1  **Modelling double strand break susceptibility to interrogate structural**

2  **variation in cancer**

3

4  **Authors**: Tracy J. Ballinger[1], Britta Bouwman[2], Reza Mirzazadeh[2], Silvano

5  Garnerone[2], Nicola Crosetto[2]*, Colin A. Semple[1]*

6  **Affiliations**:

7  1. MRC Human Genetics Unit, MRC Institute of Genetics and Molecular Medicine,

8  University of Edinburgh, Crewe Road, Edinburgh, EH4 2XU, UK

9  2. Science for Life Laboratory, Department of Medical Biochemistry and

10  Biophysics, Karolinska Institutet, Stockholm, Sweden

11  *These authors contributed equally to this work

12  **Corresponding Author**:

13  Dr Tracy Ballinger, MRC Human Genetics Unit, Institute of Genetics and

14  Molecular Medicine, University of Edinburgh, Crewe Road, Edinburgh, EH4 2XU,

15  UK; Tracy.Ballinger@igmm.ed.ac.uk

16  **Author Emails:**

17  Tracy Ballinger: tracy.ballinger@igmm.ed.ac.uk

18  Britta Bouwman: britta.bouwman@scilifelab.se

19  Reza Mirzazadeh: reza.mirzazadeh@scilifelab.se

20    Silvano Garnerone: silvano.garnerone@scilifelab.se

21    Nicola Crosetto: nicola.crosetto@scilifelab.se

22    Colin Semple: colin.semple@igmm.ed.ac.uk

23    **Running head**:

24    Modelling DNA double strand breaks

25

26    **Abstract**

27    **Background:** Structural variants (SVs) are known to play important roles in a

28    variety of cancers, but their origins and functional consequences are still poorly

29    understood. Many SVs are thought to emerge via errors in the repair processes

30    following DNA double strand breaks (DSBs) and previous studies have

31    experimentally measured DSB frequencies across the genome in cell lines.

32    **Results:** Using these data we derive the first quantitative genome-wide models

33    of DSB susceptibility, based upon underlying chromatin and sequence features.

34    These models are accurate and provide novel insights into the mutational

35    mechanisms generating DSBs. Models trained in one cell type can be successfully

36    applied to others, but a substantial proportion of DSBs appear to reflect cell type

37    specific processes. Using model predictions as a proxy for susceptibility to DSBs

38    in tumours, many SV enriched regions appear to be poorly explained by

39    selectively neutral mutational bias alone. A substantial number of these regions

40    show unexpectedly high SV breakpoint frequencies given their predicted

41    susceptibility to mutation, and are therefore credible targets of positive selection

42    in tumours. These putatively positively selected SV hotspots are enriched for

43  genes previously shown to be oncogenic. In contrast, several hundred regions

44  across the genome show unexpectedly low levels of SVs, given their relatively

45  high susceptibility to mutation. These novel 'coldspot' regions appear to be

46  subject to purifying selection in tumours and are enriched for active promoters

47  and enhancers.

48  **Conclusions:** We conclude that models of DSB susceptibility offer a rigorous

49  approach to the inference of SVs putatively subject to selection in tumours.

50

51  **Keywords**: Double strand break, cancer, structural variaton, chromatin,

52  modelling

53

54  **Background**

55

56  Structural variation (SV) in tumour genomes is known to play important roles in

57  disease progression and may be critical in driving the development of certain

58  cancer types (1–3). However, challenges remain not only in ascertaining accurate

59  SV calls, as evidenced by the compendium of SV calling algorithms used in many

60  projects (4–6), but also in predicting their functional impact. Some SVs have

61  apparently direct consequences; for example, amplification of oncogenes leading

62  to overexpression, deletion of tumor suppressors leading to dysfunction, and

63  translocations generating oncogenic fusion proteins (4). Reportedly indirect

64  consequences of SVs include changes in enhancer targeting, affecting the

65  expression of nearby genes, or "enhancer hijacking" (7). However, it remains

66     challenging to distinguish the influences of evolutionary selection versus

67     primary mutation rate in generating the SVs concerned.

68

69     A recent study of whole genome sequencing (WGS) data from breast tumours

70     identified SV hotspots and putative driver SVs, but could not discern the relative

71     contributions of mutational bias and selection underlying these hotspots (8).

72     Resolving the influences of mutational bias versus selective forces has become

73     critical given that both single nucleotide variant (SNV) and SV mutation rates

74     vary widely across the genome, in parallel with replication timing and chromatin

75     structure (9,10). In analyses of tumour SNVs, variants are routinely prioritized

76     based on algorithms including corrections for estimates of SNV mutation rate

77     variation (11), but analogous methods are not yet applied to SVs.

78

79     Variable rates of SVs observed across the genome are likely to be affected by

80     differences in the efficiency of repair of DNA double strand breaks (DSBs). DSBs

81     can be repaired by homologous recombination (HR) at the G2 and S stages of the

82     cell cycle and, more commonly, by canonical non-homologous end joining (c-

83     NHEJ) which operates throughout the cell cycle (12). The c-NHEJ process is error

84     prone and has been shown to create structural variants initiating carcinogenesis

85     (13). A third repair process, alternative NHEJ (alt-NHEJ) uses microhomology to

86     mediate repairs when the c-NHEJ pathway is unavailable, and repair by alt-NHEJ

87     appears to increase the rate of deletions, insertions, and translocations further

88     (14). The efficiency of these repair processes is often dependent upon the

89     chromatin features and nuclear organization present where the damage occurs.

90     For example, the histone modification H3K36me3, associated with active

91 transcription, recruits the HR pathway, while H4K20me1, a mark of highly

92 transcribed genes, recruits components of the NHEJ pathway (15). The

93 associations between DSB repair and the underlying chromatin landscape may

94 therefore explain the observed correlations between tumour SV rates and

95 chromatin structure (9).

96

97 Previous studies have also shown DSB formation to be influenced by underlying

98 chromatin structures and genomic sequences. It has long been known that

99 certain cytogenetically mapped loci, termed "fragile sites" undergo recurrent

100 DSBs in cells under replicative stress and in cancer (16). More recent high

101 throughput sequencing (HTS) based approaches have been developed to profile

102 DSB rates more precisely within *in vitro* populations of cells (17–25). Three of

103 these methods, BLESS (18), DSBCapture (22), and BLISS (25) are closely related

104 and have been used to generate high-resolution maps of endogenous DSBs

105 occurring in human cell lines, resulting in continuous data reflecting the

106 propensities for DSBs across all chromosomes. These studies have suggested that

107 DSBs may preferentially occur within nucleosome-depleted regions, are

108 correlated with active promoter and enhancer histone modifications, and may

109 associate with G-quadruplex sites (22,26). Certain studies have also suggested

110 DSBs to be depleted in some transposon classes and enriched in some simple

111 repeat classes, and to be unusually frequent in long, late-replicating genes

112 (18,24). Overall, previous studies have found correlations and enrichments

113 between DSBs and various inter-correlated chromatin and genomic features,

114 making it difficult to accurately assess the contribution of any particular feature

115 to DSB susceptibility. Understanding such contributions can be valuable for

5

116 understanding the underlying mutational and repair mechanisms. In addition, a

117 fuller understanding of the relative contributions of many features to DSB

118 formation can allow reliable predictions of the expected DSB frequency in a given

119 genomic region.

120

121 Random forests have been used to model a variety of biological phenomena

122 because they perform well in the presence of inter-correlated input variables

123 showing non-linear relationships. For example, they have been used to predict

124 nuclear compartments (27), cancer SNV mutational landscapes (28), and

125 enhancer-promoter interactions (29). In this study we construct random forest

126 regression models to generate quantitative measures of the relative importance

127 of a variety of matched chromatin and other features to DSB susceptibility. We

128 use multiple, high-resolution DSB profiling datasets to compare modeling

129 accuracy across several platforms and cell types. The cell types selected have

130 also been extensively profiled for a variety of chromatin features by the ENCODE

131 Project (30) and others, allowing well-matched models to be constructed for all

132 datasets. We demonstrate that these models provide accurate estimates for the

133 expected rate of DSBs in a given region and can be cross applied between DSB

134 datasets. In addition the models can be used to explore tumour SV breakpoint

135 data, to nominate novel regions putatively subject to selection in cancer.

136

137 **Results**

138

139 We uniformly processed four DSB datasets from three related platforms

140 (DSBCapture and BLISS are both based upon modifications to the BLESS

6

141  protocol) and covering three different cell types, collating matched chromatin

142  data for each. These datasets include two novel DSB mapping datasets derived

143  from the K562 erythroleukemia and MCF7 breast cancer cell lines using the

144  recently developed BLISS method (25) (see Methods) and two previously

145  published DSB mapping datasets derived from the NHEK keratinocyte cell line

146  using BLESS and DSBCapture (22) protocols. DSB frequency is defined in each

147  dataset as the number of unique reads mapping to a given 50kb region, since

148  each read in a DSBCapture, BLESS, or BLISS experiment represents an exposed

149  DNA DSB end. Replicate experiments within each dataset were strongly and

150  significantly correlated (Pearson's r = 0.905 to 0.992, p<2.2e-16) and were

151  combined to reduce noise, although random forest models generated from any

152  single one of the replicates yielded very similar results (see Methods).

153  Comparisons among DSB profiling datasets showed moderate correlations in

154  genome-wide DSB frequency between the three cell types as expected (r = 0.351

155  to 0.635, p<2.2e-16), shown in Supp Figure 1. All three cell types correspond to

156  well-characterized ENCODE cell lines, providing numerous matched chromatin

157  and genomic features exhibiting a range of correlations to DSB (Figure 1), and

158  are also inter-correlated themselves (Supp Figure 2).

159

160  **Accurate models of genome-wide DSB frequency across cell types**

161

162  We modeled DSB frequency at 50kb resolution, using the same ten matched

163  genomic features from each cell type to construct random forest models (see

164  Methods): open chromatin assayed by DNase-seq, POL2B binding, CTCF binding

165  and five histone modifications assayed by ChIP-seq, replication timing assayed

7

166    by Repli-seq, and RNA-seq. We also included G-quadruplex forming regions as an

167    additional feature, since these DNA secondary structures are associated with

168    genomic instability (31). We found strong and significant correlations between

169    predicted and observed DSB frequency for all four datasets, with Pearson's

170    coefficients ranging from 0.83 to 0.92 (Figure 2). We also generated a model for

171    the NHEK DSBCapture dataset using an extended set of 21 features, including

172    additional    histone    modifications,    histone    variants,    and    nuclear

173    compartmentalization from Hi-C data (32). This extended model resulted in

174    better predictive results for a small fraction of the genome (Supp Figure 4, Box

175    B), and a modestly increased genome-wide Pearson's coefficient between

176    predicted and observed values (11 feature model r = 0.918; 21 feature model r =

177    0.922). We conclude that models constructed using the 11 selected genomic

178    features (Figure 2) provide high predictive accuracy across cell types, with

179    additional features likely to provide only marginal gains.

180

181    Variable importance metrics for these models reveal consistent trends in the

182    most influential features in DSB frequency prediction (Figure 2,E-H). Replication

183    timing is the most important feature across all three models with early

184    replication associated with high DSB regions and late replication with low DSB

185    (Figure 3C), in agreement with previous studies (33). In addition, the histone

186    modifications H3K36me3 and H3K9me3 (demarcating active genes and gene-

187    poor heterochromatin respectively) emerge as informative features, with

188    H3K36m3 enriched in high DSB regions and H3K9me3 in low DSB regions

189    (Figure 3C). This is consistent with observations that structural variants

190    disproportionately accumulate within the early replicating, relatively gene rich

191 regions of the genome in cancer, and are relatively depleted in late replicating

192 heterochromatin (9,10). DNase-seq open chromatin ranks second in three

193 datasets and fourth in the MCF7 model and is also the most important feature for

194 predicting DSB peaks in the study of Mourad et al. (34) in which they do not

195 include replication timing. The influence of G-quadruplex forming regions is

196 notably variable, ranking as a relatively important feature in the NHEK datasets,

197 but having little and no predictive value in the K562 and MCF7 datasets. RNA-seq

198 is not a strong predictor of DSB susceptibility although DNase-seq peaks are

199 often found at the promoter regions of active genes. This suggests that open

200 chromatin at transcriptionally active genes and associated regulatory elements

201 (reflected in DNase-seq, H3K4me3 and POL2B binding), rather than

202 transcription per se, is the dominant influence on DSB frequency. CTCF binding

203 also appears to be an informative variable, genome-wide in all models, though it

204 binds at sites constituting a very small fraction of the genome. Given the critical

205 roles of CTCF in chromatin architecture and regulation (32), there has been

206 intense interest in the causes and effects of structural variants disrupting CTCF

207 binding sites (35,36).

208

209 **Influential features underlying DSB frequency differ between genomic loci**

210 **and cell types**

211

212 Beyond the general, genome-wide trends described above, we see differences in

213 the behavior of certain classes of loci. These are evident as regions departing

214 from the linear relationship between observed and predicted DSB frequency

215 seen for the majority of the genome (Figure 3A; Supp Fig 4). Deeper exploration

9

216    of the relationships between underlying genomic features and DSB frequency

217    reveals diagnostic features for these discrepant classes. One class of loci (Figure

218    3, Box A) shows unusually low values for both predicted and observed DSB

219    frequencies, and is enriched for H3K9me3 marked heterochromatin and low

220    sequence mappability (Figure 3B). These regions are likely to correspond to

221    repeat-rich regions near centromeres and on the short arms of acrocentric

222    chromosomes, which are problematic for read mapping algorithms (37). Another

223    class of H3K9me3 heterochromatin enriched loci shows higher DSB predictions

224    than observed, in spite of high mappability values (Figure 3, Box B). This class of

225    regions is absent in DSB datasets generated by the BLISS protocol (Figure 2), so

226    these aberrant predictions may reflect technical and methodological differences

227    between datasets. In any case, it is clear that model predictions may reasonably

228    be expected to be less accurate in heterochromatic regions.

229

230    The similarities in relative variable importance across datasets (Figure 2)

231    suggest that many features have a similar influence on DSB frequency in each of

232    the three cell types. Thus, a model trained in one cell type might generalize well

233    to another cell type and allow us to generate predictive DSB frequency profiles

234    for model cell lines currently lacking high resolution DSB data. We cross-applied

235    models and found models trained in one cell type often performed well in

236    another (Figure 4). For example, a model trained in NHEK cells could be used to

237    predict DSB frequencies in K562 cells (inputting K562 genomic features) with

238    high accuracy (Pearson's r = 0.85 correlation; Figure 4). This offers a substantial

239    improvement over the base correlation (r = 0.63) between NHEK and K562

240    observed DSB profiles. We measured the correlation of observed and predicted

10

241 DSB frequencies across all nine model and feature combinations and always

242 found correlations (r = 0.58 to 0.85) that improved on the base correlations (r =

243 0.38 to 0.63) seen between the observed DSB datasets (Figure 4). These

244 improvements echo the similarities in variable importance between cell types

245 (Figure 2). The moderate correlations between DSB across cell types

246 demonstrate that a substantial proportion of DSB susceptibility across the

247 genome is cell type specific, which is consistent with the established cell type

248 specific properties of many SV breakpoint regions in tumours, such as common

249 fragile sites (38). Furthermore the larger performance gap in models for cell

250 lines with altered variable rankings indicates that DSB mechanisms may differ

251 across cell types and may not be completely captured via epigenomic features.

252

253 **Tumour SV breakpoints possess variable susceptibility to DSBs**

254

255 Keratinocytes are considered to be the cell type of origin for mucosal and

256 cutaneous carcinomas, particularly squamous cell carcinomas (39), and NHEK

257 cells are often used in the literature as a model for these cancers. Similarly, MCF7

258 cells and K562 cells have been used extensively as models for breast and blood

259 cancers respectively. This motivated us to ask how the DSB models for these

260 three cell types relate to the patterns of SV breakpoints observed in squamous

261 cell carcinomas, blood cancers, and breast tumours.

262

263 A number of large structural variant (SV) collections have been established for a

264 variety of tumour types, and each possesses advantages and shortcomings. The

265 International Cancer Genome Consortium (ICGC) provides high resolution SV

11

266   calls based upon whole genome sequencing (WGS) for 2,146 patients across 17

267   cohorts (40), but sample cellularities, sequencing depths and SV calling methods

268   vary across cancer cohorts, and are expected to affect results (Supp Figure 6).

269   The Cancer Genome Atlas (TCGA) produced consistently processed copy number

270   variant (CNV) calls from SNP chip data for 23,084 patients across 33 cohorts

271   (Supp Figure 7). However, breakpoint resolution is much lower than calls based

272   upon WGS, and copy neutral SVs such as inversions and translocations are

273   absent. We analyzed ICGC and TCGA data as pancancer datasets, combining all

274   cancer types together, but also as three cancer type subgroups. TCGA subgroups

275   comprised a squamous cell carcinoma subgroup, a blood cancers subgroup

276   including two blood cancers, and breast cancer as a separate group (see

277   Methods). Similar ICGC subgroups were formed (from cohorts independent of

278   TCGA), but with the squamous cell carcinoma subgroup replaced with a

279   carcinoma subgroup, which includes seven carcinoma cancer studies excluding

280   breast cancer (see Methods).

281

282   Analogously to the DSB datasets, we determined the number of tumour SV

283   breakpoints per 50kb region for each of the ICGC and TCGA SV datasets (see

284   methods) and compared these to the DSB predictions from our models. In ICGC

285   data overall we saw low correlations between the number of SV breakpoints and

286   DSB predictions (Supp Figure 8 and Supp Figure 9). Restricting our analysis to

287   ICGC enriched SV breakpoint regions, or ESBs for the purpose of this manuscript

288   (50kb regions with SV breakpoint counts in the top 5% genome-wide, see

289   Methods), increased the agreement with DSB model predictions. Significant

290   increases in NHEK and MCF7 model predictions were seen for pancancer,

291  carcinoma, blood, and breast tumour ESBs and in K562 model predictions for all

292  cancer subsets except blood ESBs (Figure 5). The significant increase in DSB

293  model predictions seen for carcinoma ESBs indicates that DSB susceptibility

294  (captured in the models) may shape the SV landscape of these cancer types. We

295  also see a significant increase in DSB predictions for TCGA blood cancer ESBs,

296  but not for any other subgroups in TCGA data (Supp Figure 10). However, as

297  mentioned, TCGA data is of low resolution and not suitable for accurate

298  breakpoint detection.

299

300  Certain classes of relatively simple SVs (deletions, duplications, inversions,

301  translocations) are often the product of one or two DSBs, while more complex

302  intrachromosomal rearrangements can be difficult to classify accurately, and

303  may have origins in poorly understood phenomena such as chromothripsis (41).

304  Indeed, even for simple SVs there may be some ambiguity, with an unknown

305  fraction arising by mechanisms that may not involve a DSB. For example,

306  insertions can arise from transposon activity, and duplications from replication

307  slippage (42). However, even if many SV breakpoints do not arise from DSBs, we

308  might reasonably expect to see shifts to higher median DSB model prediction

309  values for many simple SV classes. We determined ESBs as above for ICGC-

310  annotated SV classes across all ICGC tumour types to examine their DSB

311  frequency predictions, compared to non-ESBs, 50kb regions that do not attain SV

312  breakpoint counts in the top 5% with at least one tumour SV breakpoint

313  detected. Overall, the models show significant elevations for ESBs covering all SV

314  classes except insertions (Figure 5). Insertions may be less influenced by DSB

315  susceptibility because they may occur via transposable element activity rather

13

316    than through DNA damage and repair pathways. Crosetto et al. (18) find an

317    enrichment of satellite repetitive elements in regions enriched for DSB in cells

318    exposed to aphidicolin. However, regions that undergo DSB under replicative

319    stress, as induced by aphidicolin, may differ from DSB regions under normal cell

320    growth conditions.

321

322    **Interrogating tumour SV data at common fragile sites with DSB models**

323

324    The predicted DSB frequencies from our models and ICGC tumour SV breakpoint

325    frequencies differ in their scaling and distributions and are not directly

326    comparable. However, it is of interest to identify outlier regions, where model

327    predictions and observed tumour SV breakpoint rates diverge most, since these

328    regions may include loci under selection in tumours. We developed a novel

329    metric, the d-score, to measure this divergence between expectations given a

330    DSB model and observed SV breakpoint rates in tumours. In brief, this metric

331    relies on fitting known distributions to the observed SV breakpoint dataset and

332    to the predicted DSB dataset. Based upon the known distributions we then

333    transform the observed SV counts and predicted DSB values to p-values,

334    reflecting the probability that each value is drawn from the fitted distribution

335    (see Methods). For each 50kb region in the genome the difference between the

336    SV breakpoint log p-value and the predicted DSB log p-value is the d-score.

337    Regions with unexpectedly high d-scores contain more SV breakpoints than

338    expected, given our model, whereas regions with unusually low d-scores contain

339    fewer SV breakpoints than expected.

340

14

341    Common fragile sites (CFSs) have long been studied for their unusual properties

342    of generating SVs, both in normal cells and in cancer (38). These regions undergo

343    frequent DSBs in tumours and have been well studied in terms of their genomic

344    context, relationship to replication timing and origins, and correlations with

345    particular chromatin states (43). They tend to occur within large genes, in G-

346    negative chromosomal bands with high DNA flexibility, are unusually late

347    replicating (44), and it is thought that their instability derives from

348    transcription-associated replication stress (38). CFSs only exist in modest

349    numbers and are defined at low resolution (by cytogenetic bands or gene loci);

350    they therefore provide an interesting, though challenging, test set of regions to

351    examine d-score performance.

352

353    We examined predicted (NHEK model) DSB frequencies at 294 50kb regions

354    coinciding with annotated CFS gene loci across the genome, in comparison to

355    regions associated with all annotated genes, and regions associated with putative

356    cancer driver genes (Figure 6C). Although significant shifts to higher frequencies

357    are seen for the driver gene sets for predicted DSB frequencies, the CFSs do not

358    show a similar increase, most likely because the model predicts DSB in early

359    replicating regions, and CFS tend to be late-replicating. Thus, the dominant

360    features influencing DSB susceptibility genome-wide do not appear to drive the

361    elevated DSB rates at CFSs, consistent with CFS instability involving replicative

362    stress (38). However, CFS d-scores show a significant shift above the distribution

363    for all genes and above the driver gene sets as well (Figure 6D).  This result is

364    replicated in the MCF7 BLISS model examined inconjunction with ICGC breast

365    cancer SV breakpoints (Sup Figure 11).  We conclude that the d-score, a measure

366    of relative DSB enrichment, offers a robust metric for the classification of regions

367    showing unusual SV breakpoint rates in tumours.

368

369    **Identification of hot and cold spots for structural variant breakpoints in**

370    **tumours**

371

372    We have developed a classification of regions of interest within ICGC tumour

373    cohorts based upon the d-score metric. We call regions with significantly more

374    SV breakpoints than expected, or SV hotspots, cancHpredL (cancer high,

375    predicted low), and regions with fewer SV breakpoints than expected, or SV

376    coldspots, cancLpredH (cancer low, predicted high) (see Methods). Figure 6

377    depicts these classes of regions in d-score plots of ICGC SV breakpoint data. Many

378    previous studies have predicted oncogenic SV hotspots simply as regions

379    repeatedly rearranged in cancers. Here we refine such predictions by assessing

380    these raw SV breakpoint frequencies relative to the predicted susceptibility of

381    each region to breakage. It is not possible to predict coldspot regions without a

382    model of expected DSB frequency, and to our knowledge SV breakpoint coldspots

383    have not been studied before.

384

385    We also define a class of regions possessing both high predicted DSB values and

386    high SV breakpoint frequencies (cancHpredH), corresponding to regions

387    showing unusually high SV frequencies on the background of high susceptibility

388    to DSBs. Finally, we define a fourth class of regions that have predicted DSB rates

389    close to zero but high SV breakpoint frequencies (cancHpredL2). In principle,

390    these regions are a class of SV hotspots but, as shown in Figure 3B, they are likely

16

391    to be repetitive, heterochromatic, and enriched for artifacts (false positives and

392    negatives in SV breakpoint) due to their association with low mappability.

393

394    We examined a range of functional annotation enrichments in the four classes of

395    regions using circular permutation to assess significance (see Methods; Figure

396    6). The annotations included two putative cancer gene sets, 260 genes from the

397    Cancer5000 dataset (45) and 561 genes from the COSMIC collection (46)). We

398    also included a set of 15,415 super enhancers (47), common fragile sites, and

399    chromatin states from ENCODE chromHMM analysis (48). Notably, the majority

400    of genes in both cancer sets are predicted to be oncogenic based on unexpectedly

401    high and functionally significant SNV (rather than SV) loads and are not

402    necessarily expected to occupy regions with higher levels of SV breakpoints. In

403    fact, both gene sets demonstrate significant enrichments in the cancHpredL class

404    of hotspot regions (Figure 6D), although RefSeq genes do not, suggesting that

405    these genes may also frequently be altered in cancer through SV. The

406    cancHpredL regions are also significantly depleted in active chromatin regions,

407    such as promoters, enchancers, and insulator regions, most likely because these

408    types of regions do not have low predicted DSB. The high susceptibility

409    cancHpredH regions occupy gene-rich areas of the genome (enriched for known

410    RefSeq genes) including both cancer genes sets, and for active promoters, strong

411    enhancers, and insulators. This is consistent with reports that CTCF bound

412    insulator elements suffer recurrent mutations in tumours. Likewise, the

413    cancLpredH class of coldspot regions occupy gene rich neighbourhoods, active

414    promoters, and strong enhancers (Figure 6), suggesting some genes and distal

415    regulatory regions may have experienced purifying selection in tumours.

416

417     Given the discrepancies mentioned above between ICGC and TCGA experimental

418     platforms, data analysis, and sample cohorts, we do not expect strong agreement

419     between ICGC and TCGA derived SV datasets. Indeed, the correlation between

420     them is low (Spearman's rho of 0.099, p<2.2e-16), and the pancancer ESBs from

421     either set do not significantly overlap (p < 0.99, see methods). However, the

422     cancLpredH class is again enriched in active promoter and strong enhancer

423     regions, in accordance with the results based upon ICGC SV data (Sup Figure 12).

424

425     We again wanted to test the utility of DSB random forest models applied to

426     different cell types by testing the accuracy of predictions made by a model

427     trained in one cell type given features for a different cell type, as in Figure 4.

428     Instead of looking at the correlation between the observed and predicted DSB

429     scores across the genome, we examined the overlap between cancHpredL,

430     cancHpredH, and cancLpredH 50kb regions for the MCF7 model versus the

431     NHEK model, using the MCF7 model as the truth set. Subsets of 50kb regions for

432     each model were derived from MCF7 features and ICGC breast cancer SV

433     breakpoints; only the training data for the models differ. We found a significant

434     overlaps between all three categories of d-score subsets, with 595/662

435     cancHpredL, 255/785 cancHpredH, and 253/594 cancLpredH regions detected

436     via the NHEK model (p<2.2e-16), demonstrating that a given model can be used

437     to detect regions of interest in various cell types.

438

439     **Functional annotation of regions of interest**

440

18

441    We closely examined the ten 50kb regions with the highest (cancHpredL) d-

442    scores to uncover genes that might be reclassified as oncogenic due to a higher

443    than expected SV breakpoint frequency in cancer. Likewise, we investigated the

444    ten regions with the lowest d-scores (cancLpredH), which we predict to be under

445    purifying selection, for signals of potential functionality. For this analysis we

446    used the NHEK model predictions paired with ICGC carcinoma SV breakpoints.

447

448    Nine out of ten regions with the highest d-scores overlap a gene, and four

449    overlap COSMIC genes. *CHEK2* and *CDKN2A* are known tumor suppressors, and

450    *TMPRSS2* and *ERG* is frequently involved in translocation events forming fusion

451    oncogenes in certain cancers. For example, it fuses with *TMPRSS2* in most

452    prostate cancers, with *EWS* in Ewing's sarcoma, and with *FUS* in AML. Two

453    adjacent 50kb regions on *chr17q12* overlap *GRB7* and *IKZF3*. *GRB7* encodes a

454    protein that interacts with epidermal growth factor receptor (*EGFR*), a well-

455    known proto-oncogene, and *IKZF3* is a zinc finger protein and transcription

456    factor involved in B lymphocyte regulation and differentiation as well as

457    chromatin remodeling. This region also corresponds to a known fragile site

458    *FRA17A* (49). Of the ten regions with the lowest d-scores, seven overlap a known

459    gene and two known oncogenes. The oncogene, *CDC27* , or cell division cycle 27,

460    encodes a component of the *APC* and has been shown to interact with other

461    mitotic checkpoint proteins. It is highly conserved and may be necessary for cell

462    survival. There is also a non-coding RNA found on chr2 in the centromeric

463    region, *LOC654342*, which overlaps an H3K27ac peak, and may be acting as a

464    regulatory element.

465

19

466 **Discussion**

467

468 Recent *in vitro* studies of DSB frequency in cell lines have suggested that a

469 variety of underlying genomic features are associated with DSB susceptibility.

470 We have shown that accurate models of genome-wide DSB frequency can be

471 built from a modest number of such features, with replication timing, open

472 chromatin, and marks of active promoter or enhancer regions associated with

473 increased DSBs. Although active regulatory regions often harbor actively

474 transcribed genes, it appears that chromatin accessibility at these sites rather

475 than transcription itself determines DSB propensity. The variable importance

476 metrics also show certain features to be more influential in particular cell types,

477 with CTCF and H3K36me3 having more predictive power in MCF7 than in NHEK

478 or K562. Not only are DSB patterns cell type specific, but the factors influencing

479 those patterns also depend on cell type, suggesting different mutational

480 mechanisms at play. As a matter of course, our models' accuracies decline when

481 applied to cell lines other than the training set, but they still generate reasonable

482 DSB frequency predictions, with correlations between 0.57 and 0.83 to the

483 observed data, which are large improvements over a simple inference. Since

484 chromatin features influence mutation patterns and are cell type specific, it will

485 be important to use mutational propensity profiles for matched cell types in

486 future cancer studies.

487

488 Our models of genome-wide DSB susceptibility predict DSB frequencies for all

489 50kb loci, and reflect the established correlations between replication timing and

490 DSB frequency (50) as well as tumour SV rates (9,10). A recent complementary

20

491    study has shown that 84,946 high confidence peaks of NHEK DSBCapture signal

492    (22), marking small (median: 391bp) sites of unusually high DSB susceptibility,

493    can be accurately classified from control sites using underlying genomic features

494    (34). Consistent with our results, this binary classifier suggested prominent roles

495    for DNase accessible regulatory sites and CTCF binding, and recapitulated many

496    of the patterns reported by Lensing et al (2016). However, the model of Mourad

497    et al (2018) omitted replication timing and does not provide quantitative

498    predictions of DSB susceptibility across the genome.

499

500    We used our genome-wide models of DSB susceptibility to interrogate the largest

501    tumour SV breakpoint collections and found surprising levels of agreement, such

502    that SV breakpoint enriched regions often show shifts to higher predicted DSB

503    susceptibility. In spite of variable sample sizes, the classes of simple SV likely to

504    arise by one or two DSBs (deletions, duplications, inversions, translocations)

505    showed significant increases in predicted DSB susceptibility. The NHEK model

506    best predicted the patterns of DSB susceptibility in tumours, showing genome-

507    wide elevations of predicted DSBs for all of these SV classes relative to control

508    regions. Thus, the chromatin-mediated DSB susceptibility captured in the model

509    may shape the landscape of SV recurrence in these classes.

510

511    There are many reasons why one might expect a much poorer agreement

512    between the predictions of in vitro DSB frequency models and the patterns of SV

513    breakpoints observed in tumour sequencing studies. The available collections of

514    SV breakpoints in tumours are far from perfect, and even the best ICGC data

515    suffer large variations in sample size, sample heterogeneity, sequencing depths

21

516 and SV calling methods across tumour cohorts. In addition, fundamental aspects

517 of tumour biology (cellular heterogeneity, disrupted repair pathways, chromatin

518 alterations etc.) are expected to place distinct limits on the agreement we can see

519 with the DSB patterns seen in cell lines. Evidence is also emerging that there are

520 important properties of the mutational landscape in tumours that are unlikely to

521 be captured by in vitro model systems. For example, a recent study of intra-

522 tumour diversification in colorectal cancer suggests that most mutations occur

523 during the final clonal expansion of these tumours, resulting from mutational

524 processes that are absent from normal colorectal cells (51). Enhanced rates of

525 DSB formation have also been observed in vitro at cryptic replication origins

526 activated by oncogene-induced replication stress, though these cryptic sites

527 seem to explain only a minority of SV breakpoints (<8%) across a variety of

528 TCGA tumour types (52). Given the many known and possible differences

529 between in vitro DSB model predictions and observed tumour SV breakpoints, it

530 is remarkable that significant agreement is found on any level.

531

532 There is great interest in 'hotspot' genomic regions harbouring recurrent SVs in

533 tumours, on the basis that such regions may be under positive selection,

534 conferring a proliferative or survival advantage to tumour cells. However,

535 rigorous inference of selection requires a proxy for the expected rate of

536 recurrence within such regions. Using model predictions as this proxy we have

537 produced refined hotspot predictions, reflecting SV breakpoint frequencies

538 relative to the predicted susceptibility of each region. Since our predictions of

539 DSB susceptibility are genome-wide it was also possible to predict coldspot

540 regions, regions possessing unexpectedly low SV breakpoint rates given model

22

541  predictions, and putatively subject to negative or purifying selection in tumours.

542  If selection in tumours is prominent in driving SV breakpoint frequencies away

543  from DSB model predictions, we might expect hotspot and coldspot regions to

544  show unusual functional enrichments. Multiple caveats apply to the annotations

545  examined but analysis using the NHEK model shows that ICGC carcinoma

546  hotspots are enriched for putative oncogenes. Coldspots occupy gene-rich

547  neighbourhoods but and are also enriched in active promoters and strong

548  enhancers, and insulators, indicating regulatory regions that may have

549  experienced purifying selection in tumours.

550

551  **Conclusions**

552

553  When inferring selection on single nucleotide variants it is standard practice to

554  make comparisons between the observed variant frequencies and the

555  frequencies expected, according to a model of single nucleotide mutation rates.

556  We have developed models of DSB mutation rates that can be used to generate

557  expected SV breakpoint frequencies and illuminate regions with significant

558  deviations from these expectations. This approach provides statistically rigorous

559  protocols to prioritize novel loci putatively under selection in tumours,

560  generating testable hypotheses for further experimental studies.

561

562

563  **Methods**

564

565  *Derivation of DSB data in the K562 and MCF7 cell lines*

23

566    DSB profiles were generated with an adapted version of the Breaks labeling *in*

567    *situ* and sequencing protocol (25), in which DSB ends are labeled with a dsDNA

568    BLISS adapter in cell suspensions of 1 million cells. Afterwards the published

569    protocol is followed with only minor modifications. Labeled DSBs are selectively

570    amplified using T7-driven linear amplification, after which sequencing libraries

571    are generated and sequenced with single-end 1x75 v2 chemistry on an Illumina

572    NextSeq 500. Raw sequencing reads were demultiplexed by Illumina's

573    BaseSpace, after which FASTQ files were downloaded and processed as

574    described in Yan et al. 2017 (SRA accession SRP150602). In brief, reads with the

575    expected prefix of 8nt UMI and 8nt sample barcode sequence were filtered using

576    SAMtools and *scan for matches*, allowing at most one mismatch per barcode.

577    Trimmed reads were then aligned to GRCh37 using bwa mem, and reads with

578    mapping scores below 30 were discarded. Next, PCR duplicates were identified

579    by searching for proximal reads (within 30bp of the reference genome) with at

580    most two mismatches in the UMI sequence, which were then grouped and

581    collapsed into a single break location. Finally, we generated .bed files with DSB

582    locations and the number of unique UMIs indicating that location.

583

584    *Generating random forest models*

585    We downloaded ten tracks from ENCODE for multiple chromatin marks,

586    replication timing, open chromatin, several DNA binding proteins, and

587    nucleosome pull-downs from the UCSC genome browser (53). We used G-

588    quadruplex data generated by Chambers et al, (GSE63874). In their study, they

589    make separate .bedgraph files available with the G-quadruplex density for each

590    strand. We used the sum of the plus and minus strands in our analysis. The list of

24

591    bigwig files used for each cell line along with their sources and graphical labels is

592    in Supplementary Table 1. We used the bigWigAverageOverBed tool from the

593    kentUtils tool library to produce average signal per 50kb in non-overlapping

594    windows across hg19 for each track. We combined the results to a single matrix

595    per cell line composed of 61,903 rows, one for each 50kb bin, and 11 columns,

596    one for each chromatin or genomic feature. These feature matrices are available

597    in supplementary data and scatter plots of each feature with the NHEK

598    DSBCapture data are shown in Supplementary Figure 3.

599

600    For the extended model in Supplementary Figure 4, we downloaded an

601    additional nine features from the UCSC genome browser (53), which were

602    processed in the same way as the ten ENCODE features used in the primary

603    feature matrix. We also downloaded .hic files for NHEK, K562, and HMEC cells

604    generated from Rao, et al. (GSE63525). We used their custom toolbox, Juicer, to

605    calculate eigenvectors per chromosome, and generated 50kb resolution

606    eigenvector profiles using the bedGraphToBigWig and bigWigAverageOverBed

607    tools from kentUtils. The figure labels and sources for these data are in

608    Supplementary Table2, and the extended feature matrices are in supplementary

609    data.

610

611    We generated DSB frequency scores from each of four HTS DSB profiling

612    datasets: two in NHEK cells, one for K562, unpublished, and one for MCF7,

613    unpublished. As mentioned in the results, two replicates for each of two DSB HTS

614    profiling methods, DSBCapture and BLESS, were available from Lensing et al.

615    (22). We took the average per 50kb of the replicates to create an NHEK

616   DSBCapture profile and an NHEK BLESS profile. We combined three replicates of

617   MCF7 BLISS data (via a sum operation) to serve as our MCF7 DSB profile. A

618   fourth MCF7 BLISS dataset is available, but we excluded it from our analysis

619   because it had a distinctly lower correlation to the other three datasets (0.90-

620   0.92 as opposed to 0.97-0.99). These scores are available as supplementary files.

621

622   We used the randomForest package in R to generate random forest models with

623   500 trees and five OOB permutations per tree (options ntree=500, nPerm=5). To

624   calculate variable importance, we used the importance command within the

625   randomForest                    package                    (https://cran.r-

626   project.org/web/packages/randomForest/index.html),    which    calculates    the

627   average prediction error rate (MSE) for each datapoint (50kb bin) across all

628   trees in the random forest. Then, for each feature variable, the values are

629   randomly permuted and the MSE for each 50kb bin is calculated again. The final

630   variable importance score is the average difference in MSE before and after the

631   permutation, normalized by the standard deviation of these differences. Because

632   many features are inter-correlated, their importance measures were very

633   similar. Therefore, in order to determine a consistent ranking of features'

634   importance values, we generated ten random forest models per dataset and

635   calculated the average and standard deviation of importance across the ten

636   models.

637   Although random forest models are not susceptible to overfitting, to confirm that

638   our models were not overfit to the DSB data, we also generated a random forest

639   model for the NHEK DSBCapture dataset, holding out one third of the data as the

640   test set and training the model on the remaining two thirds.  This model showed

26

641   0.93 Pearson's correlation between the predictions and the observed data for the

642   training set, similar to the model trained on the full dataset (Sup Figure 5).

643

644   *Determining tumour ESBs and their predicted DSB scores*

645   To determine SV DSB rates in from TCGA data, we downloaded CNV data from

646   TCGA (54), which came from Affymetrix SNP 6.0 arrays processed by the

647   DNAcopy    R-package    (https://docs.gdc.cancer.gov/Data/PDF/Data_UG.pdf).

648   DNAcopy generates a set of continuous segments, outputting regions with little

649   or no copy number change, so we filtered these, defining segments with a CN

650   ratio >1 as amplifications and ratios < -1 as deletions. The segments were lifted

651   from hg38 to hg19 using UCSC's liftOver tool. For each CNV, we counted a single

652   DSB to occur in a 50kb bin if either or both ends of the segment overlapped the

653   bin. The TCGA-BLOOD group includes the two blood cancer cohorts: acute

654   myeloid leukemia (LAML) and lymphoid neoplasm diffuse large B-cell lymphoma

655   (DLBC), while the TCGA-SCCA group includes three squamous cell carcinomas:

656   cervical squamous cell carcinoma and endocervical adenocarcinoma (CESC),

657   head and neck squamous cell carcinoma (HNSC), and lung squamous cell

658   carcinoma (LUSC). The BRCA group includes only the TCGA breast cancer cohort

659   (BRCA), and the PANC group includes all 33 cancer types, shown in

660   Supplementary Figure 7. Counts for various groups and CNV types are available

661   as Supplementary Files.

662   We downloaded available WGS SV calls from the ICGC Data Portal

663   (https://dcc.icgc.org/projects). As with the TCGA CNV, a single DSB was counted

664   per 50kb bin if either one or two ends of a SV overlapped the region. The ICGC

665   pancancer group contains SVs from 17 cancer studies, shown in Supplementary

27

666  Figure 6. The carcinoma group contains all available carcinoma cancer studies,

667  excluding breast cancer: early onset prostate cancer (EOPC-DE), liver cancer

668  (LIRI-JP), pancreatic cancer (PACA-CA, PAEN-AU, PAEN-IT), prostate cancer

669  (PRAD-CA, PRAD-UK), and skin adenocarcinoma (SKCA-BR). The ICGC blood

670  group contains chronic lymphocytic leukemia (CLLE-ES) and malignant

671  lymphoma (MALY-DE), and the breast group contains breast cancer studies

672  (BRCA-EU and BRCA-FR). A table of DSB counts per 50kb broken up by group

673  and SV type is in supplementary data.

674

675  We determined enriched SV breakpoint regions (ESBs) per cohort or SV type

676  grouping by ranking the 50kb bins by the number of DSB, excluding regions with

677  no DSB in the group, and using the number of DSB in the top 5% as the cutoff. All

678  50kb regions with a DSB count greater than or equal to the cutoff were

679  designated ESBs. We used a Wilcoxon ranked sum test (R wilcox.test command)

680  to test for significant increase in the predicted DSB values for ESBs compared to

681  all other regions, and we excluded regions in which no DSB were found in any

682  cancer study since these are likely to be unmappable or blacklisted regions.

683

684  The correlation between TCGA and ICGC pancancer SV breakpoint counts was

685  calculated using Spearman's rho and excluding 50kb regions with no SV

686  breakpoints in either the TCGA or ICGC datasets.  The top 5% ESBs were found

687  for each dataset, with 2,839 regions found in TCGA and 3,072 in ICGC, and the

688  significance of the overlap was calculated using a hypergeometric test (R

689  command phyper with q=177, m=2,839, n=61,903-2,839, and k=3,072).

690

691  *Calculating d-scores*

692  We used the R package fistdistrplus (55) to determine the distributions with the

693  best fit to the DSB prediction values and the SV breakpoint frequencies. We used

694  a likelihood maximization test (method="mle") and the BIC (Bayesian

695  Information Criterion) measure of goodness of fit to choose the best distribution.

696  We tested a lognormal, log-logistic, gamma, normal, and an exponential

697  distribution, and fitted the distributions to the bulk of the SV breakpoint or DSB

698  prediction data. We excluded 50kb regions with breakpoint frequencies greater

699  than six times the interquartile range from the median in order to exclude

700  extreme outliers. While we aimed to emphasize the fit of the tails of our data's

701  distributions, including these outliers resulted in poorly fitting distributions to

702  the bulk of the real data. Once we found the best of the three candidate model

703  distributions, we assigned a p-value to each 50kb bin from the fitted distribution

704  (using the plnorm, pllogis, or pgamma functions in R) which represent the

705  probability of seeing a given breakpoint frequency or DSB prediction or greater

706  in the known distribution. The actual and fitted distributions and quantile-

707  quantile plots are shown in Supplementary Figures 13 and 14.

708

709  Next, for each 50kb bin, we calculated the difference in log p-values between the

710  predicted DSB and the actual SV breakpoints, called d-scores. Using the

711  fistdistrplus R package again, we determined the best-fit distribution for the d-

712  scores, choosing between a t-distribution, a normal, and a Cauchy distribution.

713  Again, we used a maximum likelihood method and the BIC measurement and

714  excluded extreme outliers. In all cases, a t-distribution with four degrees of

715  freedom (df=4) was the best fit, so each 50kb bin was assigned a p-value from

716 this distribution according to its d-score. The histograms and quantile-quantile

717 plots of the d-scores and fitted distributions are shown in Supplementary Figure

718 15.

719

720 *Calculating gene set and chromatin domain enrichments*

721 We used the d-score p-values to categorize regions into informative subsets,

722 using the R command qt(p=0.01, df=4, lower.tail=FALSE) to determine the d-

723 score cutoffs. The cancHpredL class of regions have d-scores in the upper one

724 percentile (> 3.75), and the cancLpredH have d-scores in the lower one

725 percentile (< -3.75). The cancHpredH class has d-scores in the $40^{th}$ to $70^{th}$

726 percentiles and SV breakpoint frequencies or DSB predictions with p-values less

727 than 0.01, so these regions have significantly (p-value < 0.01) high SV

728 breakpoints or DSB predictions but insignificant d-scores (p-value < 0.6). The

729 cancHpredL2 class consists of regions with SV breakpoint p-values less than

730 0.01, and DSB predictions less than 0.5 for the NHEK models and less than 0.001

731 for the MCF7 model.

732 We used a binomial test to measure the significance of overlaps between sets

733 when comparing results from the MCF7 model and the NHEK model applied to

734 ICGC breast cancer data and MCF7 cell line features (R command binom.test).

735

736 We used the R package regioneR (56) to compute the overlap significance

737 between each set of regions and various genome and chromatin annotation files.

738 A list of annotation sets and their original sources are in Supplementary Table 2.

739 We matched Cancer5000 genes and Cosmic gene lists to RefSeq gene names in

740 order to get their genome coordinates, so the cancer gene lists are RefSeq gene

741 subsets. The super enhancer set (SEA) came from A549 cells, derived from a lung

742 carcinoma (47). Common fragile sites (CFS) were collected from NCBI's gene

743 archive by searching for "common fragile site" or "fragile site" within human

744 genes. Many fragile sites are annotated by chromosome band but do not have

745 exact coordinates; we filtered these out because they are low resolution. The

746 chromHMM (48) annotation came from the UCSC genome browser. We tested

747 enrichment of the NHEK states with the NHEK model d-score classes and the

748 HMEC track, from primary mammary epithelial cells, with the MCF7 model's d-

749 score classes. The regioneR package performs random circular permutation of

750 regions of interest and then computes the number of overlaps between the

751 permutated set and a second set of regions. The p-value represents how often,

752 over the course of the permutations, the two sets overlap to the same extent that

753 they do without any permutation. We used 1,000 iterations to achieve a

754 maximum p-value of 0.001.

755

756 **Declarations**

757

758 *Ethics Approval*

759 Approval for access and use of ICGC variant data was obtained from the ICGC

760 Data Access Compliance Office. Use of TCGA CNV does not require ethics

761 approval.

762

763 *Consent for Publication*

764 Not applicable

765

766 *Availability of data and materials*

767 All analysis was done using GRCh37 as the reference genome. The raw BLISS

768 sequencing data is available on SRA with accession SRP150602. All scripts and

769 commands used to do this analysis are available on github

770 (https://github.com/TracyBallinger/dsb_model). In addition, we have made

771 ipython notebooks for the figures used in this manuscript to ease reproducibility

772 and allow further exploration of the data, also available on github. All

773 supplementary files are available for download at

774 https://datashare.is.ed.ac.uk/handle/10283/3103.

775

776 *Competing Interests*

777 The authors declare they have no competing interests.

778

779 *Funding*

780 This study was funded by core funding of the UK Medical Research Council

781 (MRC) to the MRC Human Genetics Unit to C.S.; by grants from the Karolinska

782 Institutet, the Ragnar Söderberg Foundation, the Swedish Foundation for

783 Strategic Research (N.C.: BD15-0095), and the Strategic Research Programme in

784 Cancer (StratCan) at Karolinska Institutet to N.C.; and by a Rubicon fellowship

785 from the Netherlands Organisation for Scientific Research (NWO) to B.B.

786

787 *Authors' Contributions*

788 BB and RM generated the BLISS DSB profiles. SG developed the BLISS alignment

789 pipeline and generated .bed files of DSB profiles. TB performed all subsequent

790    data analysis and produced figures. TB and CS wrote the manuscript. NC and CS

791    supervised the project. TB, BB, NC, and CS edited the final manuscript.

792

796

797

798    **References**

799

800    1.  Ciriello G, Miller ML, Aksoy BA, Senbabaoglu Y, Schultz N, Sander C. Emerging
801        landscape of oncogenic signatures across human cancers. Nat Genet. 2013
802        Oct;45(10):1127–1133.

803    2.  Patch A-M, Christie EL, Etemadmoghadam D, Garsed DW, George J, Fereday S,
804        et al. Whole–genome characterization of chemoresistant ovarian cancer.
805        Nature. 2015 May;521(7553):489–494.

806    3.  Scarpa A, Chang DK, Nones K, Corbo V, Patch A-M, Bailey P, et al. Whole-
807        genome landscape of pancreatic neuroendocrine tumours. Nature. 2017
808        Mar;543(7643):65–71.

809    4.  Alaei-Mahabadi B, Bhadury J, Karlsson JW, Nilsson JA, Larsson E. Global
810        analysis of somatic structural genomic alterations and their impact on gene
811        expression in diverse human cancers. Proc Natl Acad Sci U S A. 2016
812        Nov;113(48):13768–13773.

813    5.  Li Y, Roberts N, Weischenfeldt J, Wala JA, Shapira O, Schumacher S, et al.
814        Patterns of structural variation in human cancer. bioRxiv. 2017
815        Aug;181339.

816    6.  Sudmant PH, Rausch T, Gardner EJ, Handsaker RE, Abyzov A, Huddleston J, et
817        al. An integrated map of structural variation in 2,504 human genomes.
818        Nature. 2015 Oct;526(7571):75–81.

819    7.  Weischenfeldt J, Dubash T, Drainas AP, Mardin BR, Chen Y, Stütz AM, et al.
820        Pan-cancer analysis of somatic copy-number alterations implicates IRS4
821        and IGF2 in enhancer hijacking. Nat Genet. 2017 Jan;49(1):65–74.

822    8.  Glodzik D, Morganella S, Davies H, Simpson PT, Li Y, Zou X, et al. A somatic-
823        mutational process recurrently duplicates germline susceptibility loci and
824        tissue-specific super-enhancers in breast cancers. Nat Genet. 2017
825        Jan;49(3):341–348.

826    9.  Morganella S, Alexandrov LB, Glodzik D, Zou X, Davies H, Staaf J, et al. The
827        topography of mutational processes in breast cancer genomes. Nat
828        Commun. 2016 May;7:11383.

829    10. Schuster-Böckler B, Lehner B. Chromatin organization is a major influence on
830        regional mutation rates in human cancer cells. Nature. 2012
831        Aug;488(7412):504–507.

832    11. Ding L, Getz G, Wheeler DA, Mardis ER, McLellan MD, Cibulskis K, et al.
833        Somatic mutations affect key pathways in lung adenocarcinoma. Nature.
834        2008 Oct;455(7216):1069–1075.

835    12. Jackson SP, Bartek J. The DNA-damage response in human biology and
836        disease. Nature. 2009 Oct;461(7267):1071–1078.

837    13. Biehs R, Steinlage M, Barton O, Juhász S, Künzel J, Spies J, et al. DNA Double-
838        Strand Break Resection Occurs during Non-homologous End Joining in G1
839        but Is Distinct from Resection during Homologous Recombination. Mol Cell.
840        2017 Feb;65(4):671–684.e5.

841    14. Nussenzweig A, Nussenzweig MC. A backup DNA repair pathway moves to
842        the forefront. Cell. 2007 Oct;131(2):223–225.

843    15. Clouaire T, Legube G. DNA double strand break repair pathway choice: a
844        chromatin based decision? Nucl Austin Tex. 2015;6(2):107–113.

845    16. Glover TW, Berger C, Coyle J, Echo B. DNA polymerase alpha inhibition by
846        aphidicolin induces gaps and breaks at common fragile sites in human
847        chromosomes. Hum Genet. 1984;67(2):136–142.

848    17. Canela A, Sridharan S, Sciascia N, Tubbs A, Meltzer P, Sleckman BP, et al. DNA
849        Breaks and End Resection Measured Genome-wide by End Sequencing. Mol
850        Cell. 2016 Sep;63(5):898–911.

851    18. Crosetto N, Mitra A, Silva MJ, Bienko M, Dojer N, Wang Q, et al. Nucleotide-
852        resolution DNA double-strand break mapping by next-generation
853        sequencing. Nat Methods. 2013 Mar;10(4):361–365.

854    19. Frock RL, Hu J, Meyers RM, Ho Y-J, Kii E, Alt FW. Genome-wide detection of
855        DNA double-stranded breaks induced by engineered nucleases. Nat
856        Biotechnol. 2015 Feb;33(2):179–186.

857    20. Iacovoni JS, Caron P, Lassadi I, Nicolas E, Massip L, Trouche D, et al. High-
858        resolution profiling of gammaH2AX around DNA double strand breaks in
859        the mammalian genome. EMBO J. 2010 Apr;29(8):1446–1457.

860   21. Kim D, Bae S, Park J, Kim E, Kim S, Yu HR, et al. Digenome-seq: genome-wide
861       profiling of CRISPR-Cas9 off-target effects in human cells. Nat Methods.
862       2015 Mar;12(3):237–43– 1 p following 243.

863   22. Lensing SV, Marsico G, Hänsel-Hertsch R, Lam EY, Tannahill D,
864       Balasubramanian S. DSBCapture: in situ capture and sequencing of DNA
865       breaks. Nat Methods. 2016 Aug;13(10):855–857.

866   23. Slaymaker IM, Gao L, Zetsche B, Scott DA, Yan WX, Zhang F. Rationally
867       engineered Cas9 nucleases with improved specificity. Science. 2016
868       Jan;351(6268):84–88.

869   24. Wei P-C, Chang AN, Kao J, Du Z, Meyers RM, Alt FW, et al. Long Neural Genes
870       Harbor Recurrent DNA Break Clusters in Neural Stem/Progenitor Cells. Cell.
871       2016 Feb;164(4):644–655.

872   25. Yan WX, Mirzazadeh R, Garnerone S, Scott D, Schneider MW, Kallas T, et al.
873       BLISS is a versatile and quantitative method for genome-wide profiling of
874       DNA double-strand breaks. Nat Commun. 2017;8:15058.

875   26. De S, Michor F. DNA secondary structures and epigenetic determinants of
876       cancer genome evolution. Nat Struct 38 Mol Biol. 2011 Jul;18(8):950–955.

877   27. Moore BL, Aitken S, Semple CA. Integrative modeling reveals the principles of
878       multi-scale chromatin boundary formation in human nuclear organization.
879       Genome Biol. 2015 May;16(1):1270.

880   28. Polak P, Karlic R, Koren A, Thurman R, Sandstrom R, Lawrence MS, et al. Cell-
881       of-origin chromatin organization shapes the mutational landscape of cancer.
882       Nature. 2015 Feb;518(7539):360–364.

883   29. Whalen S, Truty RM, Pollard KS. Enhancer–promoter interactions are
884       encoded by complex genomic signatures on looping chromatin. Nat Genet.
885       2016 Apr;48(5):488–496.

886   30. Consortium TEP. An integrated encyclopedia of DNA elements in the human
887       genome. Nature. 2012 Sep;489(7414):57–74.

888   31. Chambers VS, Marsico G, Boutell JM, Di Antonio M, Smith GP,
889       Balasubramanian S. High-throughput sequencing of DNA G-quadruplex
890       structures in the human genome. Nat Biotechnol. 2015 Jul;33(8):877–881.

891   32. Rao SSP, Huntley MH, Durand NC, Stamenova EK, Bochkov ID, Robinson JT, et
892       al. A 3D Map of the Human Genome at Kilobase Resolution Reveals
893       Principles of Chromatin Looping. Cell. 2014 Dec;159(7):1665–1680.

894   33. Drier Y, Lawrence MS, Carter SL, Stewart C, Gabriel SB, Lander ES, et al.
895       Somatic rearrangements across cancer reveal classes of samples with
896       distinct patterns of DNA breakage and rearrangement-induced
897       hypermutability. Genome Res. 2013 Feb;23(2):228–235.

35

898   34. Mourad R, Ginalski K, Legube G, Cuvier O. Predicting double-strand DNA
899        breaks using epigenome marks or DNA at kilobase resolution. Genome Biol.
900        2018 Dec;19(1):34.

901   35. Canela A, Maman Y, Jung S, Wong N, Callen E, Day A, et al. Genome
902        Organization Drives Chromosome Fragility. Cell. 2017 Jul;170(3):507–
903        521.e18.

904   36. Kaiser VB, Semple CA. When TADs go bad: chromatin structure and nuclear
905        organisation in human disease. F1000Research. 2017;6.

906   37. Altemose N, Miga KH, Maggioni M, Willard HF. Genomic Characterization of
907        Large Heterochromatic Gaps in the Human Genome Assembly. PLoS Comput
908        Biol. 2014 May;10(5):e1003628.

909   38. Glover TW, Wilson TE, Arlt MF. Fragile sites in cancer: more than meets the
910        eye. Nat Rev Cancer. 2017 Aug;17(8):489–501.

911   39. Quint KD, Genders RE, de Koning MN, Borgogna C, Gariglio M, Bavinck JNB, et
912        al. Human Beta-papillomavirus infection and keratinocyte carcinomas. J
913        Pathol. 2015 Jan;235(2):342–354.

914   40. Zhang J, Baran J, Cros A, Guberman JM, Haider S, Hsu J, et al. International
915        Cancer Genome Consortium Data Portal–a one-stop shop for cancer
916        genomics data. Database. 2011 Sep;2011(0):bar026–bar026.

917   41. Weckselblatt B, Rudd MK. Human Structural Variation: Mechanisms of
918        Chromosome Rearrangements. Trends Genet. 2015 Oct;31(10):587–599.

919   42. Viguera E, Canceill D, Ehrlich SD. Replication slippage involves DNA
920        polymerase pausing and dissociation. EMBO J. 2001 May;20(10):2587–
921        2595.

922   43. Fungtammasan A, Walsh E, Chiaromonte F, Eckert KA, Makova KD. A
923        genome-wide analysis of common fragile sites: what features determine
924        chromosomal instability in the human genome? Genome Res. 2012
925        Jun;22(6):993–1005.

926   44. Irony-Tur Sinai M, Kerem B. DNA replication stress drives fragile site
927        instability. Mutat Res. 2018 Mar;808:56–61.

928   45. Lawrence MS, Stojanov P, Mermel CH, Robinson JT, Garraway LA, Golub TR,
929        et al. Discovery and saturation analysis of cancer genes across 21 tumour
930        types. Nature. 2014 Jan;505(7484):495–501.

931   46. Forbes SA, Beare D, Boutselakis H, Bamford S, Bindal N, Tate J, et al. COSMIC:
932        somatic cancer genetics at high-resolution. Nucleic Acids Res. 2017
933        Jan;45(D1):D777–D783.

934   47. Wei Y, Zhang S, Shang S, Zhang B, Li S, Wang X, et al. SEA: a super-enhancer
935        archive. Nucleic Acids Res. 2016 Jan;44(D1):D172–D179.

936   48. Ernst J, Kellis M. ChromHMM: automating chromatin-state discovery and
937        characterization. Nat Methods. 2012 Feb;9(3):215–216.

938   49. Mrasek K, Schoder C, Teichmann A-C, Behr K, Franze B, Wilhelm K, et al.
939        Global screening and extended nomenclature for 230 aphidicolin-inducible
940        fragile sites, including 61 yet unreported ones. Int J Oncol. 2010
941        Apr;36(4):929–940.

942   50. Sima J, Gilbert DM. Complex correlations: replication timing and mutational
943        landscapes during cancer and genome evolution. Curr Opin Genet Dev. 2014
944        Apr;25:93–100.

945   51. Roerink SF, Sasaki N, Lee-Six H, Young MD, Alexandrov LB, Behjati S, et al.
946        Intra-tumour diversification in colorectal cancer at the single-cell level.
947        Nature. 2018 Apr;556(7702):457–462.

948   52. Macheret M, Halazonetis TD. Intragenic origins due to short G1 phases
949        underlie oncogene-induced DNA replication stress. Nature. 2018
950        Mar;555(7694):112–116.

951   53. Kent WJ, Sugnet CW, Furey TS, Roskin KM, Pringle TH, Zahler AM, et al. The
952        human genome browser at UCSC. Genome Res. 2002 Jun;12(6):996–1006.

953   54. Grossman RL, Heath AP, Ferretti V, Varmus HE, Lowy DR, Kibbe WA, et al.
954        Toward a Shared Vision for Cancer Genomic Data. N Engl J Med. 2016
955        Sep;375(12):1109–1112.

956   55. Delignette-Muller ML, Software CDJ of S, 2015. fitdistrplus: An R package for
957        fitting distributions. rdrr.io.

958   56. Gel B, Díez-Villanueva A, Serra E, Buschbeck M, 2015. regioneR: an
959        R/Bioconductor package for the association analysis of genomic regions
960        based on permutation tests. academic.oup.com.

961

962   **Figure Legends**

963

964   **Figure 1**: DSB frequency and genomic features display similar patterns. The

965   tracks show DSBCapture profiles in NHEK cells, BLESS profiles in NHEK cells,

966   BLISS in K562 cells, and BLISS in MCF7 cells.  All tracks are at 50kb resolution

967   over a representative region of chromosome 1, with a variety of chromatin and

968   sequence features to illustrate the similarities between them. Numbers in

969 parenthesis are the spearman's rho between the associated track and the NHEK

970 DSBCapture 1 dataset.

971

972 **Figure 2**: Accurate models of DSB frequency built from chromatin and sequence

973 features. Panels A-D show random forest regression model predictions built

974 upon eleven genomic features at 50kb resolution compared to observed DSB

975 frequencies for four datasets: NHEK DSBCapture, NHEK BLESS, K562 BLISS, and

976 MCF7 BLISS. The y-values reflect the sequencing depth of each dataset. The

977 models' predictions are all highly correlated with the observed data, as shown by

978 the noted Pearson's correlations ($p < 2.2e-16$ for each dataset). Panels E-H show

979 the predictive features ranked by variable importance, a measure of how useful a

980 particular feature is for the model (see methods).

981

982 **Figure 3**: Modelling accuracy and the polarity of genomic features. A) NHEK

983 DSBCapture 50kb regions data is split into three distinct groups with differing

984 modelling accuracies. Panels B and C show the values of the model features for

985 the two boxes, A and B, and for group C, which contains randomly chosen points

986 along the spectrum of DSB frequency values for the majority of the genome. The

987 columns are ordered by observed DSB frequency, shown on the top row, and the

988 rows for features used to build the model (the third to second to last row) are

989 ordered by average variable importance. The number of 50kb regions in each

990 group is shown in parenthesis above each heatmap. Each feature was

991 normalized, setting the $1^{st}$ to $99^{th}$ quantiles to values between 0 and 1, with high

992 outliers (in the top percentile) set to 1.1. B) Group A has high H3K9me3 and low

993 mappability scores, indicative of heterochromatin and repetitive sequence, while

38

994    B has feature patterns that closely match low DSB values in group C. C) For most

995    of the genome, high H3K9me3 corresponds to low DSB regions, and high, or

996    early, replication timing values and open chromatin values signify high DSB

997    regions.

998

999    **Figure 4**: DSB models improve predictions for non-model cell types. Models

1000    trained using a dataset from one cell type were used to generate predictions for a

1001    different cell type, given the matched features. The dark blue lines mark the

1002    Pearson's correlation between the two cell types. The cell type used to train the

1003    model is indicated by the colour of the bar, and the cell type on which the model

1004    is being applied is shown on the x-axis. In all cases, the random forest model

1005    greatly improves the predictions from a naïve inference, with a 1.3-1.8 fold

1006    improvement in correlation.

1007

1008    **Figure 5**: Regions enriched for cancer SV breakpoints (ESBs) display a

1009    significant increase in DSB frequency across cancer types. A-C) The regions with

1010    ICGC SV breakpoint frequencies in the top 5% are shown with their predicted

1011    DSB values as violin plots for each of the three cell type models: NHEK, K562, and

1012    MCF7. ICGC cohorts are shown all together (pancancer), and split into three

1013    cancer categories: carcinoma, blood, and breast cancers (see methods). D-F)

1014    ICGC SV breakpoint counts separated by SV type, and the top 5% of ESBs are

1015    shown with their predicted DSB values as violin plots. The numbers following

1016    the x-axis labels are SV breakpoint count cut-offs for the top 5% ESBs, and the

1017    numbers in parenthesis are the number of 50kb regions that meet the cut-off.

1018    For example, there are 225 50kb regions with more than two SV breakpoints in

39

1019    blood cancers. Stars indicate significantly higher values in DSB predictions for

1020    the ESBs relative to non-ESBs for each category, as determined by a Wilcox

1021    ranked sum test (* for $p \leq 0.05$, ** for $p \leq 0.01$, *** for $p \leq 1e-3$, and **** for
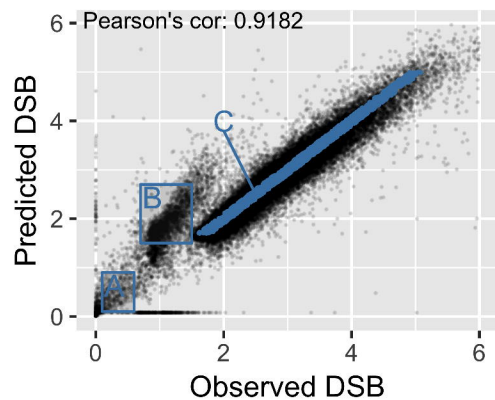
1022    $p \leq 1e-4$).

1023

1024    **Figure 6**: Inference of positively and negatively selected SV regions. A) The

1025    predicted DSB frequencies for regions overlapping RefSeq genes, two sets of

1026    cancer consensus genes, and common fragile sites (CFS) are shown as violin

1027    plots. The stars represent significantly higher values in the region subsets,

1028    compared to genomic regions that do not overlap the given annotation set, using

1029    a Wilcox ranked sum text. B) The same regions as in a), but with d-score values, a

1030    measure of the deviation of the observed breakpoint frequencies from the

1031    predicted or expected DSB frequencies. C) Observed SV breakpoint frequencies

1032    for ICGC carcinomas (excluding breast cancer) with predicted DSB frequencies

1033    from the NHEK DSBCapture model. Each point represents a 50kb region and is

1034    coloured by its d-score. Regions were split into high (cancHpredL) and low

1035    (cancLpredH) d-score categories (d-score p-value < 0.01), a cancHpredH

1036    category, representing regions with d-scores near zero, and a cancHpredL2

1037    category, representing low mappability regions (see methods). D) Each category

1038    was tested for enrichment of various annotations using circular permutation

1039    (see methods). The yellow dotted line marks p<0.01 significance, and the

1040    numbers in parenthesis indicate the number of 50kb regions in each category,
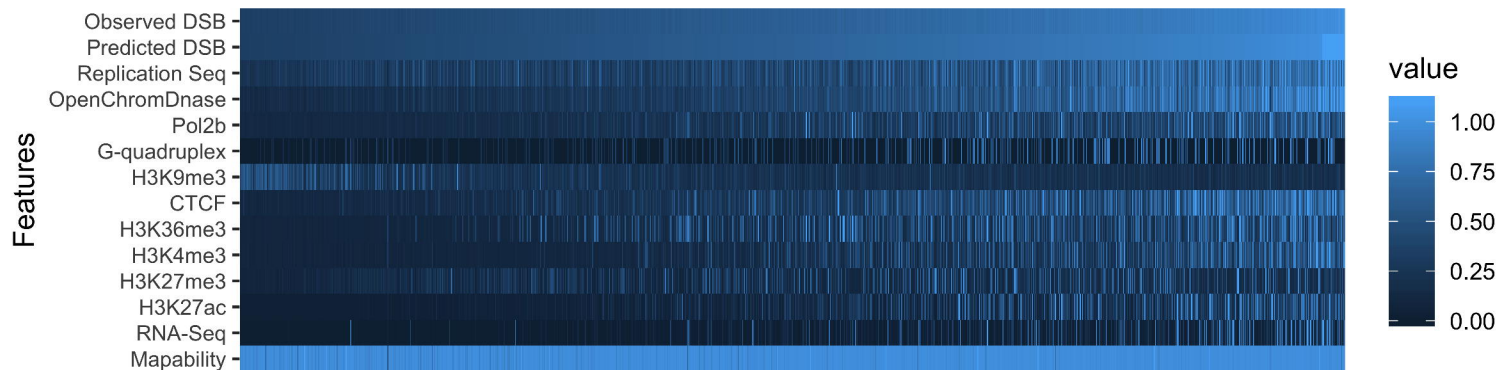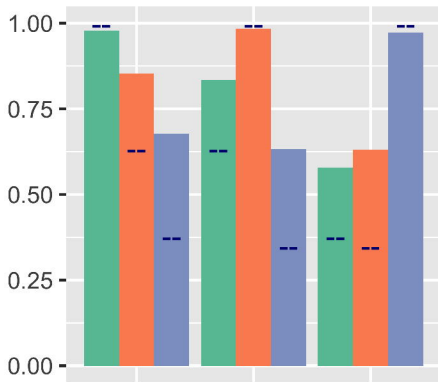
1041    out of 61,903 in total.

1042

A — NHEK DSBCapture; Pearson's cor: 0.9182
B — NHEK BLESS; Pearson's cor: 0.8337
C — K562 BLISS; Pearson's cor: 0.8986
D — MCF7 BLISS; Pearson's cor: 0.8727

A

Pearson's cor: 0.9182

B

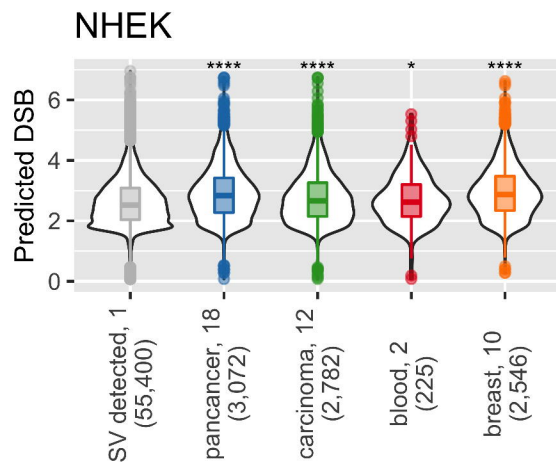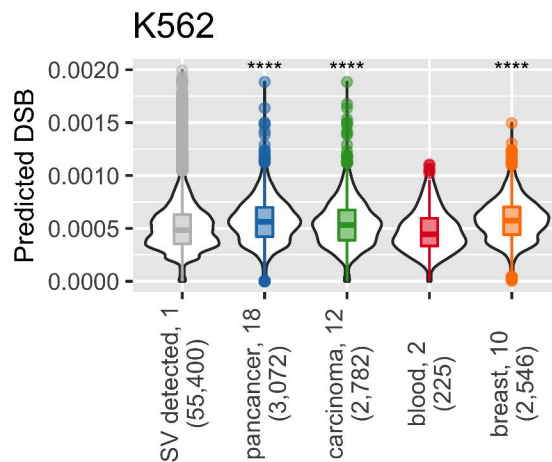Box A: 0.1-0.6  Box B: 0.7-1.5
(419)              (2063)

C

Group C: 1.7-5
(1489)

ICGC enriched SV breakpoint regions (ESBs), top 5%

# NHEK model with ICGC carcinomas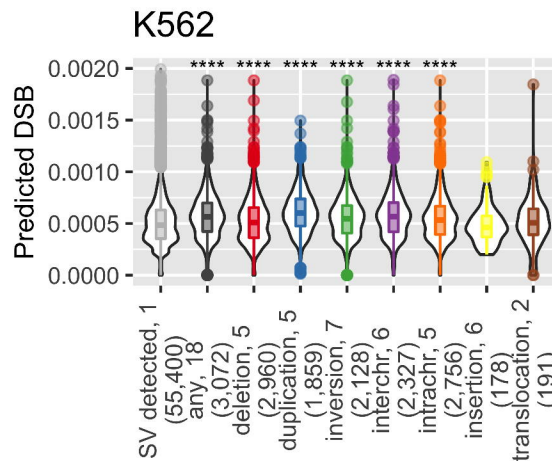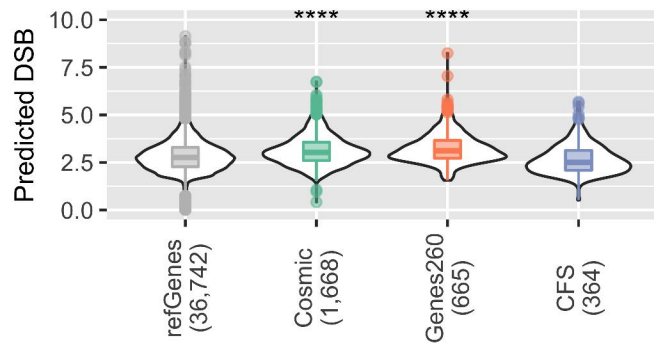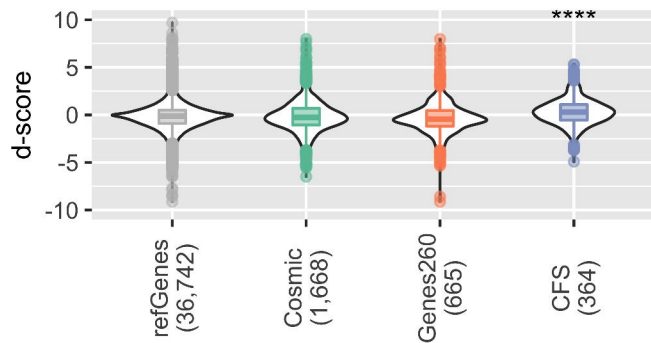