# 1 Combining accurate tumour genome simulation

# 2 with crowd-sourcing to benchmark somatic

# 3 structural variant detection

4

5 Anna Y. Lee[1,12], Adam D. Ewing[2,3,12], Kyle Ellrott[2,4,12], Yin Hu[5], Kathleen E. Houlahan[1], J.

6 Christopher Bare[5], Shadrielle Melijah G. Espiritu[1], Vincent Huang[1], Kristen Dang[5], Zechen

7 Chong[6,7,8], Cristian Caloian[1], Takafumi N. Yamaguchi[1], ICGC-TCGA DREAM Somatic Mutation

8 Calling Challenge Participants, Michael R. Kellen[5], Ken Chen[6], Thea C. Norman[5], Stephen H.

9 Friend[5], Justin Guinney[5], Gustavo Stolovitzky[9], David Haussler[2], Adam A. Margolin[4,5,13], Joshua

10 M. Stuart[2,13], Paul C. Boutros[1,10,11,13]

11

12 1 Ontario Institute for Cancer Research; Toronto, Ontario, Canada

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

14 CA, USA

15 3 Mater Research Institute; University of Queensland; Woolloongabba, QLD, Australia

16 4 Computational Biology Program; Oregon Health & Science University; Portland, OR, USA

17 5 Sage Bionetworks; Seattle, WA, USA

18 6 Department of Bioinformatics and Computational Biology; University of Texas MD Anderson

19 Cancer Center; Houston, TX, USA

20 7 Department of Genetics; University of Alabama at Birmingham; Birmingham, AL, USA

21 8 Informatics Institute; University of Alabama at Birmingham; Birmingham, AL, USA

22    9 IBM Computational Biology Center; T.J.Watson Research Center; Yorktown Heights, NY,

23    USA

24    10 Department of Medical Biophysics; University of Toronto; Toronto, Ontario, Canada

25    11 Department of Pharmacology & Toxicology; University of Toronto; Toronto, Ontario, Canada

26    12 These authors contributed equally

27    13 Corresponding authors

28

# Abstract

## Background

31    The phenotypes of cancer cells are driven in part by somatic structural variants (SVs). SVs can

32    initiate tumours, enhance their aggressiveness and provide unique therapeutic opportunities.

33    Whole-genome sequencing of tumours can allow exhaustive identification of the specific SVs

34    present in an individual cancer, facilitating both clinical diagnostics and the discovery of novel

35    mutagenic mechanisms. A plethora of somatic SV detection algorithms have been created to

36    enable these discoveries, however there are no systematic benchmarks of them. Rigorous

37    performance evaluation of somatic SV detection methods has been challenged by the lack of

38    gold-standards, extensive resource requirements and difficulties in sharing personal genomic

39    information.

## Results

41    To facilitate SV detection algorithm evaluations, we created a robust simulation framework for

42    somatic SVs by extending the BAMSurgeon algorithm. We then organized and enabled a

43    crowd-sourced benchmarking within the ICGC-TCGA DREAM Somatic Mutation Calling

44    Challenge (SMC-DNA). We report here the results of SV benchmarking on three different

45  tumours, comprising 204 submissions from 15 teams. In addition to ranking methods, we

46  identify characteristic error-profiles of individual algorithms and general trends across them.

47  Surprisingly, we find that ensembles of analysis pipelines do not always outperform the best

48  individual method, indicating a need for developing new ways to aggregate somatic SV

49  detection approaches.

50  ## Conclusions

51  The synthetic tumours and somatic SV detection leaderboards remain available as a community

52  benchmarking        resource,        and        BAMSurgeon        is        available        at

53  https://github.com/adamewing/bamsurgeon.

54  ## Keywords

55  somatic mutations, simulation, structural variants, benchmarking, cancer genomics, whole-

56  genome sequencing, crowd-sourcing

57

58  # Background

59  Somatic structural variants (SVs) are mutations that arise in tumours involving rearrangements,

60  duplications or deletions of large segments of DNA. SVs are often defined to be events larger

61  than 100 bp in size, although with significant variability in this definition. Somatic SVs are critical

62  in driving and regulating tumour biology. They can initiate tumours [1,2] and because they are

63  unique to the cancer, can serve as highly-selective avenues for therapeutic intervention [3]. The

64  overall mutation load of somatic SVs serves as a proxy for genomic instability, and can robustly

65  predict tumour aggressiveness in multiple tumour types [4,5].

66    While somatic SVs that alter copy-number can be detected using microarray assays, the

67    resolution of such studies is limited, and many other important types of SVs cannot be detected.

68    As a result, high-throughput DNA sequencing is now a standard approach for detecting SVs in

69    cancer genomes. Although RNA-based assays are useful for detecting SVs that alter protein-

70    structure, DNA-based assays are required for most others. As a result, a broad range of

71    algorithms has been developed to detect SVs from short-read sequencing data using read

72    depth analysis, split read (*i.e.* a read that maps to different parts of the reference sequence)

73    alignment, paired end mapping and de novo assembly techniques [6–9]. However, the accuracy

74    of existing methods is poorly described. There are no comprehensive benchmarks of somatic

75    SV detection approaches. Most comparison results are reported by the developers of newly

76    published methods. These developer-run benchmarks are potentially subject to several types of

77    selection biases. For example, the developers of one tool may be experts in parameterizing and

78    tuning it, but may lack the same skill in tuning methods developed by others. Further, evaluating

79    the accuracy of somatic SV detection is more challenging than evaluating the accuracy of

80    somatic single nucleotide variant (SNV) detection as validation data is more difficult to generate

81    for SVs. Even the metrics of measuring accuracy are not agreed upon, with no community-

82    accepted standards on how SV prediction accuracy should be assessed, especially when

83    predictions are close to, but not exactly at, the actual sequence breakpoints. As a result, there

84    are no robust estimates of the false positive and false negative rates of somatic SV prediction

85    tools on tumours of different characteristics.

86    To fill this gap,  we created an open challenge-based assessment of somatic SV prediction tools

87    as part of the ICGC-TCGA DREAM Somatic Mutation Calling Challenge (the Challenge). We

88    first extended BAMSurgeon [10], a tool for creating synthetic mutations, to generate somatic

89    SVs. We then created and distributed three synthetic tumours, on which 204 submissions were

90    made by 15 teams.

# Results

## Simulation of SVs with BAMSurgeon

In addition to point mutations [SNVs and short insertions or deletions (INDELs)], BAMSurgeon is capable of creating SVs through read selection, local sequence assembly, manipulation of assembled contigs, and simulation of sequence coverage over the altered contigs (Fig. 1a, Additional file 1: Figure S1). This, combined with careful tracking of read depth, yields approximations of SVs including insertions, deletions, duplication, and inversions into pre-existing backgrounds of real sequence data. The BAMSurgeon manual, available online, contains a full description of input formatting and available parameters. The input regions define where local assembly will be attempted *via* Velvet [11]. For each region, the largest assembled contig is selected and re-aligned to the reference genome using Exonerate [12]. The contig is then trimmed to the length of its longest contiguous alignment and the alignment is used to accurately track breakpoint locations within the contig in terms of reference coordinate space. The location and identity of reads from the original BAM file in the assembled contig is tracked *via* parsing of the AMOS [13] file output by Velvet [14], which also enables tracking of reads included or excluded after contig trimming. If a suitable contig is not available for a given input segment, no mutation is made for that segment. For each segment where contig assembly succeeds, the contig is rearranged according to the user specification (*e.g.* insertion, deletion, duplication, or inversion of sequence). Following rearrangement of the contig, paired reads are simulated from the rearranged contig using wgsim [15], with specific parameters controllable by the user. Because reads are simulated using the rearranged contig, breakpoint-spanning reads and reads that will be discordant versus the reference genome assembly will be created. The number of reads simulated (final coverage, $C_f$) depends on the original coverage $C_o$, the

114    difference in length between the original contig $L_o$ and the rearranged contig $L_f$, and a user-

115    specified parameter controlling variant allele fraction (VAF). Thus, $C_f = VAF*C_o*(L_f/L_o)$.

116    Duplications and insertions result in larger contigs and require new reads to be added to the

117    final BAM, and deletions yielding a smaller contig require reads to be removed from the final

118    BAM. In cases where reads must be added (duplications and insertions), additional reads are

119    added to the final BAM. Conversely, where reads need to be removed from the original BAM, a

120    list of reads to be deleted is maintained, which correspond to reads covering the deleted region

121    in the original BAM.

## Validation of simulated somatic SVs

123    To validate SVs simulated by BAMSurgeon, we performed a series of quality-control

124    experiments analogous to those performed to validate simulated SNVs [10]. Briefly, we used

125    BAMSurgeon to generate synthetic tumour-normal pairs, with the same set of target mutations,

126    that differ by the division of reads into tumour and normal sequence sets, aligner or cell line.

127    The target mutation set was designed to generate a synthetic tumour with a baseline level of

128    complexity and thus did not include insertions. We ran four SV callers using default parameters

129    on each pair: two widely used callers, CREST [16] and Delly [9], and two callers developed over

130    the course of the Challenge, Manta [17] and novoBreak [18]. We did not optimize parameters

131    for the callers since the goal of this validation was not to identify the best caller, but instead to

132    verify that the caller ranking is maintained across analogous datasets.

133    The definition of a SV suggests different scoring schemes for measuring the performance of a

134    caller. All SVs can be defined by at least one breakpoint; deletions, duplications and inversions

135    are SVs defined by a pair of breakpoints which in turn define a genomic region. Thus, we

136    compared called SVs to true-positive SVs based on i) region overlap or ii) breakpoint closeness

137    (Table 1, Additional file 1: Figure S2). The Challenge initially used a scoring scheme based on

138  region overlap (at least one or more bases in common; Additional file 1: Figure S2a). Here we

139  focus on the breakpoint closeness scheme since it is well suited for all types of SVs. A called

140  SV that is sufficiently similar to a known SV based on such criteria was considered a true

141  positive; otherwise, a false positive. We used such annotations to assess the performance of a

142  caller in terms of precision (fraction of calls that are true), recall (fraction of known SVs called)

143  and *F*-score (harmonic mean of precision and recall).

144  We performed several quality-control experiments. First, the caller ranking (by *F*-score) was

145  independent of the random division of reads: Manta > novoBreak > CREST > Delly (Additional

146  file 1: Figure S3a,b). Second, the same ranking was observed when alignments were generated

147  either using the Burrows-Wheeler Aligner (BWA) or NovoAlign with and without INDEL

148  realignment (*i.e.* local realignment to minimize mismatches across reads due to INDELs relative

149  to the reference genome), indicating that the ranking was independent of the aligner used (Fig.

150  1b, Additional file 1: Figure S3c). Lastly, when the genomic background was varied by using

151  HCC1143 BL or HCC1954 BL sequence data, the caller ranking was largely independent of the

152  cell line: Manta and novoBreak retained first and second place, respectively, while CREST and

153  Delly swapped places, although their *F*-scores were very similar to each other when HCC1954

154  BL was used (Fig. 1c, Additional file 1: Figure S3d). Overall, these results show that simulated

155  SVs are robust to changes in the read division, aligner and genomic background.

## Crowd-sourced benchmarking of somatic SV calling

157  The SV component of the Challenge consisted of the same three synthetic tumour-normal data

158  sets used in the SNV component [10]. Briefly, the data sets were derived from existing cell line

159  sequence data (thus minimizing data access restrictions) and *in silico* tumours 1-3 (IS1-IS3)

160  were generated with increasing complexity (Fig. 1d). In terms of SVs, breakpoint locations were

161  randomly selected and the tumours had increasing mutation rates (371 *vs.* 2,886 somatic SVs in

162     IS1 and IS3, respectively). Moreover, IS1 contained deletions, duplications and inversions while

163     IS2 and IS3 additionally contained insertions. Like the SNV component, the SV component of

164     the Challenge was implemented using the Dialogue for Reverse Engineering Assessments and

165     Methods (DREAM) framework. Briefly, information about the Challenge was shared on its

166     website [19], participants registered online, downloaded a data set, applied their SV calling

167     pipelines to the data set and submitted the results in Variant Call Format (VCF) v4.1. IS1-IS3

168     were released sequentially, each data set had its own competition phase and participants could

169     make multiple submissions for each data set. Each tumour genome was divided into a training

170     set and a testing set by holding out a portion of the genome. During the competition phase,

171     leaderboards showed performance measures on the training set. After the competition closed,

172     leaderboards also showed performance measures on the whole genome (training + testing

173     sets).

174     The Challenge administration team prepopulated the leaderboards with two submissions and

175     the community provided 204 submissions from 15 teams (Additional file 2). A list of all

176     submissions, and descriptions of pipelines used to generate them, can be found in Additional

177     files 3 and 4, respectively. The submissions were surprisingly discordant in format. Some

178     specified SV types that are not recognized as VCF formatted types, and between 5.5-11% of all

179     submissions were not made in valid VCF format (Additional file 5). For this reason, and the

180     ambiguity of specifying SV types (*i.e.* the same SV can be specified with a specific type, or by

181     specifying the breakpoints and break-end adjacencies), type specifications were ignored when

182     scoring submissions. Team ranking varied with the stringency of the scoring (Additional file 1:

183     Figure S2d-i). For simplicity, we focused on scoring with $f = 100$ bp due to the balance between

184     the median and variance of the resulting $F$-scores (Additional file 1: Figure S4). While the top-

185     performing teams achieved near maximal precision on the simplest tumour, IS1, their recall

186     remained less than 0.9 (Fig. 2a), and decreased further on the other tumours (Additional file 1:

187    Figure S5a,b). On all three tumours, *F*-scores on the training and testing sets were highly

188    correlated (Spearman's rank correlation coefficient ($\rho$) ≥ 0.98; Fig. 2b, Additional file 1: Figure

189    S5c,d). However, the slightly elevated *F*-scores in the training sets observed for IS1 and IS2

190    may reflect minor overfitting; overfitting occurs when a statistical model is tuned to the training

191    set, limiting generalizability. Notably, the total number of somatic SV mutations in IS3 is >4x that

192    for IS1 and IS2 (Fig. 1d). Conversely, the percentage of mutations used for training is greater for

193    IS1 (93%) and IS2 (92%) *vs.* IS3 (89%). Sampling from the IS3 mutations, we simulated training

194    and testing sets of different sizes, and computed the differences between the *F*-scores on the

195    training sets and the *F*-scores on the testing sets. We found that that the differences tend to be

196    greater when the percentage of mutations used for training is greater (Additional file 1: Figure

197    S5e). This suggests that the *F*-score differences observed for IS1 and IS2 are at least in part an

198    artefact of training set size.

## Pipeline optimization

200    The within-team variability in *F*-scores accounts for 21-43% of the total per-tumour variance in

201    *F*-scores. The large variability in submissions by certain teams highlights the impact of tuning

202    parameters during the Challenge (Fig. 3a, Additional file 1: Figure S6a,b). In comparing the

203    initial ("naive") and best ("optimized") submissions of each team, for each tumour, the maximum

204    *F*-score improvement was 0.75 (from 0.12 to 0.87 by Team 5 for IS1), and the median

205    improvements were 0.20, 0.01, and 0.07 for IS1, IS2 and IS3, respectively (Fig. 3b). At least

206    33% of teams improved their *F*-score by at least 0.05 and at least 25% of teams improved it by

207    more than 0.20, depending on the tumour. Despite these improvements by parameterization,

208    team rankings were only moderately changed: if a team's naive submission ranked in the top

209    three, their optimized submission remained in the top three 66% of the time (Fig. 3c).

210    Given the crowd-sourced nature of the Challenge, we explored "wisdom of the crowds" as an

211    approach to optimize performance [20,21]. Specifically, we aggregated SV calls into an

212    ensemble by first identifying sufficiently similar calls in the majority of the top *k* submissions.

213    Pairwise distances between calls from different submissions were computed (*i.e.* a breakpoint-

214    length distance that incorporates distances between breakpoints and differences in SV length,

215    Additional file 1: Figure S2c), and those calls with distances less than a selected threshold

216    (equal to *f*, for consistency) were considered to represent an equivalent called SV event. The

217    chromosome together with the median start and end positions of a set of similar calls would

218    then define a single ensemble SV prediction. We considered two variations of this approach: i) a

219    baseline approach with ensembles of the best submission from each team, and ii) a

220    conservative approach with ensembles of all submissions and more stringent aggregation of

221    called SVs (see Methods). The baseline ensembles were found to have *F*-scores comparable to

222    that of the best submission (*e.g.* for IS1, the best ensemble and submission have *F*-scores of

223    0.92 and 0.91, respectively; Fig. 3d, Additional file 1: Figure S7b). However, the ensembles had

224    lower *F*-scores than the best submission for IS2 (Additional file 1: Figure S7a). When *k* > 15, we

225    found that the conservative ensemble *F*-scores drop further below that of the best submission

226    (Additional file 1: Figure S7c-e; *e.g.* for IS1, the best ensemble with *k* > 15 and the best

227    submission have *F*-scores of 0.83 and 0.91, respectively). In contrast, the precision of all

228    ensembles (range: 0.993-1.00) was similar or slightly improved compared to that of the best

229    submission. Thus, any changes in the ensemble *F*-scores were mostly influenced by the

230    changes in recall as *k* varied.

231    ## Error characteristics

232    We next exploited the large number of independent analyses to identify characteristics

233    associated with false negatives (FNs) and false positives (FPs). For example, error profiles

234    differed significantly between subclonal populations in IS3, with greater FN rates for mutations

235    present at lower VAFs (Additional file 1: Figure S8; one-sided Wilcoxon signed rank $P = 0.02$ for

236    VAF = 0.2 vs. 0.33, $P = 0.04$ for VAF = 0.33 vs. 0.5, $n = 7$). We also selected the best

237    submission from each team (by $F$-score) and focused on 14 variables associated with

238    breakpoint positions, representing factors like coverage and mapping quality (Additional file 6).

239    Several of these variables were associated with false-positive rates; in particular, tumour

240    coverage ($R > 0.24$), bridging reads count (the number of reads that bridge a putative

241    breakpoint, $R > 0.21$) and mapping quality ($R < -0.29$), have stronger associations with FPs for

242    both IS2 and IS3, compared to other variables (Additional file 1: Figure S9a, S10-S25). By

243    contrast, few were associated directly with false-negative rates ($0 \leq |R| \leq 0.15$; Additional file 1:

244    Figure S9b, S10-S25).

245    To evaluate whether these variables, and additional categorical variables, contribute

246    simultaneously to somatic SV prediction error, we generated two Random Forests (non-

247    parametric learning models that can trivially merge multiple data types) [22] for each team to

248    assess variable importance for FN and FP breakpoints separately. FN breakpoints are

249    associated with variables such as high bridging reads count and strand bias (Fig. 4a,c,e,g,i;

250    Additional file 1: Figure S26a). FP breakpoints are generally associated with variables such as

251    low mapping quality (Fig. 4b,d,f,h,j; Additional file 1: Figure S26b).

252    By executing specific SV callers, CREST (Fig. 4a,b), Delly (Fig. 4c,d) and Manta (Fig. 4e,f), with

253    the same parameters on all three tumours, we identified tumour-specific error profiles. For

254    example, the distance to the nearest germline INDEL tends to have stronger associations with

255    errors in IS2 and IS3 compared to IS1 (Fig. 4a-e). Team-specific error profiles are more

256    apparent with the FP breakpoint analysis. For example, Teams 8 and 10 have distinct FP

257    profiles for the same tumour, IS2 (Fig. 4h); FPs by Teams 8 and 10 are negatively and positively

258 associated with tumour coverage, respectively. Algorithmic approaches to SV calling from

259 sequencing data include i) read depth analysis, ii) paired end mapping, iii) split read alignment,

260 and iv) *de novo* assembly [23]. Some teams submitted sufficient algorithm details to determine

261 the general approaches used, as well as the choice of aligner (Fig. 4g-j). Based on the available

262 annotations, teams using the same aligner do not have error profiles that tightly cluster for all

263 three tumours, suggesting that the aligner is not as strong a driver of those profiles, compared

264 to the caller algorithm.

265

# Discussion

266

267 Crowd-sourced benchmarking challenges are ideal for questions where significant diversity in

268 algorithmic approaches exists, particularly where individual methods are highly parameterized

269 or computationally intensive [24,25]. The detection of variants from high-throughput sequencing

270 data fits these criteria well: dozens of algorithms are in common use, many with complicated

271 sets of parameters to tune and most requiring tens to hundreds of CPU hours to execute. We

272 have quantified the critical importance of parameterization: it accounts for 21-43% of the

273 variability in performance across the 204 submissions evaluated. This is comparable to the 26%

274 of variability observed in somatic SNV detection benchmarking [10], and highlights the need for

275 algorithm developers to continue to optimize parameters, provide guidance for their tuning and

276 work toward automating their selection to make their software easier to use.

277 The "wisdom of the crowds" is the idea that an ensemble of multiple algorithms can significantly

278 outperform any individual method. Several crowd-sourced benchmarking competitions from

279 diverse fields have shown great success in combining submissions from contestants to achieve

280 a high-performing meta-predictor including challenges for somatic SNV detection [10], gene

281 regulatory network inference [21] and mRNA-based prognostic signatures for breast cancer

282 [20]. By contrast, in somatic SV detection, we do not have clear evidence that an ensemble

283 improves on the best individual method consistently across different tumours, despite testing

284 several ways of creating ensembles. This may reflect the large diversity in the biases of each

285 individual algorithm (Fig. 4), or it may represent the unique challenges of scoring SVs. While

286 some SV classes may be well-represented by overlap-based scoring methods, others benefit

287 more from breakpoint-based scoring, and the choice of scoring metric and parameter must be

288 tuned to the specific biological question of interest. For example, breakpoint identification is

289 critical when considering translocations -- especially those generating candidate fusion proteins

290 -- while overlap of the called and known regions is much more important for copy-number

291 analyses. The fact that the "wisdom of the crowds" *via* majority vote approach works very well

292 for somatic SNV detection, it appears to fail for somatic SV detection. Thus there is a need for

293 continued development of new, more complex ways to integrate multiple somatic SV detection

294 methods [26].

295 Given the paucity of gold-standard benchmarking data for somatic SVs, the creation of the

296 simulated datasets and the associated leaderboards constitutes a major contribution of this

297 Challenge. There are distinct advantages to benchmarking on simulated datasets. It is

298 dramatically easier to simulate large numbers of tumours, or to create tumours with highly

299 divergent mutational properties, leading to well-supported estimates of per-tumour caller

300 accuracy. This enables our strategy of generating synthetic tumours of increasing complexity by

301 facilitating assessment of the impact of the complexity introduced at each step. Specifically, it is

302 possible to identify strengths and weaknesses of an individual caller by comparing its tumour-

303 specific error profiles. Moreover, synthetic tumours can be designed to test the limits of callers.

304 These advantages highlight the usefulness of synthetic datasets for benchmarking callers, and

305 until synthetic datasets are completely realistic, they will serve as important complements to real

306 datasets. While 15 teams participated in the actual competitive phase of the Challenge, 8 teams

307 have exploited the IS1-3 benchmarking resources since the competition, making 73

308 submissions to benchmark their methods for pipeline evaluation and development. We hope

309 that journals will begin to expect benchmarking on these standard datasets, as well as those

310 being generated by the final phases of the ICGC-TCGA DREAM Somatic Mutation Calling

311 Challenge, as a standard part of manuscripts reporting new somatic SV detection algorithms.

312

# Conclusions

313

314 Analysis of the error profiles of the Challenge submissions showed that somatic SV calling is a

315 distinctly harder problem than somatic SNV calling, with individual pipelines having complex and

316 unique error profiles. Parameterization was a critical factor in determining the performance of

317 teams. Finally, we demonstrate that, unlike almost every past DREAM Challenge, somatic SV

318 prediction does *not* benefit from the "wisdom of the crowds" -- simple voting of multiple

319 prediction pipelines does not yield improved predictions in this instance. The synthetic tumours

320 and somatic SV detection leaderboards remain available as a community benchmarking

321 resource.

322

# Methods

323

### Simulation of SVs by BAMSurgeon

324

325 SV support in BAMSurgeon has evolved throughout the Challenge, largely as a result of

326 constructive feedback from participants. Our descriptions of BAMSurgeon's method for

327 simulating SVs is current as of commit (*i.e.*, version) b851573474 of the code available at [27].

328    As input, BAMSurgeon (addsv.py) requires an indexed reference genome, a pre-existing BAM

329    file (Additional file 1: Figure S1a), and a list of intervals (Additional file 1: Figure S1b) along with

330    the SV type and parameters (see manual [28]). The intervals should be wide enough that local

331    sequence assembly is successful in generating a contig that spans at least 2x the expected

332    library size in the input BAM file. Intervals for which a sufficiently long contig cannot be

333    generated are rejected, where the exact definition of 'sufficiently long' is an optional parameter.

334    Intervals which contain too many discordant read pairs are also rejected, subject to a

335    parameter. Following local assembly, the contig is re-arranged: the specific rearrangements for

336    each supported SV type are illustrated in Fig. 1a (step iii) and Additional file 1: Figure S1c,e,g.

337    The assembled contig is then re-aligned to the target interval in the reference genome

338    (exonerate --bestn 1 -m ungapped) and is trimmed based on the start and end coordinates of

339    this alignment. Read pairs corresponding to trimmed contig sequence are removed from further

340    consideration.

341    Read coverage is generated over the rearranged contig using a read simulator (wgsim -e 0 -R 0

342    -r 0), to achieve the same average depth as the input BAM file, which has the effect of creating

343    split reads relative to the reference genome supporting SV detection. For a deletion, the number

344    of reads required to achieve (*e.g.*) 30x coverage is fewer than the number of reads required to

345    reach 30x coverage prior to the deletion, so reads must be removed from the original BAM (Fig.

346    1a, step iv). Inversely, for duplications and insertions additional reads need to be added to the

347    original BAM (Additional file 1: Figure S1d,h). Inversions generally do not affect coverage

348    (Additional file 1: Figure S1f). To ensure any reads removed actually correspond to the deleted

349    region of the contig, the locations of reads in the assembled contig are tracked. The number of

350    reads to be replaced, added, or deleted is scaled with the desired allele fraction. Finally, any

351    read pairs in the original BAM corresponding to reads altered in the simulated SV are replaced,

352    any read pairs marked for deletion are removed from the original BAM, and any additional read

353     pairs generated are added. It is recommended that the resulting altered BAM be post-processed

354     (with postprocess.py) to ensure compliance with the SAM format specification (see manual

355     [28]).

## Synthetic tumour generation

357     Synthetic tumours were prepared by partitioning high-coverage BAMs from 'normal' cell lines

358     into two groups of reads, picking read pairs at random as described previously [10]. For the

359     three *in silico* challenges, non-overlapping regions were selected at random for SV addition,

360     resulting in 371 variants added for IS1, 655 for IS2, and 2,886 for IS3 (Fig. 1d). Variant input

361     files are available in Additional file 7. SVs were added using addsv.py with assembly

362     GRCh37/hg19 as the reference genome and default parameters except where noted. For IS3,

363     to simulate subclones a file specifying CNV fractions over SV regions was input via option -c to

364     specify the variant allele frequency (VAF) of the spiked-in variants at either 0.5, 0.33, or 0.2

365     (Additional file 7). The output BAMs were post-processed to account for any inconsistencies

366     introduced due to remapping and merging of reads supporting SVs using the script

367     postprocess.py included with BAMSurgeon. The BAMs were further adjusted with

368     RealignerTargetCreator and IndelRealigner from the Genome Analysis Toolkit (v.2.4.9). All

369     tumour-normal pairs generated via BAMSurgeon are verified for adherence to the SAM/BAM

370     format specification using the ValidateSamFile tool included in the Picard tool set [29]. Truth

371     VCF files, *i.e.* files specifying simulated mutations, for SVs were generated using the script

372     etc/makevcf_sv.py and merged with truth files for SNP and INDEL locations, where applicable.

373     SAMtools was used throughout to split, merge, sort, and index BAMs, and also index FASTA

374     files. Details on the specific BAMSurgeon commits used for generating each tumour, as well as

375     other tumour details are given at [30].

## Validation of BAMSurgeon

376

377 To validate BAMSurgeon's ability to simulate somatic SVs, we compared the output of four

378 algorithms -- two widely used SV callers, CREST [16] and Delly [9], and two callers developed

379 over the course of the Challenge, Manta [17] and novoBreak [18] -- on the IS1 tumour-normal

380 data set, and analogous datasets generated with the same spike-in set of mutations, but with an

381 alternate aligner (NovoAlign v.3.00.05 [31]), cell line (HCC1954 BL) or read division. We did not

382 optimize parameters for the callers since the goal of this validation was not to identify the best

383 caller, but instead to verify that the caller ranking is maintained across analogous datasets.

384 Each tumour-normal pair was processed by CREST (v1.0) to extract soft clipping positions for

385 each chromosome separately, using default parameters. This data was then used by CREST to

386 call somatic SVs using the default protocol, and we converted the output into VCF v4.1 format.

387 Somatic SVs were called from each tumour-normal pair using Delly (v0.5.5) with default

388 parameters. Calling was performed on each chromosome separately for all supported SV types

389 except for translocations, and we converted the translocation output into VCFv4.1 format. Calls

390 with MAPQ < 20, PE < 5, or labeled as "LowQual" or "IMPRECISE" were filtered out. Somatic

391 SVs were called from each tumour-normal pair using Manta (v0.26.3) with the following

392 parameters: -m local -j 4 -g 10. Lastly, somatics SVs were called from each dataset using

393 novoBreak (v1.04) with a modification to ensure that sequence windows around breakpoints

394 never go beyond the start of the chromosome. All sets of SV calls were scored with $f$ = 100 bp

395 and $j$ > 0, callers were ranked based on $F$-score for each tumour-normal pair, and rankings were

396 compared across pairs (Fig. 1b,c and Additional file 1: Figure S3).

## Preprocessing VCF files

397

398 We preprocess VCF files to parse out the SV-relevant details (*e.g.* the END coordinate in the

399 INFO value or from the length of the REF sequence; if the END coordinate cannot be

400    determined from those values, it is set to the POS coordinate), remove SVs that did not pass

401    filters (as indicated by the FILTER values) and ensure consistent formatting between files. To

402    ensure consistent formatting in accordance with the VCFv4.1 specification [32] we:

403        1. Add row entries to ensure that each MATEID specification has a corresponding pair of

404            entries, where only a single entry is provided

405        2. Re-assign IDs and MATEIDs to ensure unambiguous references to entries

406        3. Where possible, replace SVTYPE = BND entries with entries specifying SVTYPE =

407            {CNV, DEL, DUP, INS, INV} in accordance with REF, ALT and EVENT values

408    Testing set SVs are indicated in the truth VCF file with the addition of masked genomic regions

409    specified with CHROM, POS and END values indicating the chromosome, start and end

410    coordinates, and SVTYPE = MSK. Specifically, a SV where ≥ 50% of the corresponding region

411    overlaps a masked region is allocated to the testing set; otherwise, it is in the training set.

## Structural variant scoring

413    Our scoring approaches evaluate the accuracy of a set of called SVs and requires input VCF

414    files specifying: i) called SVs, and ii) true/known SVs. Generally, a called SV that is sufficiently

415    similar to a known SV based on specific criteria (Table 1) is considered a true positive (TP);

416    otherwise, a false positive (FP). Also, a known SV that is sufficiently similar to a called SV is

417    considered a TP; otherwise, a false negative (FN). Our scoring supports two ways of quantifying

418    similarity:

419        A. **Region overlap.** The Jaccard coefficient ($j$) is computed from the lengths (in bp) of

420            intersection and union regions (Additional file 1: Figure S2a).

421    B. **Breakpoint closeness**. The distance ($\Delta$, in bp) between called and known breakpoints

422        is computed (Additional file 1: Figure S2b). If $\Delta \leq f$ (where $f$ is a flank threshold

423        parameter), a relative closeness is computed, $c' = 1 - \Delta/f$. The overall closeness ($c$) is

424        defined as the geometric mean of the $c'$ values for the start and end breakpoints. If only

425        one of the start and end breakpoints has $\Delta \leq f$, the called and known SVs are annotated

426        as partially matching.

427    Unless otherwise specified, we scored with $f = 100$ bp. If there is an ambiguous matching of

428    called SVs to known SVs by sufficient similarity, the similarity values ($j/c$) are used to identify an

429    optimal one-to-one matching. First, we restrict the matching to the best match(es) for each

430    called and known SV. If a SV has multiple best matches by similarity, we attempt to break the

431    tie by favouring SVs with the same SVTYPE, and/or test/training set membership. If the best

432    matching is still ambiguous, we then use corresponding similarity values together with the

433    Hungarian algorithm to obtain a one-to-one matching (with the clue v0.3-48 R package [33]).

434    Finally, SVs are annotated based on this matching. SVs that have sufficient similarity but are not

435    in the final matching are annotated as partially matching. Mated breakpoints are initially

436    annotated separately. If one is annotated as partially matching or as a TP, and the other is a FP,

437    the FP annotation is replaced by a partial match annotation. Subsequently, each set of mated

438    breakpoints is treated as a single SV.

439    These annotations are used to assess the performance of a SV caller in terms of precision =

440    nTP/(nTP + nFP), recall = nTP/(nTP + nFN) and $F$-score (specifically, $F_1$-score) = 2 x precision x

441    recall/(precision + recall), where nTP, nFP and nFN represent the numbers of TPs, FPs and

442    FNs, respectively. Partial matches are not counted in these computations. Unless otherwise

443    specified, the precision, recall and $F$-score values presented here were computed on the testing

444    and training sets combined. The best submission of a given team is defined as the team's

445    submission with the greatest *F*-score computed against all known SVs.

## Execution of challenge-based benchmarking

447    The SV component of the Challenge was executed concurrently with the SNV component, and

448    the procedure has been described previously [10]. It was implemented using the Dialogue for

449    Reverse Engineering Assessments and Methods (DREAM) framework. Briefly, information

450    about the Challenge was shared on its website [19], participants registered online, downloaded

451    a data set, applied their SV calling pipelines to the data set and submitted the results in

452    VCFv4.1 format. IS1-IS3 were released sequentially, each data set had its own competition

453    phase and participants could make multiple submissions for each data set. Each tumour

454    genome was divided into a training set and a testing set. During the competition phase,

455    leaderboards showed performance measures on the training set. After the competition closed,

456    leaderboards also showed performance measures on the whole genome (training + testing

457    sets), thus benchmarking the SV calling pipelines. The SV leaderboards for IS1 and IS2 were

458    pre-populated with results from BreakDancer (v1.1.2_2013_03_08 [7]) run with default

459    parameters; a reference point submission indicated labeled as "Standard" in figures and tables.

460    Due to our exploration of multiple SV scoring methods in this manuscript, the leaderboard

461    results are not completely consistent with the results presented here, but all raw and

462    leaderboard data are available.

## Overfitting artefact analysis

464    Due to the order of magnitude greater number of SVs spiked into IS3, we simulated training and

465    testing sets of different sizes by sampling from the IS3 training set. Specifically, we assessed

466    mutation totals of 100 to 1000 (by increments of 100), and training sets that were 80-95% (by

467    increments of 1%) of the total, by sampling each {mutation-total, training-set%} combination 100

468     times. For each sample, we computed $F_{train}$ - $F_{test}$ for each IS3 submission where $F_{train}$ and $F_{test}$

469     are $F$-scores computed on the simulated training and testing sets, respectively. We then

470     computed the median difference across samples to obtain a summary value for each

471     submission, and finally show the median across submissions in Additional file 1: Figure S5e.

472     ($F_{train}$ - $F_{test}$) > 0 suggests overfitting but such values are an artefact of testing set size since no

473     fitting/training was done in this analysis.

## Team variation

475     For each tumour-normal pair, we computed the percentage of variation in $F$-score, across all

476     submissions, that is accounted for by within-team variation. Specifically, we computed the

477     within-team sum of squares as a percentage of the total sum of squares.

## Definition of ensembles

479     We aggregated SV calls from $k$ submissions into an ensemble set with the following general

480     approach:

481        1. **BND filter**. Calls defined with SVTYPE = BND were excluded for simplicity.

482        2. **Compute call distances**. Pairwise distances ($d$, in bp) between remaining predictions

483            were computed (*i.e.* a breakpoint-length distance that incorporates distances between

484            breakpoints and differences in predicted SV length, Additional file 1: Figure S2c).

485            Distances were only computed between predictions from different submissions.

486        3. **Generate sets of similar calls**. A distance less than a selected threshold (100 for

487            consistency with $f$, see **Structural variant scoring**) indicated sufficiently similar calls.

488            We assessed two variations:

489               a. **Baseline**. We defined a graph such that vertices represented calls and edges

490                 connected sufficiently similar calls. We identified the connected components to

491        define the sets of similar calls. Sets with median intra-set distances > $f$ were

492        refined. Specifically, the call with the greatest median distance to the other set

493        members was iteratively removed until the median intra-set distance dropped

494        below $f$, or the set became empty.

495        b.  **Conservative**. We used the added constraint that called SVs overlap by ≥ 1 bp

496        to be treated as sufficiently similar. Sets of similar calls were constructed by

497        iterating over the sufficient similarity pairs from least to most distant. If a pair did

498        not contain a call in an existing call set, the pair was used to define a new call

499        set. Otherwise, one call was already in a set, and the other was a candidate for

500        addition to the same set via guilt-by-association. If the candidate came from a

501        submission that was not already covered by the set, and its median distance to

502        the existing set members ≤ $f$, it was added to the set. Any unprocessed pairs

503        within or between the prediction sets at that stage were excluded from

504        consideration.

505    4.  **Majority vote filter**. Sets with calls from ≤ $k/2$ submissions were excluded; each

506        remaining set covered the majority of submissions.

507    5.  **Aggregate sets to define ensemble calls**. The chromosome together with the median

508        start and end positions of each set of calls defined a single ensemble SV call.

509    An additional distinction between the baseline and conservative approaches is that the baseline

510    approach only involved the best submission from each team whereas all submissions were

511    used with the conservative approach. To investigate different ensembles of $N$ submissions for

512    the same tumour-normal pair, we first ordered the submissions by overall $F$-score, computed

513    after excluding calls with SVTYPE = BND. We then generated an ensemble call set with the top

514    $k$ submissions, for $k = 2..N$. The performance of ensembles was compared to that of the

515    individual submissions, after excluding calls with SVTYPE = BND (*e.g.* Fig. 3d).

## Error characterization

517    To characterize the errors made by a team, we assessed the team's best submission for a given

518    tumour-normal pair. We also assessed errors made by CREST, Delly and Manta when run, with

519    the same protocols described in the **Validation of BAMSurgeon** section, on all three tumour-

520    normal pairs. Characterizing FNs and FPs involved comparisons to TPs and true negatives

521    (TNs), respectively. Moreover, we characterized errors at the level of breakpoints.

522    ***Sampling true negatives.*** Given a set of submissions for the same tumour-normal pair, we

523    identified the maximum number of FPs from a single submission, $m$. We then sampled $\geq m$ TNs

524    for each submission, by sampling regions from the reference genome that satisfied these

525    criteria:

526        1.  length sampled from a log-normal distribution with mean and standard deviation equal to

527            that of the logged lengths of the known SVs

528        2.  start position is not in known gap and repeat regions

529        3.  region does not overlap with any known SVs

530        4.  region does not overlap with any SVs called in the submission

531    Some sampled regions qualified as TNs for multiple submissions. For IS2, we excluded Team

532    14's submission because it had a very large number (17,806) of FPs, and thus was

533    computationally problematic for the subsequent Random Forest analysis.

534     ***Breakpoint annotations based on scoring.*** A single breakpoint may be associated with

535     multiple (called/known) SVs, and therefore may be associated with multiple annotations

536     depending on the scoring approach used, *i.e.* > 1 of {TP, FN, FP}. To remove ambiguity, we

537     choose a single annotation for each breakpoint by prioritizing as follows: TP > FN > FP. This

538     prioritization favours good performance (*i.e.* TP has highest priority) and then recall (*i.e.* FN >

539     FP) since it appears to be a greater challenge than precision for SV calling (Fig. 2a, Additional

540     file 1: Figure S5a,b). TN breakpoints should be unambiguous due to the way in which they were

541     sampled (see above).

542     ***Genomic variables.*** For each breakpoint position, we computed 16 genomic factors, 12 of

543     which were previously described [10]. The additional genomic variables were computed as

544     follows:

545     A. **Bridging reads count.** We used samtools v0.1.19 to identify reads in the tumour BAM

546         mapped to a genomic region containing the window defined by the breakpoint position

547         +/- 1 bp. The bridging read count was defined as the number of identified reads. Note

548         that a bridging read does not necessarily have a secondary mapping for part of the read,

549         as one might expect for a split read.

550     B. **Distance to nearest germline INDEL.** Germline calls were obtained as previously

551         described [10] and INDELs were parsed out. The distance of a breakpoint to the nearest

552         germline INDEL was computed using BEDTools closest v2.18.2.

553     C. **Nucleotide complexity.** The sequence for the window defined by the breakpoint

554         position +/- 50 bp was extracted from the reference fasta file. The nucleotide complexity

555         was defined as the entropy of the sequence: $-\Sigma p_x \log_2(p_x)$ over $x \in \{A, G, C, T\}$ where $p_x$

556         is the proportion of the sequence with *x* (case-insensitive).

557    D. **Strand bias.** We used samtools v0.1.19 to identify reads in the tumour BAM mapped to

558        a genomic region containing the breakpoint position. The strand bias was defined as the

559        proportion of these reads mapped to the + strand.

560    *Univariate analysis.* To assess the relationship between each non-categorical variable and

561    prediction error rates, we computed the Pearson correlation coefficient between the variable

562    values and the proportion of teams with a FN/FP at the breakpoints, as well as the

563    corresponding *P* value. Reference and alternative allele counts, base quality, tumour and

564    normal coverages, bridging reads counts and distances to the germline SNP and INDEL were

565    logged (base 10) prior to computing correlations (zeros were replaced with -1 instead of

566    logged). For the categorical variables, trinucleotide and genomic location, the *P* value measured

567    the significance of the variable in a fitted binomial model predicting the FN/FP rate at a

568    breakpoint. A binomial model was fitted because it is a relatively simple model (and thus less

569    prone to overfitting) to test the relationship between a categorical variable and a proportion

570    variable (*i.e.* an error rate).

571    *Multivariate analysis.* Random Forests were generated as previously described [10] with a few

572    alterations. Here, a total of 16 genomic variables (Fig. 4) were used to build: i) a classifier to

573    distinguish FN and TP breakpoints, and ii) a classifier to distinguish FP and TN breakpoints. A

574    FP classifier was not generated for Team 7 with respect to IS1 since the team produced only

575    one FP, and thus there was insufficient data to generate an accurate model. Conversely, a FP

576    classifier was not generated for Team 14 with respect to IS2 since the team produced a very

577    large number of FPs (17,806) that caused a failure to converge. Computation of the directional

578    effect of variables was also as previously described [10].

579    Non-parametric tests (*i.e.* Wilcoxon and Mann-Whitney tests) were used throughout to avoid

580    assumptions about the distributions of the tested populations; all tested populations had $n \geq 7$.

581  The BEDTools suite (v2.18.2 [34]) was used with the bedR R package (v0.5.3 [35]) throughout.

582  Plots were generated with the BPG (v5.3.9), lattice (v0.20-33) and latticeExtra (v0.6-26) R

583  packages and R (v3.2.1) was used throughout.

584

# Declarations

585

## Availability of data and materials

586

587  Sequences files are available at the Sequence Read Archive (SRA) under accession number

588  SRP042948. BAMSurgeon [commit (*i.e.*, version) b851573474] is available at [27]. Submission

589  and known mutation (*i.e.* ground truth) VCF files are available from the Challenge website [19]

590  following registration and subsequent login.

## Acknowledgments

591

## Funding

## Author's contributions

623   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.,

624   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

625   DREAM Somatic Mutation Calling Challenge. A.Y.L., A.D.E., Y.H., K.E.H., S.M.G.E., V.H., K.D.,

626   Z.C., C.C., and T.N.Y. created datasets and analyzed sequencing data. A.Y.L., Y.H., K.E.H, and

627    P.C.B. were responsible for statistical modelling. Research was supervised by K.C., S.H.F.,

628    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

629    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

630    authors.

631    **Ethics approval and consent to participate**

632    Not applicable.

633    **Consent for publication**

634    Not applicable.

635    **Competing interests**

636    All authors declare that they have no competing interests.

637    **Additional files**

638    Additional file 1: Figures S1-S26. (PDF)

639    Additional file 2: Table S1. Challenge participation. (XLS)

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

641    (XLS)

642    Additional file 4: Descriptions of pipelines used to generate submissions. (PDF)

643    Additional file 5: Table S3. Invalid SV types. (XLS)

644    Additional file 6: Table S4. Univariate error analysis. (XLS)

645    Additional file 7: BAMSurgeon input files used to generate the three *in silico* tumour-normal

646    pairs (IS1-IS3). (TAR.GZ)

647

648 **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.e.* breakpoints are within $f$ bp of one another where $f$ is a flank parameter |
| Strengths | <ul><li>identifies genomic regions affected by the known SVs</li></ul> | <ul><li>suited to all types of SVs</li><li>evaluates precision of breakpoint predictions, facilitating subsequent breakpoint validation</li></ul> |
| Weaknesses | <ul><li>some SVs are not accurately defined by genomic regions, e.g. an insertion may be characterized by a single breakpoint</li><li>need criteria to define sufficient overlap</li></ul> | <ul><li>need criteria to define sufficient closeness</li></ul> |

649


650

# 651 Figure Legends

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

653 Method for adding SVs to existing BAM alignments. **a** Overview of SV (*e.g.* deletion) spike-in:

654 Starting with an original BAM (i), a region (ii) is selected where a deletion is desired. iii) Contigs

655 are assembled from reads in the selected region, and the contig is rearranged by deleting the

656 middle. The amount of contig deleted is a user-definable parameter. Read coverage is

657 generated over the contig using wgsim to match the number of reads per base in the original

658 BAM. Since the deletion contig is shorter than the original, fewer reads will be required to

659 achieve the equivalent coverage. iv) Generated read pairs include discordant pairs (*i.e.* paired

660 reads that do not align to the reference genome with the expected relative orientation and inner

661 distance) spanning the deletion and clipped reads (*i.e.* reads that are only partially aligned to the

662 reference). Reads mapping to the deleted region of the contig are not included in the final BAM.

663 **b,c** To test the robustness of BAMSurgeon with respect to changes in (**b**) aligner and (**c**) cell

664 line, we compared the ranks of CREST, Delly, Manta and novoBreak on two new tumour-normal

665 data sets: one with an alternative aligner, NovoAlign, and the other on an alternative cell line,

666 HCC1954 BL. Callers were scored with $f$ = 100 bp (Additional file 1: Figure S2b); Manta retained

667 the top position, independent of aligner and cell line. **d** Summary of the three *in silico* (IS)

668 tumours described here. Abbreviations: DEL, deletion; DUP, duplication; INV, inversion; INS,

669 insertion.

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

671 **a** Precision-recall plot of IS1 submissions. Each point represents a submission, each colour

672 represent a team and the best submission from each team (top *F*-score) is circled. The

673 "Standard" point corresponds to the reference point submission provided by Challenge

674     organizers. **b** The *F*-scores of submissions on the training and testing sets are highly correlated

675     for IS1 (Spearman's $\rho$ = 0.98), falling near the plotted *y* = *x* line.

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

677     **a** Recall, precision and *F*-score of all IS1 submissions plotted by team, then submission order.

678     Teams were ranked by the *F*-score of their best submission, colour coding (top bar) as in Fig. 2.

679     The "Standard'" lines correspond to the reference point submission provided by Challenge

680     organizers. **b** For each tumour, the improvement in *F*-score from the initial ("naive") to the best

681     ("optimized") submissions of each team. Darker shades of blue indicate greater improvement. **c**

682     For each tumour, team rankings based on their naive or optimized submissions. Larger dot

683     sizes indicate better ranks by *F*-score. **b,c** An "X" indicates that the team did not make a

684     submission for the specific tumour (or changed team name). **d** Recall, precision and *F*-score of

685     ensembles versus individual submissions for IS1. At the *k*th rank, the triangles indicate

686     performance of the ensemble of the top *k* submissions, and the circles indicate performance of

687     the *k*th ranked submission. The ensemble analysis focused on the best submission from each

688     team.

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

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

691     (**a,c,e,g,i**) and FP (**b,d,f,h,j**) breakpoints. Each panel shows variable importance on the left,

692     where each row represents the best performing set of predictions by the given team/caller (on

693     the given *in silico* tumour), and each column represents the indicated variable. Dot size reflects

694     variable importance, *i.e.* the mean change in accuracy caused by removing the variable from

695     the model (generated to predict erroneous breakpoints). Colour reflects the directional effect of

696     each variable (red and blue for greater and lower variable values, respectively, associated with

697     erroneous breakpoints; black for categorical variables or insignificant directional associations,

698     two-sided Mann-Whitney $P > 0.01$). Background shading indicates the accuracy of the model

699     (see colour bar). Variable importance for FN and FP breakpoints in each of the three tumours is

700     shown for the following SV callers: CREST (**a,b**), Delly (**c,d**) and Manta (**e,f**). Manta only called

701     two FPs in IS1; thus, variable importance for FP breakpoints could not be computed (indicated

702     by Xs in the plot). Variable importance for FN and FP breakpoints in IS2 (**g,h**) and IS3 (**i,j**) is

703     shown for each team. In the right plot (**g-j**), the first four columns indicate usage of the indicated

704     algorithmic approaches by each team, and the last column indicates the aligner used. Grey

705     indicates that algorithmic approaches and aligner are unknown for the given team.

706     Abbreviations: Algm, algorithm; SNP, single-nucleotide polymorphism; INDEL, short insertion or

707     deletion.

708

# 709 References

710   1.     Northcott PA, Lee C, Zichner T, Stütz AM, Erkek S, Kawauchi D, et al. Enhancer

711         hijacking activates GFI1 family oncogenes in medulloblastoma. Nature. 2014;511:428–

712         34.

713   2.     Taub R, Kirsch I, Morton C, Lenoir G, Swan D, Tronick S, et al. Translocation of the c-

714         myc gene into the immunoglobulin heavy chain locus in human Burkitt lymphoma and

715         murine plasmacytoma cells. Proc. Natl. Acad. Sci. U. S. A. 1982;79:7837–41.

716   3.     Huang M, Ye Y, Chen S, Chai J, Lu J, Zhoa L, et al. Use of all-trans retinoic acid in the

717         treatment of acute promyelocytic leukemia. Blood. 1988;72.

718   4.     Lalonde E, Ishkanian AS, Sykes J, Fraser M, Ross-Adams H, Erho N, et al. Tumour

719         genomic and microenvironmental heterogeneity for integrated prediction of 5-year

720         biochemical recurrence of prostate cancer: a retrospective cohort study. Lancet. Oncol.
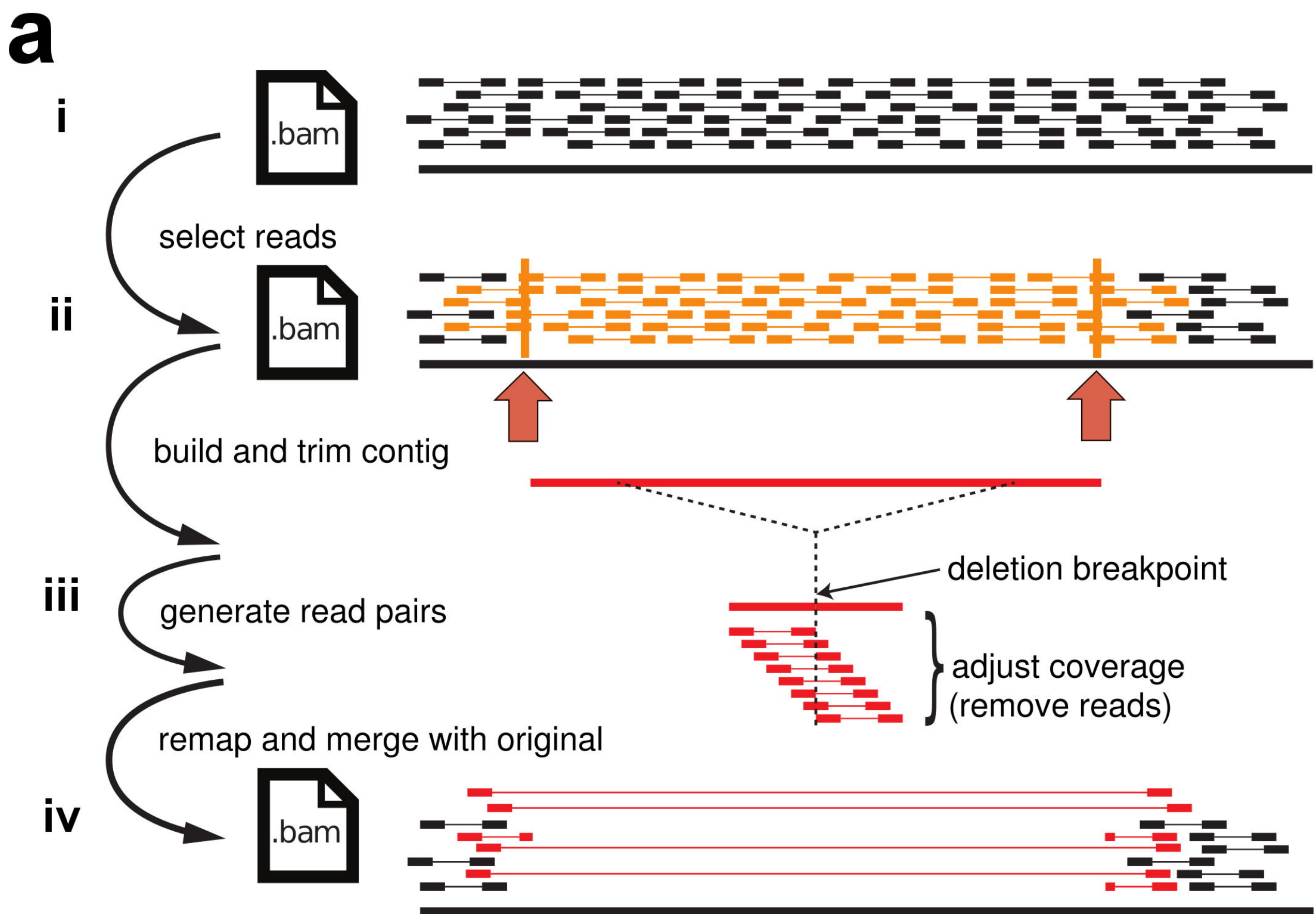
721         2014;15:1521–32.

722  5.   Vollan HKM, Rueda OM, Chin S-F, Curtis C, Turashvili G, Shah S, et al. A tumor DNA

723       complex aberration index is an independent predictor of survival in breast and ovarian

724       cancer. Mol. Oncol. 2015;9:115–27.

725  6.   Medvedev P, Stanciu M, Brudno M. Computational methods for discovering structural

726       variation with next-generation sequencing. Nat. Methods. 2009;6:S13–20.

727  7.   Chen K, Wallis JW, McLellan MD, Larson DE, Kalicki JM, Pohl CS, et al. BreakDancer:

728       an algorithm for high-resolution mapping of genomic structural variation. Nat. Methods.

729       2009;6:677–81.

730  8.   Hormozdiari F, Hajirasouliha I, Dao P, Hach F, Yorukoglu D, Alkan C, et al. Next-

731       generation VariationHunter: combinatorial algorithms for transposon insertion discovery.

732       Bioinformatics. 2010;26:i350-7.

733  9.   Rausch T, Zichner T, Schlattl A, Stütz AM, Benes V, Korbel JO. DELLY: structural variant

734       discovery by integrated paired-end and split-read analysis. Bioinformatics. 2012;28:i333–

735       9.

736  10.  Ewing AD, Houlahan KE, Hu Y, Ellrott K, Caloian C, Yamaguchi TN, et al. Combining

737       tumor genome simulation with crowdsourcing to benchmark somatic single-nucleotide-

738       variant detection. Nat. Methods. 2015;12:623–30.

739  11.  Zerbino DR, Birney E. Velvet: algorithms for de novo short read assembly using de Bruijn

740       graphs. Genome Res. 2008;18:821–9.

741  12.  Slater GSC, Birney E. Automated generation of heuristics for biological sequence

742       comparison. BMC Bioinformatics. 2005;6:31.

743  13.  Pop M, Phillippy A, Delcher AL, Salzberg SL. Comparative genome assembly. Brief.

744       Bioinform. 2004;5:237–48.

745  14.  Zerbino DR, McEwen GK, Margulies EH, Birney E. Pebble and rock band: heuristic

746       resolution of repeats and scaffolding in the velvet short-read de novo assembler.

747         Salzberg SL, editor. PLoS One. 2009;4:e8407.

748   15.   GitHub Code Repository: wgsim. https://github.com/lh3/wgsim. Accessed 22 November

749         2017.

750   16.   Wang J, Mullighan CG, Easton J, Roberts S, Heatley SL, Ma J, et al. CREST maps

751         somatic structural variation in cancer genomes with base-pair resolution. Nat. Methods.

752         2011;8:652–4.

753   17.   Chen X, Schulz-Trieglaff O, Shaw R, Barnes B, Schlesinger F, Källberg M, et al. Manta:

754         rapid detection of structural variants and indels for germline and cancer sequencing

755         applications. Bioinformatics. 2016;32:1220–2.

756   18.   Chong Z, Ruan J, Gao M, Zhou W, Chen T, Fan X, et al. novoBreak: local assembly for

757         breakpoint detection in cancer genomes. Nat. Methods. 2017;14:65–7.

758   19.   ICGC-TCGA DREAM Mutation Calling challenge.

759         https://www.synapse.org/#!Synapse:syn312572/wiki/58893. Accessed 22 November

760         2017.

761   20.   Margolin AA, Bilal E, Huang E, Norman TC, Ottestad L, Mecham BH, et al. Systematic

762         analysis of challenge-driven improvements in molecular prognostic models for breast

763         cancer. Sci. Transl. Med. 2013;5:181re1.

764   21.   Marbach D, Costello JC, Küffner R, Vega NM, Prill RJ, Camacho DM, et al. Wisdom of

765         crowds for robust gene network inference. Nat. Methods. 2012;9:796–804.

766   22.   Strobl C, Boulesteix A-L, Zeileis A, Hothorn T. Bias in random forest variable importance

767         measures: illustrations, sources and a solution. BMC Bioinformatics. 2007;8:25.

768   23.   Tattini L, D'Aurizio R, Magi A. Detection of Genomic Structural Variants from Next-

769         Generation Sequencing Data. Front. Bioeng. Biotechnol. 2015;3:92.

770   24.   Boutros PC, Margolin AA, Stuart JM, Califano A, Stolovitzky G. Toward better

771         benchmarking: challenge-based methods assessment in cancer genomics. Genome Biol.

772      2014;15:462.

773   25.   Meyer P, Alexopoulos LG, Bonk T, Califano A, Cho CR, de la Fuente A, et al. Verification

774         of systems biology research in the age of collaborative competition. Nat. Biotechnol.

775         2011;29:811–5.

776   26.   Mohiyuddin M, Mu JC, Li J, Bani Asadi N, Gerstein MB, Abyzov A, et al. MetaSV: an

777         accurate and integrative structural-variant caller for next generation sequencing.

778         Bioinformatics. 2015;31:2741–4.

779   27.   GitHub Code Repository: BAMSurgeon. https://github.com/adamewing/bamsurgeon.

780         Accessed 22 November 2017.

781   28.   BAMSurgeon Manual.

782         https://github.com/adamewing/bamsurgeon/blob/master/doc/Manual.pdf. Accessed 22

783         November 2017.

784   29.   Picard Tools - By Broad Institute. http://broadinstitute.github.io/picard/. Accessed 22

785         November 2017.

786   30.   ICGC-TCGA DREAM Mutation Calling challenge: Synthetic Tumours.

787         https://www.synapse.org/#!Synapse:syn312572/wiki/62018. Accessed 22 November

788         2017.

789   31.   Novocraft. http://www.novocraft.com/. Accessed 22 November 2017.

790   32.   The Variant Call Format (VCF) Version 4.1 Specification. https://samtools.github.io/hts-

791         specs/VCFv4.1.pdf. Accessed 22 November 2017.

792   33.   Kuhn HW. The Hungarian method for the assignment problem. Nav. Res. Logist. Q.

793         Wiley Subscription Services, Inc., A Wiley Company; 1955;2:83–97.

794   34.   Quinlan AR, Hall IM. BEDTools: a flexible suite of utilities for comparing genomic

795         features. Bioinformatics. 2010;26:841–2.

796   35.   Haider S, Waggott D, Lalonde E, Fung C, Liu F-F, Boutros PC. A bedr way of genomic

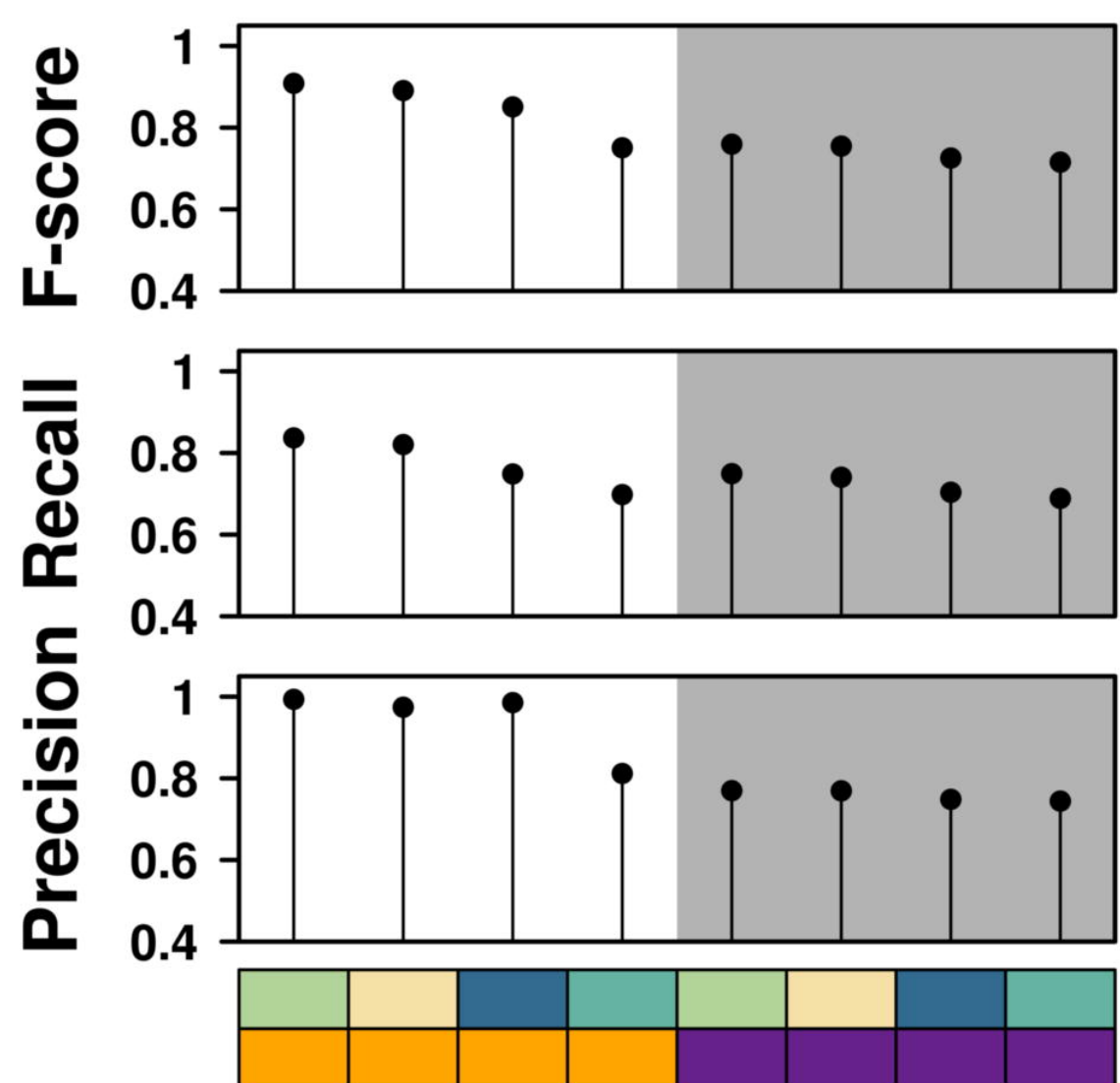797          interval processing. Source Code Biol. Med. 2016;11:14.
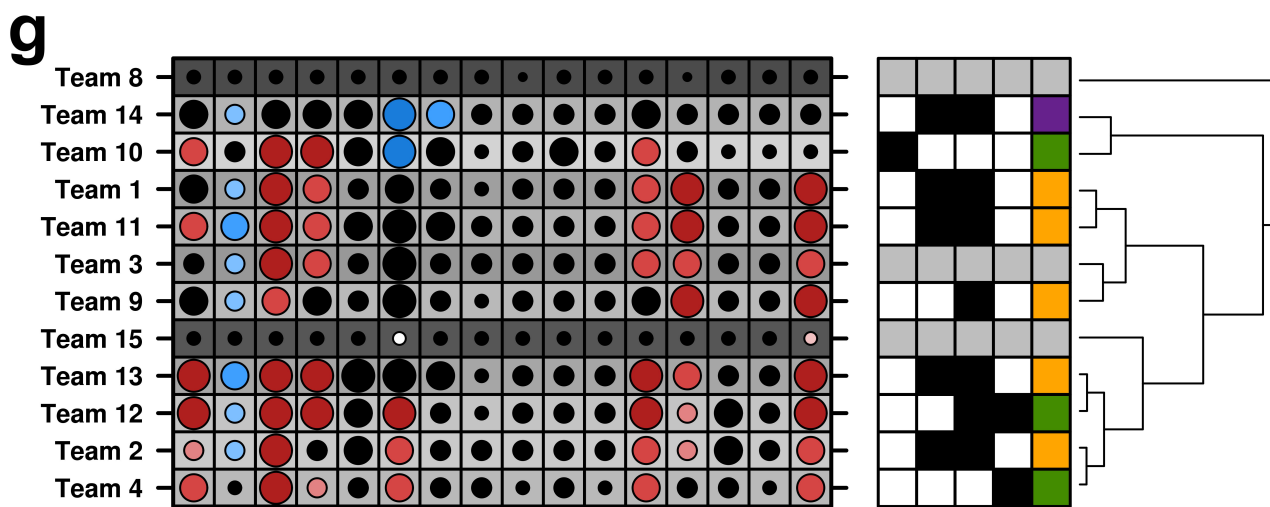
798

**d**

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

**False Negatives**

**False Positives**

**Importance**
- > 0.05
- (0.01, 0.05]
- (0.005, 0.01]
- (0.001, 0.005]
- (0, 0.001]
- 0
- [-0.001, 0)
- [-0.005, -0.001)
- [-0.01, -0.005)
- [-0.05, -0.01)
- < -0.05

**Algm Feature**
- Yes
- No

**Aligner**
- BWA-backtrack
- BWA-MEM
- Subread

Accuracy of model

Feature labels (columns): Reference allele count, Alternate allele count, Base quality, Tumour coverage, Normal coverage, Mapping quality, Read position, Trinucleotide, Homopolymer rate, GC content, Genomic location, Bridging reads count, Distance to germline INDEL, Distance to germline SNP, Nucleotide complexity, Strand bias

Algm feature labels: Read depth analysis, Paired end mapping, Split read alignment, De novo assembly, Aligner