accuracy of the model
864 (see colour bar). Variable importance for FN and FP breakpoints in each of the three tumours is
865 shown for the following SV callers: CREST (a,b), Delly (c,d) and Manta (e,f). Manta only called
866 two FPs in IS1; thus, variable importance for FP breakpoints could not be computed (indicated
867 by Xs in the plot). Variable importance for FN and FP breakpoints in IS2 (g,h) and IS3 (i,j) is
868 shown for each team. In the right plot (g-j), the first four columns indicate usage of the indicated
869 algorithmic approaches by each team, and the last column indicates the aligner used. Grey
870 indicates that algorithmic approaches and aligner are unknown for the given team.
871 Abbreviations: Algm, algorithm; SNP, single-nucleotide polymorphism; INDEL, short insertion or
872 deletion.

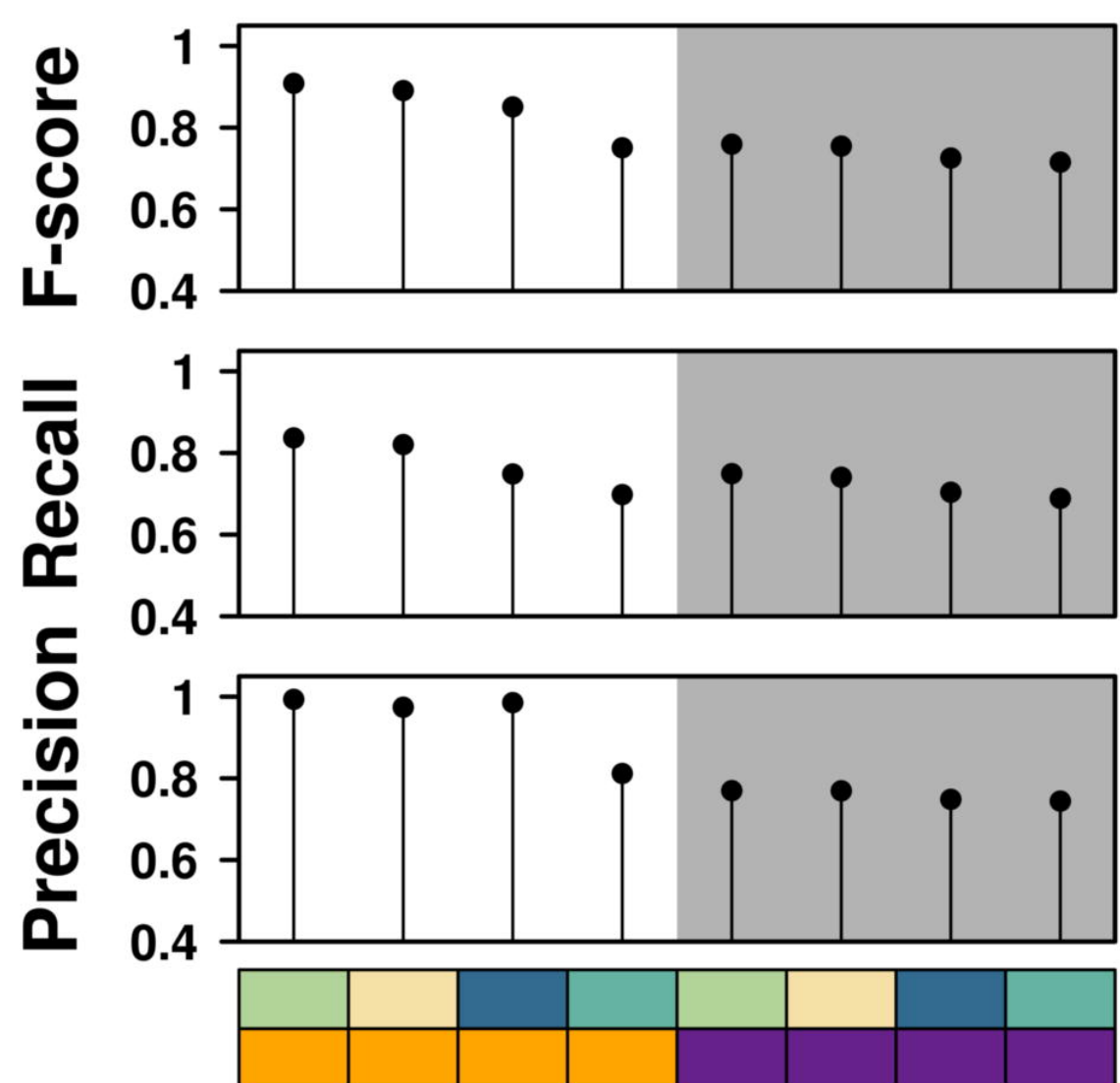


b

bioRxiv preprint doi: <https://doi.org/10.1101/224733>; this version posted August 28, 2018. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

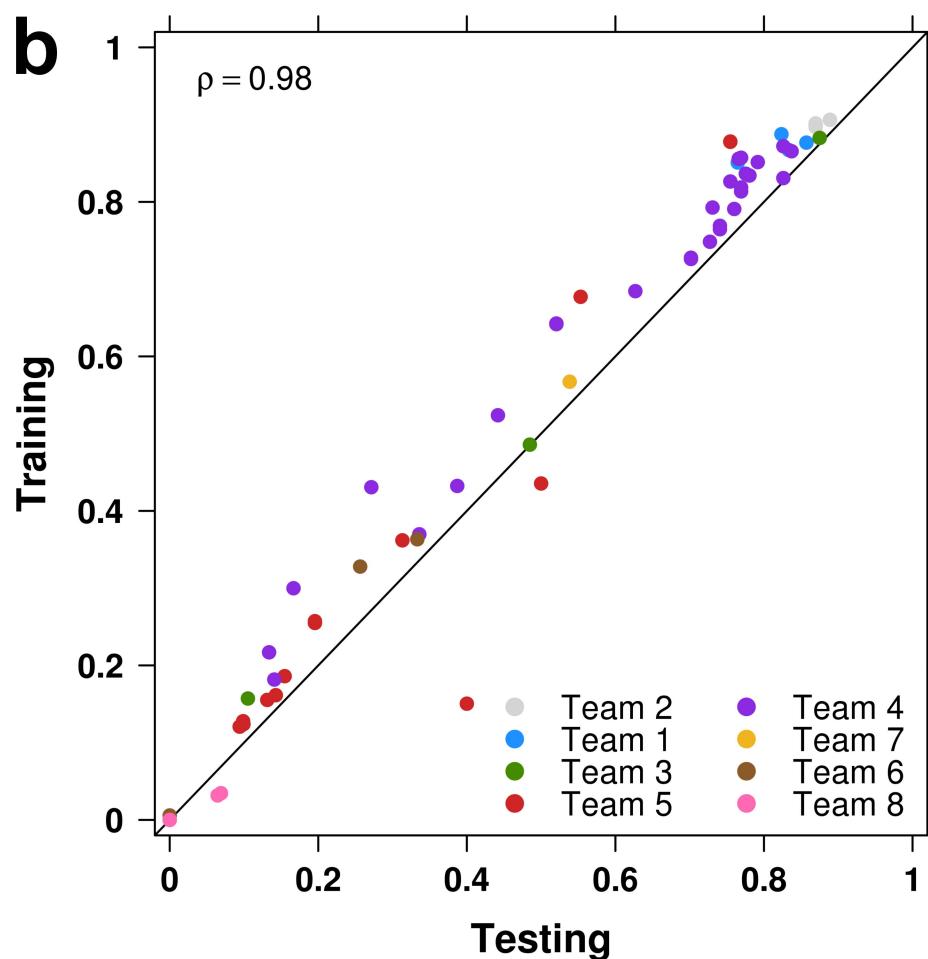
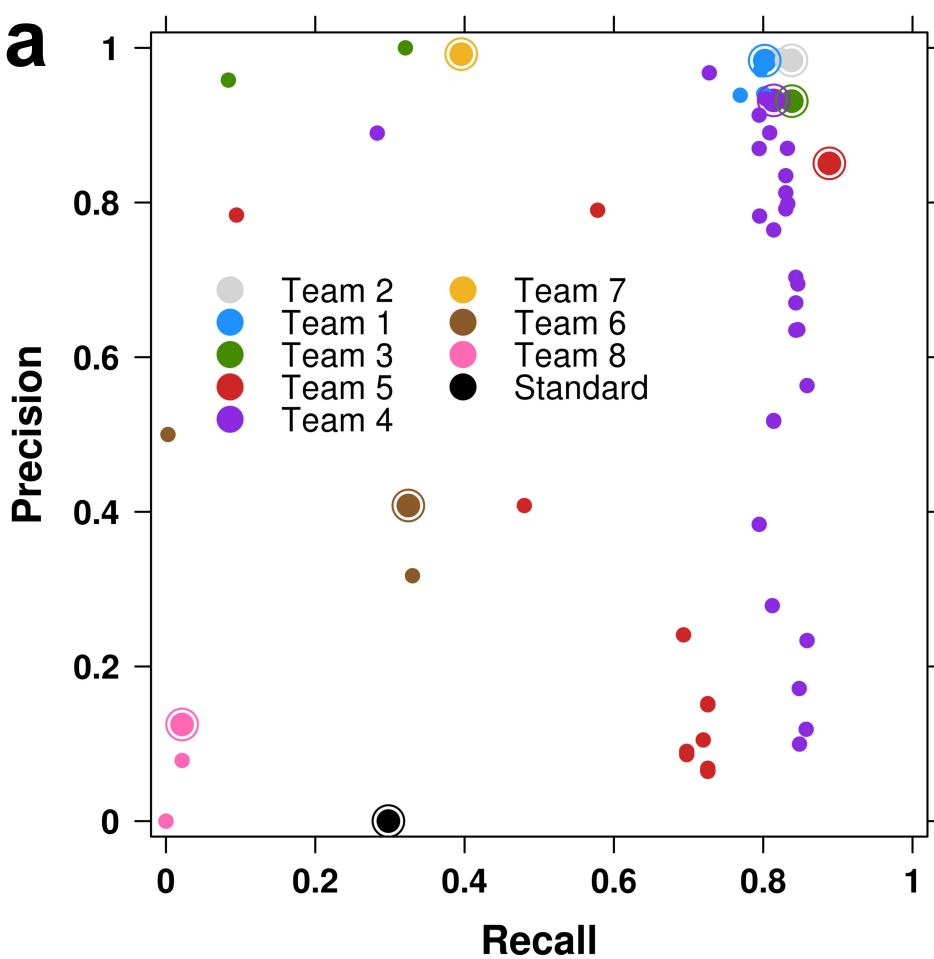


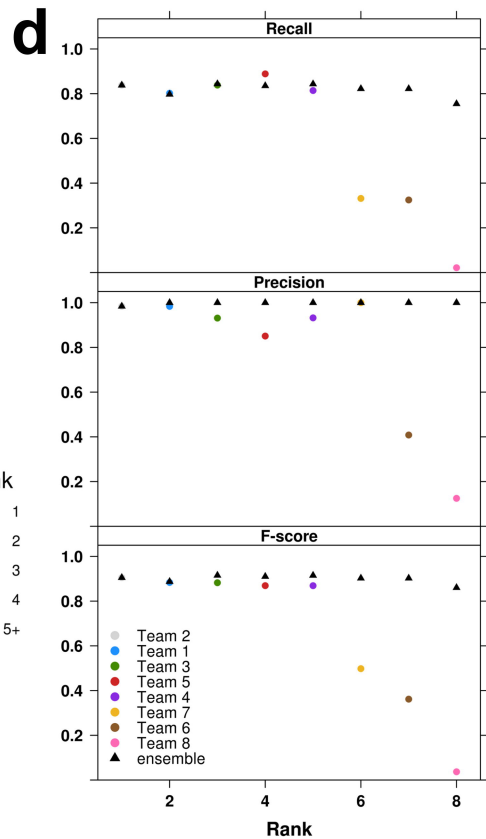
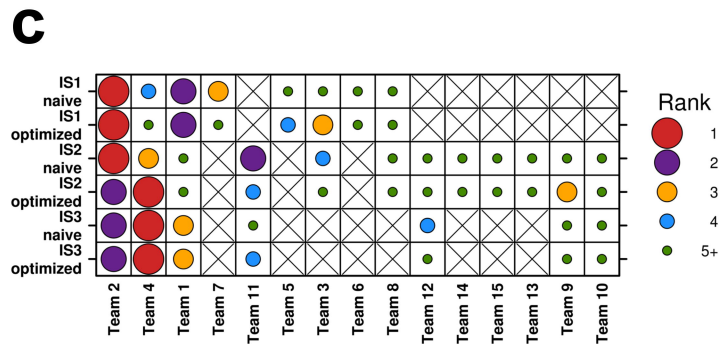
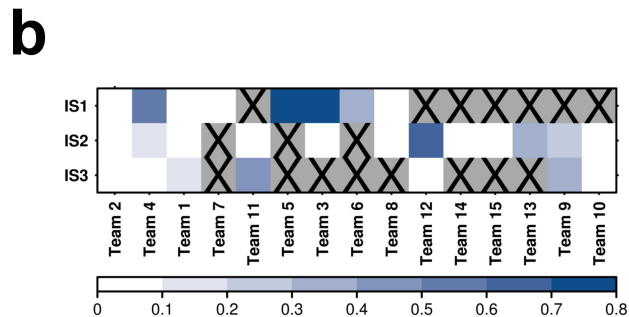
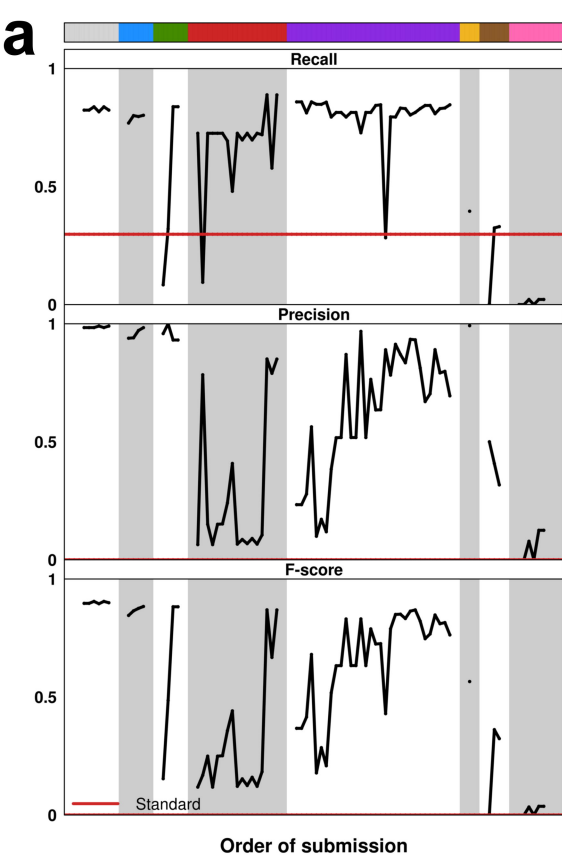
c



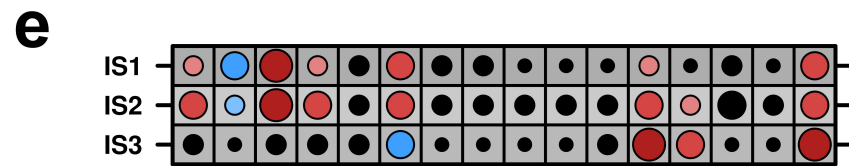
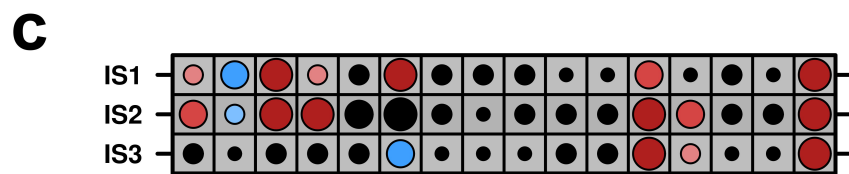
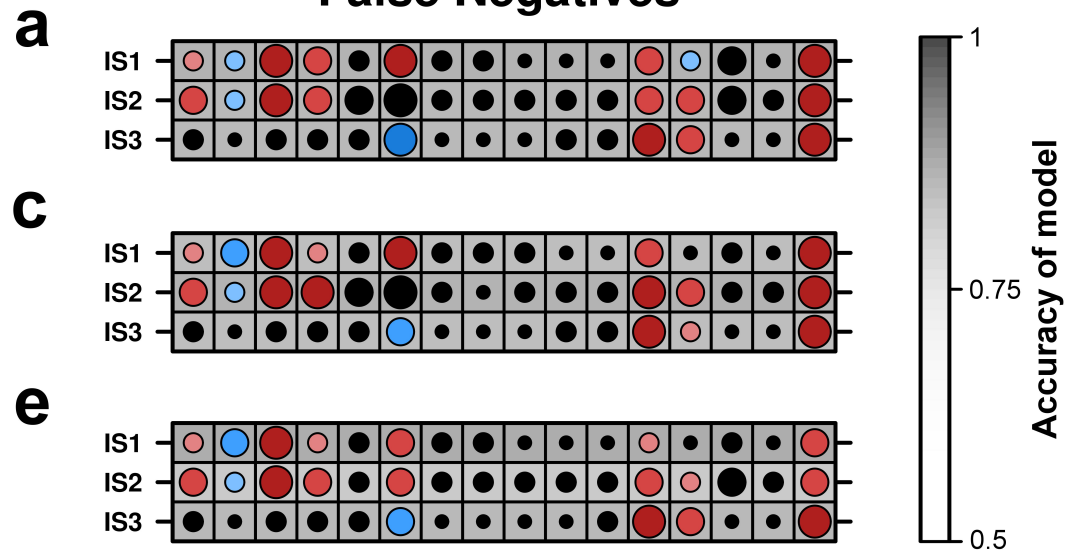
d

Tumour	Cell line	Number of somatic SVs	SV types	Cellularity (%)
<i>in silico</i> 1	HCC1143 BL	371	DEL, DUP, INV	100
<i>in silico</i> 2	HCC1954 BL	655	DEL, DUP, INV, INS	80
<i>in silico</i> 3	HCC1143 BL	2,886	DEL, DUP, INV, INS	100





False Negatives



False Positives

