

1 **DSBCapture: *in situ* capture and direct sequencing of dsDNA breaks**

2
3 Stefanie V. Lensing¹, Giovanni Marsico¹, Robert Hänsel-Hertsch¹, Enid Y. Lam^{1,4}, David Tannahill¹ &
4 Shankar Balasubramanian^{1,2,3}

5
6 1. Cancer Research UK Cambridge Institute, University of Cambridge, Cambridge, UK.

7 2. Department of Chemistry, University of Cambridge, Cambridge, UK.

8 3. School of Clinical Medicine, University of Cambridge, Cambridge, UK

9 4. Present address: Cancer Research Division, Peter MacCallum Cancer Centre, East Melbourne,
10 Victoria, Australia

11 Correspondence should be addressed to S.B.: sb10031@cam.ac.uk

12
13 Double-strand DNA breaks (DSBs) continuously arise and cause mutations and chromosomal
14 rearrangements. Here, we present DSBCapture, a sequencing-based method that captures DSBs *in situ*
15 and directly maps these at single nucleotide resolution enabling the study of DSB origin. DSBCapture
16 shows substantially increased sensitivity and data yield compared to other methods. Employing
17 DSBCapture, we uncovered a striking relationship between DSBs and elevated transcription within
18 nucleosome-depleted chromatin.

19
20 Due to their mutagenic potential, it is important to identify the precise location of endogenous DSBs
21 and those induced by therapeutic drugs and environmental insults^{1,2}. ChIP-seq³, GUIDE-seq⁴ and
22 BLESS⁵, are currently used for the genome-wide investigation of DSBs *in situ*. ChIP-seq, however, is
23 indirect and relies on the capture of specific proteins that mark DSBs by proxy, thus compromising
24 both accuracy and resolution. GUIDE-seq profiles off-target cleavage by CRISPR-Cas nucleases and
25 relies on the integration of a double-stranded oligonucleotide at DSBs by non-homologous end joining
26 (NHEJ)⁴. BLESS captures DSBs *via* blunt-ended ligation of a barcoded adapter to broken DNA ends⁵,
27 a process less efficient than cohesive-end ligation⁶. Furthermore, sequential addition of the capture and
28 sequencing adapters necessitates two rounds of PCR, a technique known to introduce biases⁷. The
29 resulting libraries suffer from low sequence diversity (i.e. low variation in the initial bases of the
30 sequences) leading to problems during Illumina sequencing. This is normally remedied by diluting the
31 sample DNA with an unrelated library (typically phiX) to artificially increase the sequence diversity
32 however, this reduces data yield of the sample under investigation by over 50 %⁸. Finally, as
33 sequencing is not only initiated from the captured DSB (proximal end) but also from the end generated
34 through fragmentation (distal end), the number of sequencing reads that directly identify the site of
35 DNA damage in single-end sequencing is halved. To overcome these limitations we developed
36 DSBCapture, a substantially improved method derived from BLESS that enables direct *in situ* capture
37 of DSBs using a modified P5 Illumina adapter to facilitate sequencing without additional library
38 preparation steps (**Fig. 1a** and **Supplementary Fig. 1 a,b,c,d**). Ligation of a T-tailed modified P5
39 Illumina adapter to break sites identifies solely sites of *in situ* DSB formation in single-end sequencing
40 and, as DSBCapture libraries have high sequence diversity, no spike-in of another library is required
41 for sequencing. DSBCapture displays enhanced data yield and quality, higher reproducibility and
42 superior sensitivity compared to BLESS. In a head-to-head comparison DSBCapture identified 4.5-fold

43 more DSBs (84,946 compared to 18,816) in normal human epidermal keratinocytes (NHEK), providing
44 the most comprehensive DSB landscape in a normal human cell line.

45 We initially validated DSBCapture using controlled double-strand cleavage by EcoRV in
46 fixed HeLa nuclei. DSBCapture identified 93.7 % of the 430,897 EcoRV restriction sites in the human
47 genome as cleaved (**Supplementary Fig. 2a** and **Supplementary Table 1**). Furthermore, DSB
48 detection is not influenced by the chromatin state or the DNA sequence content. 0.5 % of the detected
49 EcoRV sites lie within open chromatin whilst 99.5 % lie in heterochromatin which corresponds to the
50 distribution of the predicted sites. The average GC content 100 bp either side of cleaved and uncleaved
51 EcoRV sites is comparable (37 % and 36 %, respectively) demonstrating no evident GC bias for DSB
52 detection.

53 To verify that DSBCapture can detect *in situ* generated DSBs, we performed DSBCapture on
54 AID-DIVa cells (**Supplementary Table 2**) - a U2OS cell line expressing the AsiSI restriction enzyme
55 fused to a modified estrogen receptor ligand-binding domain⁹. Upon 4-hydroxy tamoxifen (4OHT)
56 treatment, the restriction enzyme translocates to the nucleus causing sequence-specific DSBs at
57 GCGATCGC sites. The locations of AsiSI-induced DSBs have been evaluated by ChIP-seq using the
58 DNA damage markers γ H2AX, RAD51 and XRCC4⁹. A total of 121 AsiSI sites were identified by
59 DSBCapture, recovering 74 of the 100 ‘most cleaved’ sites previously identified by γ H2AX ChIP-seq⁹
60 (**Fig. 1b** and **Supplementary Fig. 2b**). DSBCapture also identified 47 additional AsiSI sites of which
61 28 (60 %) overlap with XRCC4 or RAD51 ChIP-seq peaks⁹, illustrating that DNA repair proteins
62 localize to these sites. The ability of DSBCapture to precisely map nuclease-induced DSBs could make
63 it useful for the study of CRISPR off-target cleavage sites.

64 In AsiSI-expressing U2OS cells, DSBCapture also recovered 2,372 endogenous DSBs
65 (lacking AsiSI sites). By comparison, GUIDE-seq identified 25 endogenous DSBs in U2OS cells⁴
66 (**Supplementary Fig. 3**). The substantially higher number of DSBs revealed by DSBCapture probably
67 results from the direct mechanism of DSB detection since DSBCapture directly captures DSBs *via*
68 ligation of an oligonucleotide to DNA ends, whereas GUIDE-seq relies on NHEJ to integrate a double-
69 stranded oligonucleotide at DSBs during DNA repair. Overall, this supports improved sensitivity and
70 genome-wide coverage for DSB detection by DSBCapture.

71 DSBCapture and BLESS were compared side-by-side to study endogenous DSBs in two
72 biological replicates of NHEKs. No DNA was recovered in three negative controls (DSBCapture
73 without T4 DNA ligase during the first or second ligation or using a non-biotinylated modified P5
74 Illumina adapter) (**Supplementary Fig. 1e**). After artificially enhancing BLESS data yield by repeated
75 sequencing to obtain similar read numbers to DSBCapture, we averaged both replicates and compared
76 the data quality of the two methodologies: 69 % of the DSBCapture data passed filtering versus 26 %
77 of the BLESS data (**Supplementary Table 3**). Overall, 84,946 common high confidence peaks were
78 observed for DSBCapture replicates versus 18,816 for BLESS (**Supplementary Fig. 4a**). **Fig. 1c**
79 exemplifies the increased signal of DSBCapture for two representative genes: *MAP2K3* and *MYC*.
80 Moreover, DSBCapture showed higher reproducibility: 83 % overlapping peaks between two
81 independent experiments versus 63 % in BLESS (**Supplementary Fig. 4a**). Of the DSBs identified
82 using BLESS, 99 % were observed by DSBCapture whereas 78 % of DSBs identified by DSBCapture

83 were missed by BLESS (**Fig. 1d, Supplementary Fig. 4b**). 73 % of the original peaks were also
84 identified when performing DSBCapture with reduced DNA input (20 μ g instead of 50 μ g), illustrating
85 good overlap irrespective of input material (**Supplementary Fig. 4c**). These results indicate that
86 DSBCapture has improved sensitivity and higher statistical confidence for DSB detection. On
87 comparing the sequence context of DSBs detected uniquely by DSBCapture to those shared between
88 BLESS and DSBCapture we found that DSBCapture detected a significantly greater proportion of
89 DSBs in genomic regions where the GC content surpassed 70 %. When the GC content of DSBs
90 exceeded 80 %, DSBCapture unique peaks were more than 4.5-fold enriched compared to those shared
91 with the BLESS dataset (**Supplementary Fig. 4d**). G-quadruplex DNA secondary structures form in
92 G-rich DNA in human cells¹⁰. G-quadruplexes have been implicated as fragile sites during transcription
93 and replication^{3,11,12}, and associate with somatic copy number alterations in human cells¹³. We found
94 that G-quadruplex-sites, previously mapped in genomic DNA (OQs)¹³, were 3-fold enriched over
95 random within high confidence DSBCapture peaks (**Supplementary Fig. 4e**), suggesting that G-
96 quadruplexes are intrinsic sites of DSB formation in NHEKs.

97 We next analysed high confidence DSBCapture data for NHEKs in the context of chromatin
98 features using ChIP-seq and DNase-seq datasets from ENCODE¹⁴. Notably, 58 % (11.4-fold
99 enrichment) of DSBs localized with histone H2A.Z, a histone variant transiently incorporated at DSBs
100 and known to have a role in DSB repair¹⁵ (**Fig. 2a, Supplementary Fig. 5a and Supplementary Table**
101 **4**). Over 76 % (33.3-fold enrichment) of the DSBCapture peaks overlapped with regulatory,
102 nucleosome-depleted regions (DNase-seq), revealing a relationship between regulatory chromatin and
103 genome instability (**Fig. 2a and Supplementary Table 4**). This is consistent with the notion that
104 nucleosome density mediates the susceptibility of DNA to genotoxic insults, with euchromatin showing
105 enhanced DNA repair signaling^{16,17}. DSBs also correlate with markers of active genes (mono-, di- and
106 tri-methylated H3K4); enhancer regions (H3K27ac and H3K4me1; **Supplementary Fig. 5b**), the
107 architectural protein CTCF and the transcription factor P63 (**Fig. 2a and Supplementary Table 4**). 38
108 % (12.2-fold enrichment) of DSBs overlapped with RNA polymerase II (POL2B) peaks, linking DSBs
109 to transcription (**Fig. 2a and Supplementary Table 4**). No enrichment of DSBs was found in
110 heterochromatic regions of the genome (H3K9me3 and H3K27me3), or in regions featuring the
111 transcriptional repressor EZH2 (**Fig. 2a and Supplementary Table 4**). Transcription has been linked to
112 DNA damage¹⁸ and DSBs have been found to occur in the proximity of transcription start sites (TSSs)
113 of highly expressed genes^{19,20}. Indeed, we found DSBs to be enriched in genic compared to intergenic
114 regions (**Supplementary Fig. 5c**), particularly at 5'UTRs (21.5-fold) and promoters (12.8-fold) (**Fig.**
115 **2b**). We analyzed the average number of DSBs present at TSSs (\pm 1 kb) and within gene bodies of
116 genes at different transcriptional levels. Increased gene expression correlated with increased DNA
117 damage around TSSs, whereas damage within gene bodies showed little association with gene
118 expression (**Fig. 2c**). Furthermore, we found that genes that do not contain DSBs at the TSS are
119 generally not expressed (**Supplementary Fig. 5d**). Four of the five most highly expressed genes in
120 NHEKs are keratin genes (*KRT5*, *KRT14*, *KRT6A* and *KRT17*); DSBCapture identified DSBs at the
121 TSS in all of these genes. In fact, DSBCapture identified DSBs at the TSS in 93 % of the most highly

122 expressed genes (top 5 %). Taken together, these data clearly link elevated gene expression and sites of
123 regulatory, nucleosome-depleted chromatin to DSB formation in NHEKs.
124
125 DSBCapture will be a valuable methodology for the study of DNA damage and repair, where
126 it will enable the elucidation of sites of endogenous and drug-induced DNA damage.

127 **Accession code**

128 The data reported in this paper are available at the NCBI's GEO repository, accession number
129 GSE78172 (<http://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE78172>)

130

131 **Acknowledgements**

132 We thank G. Legube, LBCMCP, Center for Integrative Biology (CBI), Université de Toulouse,
133 Toulouse, France for providing U2OS AID-DivA cells. We thank the genomic core facility at Cancer
134 Research UK Cambridge Institute. R.H.H. acknowledges EMBO for support (EMBO Long-Term
135 Fellowship to R.H.H.). We acknowledge support from the University of Cambridge and the Cancer
136 Research UK program. The Balasubramanian laboratory is supported by core funding from Cancer
137 Research UK (C14303/A17197 to S.B.) and by an ERC Advanced Grant (S.B.). S.B. is a Senior
138 Investigator of the Wellcome Trust.

139

140 **Author Contributions**

141 S.V.L. developed the DSBCapture method, conceived the study, conducted experiments, interpreted
142 results and wrote the manuscript. G.M. conceived the study, performed bioinformatics analyses,
143 interpreted results and wrote the manuscript. R.H.H. contributed to the development of the
144 DSBCapture method, conceived the study, conducted experiments, contributed to bioinformatics
145 analyses, interpreted results and wrote the manuscript. E.Y.L. contributed to the development of the
146 DSBCapture method, conceived the study and conducted experiments. D.T. conceived the study,
147 interpreted results and wrote the manuscript. S.B. conceived the study, interpreted results and wrote the
148 manuscript.

149

150 **Competing financial interests**

151 The authors declare no competing financial interests.

152 1 Srivastava, M. & Raghavan, S. C. DNA double-strand break repair inhibitors as cancer
153 therapeutics. *Chem. Biol.* **22**, 17-29 (2015).

154 2 Jackson, S. P. & Bartek, J. The DNA-damage response in human biology and disease. *Nature*
155 **461**, 1071-1078 (2009).

156 3 Rodriguez, R. *et al.* Small-molecule-induced DNA damage identifies alternative DNA
157 structures in human genes. *Nat. Chem. Biol.* **8**, 301-310 (2012).

158 4 Tsai, S. Q. *et al.* GUIDE-seq enables genome-wide profiling of off-target cleavage by
159 CRISPR-Cas nucleases. *Nat. Biotechnol.* **33**, 187-197 (2015).

160 5 Crosetto, N. *et al.* Nucleotide-resolution DNA double-strand break mapping by next-
161 generation sequencing. *Nat. Methods* **10**, 361-365 (2013).

162 6 Marchuk, D., Drumm, M., Saulino, A. & Collins, F.S. Construction of T-vectors, a rapid and
163 general system for direct cloning of unmodified PCR products. *Nucleic Acids Res.* **19**, 1154
164 (1990).

165 7 Aird, D. *et al.* Analyzing and minimizing PCR amplification bias in Illumina sequencing
166 libraries. *Genome Biol.* **12**, R18 (2011).

167 8 Mitra, A., Skrzypczak, M., Ginalski, K. & Rowicka, M. Strategies for achieving high
168 sequencing accuracy for low diversity samples and avoiding sample bleeding using illumina
169 platform. *PLoS One* **10**, e0120520 (2015).

170 9 Aymard, F. *et al.* Transcriptionally active chromatin recruits homologous recombination at
171 DNA double-strand breaks. *Nat. Struct. Mol. Biol.* **21**, 366-374 (2014).

172 10 Biffi, G., Tannahill, D., McCafferty, J. & Balasubramanian, S. Quantitative visualization of
173 DNA G-quadruplex structures in human cells. *Nat. Chem.* **5**, 182-186 (2013).

174 11 Ribeyre, C. *et al.* The yeast Pif1 helicase prevents genomic instability caused by G-
175 quadruplex-forming CEB1 sequences in vivo. *PLoS Genet.* **5**, e1000475 (2009).

176 12 Paeschke, K., Capra, J. A. & Zakian, V. A. DNA replication through G-quadruplex motifs is
177 promoted by the *Saccharomyces cerevisiae* Pif1 DNA helicase. *Cell* **145**, 678-691 (2011).

178 13 Chambers, V. S. *et al.* High-throughput sequencing of DNA G-quadruplex structures in the
179 human genome. *Nat. Biotechnol.* **33**, 877-881 (2015).

180 14 Consortium, E. P. An integrated encyclopedia of DNA elements in the human genome. *Nature*
181 **489**, 57-74 (2012).

182 15 Gursoy-Yuzugullu, O., Ayrapetov, M. K. & Price, B. D. Histone chaperone Anp32e removes
183 H2A.Z from DNA double-strand breaks and promotes nucleosome reorganization and DNA
184 repair. *Proc. Natl. Acad. Sci. USA* **112**, 7507-7512 (2015).

185 16 Storch, K. *et al.* Three-dimensional cell growth confers radioresistance by chromatin density
186 modification. *Cancer Res.* **70**, 3925-3934 (2010).

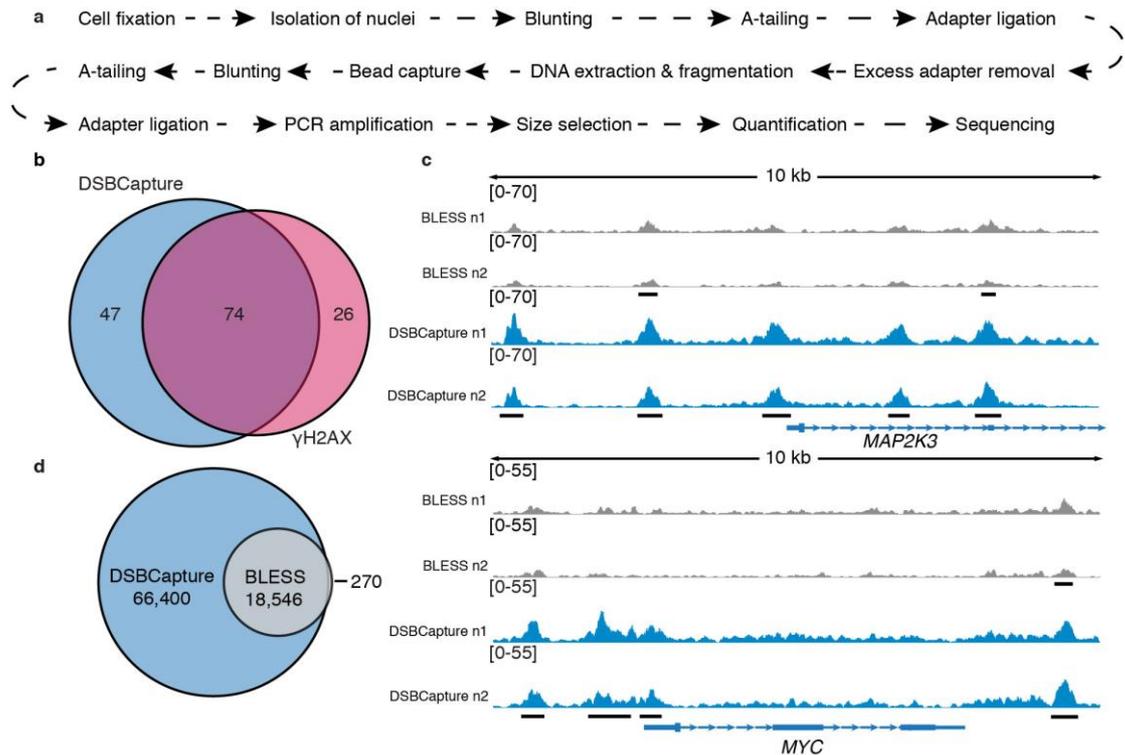
187 17 Misteli, T. & Soutoglou, E. The emerging role of nuclear architecture in DNA repair and
188 genome maintenance. *Nat. Rev. Mol. Cell. Biol.* **10**, 243-254 (2009).

189 18 Fong, Y. W., Cattoglio, C. & Tjian, R. The intertwined roles of transcription and repair
190 proteins. *Mol. Cell* **52**, 291-302 (2013).

191 19 Yang, F., Kemp, C. J. & Henikoff, S. Anthracyclines induce double-strand DNA breaks at
192 active gene promoters. *Mutat. Res.* **773**, 9-15 (2015).

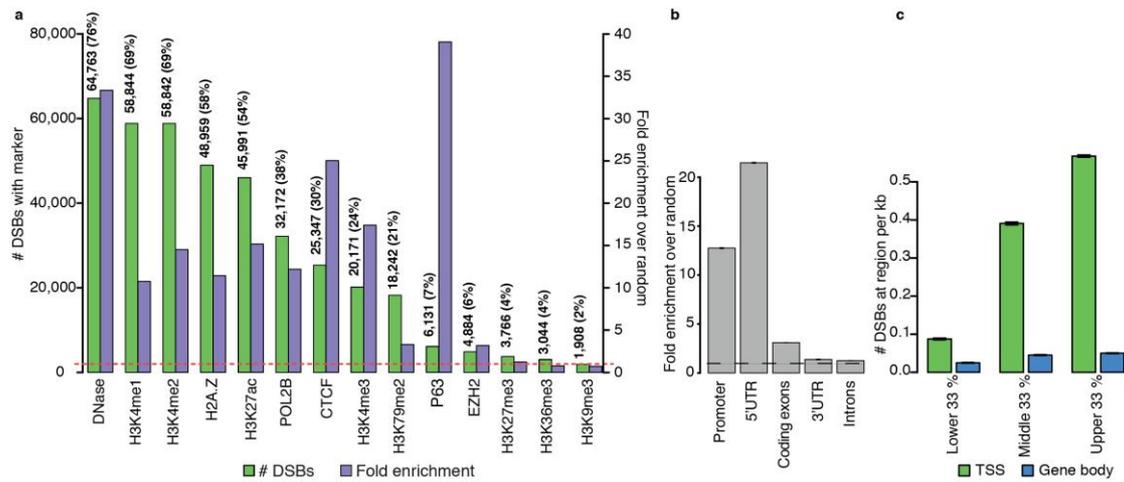
193 20 Schwer, B. *et al.* Transcription-associated processes cause DNA double-strand breaks and
194 translocations in neural stem/progenitor cells. *Proc. Natl. Acad. Sci. USA* **113**, 2258-2263
195 (2016).

196
197
198
199
200



201
 202
 203
 204
 205
 206
 207
 208
 209
 210

Figure 1 DSBCapture methodology and comparison to BLESS (a) DSBCapture workflow. (b) Venn diagram illustrating the overlap of DSBs detected at AsiSI sites by DSBCapture and γ H2AX ChIP-seq⁹. (c) Representative genomic view of DSBs detected by BLESS (grey) and DSBCapture (blue) in two biological replicates. Peaks in common between the two replicates for each method are underlined with black bars; the data range (absolute read counts) is shown in square brackets. Two 10 kb genomic regions in proximity of the *MAP2K3* and *MYC* genes are shown. (d) Venn diagram depicting the overlap of high confidence DSBCapture and BLESS peaks.



211
 212
 213
 214
 215
 216
 217
 218
 219
 220
 221
 222
 223

Figure 2 Genomic location and epigenetic context of endogenous DSBs in NHEK cells **(a)** Composite bar plot showing the overlap of high confidence DSBs detected by DSBCapture with epigenetic marks and DNA-binding proteins taken from ENCODE¹⁴. Green bars: number of DSBs overlapping each mark, ordered from highest (left) to lowest (right). Purple bars: fold change of DSBs overlapping each mark over random, a signal above 1 (dashed red line) is indicative of enrichment. The number and percentages of DSBs that overlap with each mark are shown above their respective bars. **(b)** Fold enrichment of high confidence DSBs in different genomic regions, calculated as the number of peaks divided by the number of randomly shuffled peaks overlapping to each region. Error bars: standard deviation of fold enrichment over random. **(c)** Number of high confidence DSBs detected at TSSs (green bars) or within gene bodies (blue bars) for the lower, middle and upper third of gene expression values, split according to the rpkm values. Error bars: SEM; n = 7,430 for each category.

224 **Online Methods**

225

226 **Cell culture**

227 HeLa cells (CRUK-CI Biorepository and Cell Services Core) were cultured in DMEM (Sigma, D6429)
228 supplemented with 10 % heat inactivated FBS. U2OS AID-DIVa cells were cultured in DMEM
229 supplemented with 10 % heat inactivated FBS and 800 µg/mL G418 Sulfate (Gibco, 10131). To induce
230 nuclear localization of AsiSI, AID-DIVa cells were treated with 300 nM 4OHT (Sigma, H7904) for 4
231 h. Normal human epidermal keratinocytes (NHEK), pooled from multiple donors, were purchased from
232 Thermofisher (A13401) and cultured in EpiLife medium (Thermofisher, M-EPI-500-CA)
233 supplemented with human keratinocyte growth supplement (HKGS) (Thermofisher, S-001-5). NHEKs
234 were detached using accutase (Sigma, A6964). HeLa and U2OS AID-DIVa cells were STR genotyped
235 and mycoplasma tested. NHEK cells were obtained from Thermofisher and were certified as
236 mycoplasma negative.

237

238 **BLESS**

239 The method published on <http://breakome.utmb.edu> was followed; with the exception that DNA was
240 solely fragmented by sonication and amplified using Phusion HF polymerase (NEB, M0530S) as
241 published by Crosetto, *et al*⁵. The TruSeq Nano DNA Library Prep Kit (Illumina, FC-121-4001) was
242 used to generate libraries for sequencing (size selection for 550 bp) and all centrifugation steps were
243 carried out at the speeds described for DSBCapture in the step by step method on the protocol
244 exchange (DOI: 10.1038/protex.2016.52).

245

246 **Annealing of oligonucleotides**

247 To anneal the modified P5 and P7 Illumina adapters; oligonucleotides (modified P5 and modified P5
248 complement / modified P7 and modified P7 complement) were made up to 10 µM in 1 × T4 DNA
249 ligase reaction buffer (NEB, B0202S). Oligonucleotides were heated to 95 °C for 10 minutes. Tubes
250 were removed from the heat source and gradually cooled to room temperature.

251

252 **DSBCapture (EcoRV cleavage in HeLa cells)**

253 During the development of the DSBCapture method, pilot experiments were performed in which
254 permeabilized nuclei were treated with the restriction enzyme EcoRV to generate DSBs. This pilot
255 method was essentially comparable to the final DSBCapture protocol described below with the
256 following significant differences: Cells were fixed in formaldehyde for 10 min. After re-suspension in
257 nucleus break buffer the nuclei were washed 2 × with 1 × NEBuffer 4 (NEB, B7004S) + 1 % TritonX-
258 100 followed by re-suspension in 500 µL 1 × CutSmart buffer (NEB, B7204S). 500 units of EcoRV
259 (NEB, R0195T) were added and the nuclei were incubated over night at 37 °C with shaking at 950 rpm.
260 The 8 min Proteinase K digest was performed at a final concentration of 50 µg/mL. The first blunting
261 step was missed out before A-tailing, as EcoRV is a blunt cutting enzyme. No treatment with Lambda
262 Exonuclease was performed before subsequent DNA isolation. 20 µg of extracted DNA was bound to
263 beads. PCR reactions was performed each using 5 µL beads, 1 µL 20 µM PCR F primer, 1 µL 20 µM

264 PCR R primer, 10 μ L 5 \times Phusion HF buffer, 1 μ L 10 mM dNTPs and 0.5 μ L Phusion HF DNA
265 polymerase (NEB, M0530S) in a final reaction volume of 50 μ L with the following cycling parameters:
266 94 $^{\circ}$ C 5 min, 94 $^{\circ}$ C 1 min, 60 $^{\circ}$ C 45 s, 72 $^{\circ}$ C 1 min, 72 $^{\circ}$ C 10 min, 4 $^{\circ}$ C hold, repeating steps 2-4, for a
267 total of 18 cycles. The DNA was size selected (400-500 bp) by gel electrophoresis and was
268 subsequently extracted using the MinElute gel extraction kit (Qiagen, 28604). The DNA library was
269 subjected to paired-end sequencing on an Illumina MiSeq platform.

270

271 **DSBCapture (U2OS and NHEK):** A detailed protocol can be found on the protocol exchange (DOI:
272 10.1038/protex.2016.52). Briefly: cells were detached, counted and fixed in complete medium with 2
273 % formaldehyde (Pierce, 28908) at a density of 1 million cells/1.5 mL for 30 min at room temperature
274 whilst gently rotating. Formaldehyde was quenched by adding glycine to a final concentration of 125
275 mM. Cells were washed twice in ice cold PBS and lysed in lysis buffer (10 mM Tris-HCl pH 8, 10
276 mM NaCl, 1 mM EDTA, 1 mM EGTA, 0.2 % NP40 substitute (Sigma, 74385), cOmplete Roche
277 proteinase inhibitors (REF 11873580001), 1 mM DTT) by gently rotating at 4 $^{\circ}$ C for 90 min. Nuclei
278 were subsequently re-suspended in nucleus break buffer (10 mM Tris-HCl pH 8, 150 mM NaCl, 1 mM
279 EDTA, 1 mM EGTA, 0.3 % SDS, 1 mM DTT) and incubated at 37 $^{\circ}$ C for 45 min whilst gently
280 rotating. Nuclei were then re-suspended in 1 \times NEBuffer 2 + 0.1 % TritonX-100 (10 mM Tris-HCl, 50
281 mM NaCl, 10 mM MgCl₂, 1 mM DTT, pH 7.9 at 25 $^{\circ}$ C) and transferred into 2 mL eppendorf tubes (10
282 million nuclei/tube). Proteinase K (Ambion, AM2546) was added to a final concentration of 100
283 μ g/mL on ice and nuclei were incubated at 37 $^{\circ}$ C for 8 min before adding an equal volume of 1 \times
284 NEBuffer 2 + 0.1 % TritonX-100 + 1:50 PMSF (Sigma, 93482) to kill the reaction. Nuclei were then
285 washed twice in 1 \times NEBuffer 2 + 0.1 % TritonX-100 at 4 $^{\circ}$ C. To repair DNA ends, nuclei were
286 washed once in 1 \times Blunting Buffer (NEB, E1201L) + 100 μ g/mL BSA (NEB, B9000S) and were
287 blunt-end repaired using the NEB Quick Blunting Kit (NEB, E1201L) + 100 μ g/mL BSA in a final
288 volume of 100 μ L at 25 $^{\circ}$ C for 45 min, shaking for 10 s at 800 rpm, every 5 min. Nuclei were then
289 washed three times with 1 \times NEBuffer 2 + 0.1 % TritonX-100 at 4 $^{\circ}$ C and were subsequently A-tailed
290 using Klenow Fragment 3'-5' exo- (NEB, M0212L) and dATP (Promega, U120D) in a final volume of
291 50 μ L at 37 $^{\circ}$ C for 45 min, shaking for 10 s at 800 rpm, every 10 min. Following A-tailing, nuclei were
292 washed 3 \times with 1 \times NEBuffer 2 + TritonX-100 at 4 $^{\circ}$ C, once with T4 DNA ligase buffer (NEB,
293 B0202S) + 0.1 % TritonX-100 at 4 $^{\circ}$ C and once with T4 DNA ligase buffer at 4 $^{\circ}$ C. The modified P5
294 Illumina adapter, previously annealed in 1 \times T4 DNA ligase buffer, was then ligated to DNA ends
295 using T4 DNA ligase (NEB, M0202M) for 15-20 h at 16 $^{\circ}$ C in a final reaction volume of 50 μ L whilst
296 shaking at 350 rpm for 15 s every 45 min during the ligation. To remove excess adapters, nuclei were
297 washed 2 \times in WB buffer + 0.1 % TritonX-100 (5 mM Tris-HCl pH 7.5, 1 mM EDTA, 1 M NaCl) and
298 once in 1 \times Lambda Exonuclease Reaction Buffer (NEB, M0262L) and were then treated with 50 units
299 Lambda Exonuclease (NEB, M0262L) in a final reaction volume of 50 μ L for 30 min at 37 $^{\circ}$ C.
300 Following washing in 1 \times NEBuffer 2 + 0.1 % TritonX-100 nuclei were then re-suspended in 1 \times
301 NEBuffer 2 + 0.5 % TritonX-100 at 4 $^{\circ}$ C. To extract genomic DNA, nuclei were lysed with Proteinase
302 K (200 μ g/ml) for 30 min at 55 $^{\circ}$ C, shaking at 800 rpm followed by 30 min at 65 $^{\circ}$ C, shaking at 800
303 rpm and the DNA was precipitated using isopropanol. The DNA was then re-suspended in 90 μ L H₂O

304 and incubated at 55 °C shaking at 800 rpm for 1 h. The isolated DNA was fragmented by sonication to
305 obtain an average fragment size of 200-500 bp and then 50 µg (20 µg only for data depicted in
306 **Supplementary Fig. 4c**) of DNA was bound to 5 µl MyOne streptavidin C1 Dynabeads (Invitrogen,
307 65001) at 4 °C for 45 min, whilst gently rotating. Beads were washed 3 × with WB buffer + 0.1 %
308 TritonX-100 and captured DNA was blunt end-repaired using the Quick Blunting Kit in a final reaction
309 volume of 50 µL for 45 minutes at 25 °C, shaking at 800 rpm every 5 min for 10 s. The beads were
310 then washed 3 × in WB buffer + 0.1 % TritonX-100 and DNA ends were once more A-tailed using
311 Klenow Fragment 3'-5' exo- in a final reaction volume of 25 µL at 37 °C for 45 min, shaking at 800
312 rpm for 10 s every 10 min. After washing beads 3 × with WB buffer + 0.1 % TritonX-100, the
313 modified P7 Illumina adapter, previously annealed in 1 × T4 DNA ligase buffer, was ligated to DNA
314 ends in a final reaction volume of 50 µL for 15-20 h at 16 °C, every 45 minutes samples were mixed
315 for 1 min at 1200 rpm. Beads were then washed 3 × in WB buffer + 0.1 % TritonX-100 before re-
316 suspending in 25 µL nuclease free water. The DNA DSBCapture library was then amplified by PCR
317 with PCR F and PCR R primers and NEB Next PCR mix for 15 cycles, following the manufacturers
318 recommendations (annealing temperature: 65 °C; NEB, M0541L). The amplified DNA library was
319 then purified using a MinElute PCR Purification kit (Qiagen, 28004) and the library was size-selected
320 (250-1200 bp) using a BluePippin (Sage Science). Samples were run on the Bioanalyzer (Agilent,
321 5067-4626) to determine the average library size, and the library was quantified using the KAPA
322 library quantification kit (Kapa Biosystems, kk4824), following the manufacturers recommendations.
323 Libraries were subsequently sequenced paired-end on an Illumina NextSeq 500 platform.

324

325 **RNA-seq**

326 Total RNA for RNA-seq experiments was extracted using the RNeasy kit (Qiagen, cat. no. 74104),
327 following the manufacturer's instructions. RNA-seq libraries were generated using the Illumina Truseq
328 RNA HT (stranded mRNA) kit (RS-122-2103).

329

330 **Epigenome analysis**

331 The ENCODE project NHEK epigenome ChIP-seq data¹⁴ was retrieved from the NCBI's GEO
332 repository as follows: H3K4me3 (GSM945175); H3K4me2 (GSM733686); H3K27ac (GSM733674);
333 H3K4me1 (GSM733698); H3K79me2 (GSM1003527); H3K36me3 (GSM945174); H2A.Z
334 (GSM1003488); POLR2B (GSM733671); H3K9me3 (GSM1003528); H3K27me3 (GSM733701);
335 EZH2 (GSM1003489); CTCF (GSM822271); TP63 (GSE32061), DNase-seq (GSM736556).

336

337 **Sequencing**

338 For the EcoRV enzyme cleavage experiment, 50 bp paired-end sequencing was conducted on the
339 Illumina MiSeq instrument. For the AID-DivA U2OS and NHEK cell experiments, 75 bp paired-end
340 sequencing was performed on the Illumina NextSeq 500 instrument. Despite DSB detection with
341 DSBCapture requiring only single-end sequencing, we sequenced all DSBCapture libraries in paired-
342 end mode in order to be able to distinguish DSBs at restriction sites from PCR duplicates and to be able
343 to compare our method to BLESS in NHEK cells.

344

345 **Sequencing analysis**

346 DSBCapture libraries: fastq files containing sequencing reads were pre-processed to remove the
347 Illumina adapters and to trim low quality tails using trim_galore
348 (http://www.bioinformatics.babraham.ac.uk/projects/trim_galore/); reads were aligned to the human
349 reference genome (*hg19*) using the bwa mem aligner (<http://sourceforge.net/projects/bio-bwa/files/>)
350 and cleaned for low quality alignments (mapQ < 10) using samtools (clean reads)
351 (<http://samtools.sourceforge.net>). Duplicates were identified using picard tools
352 (<https://github.com/broadinstitute/picard/>) and removed (unduplicated reads). Only reads from read 1,
353 which are proximal to the break site, were retained for downstream analysis and are shown in the
354 coverage plots.

355 BLESS libraries: when performing the Illumina adapter trimming, the BLESS linker sequences are
356 trimmed (TCGAGGTAGTA and TCGAGACGACG for proximal and distal linkers, respectively), and
357 only read pairs containing both linkers are retained. Trimmed fastq files then undergo the same
358 processing pipeline as for DSBCapture libraries. Finally, only sequencing reads that originally
359 contained the proximal linker (originating either from read 1 or read 2) are retained for subsequent
360 analysis and are shown in the coverage plots.

361

362 **Peak calling analysis**

363 Regions with enriched reads over the background were called as peaks using the MACS2 software
364 (<https://github.com/taoliu/MACS/>) using default parameters for the statistical threshold (i.e., corrected
365 p-value $q \leq 0.05$) and the options “no-model” and “no control”. For the NHEK cell analysis, peaks
366 were called independently for each replicate on the clean bam files containing only the proximal
367 linker, and the two peaks files were then intersected using the bedtools package
368 (<http://bedtools.readthedocs.org/>) to generate the high-confidence peak files for both BLESS and
369 DSBCapture. The sample size was constrained by experimental considerations. For the endogenous
370 DSBs mapping in primary NHEK cells, we observed good reproducibility with 2 biological replicates
371 for both BLESS and DSBCapture.

372

373 **Comparison of DSBs identified by DSBCapture with 50 µg and 20 µg of input material**

374 DSBCapture was performed using 50 µg and 20 µg of input material from the same biological sample
375 ($n = 1$). The two DSBCapture experiments were processed in the same way, as previously described in
376 the paragraphs "Sequencing analysis" and "Peak calling analysis". The 20 µg library was sub-sampled
377 in order to have the same number of unduplicated, proximal reads as the 50 µg library, i.e. 124 million
378 reads. Peak calling was performed on the subsampled library. Peak intersection was calculated by using
379 the bedtools, as previously explained.

380

381 **GC content analysis of the EcoRV data set**

382 The 430,897 predicted EcoRV sites were split into two subsets: 403,905 (94 % of total) sites displaying
383 a coverage of at least 5 reads, marked as detected, and the remaining 26,992 sites (6 % of the total)

384 displaying coverage below 5, marked as not detected. The GC content around the EcoRV restriction
385 site was calculated, by considering 100 base pairs either side of the restriction site.

386

387 **Analysis of EcoRV restriction sites at DNase hypersensitivity sites**

388 DNase hypersensitivity sites in HeLa cells were downloaded from the ENCODE project
389 (<https://www.encodeproject.org/>; GEO sample accession numbers GSM736564 and GSM736510),
390 high confidence sites were calculated as the intersection of the genomic regions from the two replicates
391 (bedtools intersect), resulting in 96,541 DNase sites. The overlap of the 430,897 predicted EcoRV sites
392 and the cleaved subset (403,905; coverage ≥ 5 reads) with the high confidence DNase sites was
393 calculated. The percentage of total and cleaved sites within DNase sites is reported. The non-
394 parametric Chi-squared test for proportions was used for statistical testing.

395

396 **AsiSI-induced DSB analysis in AID-DIV A U2OS cells**

397 AsiSI restriction sites were identified as DSBs by peak calling as described above. The list of 100 most
398 cleaved AsiSI sites identified by γ H2AX ChIP-seq signal, and the RAD51 and XRCC4 ChIP-seq data,
399 were obtained from Aymard, *et al*⁹. Peaks on RAD51 and XRCC4 were called using the MACS2 peak
400 caller with default parameters ($q \leq 0.05$) using the respective input files.

401

402 **Comparison of DSBCapture to GUIDE-seq**

403 The 25 endogenous DSB hotspots identified in U2OS cells by GUIDE-seq⁴ were compared (bedtools
404 intersect) to the 2,372 endogenous DSBCapture peaks identified by MACS2 in AID-DIV A U2OS cells.

405

406 **GC content and G-quadruplex analysis**

407 The high confidence ($\sim 85,000$) peaks detected in both replicates of DSBCapture were split into two
408 subsets: DSBs uniquely identified by DSBCapture ($\sim 66,000$ peaks) and peaks common to BLESS (\sim
409 19,000). GC content was measured within the peak regions and each peak was assigned to a GC %
410 category. Independently, for both the DSBCapture unique and BLESS-DSBCapture common peak
411 subsets, the number of peaks in each GC % category were divided by the total number of peaks in the
412 subset, resulting in a fraction of peaks within each GC % category for both subsets. The two fractions
413 of each GC % category were divided one by each other to obtain the fold enrichment of DSBCapture
414 unique peaks over the BLESS-DSBCapture common peaks. Values > 1 indicate that a given GC %
415 category is enriched among the DSBCapture unique peaks, or otherwise is depleted in the BLESS-
416 DSBCapture common peaks. Similarly for the G-quadruplex analysis, the fraction of reads from the
417 DSBCapture high confidence peaks overlapping to observed G-quadruplex-forming sequences (OQs¹³)
418 was calculated for each GC % content group and compared to random (fold enrichment). The random
419 sets were generated by random shuffling of DSB genomic intervals throughout the entire genome and
420 calculating the random overlap to OQs ($n = 3$).

421

422 **Comparison to the ENCODE data sets in NHEK cells**

423 Data sets for epigenetic marks (histone modifications), DNase hypersensitivity sites (DNase) and DNA
424 binding proteins (CTCF, P63, POL2B) were downloaded from ENCODE (see Epigenome analysis,
425 Methods). The number of high confidence DSBCapture peaks overlapping each mark (command
426 intersectBed of the bedtools package) was calculated for each mark (independently of all other marks)
427 and visualized on a bar plot (green bars in **Fig. 2a** and **Supplementary Fig. 5b**). Additionally, the
428 DSBCapture peak files were randomly shuffled across the genome (command shuffleBed of the
429 bedtools package) three independent times and the overlap of the random sets with each mark was
430 calculated. The ratio of the DSB overlap divided by the average ($n = 3$) random overlap for each mark
431 was computed (fold enrichment over random) and visualized (purple bars in **Fig. 2a** and
432 **Supplementary Fig. 5b**). Fold enrichment values > 1 indicates that a mark overlaps to DSBs more
433 often than random. Standard deviations from the 3 independent randomizations are not indicated in the
434 bar plot as they lay in the range 0.004-0.12 for all marks, and are therefore negligible when compared
435 to the fold change values.

436

437 **Genomic regions analysis**

438 Gene annotation files for the human genome (*hg19*) were downloaded from the Illumina iGenomes
439 support website (https://support.illumina.com/sequencing/sequencing_software/igenome.html, and
440 different gene features were calculated as follows: Promoter = 1 kb upstream of the transcription start
441 site (TSS); 5'UTR = sequence from the TSS to the annotated translation start codon; Coding exons =
442 all exons that are translated (coding sequences) from the start until the stop codon; 3'UTR = sequence
443 from the annotated end of translation (stop codon) to the end of the last exon; introns = all the regions
444 spliced out during transcript processing. The number of high confidence DSB sites overlapping each
445 region was calculated by intersecting DSBCapture high confidence peaks to each feature separately.
446 The DSBCapture peak files were then randomly shuffled across the genome (command shuffleBed of
447 the bedtools package) three independent times and the overlap of the random sets to each feature was
448 calculated. The ratio of the DSB overlap divided by the average ($n = 3$) random overlap for each region
449 was calculated (fold enrichment over random). Similarly, the fold enrichment over random for genic
450 regions (i.e. within gene bodies) versus intergenic features was assessed.

451

452 **Gene expression analysis**

453 RNA-seq sequencing reads were pre-processed by the trim galore software
454 (http://www.bioinformatics.babraham.ac.uk/projects/trim_galore/) for removal of Illumina sequencing
455 adapters and low quality read tails. Trimmed data were then aligned to the human reference genome
456 (*hg19*) using tophat2 (<https://ccb.jhu.edu/software/tophat/index.shtml/>). Reads overlapping
457 unambiguously to each gene feature were assessed and counted by the htseq-count ([http://www-
458 huber.embl.de/users/anders/HTSeq/](http://www-huber.embl.de/users/anders/HTSeq/)). Gene expression values were split into three groups containing
459 the same number of genes (lower, middle and upper third) according to their rpkm gene expression
460 values (rpkm < 0.041 ; rpkm 0.041-6.08 and rpkm > 6.08), where rpkm is the reads per kilo base (of the
461 total exon length) per million sequenced reads. The average number of high confidence DSBs per kb
462 overlapping ± 1 kb of the TSS (green bars in **Fig. 2c**) and to gene bodies, i.e. all exons and introns

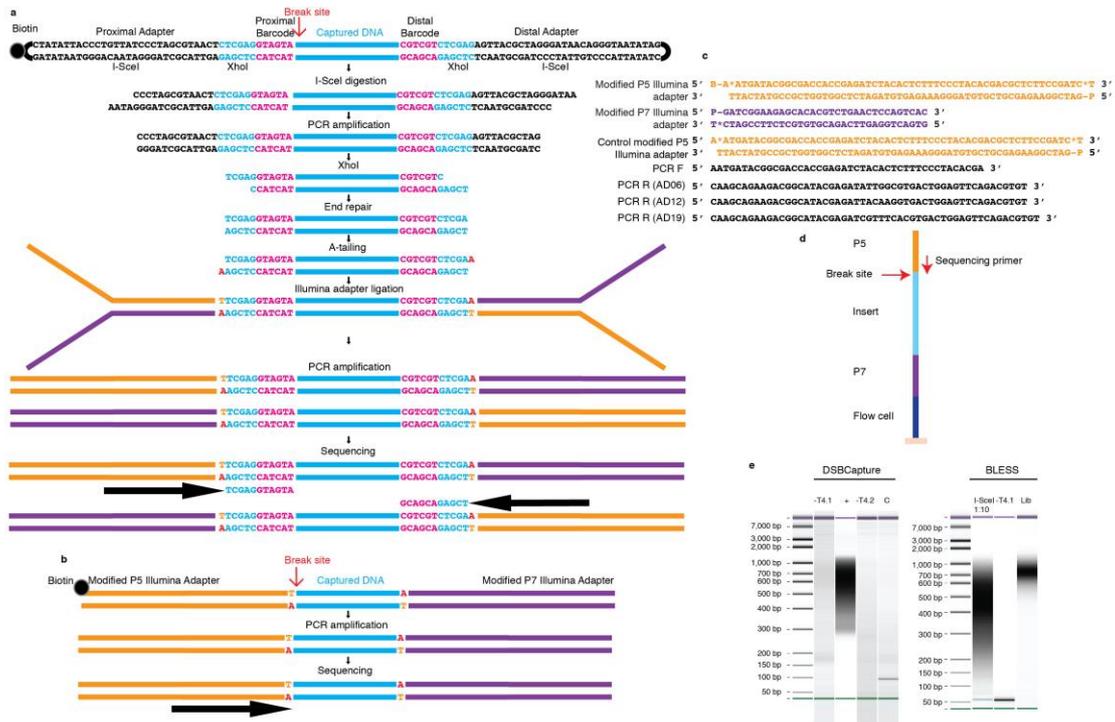
463 excluding ± 1 kb around the TSS (Gene body, blue bars in **Fig. 2c**), were calculated. Genes with length
464 less than 2 kb were excluded from the gene body category. Additionally, all genes were split into two
465 categories according to the presence ($n = 12,984$) / absence ($n = 9,272$) of DSBs at ± 1 kb around the
466 TSS and the rpkm for all genes in the two groups was inspected and visualized as box plots.

467

468 **Code availability**

469 The computer code used to analyze all data in this work, including sequencing processing and data
470 comparisons, can be obtained on request by contacting the corresponding author.

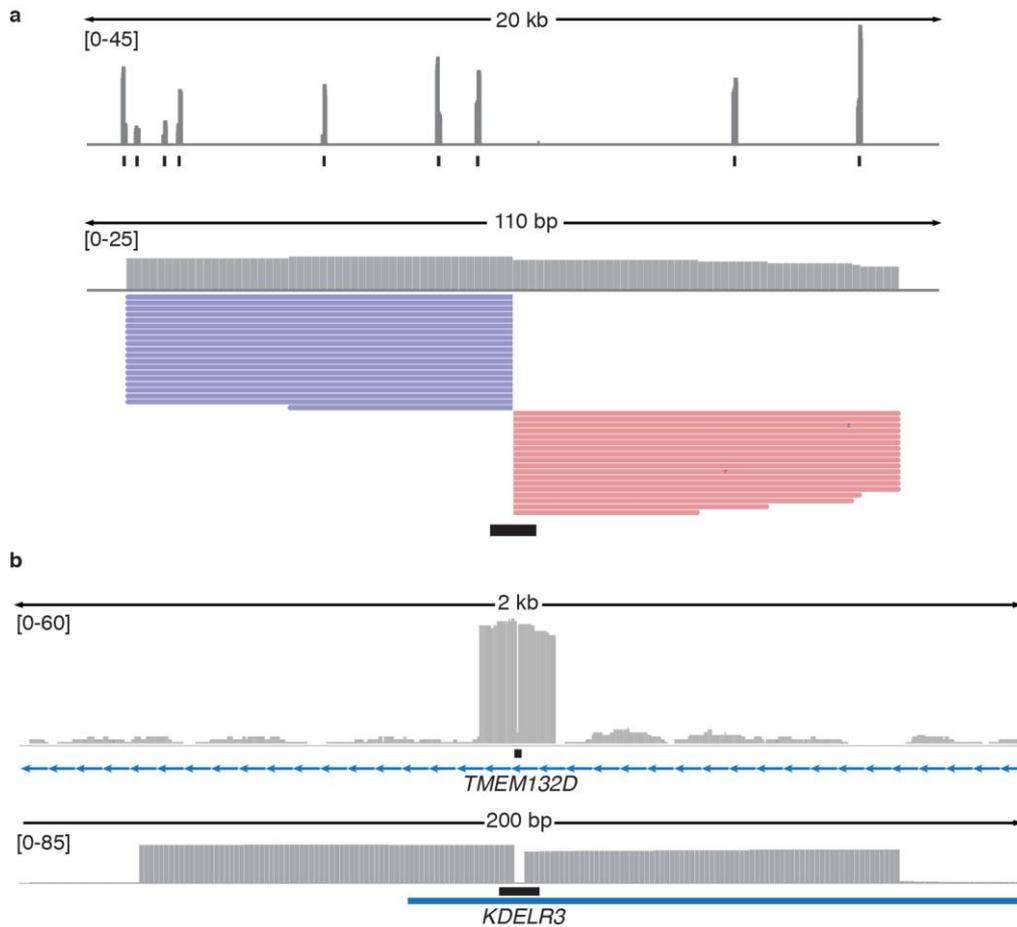
471



472

473 **Supplementary Figure 1**

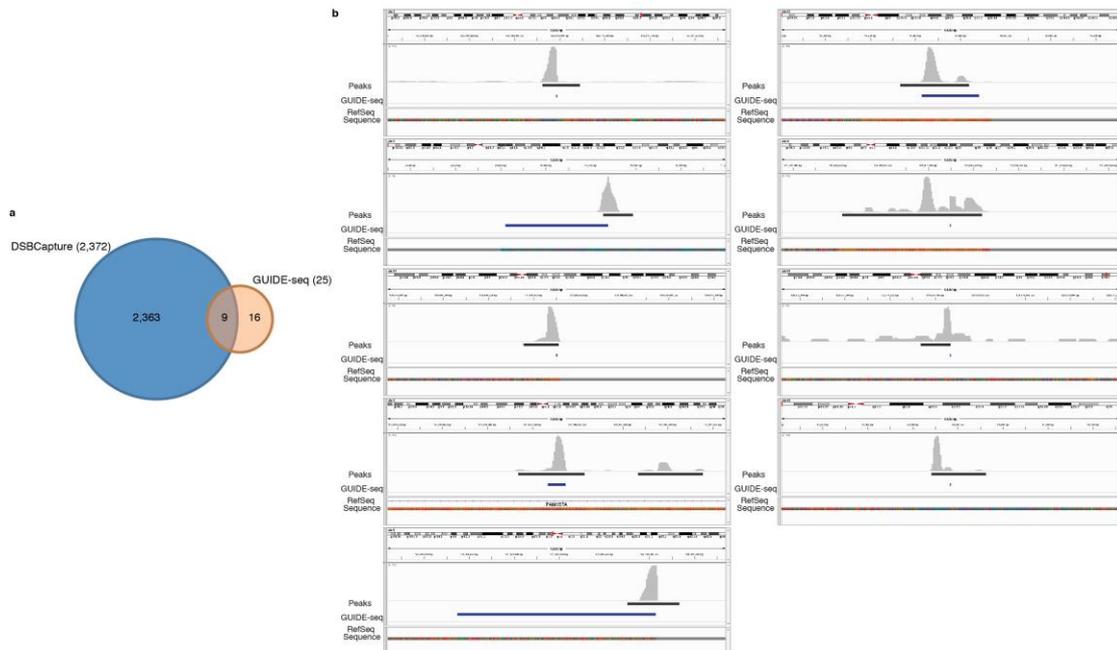
474 DNA processing workflows, adapter sequences and controls. **(a)** BLESS DNA processing workflow
 475 following the ligation of both proximal and distal adapters. DNA is digested by I-SceI, PCR amplified,
 476 digested by XhoI and subsequently subjected to Illumina library preparation, consisting of end repair,
 477 size selection (not shown), A-tailing, Illumina adapter ligation and PCR amplification. Large black
 478 arrows indicate the site at which sequencing is initiated; the first 11 bases sequenced are shown. **(b)**
 479 DSBCapture DNA processing workflow following the ligation of both modified P5 and P7 Illumina
 480 adapters. The DNA is PCR amplified, size selected (not shown), and sequenced. Large black arrow
 481 indicates the site at which sequencing is initiated; the first base sequenced is the site of *in situ* DSB
 482 formation. **(c)** Sequences of the modified P5, modified P7 and control modified P5 Illumina adapters as
 483 well as DSBCapture PCR primers (forward (PCR F) and reverse (PCR R)). AD identifies the Illumina
 484 adapter barcode sequence, three example reverse primers are shown; further primers can be created by
 485 substituting the barcode sequence. B = biotin; P = phosphorylated; * = phosphorothioate bond. **(d)**
 486 Orientation of the DSBCapture library on the Illumina flow cell. The first sequencing primer has
 487 complementarity to the P5 Illumina adapter and therefore sequencing is initiated from the P5 end. The
 488 ligation of the modified P5 Illumina adapter to the DSB *in situ* enables direct sequencing of the break
 489 site in single-end sequencing. The first base sequenced directly identifies the DSB. **(e)** Bioanalyser
 490 profiles of the DNA products from DSBCapture and BLESS NHEK libraries. DSBCapture: no product
 491 is present in the controls performed without T4 DNA ligase during the first (-T4.1) or second (-T4.2)
 492 ligation reactions, or in the control performed with the non-biotinylated control modified P5 Illumina
 493 adapter (C). A DSBCapture library was only generated when the complete procedure was carried out
 494 (+). BLESS: No product is present in the control performed without T4 DNA ligase during the first
 495 ligation reaction (-T4.1). The product of BLESS is shown before Illumina library preparation (I-SceI;
 496 diluted 1:10) and after Illumina library preparation (Lib).



497
 498
 499
 500
 501
 502
 503
 504
 505
 506
 507
 508
 509
 510

Supplementary Figure 2

DSBs mapped by DSBCapture at EcoRV and AsiSI restriction sites. **(a)** DSBs created by EcoRV cleavage in fixed nuclei, mapped by DSBCapture (n = 1). PCR duplicates have been removed. Data range is shown in square brackets and black boxes illustrate the genomic location of EcoRV sites. A 20 kb region and a 110 bp region are shown. Pink and purple lines: reads from the sense and antisense strand, respectively. As EcoRV is a blunt cutter, reads originate directly from the cleavage site. **(b)** AsiSI cleavage sites (black boxes) detected by DSBCapture (n = 1). Cleavage by AsiSI generates a 2 bp 3' overhang; end processing removes this overhang generating the 2 bp gap in the center of the peak. A 2 kb and a 200 bp region are shown.



511

512

513 **Supplementary Figure 3**

514 Overlap of DSBs detected by DSBCapture and GUIDE-seq in U2OS cells. **(a)** Venn diagram showing

515 the overlap between the DSBCapture peaks and the 25 sites detected by GUIDE-seq⁴. **(b)** Genomic

516 tracts showing the 9 DSBs detected by GUIDE-seq that are also detected by DSBCapture. Each panel

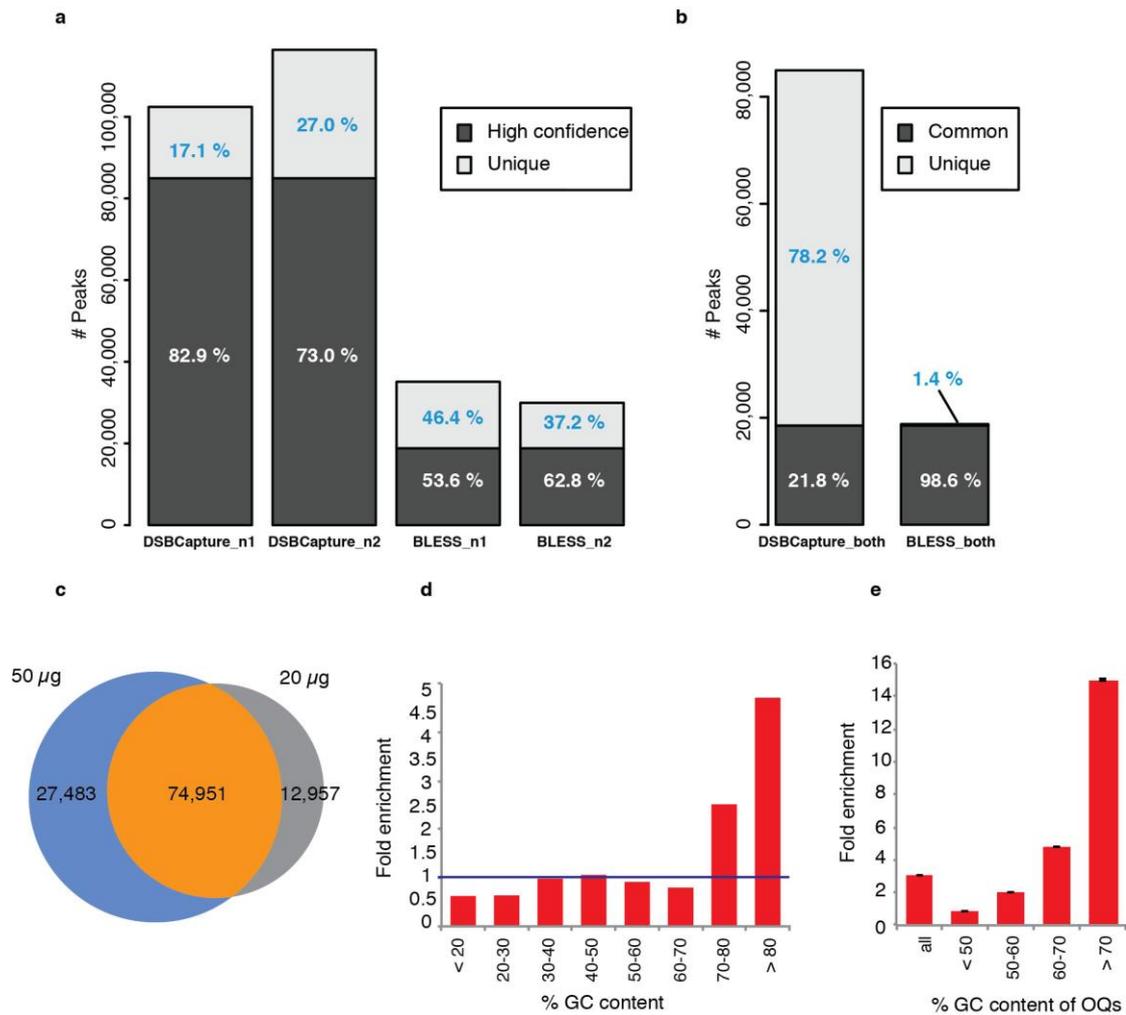
517 shows a genomic view of 2,000 bp around the GUIDE-seq detected DSB hotspot. In each panel, from

518 top to bottom: DSBCapture coverage (grey track); peaks detected in DSBCapture by peak calling

519 (black track); GUIDE-seq sites (dark blue track); RefSeq gene track; reference genome sequence

520 (*hg19*).

521



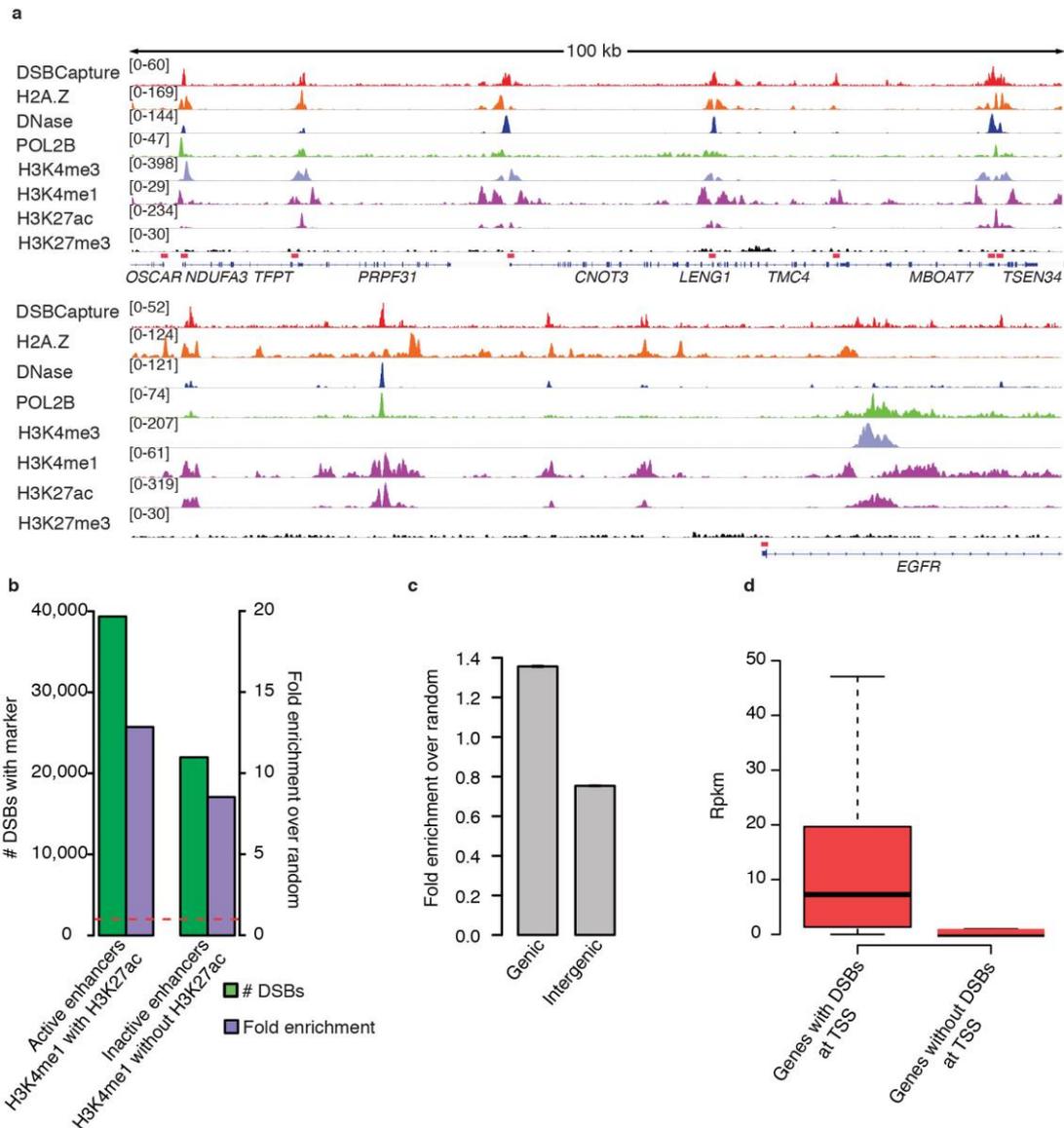
522

523

524 **Supplementary Figure 4**

525 Analysis of DSBs detected by DSBCapture and BLESS. **(a)** Overlap of peaks between two biological
 526 replicate experiments for DSBCapture and BLESS. Peaks called in both replicates (high confidence
 527 peaks) were used for data analysis. 84,946 and 18,816 high confidence peaks were identified in
 528 DSBCapture and BLESS, respectively. **(b)** Overlap between the high confidence peaks (shown in **a**)
 529 from the BLESS and DSBCapture experiments. The vast majority (98.6 %) of the BLESS peaks are
 530 also identified by DSBCapture, whereas 78.2 % of DSBCapture peaks are unique to this method. **(c)**
 531 Venn diagram showing the overlap of DSBs detected as peaks by DSBCapture performed with 50 µg
 532 and 20 µg input material. 74,951 peaks are commonly identified by the two conditions (n = 1). **(d)**
 533 Fraction of DSBs with different GC content in the DSBCapture unique peaks divided by the fraction of
 534 peaks shared between BLESS and DSBCapture within the same GC content range (fold enrichment). A
 535 fold change greater than one represents an increase in DSBs with that particular GC content in the
 536 DSBCapture unique peaks. **(e)** Fold enrichment of DSBCapture peaks with OQs¹³, calculated as the
 537 number of DSBs overlapping to OQs at each indicated % GC sequence content category (x-axis labels)
 538 divided by random overlap. All = all 716,311 OQs, irrespectively of GC content; error bars: standard
 539 deviation of the fold enrichment over random.

540



541

542 **Supplementary Figure 5**

543 Correlation of DSBs with chromatin marks, genic regions and transcription. **(a)** Genomic location of
 544 DSBs detected by DSBCapture with respect to histone marks H2A.Z, H3K4me3, H3K4me1, H3K27ac,
 545 H3K27me3 as well as DNase and POL2B. Two 100 kb genomic regions are shown; upper: a gene
 546 dense region on chromosome 19; lower: a region upstream of the *EGFR* gene. The data range is shown
 547 in square brackets, 5'UTRs are highlighted with red boxes. **(b)** Number and fold enrichment of DSBs
 548 at active and inactive enhancers. **(c)** Fold enrichment of DSBs in genic and intergenic regions over
 549 random. Values > 1 indicate that DSBs are preferentially found within that genomic location. Error
 550 bars: standard deviation of the fold enrichment over random. **(d)** Gene expression values measured as
 551 rpkM for genes with (left box) or without (right box) DSBs within ± 1 kb of the TSS. Boxes span from
 552 the 25th to the 75th percentile with the median marked by a solid bar. All whiskers extend from the 5th
 553 to the 95th percentile.

554

555 **Supplementary Table 1.** EcoRV (HeLa) DSBCapture sequencing data. Alignment statistics for a
 556 paired-end MiSeq sequencing run of the EcoRV DSBCapture library.

557
 558
 559
 560
 561
 562
 563
 564
 565
 566
 567
 568
 569
 570
 571
 572
 573
 574
 575
 576
 577
 578
 579

EcoRV (HeLa) DSBCapture experiment	Number
Reads sequenced	43,794,614
Clean reads aligned	31,267,705
Clean reads from read 1	16,366,684
Unduplicated clean reads from read 1	13,449,911
Unduplicated clean reads from read 1 (forward strand)	6,718,566
Unduplicated clean reads from read 1 (reverse strand)	6,731,345

Supplementary Table 2. AID-DIV A U2OS DSBCapture sequencing data. Alignment statistics for a
 paired-end NextSeq sequencing run of the AID-DIV A U2OS DSBCapture library. The % of proximal
 reads used was calculated as the number of unduplicated proximal reads divided by the total number of
 proximal reads sequenced (total number of proximal reads = # Reads/2).

Library	# Reads	# Clean Aligned Reads	# Unduplicated Reads	% Duplication	# Unduplicated Proximal Reads	% Proximal Reads Used
U2OS DSB-Capture	267,536,292	230,074,135	83,734,433	63.6	37,347,231	27.9

580
 581
 582

583 **Supplementary Table 3.** NHEK DSB-Capture and BLESS sequencing data. Combined data from
 584 multiple paired-end NextSeq sequencing runs of the NHEK BLESS and DSB-Capture libraries, where
 585 two independent biological replicates were performed (N1 and N2). The % of proximal reads used was
 586 calculated as the number of unduplicated proximal reads divided by the total number of proximal reads
 587 sequenced (total number of proximal reads = # Reads/2).

Library	# Reads	# Clean Aligned Reads	% Clean Reads	# Unduplicated Reads	% Duplication	# Unduplicated Proximal Reads	% Proximal Reads Used
DSB-Capture N1	353,056,932	314,223,106	89.0	256,152,484	18.5	124,192,457	70.4
DSB-Capture N2	229,658,364	204,047,616	88.8	160,803,701	21.2	77,523,430	67.5
BLESS N1	324,872,812	228,953,681	70.5	147,477,098	35.6	71,305,629	43.9
BLESS N2	585,674,924	317,449,269	54.2	51,725,762	83.7	24,766,932	8.5

588
 589
 590
 591

Supplementary Table 4. Overlap of high confidence DSB-Capture peaks in NHEKs with various marks.

Mark	# DSBs overlapping mark	Fold enrichment of DSB overlapping mark/random	% DSBs overlapping mark
DNase	64,763	33.3	76.2
H3K4me1	58,844	10.7	69.3
H3K4me2	58,842	14.5	69.3
H2AZ	48,959	11.4	57.6
H3K27ac	45,991	15.2	54.1
POL2B	32,172	12.2	37.9
CTCF	25,347	25.0	29.8
H3K4me3	20,171	17.4	23.7
H3K79me2	18,242	3.3	21.5
P63	6,131	39.1	7.2
EZH2	4,884	3.2	5.7
H3K27me3	3,766	1.2	4.4
H3K36me3	3,044	0.8	3.6
H3K9me3	1,908	0.7	2.2

592
 593