

1 **Mutational bias in spermatogonia impacts the anatomy of regulatory sites in the**  
2 **human genome**

3

4 Vera B. Kaiser<sup>1</sup>, Lana Talmane<sup>1</sup>, Yatendra Kumar<sup>1</sup>, Fiona Semple<sup>1</sup>, Marie  
5 MacLennan<sup>1</sup>, Deciphering Developmental Disorders Study<sup>1,2</sup>, David R. FitzPatrick<sup>1</sup>,  
6 Martin S. Taylor<sup>1\*</sup>, Colin A. Semple<sup>1\*</sup>

7 <sup>1</sup>MRC Human Genetics Unit, MRC Institute of Genetics and Cancer, The University  
8 of Edinburgh, Western General Hospital, Crewe Road South, Edinburgh EH4 2XU,  
9 UK

10 <sup>2</sup>The Wellcome Trust Sanger Institute, Wellcome Trust Genome Campus, Hinxton,  
11 Cambridge, CB10 1SA, UK

12 \*Equal contribution

13

14 Corresponding author: Vera B Kaiser vera.kaiser@ed.ac.uk

15

16 Keywords: germline structural variation, ATAC-seq, regulatory genomics,  
17 spermatogonia, PRDM9, NRF1

18

19

20 **Abstract**

21

22 Mutation in the germline is the ultimate source of genetic variation, but little is known  
23 about the influence of germline chromatin structure on mutational processes. Using  
24 ATAC-seq, we profile the open chromatin landscape of human spermatogonia, the  
25 most proliferative cell-type of the germline, identifying transcription factor binding  
26 sites (TFBSs) and PRDM9-binding sites, a subset of which will initiate meiotic  
27 recombination. We observe an increase in rare structural variant (SV) breakpoints at  
28 PRDM9-bound sites, implicating meiotic recombination in the generation of  
29 structural variation. Many germline TFBSs, such as NRF, are also associated with  
30 increased rates of SV breakpoints, apparently independent of recombination.  
31 Singleton short insertions ( $\geq 5$  bp) are highly enriched at TFBSs, particularly at sites  
32 bound by testis active TFs, and their rates correlate with those of structural variant  
33 breakpoints. Short insertions often duplicate the TFBS motif, leading to clustering of  
34 motif sites near regulatory regions in this male-driven evolutionary process. Increased  
35 mutation loads at germline TFBSs disproportionately affect neural enhancers with  
36 activity in spermatogonia, potentially altering neurodevelopmental regulatory  
37 architecture. Local chromatin structure in spermatogonia is thus pervasive in shaping  
38 both evolution and disease.

39

40 **Introduction**

41

42 Mutation is the ultimate source of genetic variation, and inherited variation must  
43 invariably arise in the germline. It is well established from cross-species comparisons  
44 that the rate of nucleotide substitution mutations fluctuates at the multi-megabase

45 (>10<sup>6</sup> bp) scale across the genome (Wolfe et al. 1989; Hodgkinson and Eyre-Walker  
46 2011), with early replicating regions subject to reduced rates of mutation. These  
47 patterns similarly manifest in the rate of human single nucleotide polymorphisms  
48 (SNPs) (Stamatoyannopoulos et al. 2009). Germline structural variation in the human  
49 genome is also associated with replication timing, such that copy number variants  
50 (CNVs) emerging from homologous recombination-based mechanisms are enriched  
51 in early replicating regions, while CNVs arising from non-homologous mechanisms  
52 are enriched in late replicating regions (Koren et al. 2012). Local chromatin structure  
53 also influences the mutation rate. However, finer-scale variation (<1Mb) in the  
54 germline mutation rate has so far only been related to genomic features derived from  
55 somatic cells (Gonzalez-Perez et al. 2019) because human germline-derived measures  
56 of chromatin structure have only recently become available (Guo et al. 2017; Guo et  
57 al. 2018). Transcription factor binding sites (TFBSs) are particularly prone to point  
58 mutations in cancer (Kaiser et al. 2016), probably due to interference between TF  
59 binding and the replication and repair machinery (Reijns et al. 2015; Sabarinathan et  
60 al. 2016; Afek et al. 2020), but the mutational consequences of binding at these sites  
61 in the germline is unknown.

62

63 During meiosis, homologous recombination may introduce short mutations or render  
64 genomic regions prone to rearrangements (Pratto et al. 2014; Halldorsson et al. 2019).  
65 A key player in this process is PRDM9, which binds its cognate sequence motif and  
66 directs double-strand break (DSB) formation in meiotic prophase (Baudat et al. 2010;  
67 Myers et al. 2010). In humans, PRDM9 binding site occupancy has only been directly  
68 assayed in a somatic cell line (Altemose et al. 2017), whereas indirect measures of  
69 PRDM9 activity include a proxy for DSBs (DMC1-bound single stranded DNA

70 (ssDNA)) in testis (Pratto et al. 2014), and population genetic based measures of  
71 recombination hotspots (HSs) (Myers et al. 2005; The 1000 Genomes Project  
72 Consortium et al. 2015). The method ATAC-Seq (Buenrostro et al. 2013) reports  
73 chromatin accessibility and provides a snapshot of all active regulatory regions and  
74 occupied binding sites in a given tissue. In particular, ATAC-Seq footprinting  
75 (Sherwood et al. 2014; Li et al. 2019), when applied to spermatogonia, has the  
76 potential to reveal the binding of hundreds of TFs, as well as PRDM9, in the male  
77 germline. In addition, large human genome sequencing projects can be used to reveal  
78 patterns of mutation rates, by focussing on extremely rare variants (Messer 2009;  
79 Carlson et al. 2018; Li and Luscombe 2020). Making use of such variant datasets as  
80 well as novel ATAC-Seq data in spermatogonia, we study the mutational landscape at  
81 transcription factor binding sites (TFBSs) in accessible human spermatogonial  
82 chromatin.

83

## 84 **Results**

85

### 86 **Spermatogonial regulatory regions are enriched for rare deletion breakpoints**

87

88 ATAC-Seq (Buenrostro et al. 2013) reports local chromatin accessibility and provides  
89 a snapshot of active regulatory regions and genomic regions occupied by DNA-  
90 binding proteins in a given tissue. We used ATAC-Seq to identify open chromatin  
91 sites in FGFR3-positive spermatogonial cells isolated from dissociated human  
92 testicular samples. FGFR3 is most highly expressed in self-renewing spermatogonial  
93 stem cells, with low expression also being detected in early differentiating  
94 spermatogonia (Guo et al. 2018; Sohni et al. 2019); its expression thus overlaps with

95 the onset of PRDM9 expression in pre-meiotic spermatogonia (Human Protein Atlas:  
96 <https://www.proteinatlas.org/ENSG00000164256-PRDM9/celltype/testis> and  
97 <https://www.proteinatlas.org/ENSG00000068078-FGFR3/celltype/testis>) (Guo et al.  
98 2018). Open chromatin in FGFR3-positive cells was identified using standard peak  
99 detection analysis (Methods) and multiple metrics (Supplemental Fig. S1) indicated  
100 high data quality (Yan et al. 2020). Hierarchical clustering (Ramirez et al. 2016)  
101 showed that this novel spermatogonial ATAC-Seq dataset displays a genome-wide  
102 distribution of peaks consistent with other spermatogonial derived data, and is distinct  
103 from ES cell and somatic tissue datasets (Supplemental Fig. S2).  
104  
105 Next, we assessed the enrichments of different classes of sequence variants at  
106 spermatogonial active sites, including singleton SV breakpoint frequencies as a proxy  
107 for the mutation rate of such variants. We made use of ultra-rare genomic variants  
108 from a variety of human sequencing studies: the Deciphering Developmental  
109 Disorders (DDD) study (Deciphering Developmental Disorders Study 2015; Mcrae et  
110 al. 2017) of severe and undiagnosed developmental disorders  
111 (<https://www.ddduk.org/>), a large collection of variants from an aggregated database  
112 (gnomAD; <http://gnomad.broadinstitute.org/>), and *de novo* variants from trio  
113 sequencing studies (<http://denovo-db.gs.washington.edu/>, <https://research.mss.ng/>, An  
114 et al. (2018)). Based on the DDD dataset - a combination of high-density arrayCGH  
115 and exome sequencing (Deciphering Developmental Disorders Study 2015) - we  
116 identified 6,704 singleton deletion variants among 9,625 DDD probands (carrier  
117 frequency of ~ 0.002) (Supplemental Table S1). Permutation analysis demonstrates  
118 that DDD singleton breakpoints are enriched at spermatogonial ATAC-Seq sites, their  
119 overlap being > 4-times the expected genome-wide rate (Supplemental Table S2), and

120 shifted permutation Z-scores reveal that the enrichment is specific to the ATAC-Seq  
121 peaks as opposed to wider genomic regions (Figure 1). We also considered 6,013  
122 deletions (with 7,365 unique breakpoints) that were present in the DDD consensus  
123 dataset (Deciphering Developmental Disorders Study 2015) (Methods) at a frequency  
124 of at least 1%, representing variants expected to be relatively common in human  
125 populations (Supplemental Table S1). These variants show a dip in frequency and  
126 downward trend near active sites (Figure 1a). However, we note that the overlap  
127 between common variant breakpoints and ATAC-Seq peaks is still ~ 2-fold higher  
128 than the expected genome-wide rate ( $p < 10^{-4}$ ). We conclude that singleton deletion  
129 breakpoints often occur at TFBSs in spermatogonia, suggesting a higher mutational  
130 input or less accurate repair at these sites compared to neighbouring regions. The  
131 breakpoints of more common variants are observed less frequently at the same  
132 binding sites, which may indicate the action of purifying selection in the removal of  
133 deleterious mutations at these active regulatory sites.

134         Similar trends are also observed for singleton deletion breakpoints from an  
135 independent large-scale aggregated dataset of human variants (Figure 1e, 1f) from  
136 whole genome sequence (WGS) analysis (Collins et al. 2020) (Supplemental Table  
137 S1). We again find a significant enrichment of singleton variant breakpoints at  
138 ATAC-Seq peaks, and this enrichment is not seen for common variants (Figure 1).

139

#### 140 **Locally elevated mutation at spermatogonial TFBSs**

141

142 Compared to larger structural variants, such as those (up to megabase sized) deletions  
143 examined above, indels have been shown to occur at a higher rate of about 6 new  
144 variants per genome and generation (Besenbacher et al. 2016). Short indels ( $\leq 4$  bp)

145 are thought to arise due to replication slippage (Levinson and Gutman 1987;  
146 Montgomery et al. 2013), whereas longer variants have been considered a hallmark of  
147 inaccurate DNA repair after DSBs (Rodgers and McVey 2016). Here, we focus on  
148 gnomAD singleton indels  $\leq 20$  bp as these variants are expected to be well resolved  
149 using short read sequencing. To enable higher spatial resolution of the mutation  
150 patterns at ATAC-seq defined accessible chromatin regions, and for the subsequent  
151 inference of the associated DNA binding proteins, we identified 706,008 protein  
152 binding sites using ATAC-Seq footprinting analysis (Li et al. 2019) (Methods;  
153 Supplemental Tables S3 and S4). The rate of singleton 5-20 bp insertions at  
154 footprinted spermatogonial protein binding sites approximately doubles from  
155 background expectation and is highly concentrated to within 1 kb of the binding site  
156 (Figure 2); shifted Z-scores based on genome-wide circular permutations similarly  
157 show a highly localized spike of insertions around TFBSs (Figure 2). This pattern  
158 starkly contrasts the localised depletion of common variants of the same mutation  
159 class at the same binding sites (Figure 2), again implicating a locally elevated  
160 mutation rate and purifying selection. In fact, most classes of rare mutation (singleton  
161 SVs, smaller and longer indels, SNPs) are significantly enriched at spermatogonial  
162 TFBSs (Figure 3), and in the gnomAD dataset, where all singleton classes have been  
163 ascertained by WGS, the enrichment is strongest for insertions  $\geq 5$  bp. We  
164 confirmed the enrichment of singleton short insertions and SV deletion breakpoints at  
165 spermatogonial TFBSs, using an independent permutation approach with bedtools  
166 shuffle (Quinlan and Hall 2010) (Supplemental Table S5).

167 In addition to singleton variants from large population samples, we also  
168 compiled a set of “gold standard” *de novo* short variants from a range of trio  
169 sequencing studies (see Methods). The *de novo* variants show a strikingly similar

170 trend to the gnomAD singleton variants, with a moderate (~10-60%) increase of  
171 mutation rates at TFBSs for all categories of short 1-2bp sequence variants, a but  
172 larger increase of ~130% for insertions of 5-20 bp (Figure 3). These results were  
173 confirmed using a set of independent positive and negative control sites  
174 (Supplemental Fig. S3). We conclude that regulatory sites that are active in  
175 spermatogonia show unusual parallel enrichments for both large SV breakpoints and  
176 5-20 bp insertions, consistent with localised DNA damage or error-prone repair.

177

### 178 **Germline PRDM9 and NRF1 binding generate hotspots for structural variation**

179

180 To examine any differences in mutational loads associated with different binding  
181 factors, we analysed mutational patterns stratified by the binding factors included in  
182 the JASPAR database (Sandelin et al. 2004). We accounted for redundancy caused by  
183 multiple factors binding to a single motif by considering 167 motif families  
184 (Supplemental Table S6). Furthermore, using the reported binding site motif for  
185 PRDM9 (Myers et al. 2008), we defined 9,778 putative PRDM9-bound sites  
186 corroborated by evidence for H3K4me3 enrichment in testis (Methods).

187 The spermatogonial binding sites of 11% (19/167) of motif families overlapped DDD  
188 singleton deletion breakpoints more often than expected (Bonferroni corrected  $p =$   
189 0.017) and no motif family was found to be depleted for breakpoints (Supplemental  
190 Table S3), suggesting that increased load is a common feature of TFBSs bound by  
191 different transcription factors in the germline. Similarly, singleton 5-20bp insertions  
192 from the gnomAD database were found to be significantly enriched at 29% (48/167)  
193 of families (Bonferroni corrected  $p = 0.017$ ) and, nominally, 84% (140/167) of  
194 families showed enrichment for these insertions (Supplemental Table S4). Again, no



195 TFBS family was found to be depleted for these rare variants. Collectively, these  
196 results suggest that TFBSs active in spermatogonia incur locally elevated burdens of  
197 short insertions and large structural variants across many different binding motifs.

198

199 Certain motif families appear to carry notably higher mutational loads than the  
200 general disruption seen across all TFBSs. Based on the insertion fold enrichment  
201 (IFE), i.e. the ratio of the observed to expected numbers of insertions (5-20 bp),  
202 PRDM9 binding sites are among the most disrupted sites in the genome (IFE = 6.3),  
203 and this also holds for PRDM9 sites outside known sites of recombination (IFE = 6.7  
204 for 8,139 PRDM9 sites with a distance of at least 500bp from HSs and ssDNA sites,  
205 respectively). PRDM9 sites are similarly associated with higher rates of singleton  
206 deletion breakpoints (Figure 4a, 4c), in line with the roles of PRDM9 during  
207 recombination, though PRDM9 sites outside known sites of recombination also show  
208 this trend (observed overlaps with deletion breakpoints = 9; expected = 1;  $p < 10^{-4}$ ).  
209 Two other TFBS families, corresponding to NRF1 (Nuclear Respiratory Factor 1;  
210 IFE=7.0) and HINFP (IFE=6.6) exceed the disruption seen at PRDM9 sites, and  
211 remarkably, NRF1 sites are disrupted at high rates according to DDD breakpoint data  
212 (Supplemental Table S3). Shifted Z-scores for the enrichment of short insertions (5-  
213 20 bp) at both NRF1 and PRDM9 binding sites are in the top four, next to SP/KLF  
214 transcription factors (motif families 938 and 992), suggesting strong focal  
215 enrichments at these sites (Supplemental Tables S4 and S6). NRF1 has been shown to  
216 be an important testis-expressed gene with meiosis-specific functions (Wang et al.  
217 2017; Palmer et al. 2019), but NRF1 binding sites have, to our knowledge, not been  
218 reported to be foci for genomic instability. We find strikingly similar enrichments of  
219 short insertions (5-20 bp) at TFBSs in SSEA4- and KIT-marked spermatogonial

220 samples produced in previous ATAC-seq studies (Guo et al. 2017; Guo et al. 2018).  
221 Reprocessing these previous datasets identically to our own reveals that PRDM9,  
222 NRF1 and HINFP sites are again among the top 5 disrupted motif families  
223 (Supplemental Tables S7 and S8).  
224  
225 Although both PRDM9 and NRF1 binding sites are GC-rich, their modest motif  
226 similarity suggests that the two factors occupy distinct binding motifs (PWMclus:  
227 Pearson's correlation distance  $r = 0.35$  for PRDM9 *versus* and NRF1) and should not  
228 converge on the same sites. However, in practice, PRDM9 and NRF1 binding sites  
229 were often found within the same regulatory regions, such that many (1,199) ATAC-  
230 Seq peaks contained both the NRF1 and PRDM9 binding motifs. The disruption of  
231 motifs within these co-bound peaks was notably higher, with NRF1-motifs being  
232 disrupted by short insertions 10.8-fold the expected rate (observed: 108; expected:  
233 10), and PRDM9-motifs 11.2-fold the expected rate (observed: 146; expected: 13)  
234 when co-occurring with the other factor ( $p < 10^{-4}$  in each case). Similarly, 1,311  
235 ATAC-Seq peaks contained a motif for both CTCF and PRDM9, and CTCF motifs in  
236 these peaks were more highly disrupted by short insertions (ratio = 6.3; observed: 69;  
237 expected: 11) compared to all CTCF motifs (Supplemental Table S4), as was PRDM9  
238 (ratio = 8.2; observed: 115; expected: 14) ( $p < 10^{-4}$  in each case).  
239 Importantly, the excess of insertions observed at particular motif sites is not a trivial  
240 consequence of statistical power (i.e. the number of TFBSs in the genome); for  
241 example, the number of binding sites identified for PRDM9 and NRF1 is fewer than  
242 many other factors (< 10,000 sites each; Supplemental Tables S3 and S4).  
243 In general, mutational loads appear to be dependent on the level of chromatin  
244 accessibility (MACS2 peak scores) and the number of factors predicted to bind at

245 ATAC-Seq defined regulatory regions, such that regions in the upper quartile of  
246 accessibility that are also occupied by more than 4 factors incur the highest indel  
247 loads (Supplemental Fig. S4). The significant positive correlation between the rates of  
248 binding site disruption via singleton insertions and deletion breakpoints across all  
249 motif families (Supplemental Fig. S5; Spearman's  $R = 0.52$ ,  $p < 10^{-5}$ ) suggests that  
250 the two types of damage may be mechanistically linked. In support of this idea,  
251 singleton short insertions (5-20 bp) and singleton SV deletion breakpoints overlap at  
252 the exact nucleotide position more often than expected (genome-wide  $Z$ -score =  
253 26.31;  $p < 10^{-4}$ ; see also Supplemental Fig. S6). This overlap is unlikely to be due to  
254 erroneous variant calling in the singleton dataset since we observe similar patterns for  
255 common variants of the same variant categories (genome-wide  $Z$ -score = 62.9,  $p < 10^{-4}$ ).  
256

257

### 258 **Short insertions generate clustered binding sites within regulatory regions**

259

260 Intriguingly, the 5-20 bp insertions observed at TFBSs frequently occur within only a  
261 few nucleotides of the binding motifs, whereas other classes of short variants do not  
262 show such a precisely localized increase (Figure 5 and Supplemental Fig. S7). Despite  
263 a moderate genome-wide enrichment (Figure 3), the 1-2 bp insertions characteristic of  
264 polymerase slippage, do not peak in the immediate neighbourhood of TFBSs (Figure  
265 5 and Supplemental Fig. S7). We examined the consequences of elevated 5-20 bp  
266 insertion rates at TFBSs using an exhaustive motif search algorithm (Bailey et al.  
267 2009), which finds overrepresented sequence motifs among a set of input sequences.  
268 We found that the inserted sequences at a mutated TFBS often contain additional  
269 copies of the sequence motif corresponding to the original TFBS (Figure 6 and

270 Supplemental Fig. S8), suggesting that many insertions at TFBSs are tandem  
271 duplication events, including events at CTCF, NRF1 and PRDM9 sites. The presence  
272 of these motif-containing singleton insertions appears to reveal a novel mutational  
273 mechanism expected to increase the number of binding sites for a binding factor and  
274 to lead to the expansion of TFBS clusters. CTCF-binding sites are known to occur in  
275 clusters (Kentepozidou et al. 2020) and are often affected by singleton insertions in  
276 our dataset (ranked 12<sup>th</sup> out of 167 motif families, based on the number of insertions  
277 per TFBS; Supplemental Table S4). We find that spermatogonial active sites exhibit  
278 greater homotypic clustering of TFBS than ATAC-seq defined binding sites from  
279 somatic tissues (Figure 6). Combined with a positive correlation between homotypic  
280 motif clustering and insertion rate, this suggests that spermatogonia binding sites are  
281 progressively accruing motif clusters.

282 These unusual patterns of clustered TFBSs at indel breakpoints appear to be specific  
283 to spermatogonial ATAC-seq peaks, and do not reflect genome-wide trends. Applying  
284 the MEME-Chip algorithm on 50bp regions flanking singleton insertion and deletion  
285 breakpoints, we were able to re-discover the sequence motifs of commonly disrupted  
286 binding sites, including the motifs of PRDM9 and NRF1 (Supplemental Table S9). In  
287 contrast, genome-wide, the motifs discovered flanking these variants were more likely  
288 to be simple repeats and other low complexity sequences that did not match known  
289 TFBS motifs, suggesting that processes other than transcription factor binding drive  
290 DNA breakage outside of active regulatory sites.

291

292 **Genomic instability at spermatogonial TFBSs impacts enhancers active in neural**  
293 **development**

294

295 Since many regulatory regions of the genome are active across a variety of cell types  
296 (Andersson et al. 2014), mutation at TFBSs in spermatogonia might disrupt gene  
297 regulation in other tissues. The developing brain is of particular interest, given reports  
298 of increased SV burdens in neurodevelopmental disorders (Girirajan et al. 2011;  
299 Leppa et al. 2016; Collins et al. 2017). We classified developmentally active human  
300 brain enhancers (distal regulatory elements) supported by neocortical ATAC-Seq data  
301 (de la Torre-Ubieta et al. 2018) according to whether they were either active (10,888  
302 brain enhancers) or inactive in the male germline (26,162 brain enhancers). We then  
303 calculated the odds ratio of a singleton mutation affecting a brain enhancer which is  
304 also *active* in spermatogonia, relative to a brain enhancer which is *inactive* in  
305 spermatogonia. For DDD singleton deletion breakpoints, the odds ratio was 6.82  
306 (95% CI = [5.34,8.71]), and for a singleton gnomAD insertion (5-20 bp), it was 4.69  
307 (95% CI = [4.46,4.93]). This suggests that activity in spermatogonia greatly  
308 predisposes a brain enhancer to DNA damage, and this damage manifests in  
309 enhancers that share activity with the male germline (Figure 7). Brain enhancers that  
310 are shared with spermatogonia are, on average, more accessible in the developing  
311 brain than those that are inactive in the germline (the median “mean of normalized  
312 counts” for the two types of brain enhancers were 104.8 and 54.1, respectively;  
313 Wilcoxon test  $W = 197340000$ ,  $p\text{-value} < 2.2e\text{-}16$ ), suggesting a link between  
314 enhancer activity, the sharing of enhancers across tissues and propensity to mutation.  
315 The subset of brain enhancers which overlapped spermatogonial active sites were not  
316 enriched for specific motifs, and the number of motif sites for each motif family were  
317 highly correlated between brain and spermatogonia (Spearman's  $\rho = 0.95$ ,  $p < 10$

318 <sup>15</sup>). That is, the propensity to mutation does not appear to be driven by an enrichment  
319 of specific motif families in brain enhancers.

320

321 **Spermatogonia accessible TFBS motifs incur increased rates of disruption**

322

323 We cannot exclude a small contribution of the TFBS sequence itself on the  
324 predisposition to mutation (Kondrashov and Rogozin 2004), but our data suggest that  
325 TF binding is a major driver of insertion and deletion mutation in the human  
326 germline. This is supported by the fact that we see an increase of disruption of brain  
327 enhancers if they are active in spermatogonia (Figure 7) and, more generally, an  
328 increase in the mutational load for sites that are active across other somatic tissue if  
329 binding also occurs in the germline (Supplemental Table S10). In addition, control  
330 motif sites (representing the same TFBS but located outside of ATAC-Seq peaks) are  
331 subject to lower rates of mutation compared to motifs within spermatogonial ATAC-  
332 Seq peaks (Figure 6c). Motifs within peaks carry, on average, 73% more mutations  
333 than their control counterparts, and for the most highly disrupted motifs, the  
334 discrepancy between active and control motifs is even larger. For example, PRDM9  
335 motifs are 3.4-fold, HINFP 2.9-fold and NRF1 motifs 2.6-fold more disrupted if they  
336 are active in spermatogonia, relative to spermatogonia inactive motifs. We note that  
337 this increase in disruption is likely to be a conservative estimate since some control  
338 sites may be bound at time points in the germline that our ATAC-Seq data cannot  
339 ascertain.

340 Since the X chromosome spends only one third of its time in males - the sex with the  
341 higher number of germ cell divisions - a depletion of mutations on the X chromosome  
342 is expected for a male-biased mutational process. We find the X chromosome to be

343 strongly depleted for short singleton gnomAD insertions (5-20 bp), with a ratio of X  
344 to autosome variants per uniquely mappable site of 0.78 (Supplemental Table S11).  
345 However, we note that, despite the overall reduced rate of insertions on the X, ATAC-  
346 Seq peaks on the X are still subject to increased rates of insertions compared to  
347 genome-wide expectations, suggesting that the inferred effects of protein-binding on  
348 mutation are larger than the reduction in mutation due to X-linkage (38 observed  
349 insertions in X-linked ATAC-Seq peaks, whereas 11 were expected;  $p < 10^{-4}$ ).  
350  
351 To test which candidate genomic feature most reliably predicts DNA damage, we  
352 used random forest regression to model the rate of singleton variants within 5 kb  
353 genomic windows, based on their overlap with spermatogonial TFBSs, ssDNA sites,  
354 LD-based hotspots, average GC content, mappability, gene density, replication time  
355 as well as various repeat families (LTRs, SINEs, LINEs and simple repeats). In  
356 models of genome-wide short insertion rates or deletion breakpoint rates, measures of  
357 replication timing and GC content were important predictors of mutation load as  
358 expected (Supplemental Fig. S9). Mappability was an important factor for predicting  
359 mutation rates genome-wide, perhaps reflecting the association between segmentally  
360 duplicated (low mappability) regions and rapid structural evolution, or perhaps  
361 suggesting that a fraction of variants may be erroneously called in the gnomAD  
362 dataset. (Only regions with high mappability were included in our more detailed  
363 analyses of TFBSs (Figures 3-7 and Supplemental Fig. S7)). However,  
364 spermatogonial ATAC-Seq derived TFBSs contributed additional predictive power to  
365 the models, even at the scale of the entire genome. The same TFBSs appear to be  
366 somewhat more important features in models that specifically predict damage at  
367 active brain enhancers (Supplemental Fig. S9). Genome-wide, deletion breakpoints

368 and 5-20 bp insertions were enriched in early replicating DNA (Spearman's rank  
369 correlation with replication timing:  $\rho = 0.08$ ,  $p < 10^{-15}$  and  $\rho = 0.07$ ,  $p < 10^{-15}$ ,  
370 respectively). In contrast, the presence of repeat elements had almost no impact in  
371 predicting either short insertion or deletion breakpoint rates (Supplemental Fig. S9).  
372 We conclude that germline active regulatory sites, through their occupancy by DNA  
373 binding factors, make a substantial contribution to genome-wide *de novo* structural  
374 variant rates, independent of other genomic features.

375

## 376 **Discussion**

377

378 We have demonstrated enrichments of rare and *de novo* SV breakpoints at  
379 spermatogonial regulatory sites defined by ATAC-Seq, suggesting that these sites  
380 suffer high rates of DSBs in the male germline. The same sites show unusual parallel  
381 enrichments for short variants, and particularly 5-20bp insertions. No TFBS family  
382 examined was found to be depleted for these rare variants, suggesting that many  
383 different TFBSs active in spermatogonia are prone to higher mutational loads. These  
384 loads appear to be positively correlated with the levels of chromatin  
385 accessibility/nucleosome disruption (ATAC-Seq peak binding strength) and the  
386 number of factors predicted to bind within the region. Sites bound by PRDM9, NRF1  
387 and HINFP incur the highest levels of disruption, but 11% of 167 TF families  
388 examined showed evidence for significantly elevated mutation rates. These results  
389 have implications for the evolution of binding site patterns within regulatory regions,  
390 and for disrupted regulation in somatic tissues.

391



392 Homotypic clusters of TFBSs are a pervasive feature of both invertebrate and  
393 vertebrate genomes, and have long been known to be a common feature of human  
394 promoter and enhancer regions (Gotea et al. 2010). Various adaptive hypotheses have  
395 been proposed for the presence of such clusters such that they provide functional  
396 redundancy within a regulatory region, enable the diffusion of TF binding across a  
397 region, and allow cooperative DNA binding of TF molecules (Gotea et al. 2010).  
398 More recently it has been suggested that homotypic TFBS clusters may also  
399 contribute to phase separation and the compartmentalisation of the nucleus  
400 (Kribelbauer et al. 2019). Similarly, the clustered patterns of CTCF sites in the  
401 genome have been ascribed critical roles in chromatin architecture and regulation,  
402 particularly at regulatory domain boundaries. However, these boundary regions have  
403 been shown to exhibit genome instability (Kaiser and Semple 2018) and recurrently  
404 acquire new CTCF binding sites in dynamically evolving clusters (Kentepozidou et  
405 al. 2020). The data presented here suggest that binding site clusters may arise solely  
406 as a selectively neutral consequence of the unusual mutational loads at germline  
407 TFBSs, with clusters maintained by recurrent DNA damage and misrepair.  
408  
409 We observe significant enrichments of both large SV breakpoints and small insertions  
410 together at spermatogonial TFBSs. This parallel enrichment of both types of mutation  
411 may originate from DNA breakage, followed by misrepair, conceivably via a pathway  
412 such as non-allelic homologous recombination (NAHR). It is known that NAHR can  
413 create large insertions and deletions (Kim et al. 2016), and PRDM9 activity is  
414 implicated in certain developmental disorders arising via NAHR (McVean 2007;  
415 Myers et al. 2008; Berg et al. 2010). For example, the locations of PRDM9 binding  
416 hotspots coincide with recurrent SV breakpoints causing Charcot-Marie-Tooth

417 disease, and Hunter and Potocki-Lupski/Smith-Magenis syndromes (Pratto et al.  
418 2014). It is possible that the sequence similarity at TFBSs scattered across the genome  
419 may make them particularly prone to NAHR. However, we note that the sequence  
420 similarity between the low copy repeat units, known to be involved in NAHR, is  
421 usually of the size of several kb (Gu et al. 2008), rather than sequences on the scale of  
422 TFBSs. The NHEJ pathway can also lead to short insertions after DNA breakage,  
423 usually in G0 and G1 phases of the cell cycle. Indeed, NHEJ is the most common  
424 repair pathway of DSBs in mammals and it is typically error prone (van Gent et al.  
425 2001; Lieber et al. 2003). During NHEJ, double-strand break ends are resected to  
426 form single-stranded overhangs, but when pairing occurs between the tips of the  
427 overhangs, sequences near the breakpoints will often be duplicated (Rodgers and  
428 McVey 2016). Interestingly and consistent with our results based on ultra-rare  
429 sequence variants, two previous studies using human–chimpanzee–macaque multiple  
430 alignments have shown that high numbers of short insertions have occurred in the  
431 human lineage (Kvikstad et al. 2007; Messer and Arndt 2007), and both conclude that  
432 these insertions preferentially take place in the male germline, evidenced by  
433 decreased mutation rates on the X chromosome, with similar observations in rodents  
434 (Makova et al. 2004).

435

436 The data presented here suggest that different DNA binding proteins differ widely in  
437 their impact on mutation rates. The two proteins with the largest impacts, NRF1 and  
438 PRDM9, are both highly expressed in testis, revealing a possible link between the  
439 expression level of a gene encoding a DNA binding protein and the propensity for  
440 breakage or inefficient repair at the sites the protein binds. Incidentally, NRF1 has a  
441 pLI score of 0.999, indicating that it is extremely loss-of function intolerant and

442 crucial for the organism's functioning (Karczewski et al. 2020). A previous study  
443 (Montgomery et al. 2013), using 1K genomes polymorphism data, failed to find an  
444 increase in indels at PRDM9 motifs genome-wide. This highlights the importance of  
445 using ATAC-Seq data to confine the search for motifs to germline active sites only,  
446 combined with singleton variants from large-scale sequencing studies as a more  
447 powerful strategy to explore fine scale mutational patterns.

448

449 Although studies of coding sequences, such as the DDD (Deciphering Developmental  
450 Disorders Study 2015), have revealed many of the genes disrupted in developmental  
451 disorders, more than half of cases lack a putatively causal variant (Mcrae et al. 2017),  
452 stimulating interest in the noncoding remainder of the genome, and particularly  
453 regulatory regions active in development. Limited sequencing data, covering a  
454 fraction of human regulatory regions, suggests that *de novo* mutations are enriched in  
455 these regions and are therefore likely to contribute to neurodevelopmental disorders at  
456 some level (Short et al. 2018; Gerrard et al. 2020). However, there appear to be very  
457 few, if any, individual regulatory elements recurrently mutated across multiple cases  
458 to cause neurodevelopmental disorders with a dominant mechanism (Short et al.  
459 2018). The data presented here suggests a potential solution to this paradox, where  
460 combinations of mutations at multiple regulatory regions may underlie a disease  
461 phenotype. The frequency of such combinations is expected to be many times higher  
462 if they involve regulatory regions bound by factors such as NRF1. In such cases, an  
463 entire class of sites, rather than an individual site, is subject to recurrent mutation.

464

465 **Methods**

466

## 467 **Identification of spermatogonial binding sites**

468 Samples of testicular tissue were obtained from three patients undergoing  
469 orchiectomy with total processing completed within ~5-7 hours of explant. Tissue  
470 was obtained after informed consent through the Lothian NRS BioResource, and the  
471 study was approved by NHS Lothian (Lothian R&D Project Number 2015/0370TB).  
472 Tissue samples were disaggregated into cells, and cells were labelled with  
473 phycoerythrin (PE)-conjugated antibody against the cell surface marker FGFR3  
474 (FAB766P, clone 136334, R&D systems). Spermatogonial cells were isolated using a  
475 FACSaria II cell sorter (BD bio- sciences) based on PE fluorescence and cell shape,  
476 according to Forward/Side Scatter. Isolated cells were subjected to ATAC-seq using  
477 the protocol and reagents described in (Buenrostro et al. 2013), followed by paired-  
478 end sequencing on Illumina HiSeq4000 (75 bp read length). We combined reads from  
479 separate sequencing runs into three biological replicates, based on origin and  
480 morphological appearance of the FACS sorted cells. Replicate 1: combined  
481 sequencing runs H.5.1 and H.5.4; a non-cancer patient; large cells, high side scatter;  
482 58,000 and 42,000 cells, respectively. Replicate 2: combined sequencing runs H.5.2  
483 and H.5.5; the same non-cancer patient as Replicate 1; large cells; 36,000 and 23,000  
484 cells, respectively. Replicate 3: combined sequencing runs H.7.3 and H.10.2; normal  
485 tissue from cancer patients; large cells; 69,000 and 24,000 cells, respectively. ATAC-  
486 seq raw reads were trimmed to remove any retained adaptor sequences using cutadapt  
487 (Martin 2011) with parameters *-n3 --format=fastq --overlap=3 -g*  
488 *GAGATGTGTATAAGAGACAG -g CAGATGTGTATAAGAGACAG -a*  
489 *CTGTCTCTTATAACACATCTG -a CTGTCTCTTATAACACATCTC*. Reads were aligned  
490 to the GRCh38/hg38 genome assembly with Bowtie2 (Langmead and Salzberg 2012)  
491 in paired-end mode, limiting the insert size to 4 kb. Any reads with quality score < 30

492 were discarded. Paired end reads were converted to fragments using bedtools  
493 bamtobed -bedpe, followed by extraction of the most 5' and 3' coordinates of each  
494 pair. PCR duplicates were removed by retaining only one instance of a fragment with  
495 identical coordinates within a sample. Fragments overlapping with the regions  
496 previously blacklisted as mitochondrial homologs (Buenrostro et al. 2013) were  
497 discarded. Peaks were identified from short fragments of  $\leq 100$  bp (Supplemental  
498 Fig. S1), thought to arise due to transposition events around transcription factor  
499 binding sites – and distinct from fragments spanning the larger nucleosomes  
500 (characterized by a  $\sim 200$  bp periodicity) (Buenrostro et al. 2013). Peaks were called  
501 from short fragments using macs2 callpeak (Zhang et al. 2008) with the following  
502 parameters: -B -q 0.01 -f BAMPE --nomodel --nolambda --keep-dup auto --call-  
503 summits. The clustering of ATAC-Seq peaks near transcription start sites and  
504 promoters was assessed using the ChIPseeker R package (Yu et al. 2015).

505

506 For the downstream mutation analyses, ATAC-Seq peaks from Replicates 1 and 2  
507 (the non-cancer patient) were merged, creating a single peak set. This dataset also  
508 formed the basis for the footprinting analysis, which used, as input, the combined  
509 short sequencing fragments of Replicates 1 and 2, running “rgt-hint footprinting” with  
510 --atac-seq and --bias-correction, followed by “rgt-motifanalysis matching” with the  
511 option --remove-strand-duplicates (Li et al. 2019). Input motifs were the 579 position  
512 weight matrixes (PWMs) of the japspar vertebrate database (Sandelin et al. 2004) as  
513 well as the 13-mer PRDM9 motif “CCNCCNTNNCCNC” (Myers et al. 2010) which  
514 was also provided as a PWM. The tissue donor for Replicates 1 and 2 was a carrier of  
515 the most common (European) alleles of PRDM9, which was confirmed by  
516 investigating his allelic state at the SNP (rs6889665) identified by Hinch et al. (2011);

517 this SNP was covered by our ATAC-Seq by 10 reads, all of which were “T”.

518 Accordingly, we assume that the donor is a carrier of the A and/or B allele of PRDM9

519 (both of which bind the same DNA motif), and the search for the 13-mer PRDM9

520 motif in this patient’s ATAC-Seq data can be used as a proxy for PRDM9 binding in

521 European populations. In addition, Replicate 3 was processed in the same way as the

522 combined Replicates 1 and 2 and served as a positive control to assess the genome-

523 wide enrichment of mutations at spermatogonial accessible sites (Supplemental Fig.

524 S3).

525

526 Jaspar input motifs are often highly similar, resulting in multiple binding proteins

527 being identified by the rgt-hint pipeline to bind at the same ATAC-Seq footprint; this

528 is biologically implausible (since only one protein is likely to occupy a given site),

529 and we clustered motifs by similarity, using the default parameters of the PWMclus

530 CCAT package (Jiang and Singh 2014). This resulted in a set of 167 motif families of

531 similar binding motifs (Supplemental Table S6). Using bedtools (Quinlan and Hall

532 2010), we merged overlapping binding sites that belonged to motifs of the same

533 family (thus calling them only once), and we also merged palindromic binding sites

534 called on both strands. Since PRDM9 is known to leave a characteristic histone

535 methylation mark on bound DNA (Grey et al. 2011; Powers et al. 2016), we

536 intersected the PRDM9 motif sites with testis-derived H3K4me3 marks (called in an

537 PRDM9 A/B heterozygous individuals) from Pratto et al. (2014). This resulted in a

538 stringent set of PRDM9 sites, which were both located in ATAC-Seq footprints and

539 also carried the H3K4me3 mark in human testis. ATAC-Seq-defined PRDM9 sites

540 showed moderate overlap with DMC1-bound ssDNA sites (Pratto et al. 2014) as well

541 as recombination HSs (Myers et al. 2005), which may reflect the fact that most cells

542 in our experiments are likely to be pre-meiotic: only 10 and 11% of PRDM9 sites  
543 were within 500 bp of a ssDNA peak and a recombination HS, respectively, whereas  
544 44% of DMC1-bound sites overlap with LD-defined HSs. However, we find that  
545 stronger ssDNA peaks are more likely to be near a PRDM9-binding site  
546 (Supplemental Fig. S10).

547

#### 548 **Comparisons between ATAC-Seq datasets**

549

550 Using the same procedure as described above, we processed raw ATAC-Seq reads  
551 from previously published datasets in order to call MACS2 peaks from short  
552 sequencing fragments. Datasets included ATAC-Seq reads from the germinal zone  
553 and cortical plate of the developing brain (SRR6208926, SRR6208927, SRR6208938,  
554 SRR6208943) (de la Torre-Ubieta et al. 2018), ATAC-Seq experiments of KIT+  
555 spermatogonia (sra accessions SRR7905001 and SRR7905002) (Guo et al. 2018),  
556 SSEA4+ spermatogonia (SRR5099531, SRR5099532, SRR5099533, SRR5099534)  
557 (Guo et al. 2017) and ESC cells (SRR5099535 and SRR5099536) (Guo et al. 2017).  
558 Adapter sequences within raw sequencing data were identified using bbmerge.sh of  
559 bbmap (<https://sourceforge.net/projects/bbmap/>) and removed using cutadapt (Martin  
560 2011), as above. Encode ATAC-Seq datasets (Encode Project Consortium 2012;  
561 Davis et al. 2018) (Liver: ENCF628MCV, Ovary: ENCF780JBA, Spleen:  
562 ENCF294ZCT, Testis: ENCF048IOT, Transverse Colon: ENCF377DAO) were  
563 downloaded as bam files, converted to BEDPE format, and short fragments were  
564 identified for peak calling.

565

566 BedGraph files (from the MACS2 output), describing the fragment pileup, were  
567 converted to bigwig format using bedGraphToBigWig and uploaded to the Galaxy  
568 server at <https://usegalaxy.eu/> (Afgan et al. 2018). DeepTools2's  
569 multiBigwigSummary (with default parameters) and plotCorrelation (with parameters  
570 `-skipZeros -removeOutliers`) (Ramirez et al. 2016) were used to create a heatmap of  
571 ATAC-Seq signals in the different tissues and datasets. The KIT+ and SSEA4+  
572 spermatogonial ATAC-Seq datasets of Guo et al. (2017) and Guo et al. (2018) were  
573 further used to perform footprinting and motif matching analyses (Li et al. 2019) as  
574 described above for the FGFR3-positive cells. Peaks and motif sites that are  
575 accessible in the developing brain but not in spermatogonia were identified using  
576 bedtools intersect (Quinlan and Hall 2010).

577

#### 578 **Structural Variant Breakpoint data**

579

580 Large SVs, identified by high-density arrayCGH, or a combination of arrayCGH +  
581 exome sequencing, were extracted from a cohort of 9,625 DDD patients, using variant  
582 calling procedures as described in (Deciphering Developmental Disorders Study  
583 2015). We filtered the DDD variants to only keep variants which fulfilled the  
584 following criteria: a CNsolidate wscore  $\geq 0.468$ , a callp  $< 0.01$  and a mean log2 ratio  
585 of  $< -0.41$  for deletions and  $0.36$  for duplications; CIPHER “false positives” were  
586 removed. Singleton variants were identified as being annotated as “novel” by the  
587 DDD release, only seen once among the DDD patients, and not seen in the dgv  
588 (MacDonald et al. 2014) and gnomAD V.2 (Collins et al. 2020) structural variant  
589 datasets (80% reciprocal overlap criterion). Since there are 9,625 patients in the DDD  
590 dataset, the gnomAD V.2 dataset contains SVs from 10,738 genomes and the dgv



591 contains SVs from 29,084 individuals, this puts an upper limit of the frequency of  
592 carriers of a singleton variant at ~ 0.002%. Breakpoints were identified as the 5' and  
593 3' coordinates of SVs, resulting in 13,406 singleton deletion and 3,406 duplication  
594 breakpoints; the resolution of the breakpoints was such that the median and mean  
595 confidence intervals were 300 bp and 12 kb, respectively. Further, we identified  
596 11,962 “common” deletion variants in the DDD dataset, which had a minimum  
597 variant frequency of 1% in the consensus CNV dataset as described by the DDD  
598 study (2015), i.e. pooled CNV datasets of Conrad et al. (2010), the Genomes Project  
599 Consortium (2010), the Wellcome Trust Case Control (2010) and the DDD normal  
600 controls. We used the 80% reciprocal overlap criterion and grouped common variants  
601 using the bedmap options --echo-map --fraction-both 0.8, followed by bedops --merge  
602 (Neph et al. 2012). The breakpoints of common variants are thus the outermost  
603 coordinates of all SVs that are collapsed into a given variant. The overlap of such  
604 “common” breakpoints with ATAC-Seq peaks was assessed independently of SV  
605 allele frequencies, i.e. a group of common SVs contributed two breakpoints to the  
606 analysis.

607 We also identified a set of singleton CNVs called with the Manta algorithm (Chen et  
608 al. 2016) from the gnomAD V.2 database (Collins et al. 2020) (80% reciprocal  
609 overlap criterion with gnomAD V.2, dgv and DDD variants), resulting in a set of  
610 73,063 deletion and 15,419 duplication breakpoints seen in ~ 0.002% of individuals  
611 but called with a different approach compared to the DDD. Common deletions and  
612 duplications ( $p \geq 0.05$ ) were also extracted from the gnomAD V.2 dataset; these  
613 variants had also been called with the Manta algorithm and included 5,954 deletion  
614 and 1,586 duplication breakpoint sites.

615

616 **Indels and SNP data**

617

618 The recently released gnomAD V.3 variants (indels and SNPs) were downloaded  
619 from <https://gnomad.broadinstitute.org/>. Only variants that passed all filters were kept  
620 (filtering using VCFtools --remove-filtered-all). Multiallelic variants were split using  
621 bcftools norm, and bcftools norm --IndelGap 2 was applied to indels, to allow only  
622 variants to pass that were separated by at least 2 bp. Singleton variants were defined  
623 as having an allele count of one, and the allele number was  $\geq 100,000$ , i.e. the allele  
624 frequency of singletons was  $p \leq 0.001\%$ .

625 We subdivided gnomAD indels into singleton insertions and deletions of different  
626 sizes: 1-2 bp (most commonly arising due to replication slippage) and those 5-20 bp  
627 (arising due to other mechanisms of DNA instability and within the size range reliably  
628 detected by short-read sequencing). To speed up simulations and allow for easy  
629 comparison between categories of variants, all classes of InDels and single nucleotide  
630 variants were down-sampled to 650,000 variants each.

631

632 A total of 854,409 *de novo* SNPs and indels were compiled from three different  
633 sources, lifted over to the hg38 assembly using the UCSC liftOver tool as required.  
634 First, we downloaded variants from <http://denovo-db.gs.washington.edu/>, including  
635 only samples from whole genome sequencing studies (Michaelson et al. 2012; Ramu  
636 et al. 2013; Genome of the Netherlands 2014; Besenbacher et al. 2015; Turner et al.  
637 2016; Yuen et al. 2016; Jonsson et al. 2017; RK et al. 2017; Turner et al. 2017;  
638 Werling et al. 2018), which included a total of 404,238 variants from 4,560 samples.  
639 Additional samples, which were not already included in the denovo-db dataset, were  
640 downloaded from the MSSNG database (<https://research.mss.ng/>), version

641 2019/10/16, which added 2,243 samples and 215,044 *de novo* mutations. A third  
642 source of *de novo* variants came from (An et al. 2018) - 3,805 samples and 255,107  
643 mutations.

644

### 645 **Circular Permutation**

646

647 To obtain a genome-wide estimate of enrichment of overlap between genomic  
648 features (e.g. TFBSs and mutations), we performed circular permutations using the  
649 Bioconductor regioneR package in R (Gel et al. 2016). We used the permTest()  
650 function with parameters ntimes=10000,  
651 randomize.function=circularRandomizeRegions, evaluate.function=numOverlaps,  
652 genome=hg38\_masked, alternative="auto", where hg38\_masked =  
653 getBSgenome("BSgenome.Hsapiens.UCSC.hg38.masked"). This test evaluates the  
654 number of overlaps observed between two sets of genomic features, given their order  
655 on the chromosome and the distance between features, i.e. taking their degree of  
656 clustering into account. At each iteration, all positions are shifted by the same  
657 randomly generated distance on conceptually "circularized" chromosomes, in effect  
658 "spinning" the position of features on each chromosome, while excluding masked  
659 regions (i.e. unmappable, repetitive and low-complexity segments). The statistical  
660 output is a z-score, which is defined as the distance between the expected number of  
661 overlaps (the distribution over 10,000 permutations) and the observed one, measured  
662 in terms of standard deviations. Shifted Z-score analysis can further indicate whether  
663 a given overlap pattern is caused by the precise locations of two feature sets or by  
664 broader, regional effects. One set of features is shifted in either direction from their  
665 original positions by a number of bases, so that the numbers of overlaps and the

666 degree to which the z-score changes in response to the shift can be tested. A sharp  
667 peak of shifted z-scores around the zero-coordinate indicates a precise overlap  
668 between features, whereas a flat profile may indicate overlaps attributable to regional  
669 effects, as is seen for features that tend to co-occur in regions with similar base  
670 composition or gene density.

671

672 For permutations involving SVs, we used the two breakpoints of each SV, and  
673 assessed the overlap of breakpoints with another feature of interest (i.e. ATAC-Seq  
674 sites), treating each breakpoint separately.

675 Circular permutations in regioneR (Gel et al. 2016) were also used assess the mean  
676 distance between ATAC-Seq peaks and deletion breakpoints, for common and  
677 singleton variants separately.

678

### 679 **Simple permutations**

680

681 Spermatogonial binding sites were randomly permuted, using the “bedtools shuffle”  
682 command (Quinlan and Hall 2010), with the parameter “-noOverlapping” and  
683 excluding assembly gaps with “-excl hg38.gap.bed” (downloaded from  
684 <https://genome.ucsc.edu/cgi-bin/hgTables>). In each of 10,000 simulations, we  
685 assessed the overlap of the permuted binding sites with short insertions and deletion  
686 breakpoints, respectively. This resulted in an expected distribution of overlap, given a  
687 random positioning of binding sites across the genome. This distribution was  
688 compared to the observed number of overlaps (using bedtools intersect), and the p-  
689 value was defined as the percentage of simulated overlaps that was larger than the  
690 observed overlap. This permutation framework does not take into account the spacing

691 and clustered nature of binding sites, and does not allow for an assessment of the  
692 precision of overlap between features. In this study, we favor the more conservative  
693 measures of significance provided by the circular permutation strategy.

694

#### 695 **Brain enhancer data**

696

697 Active brain enhancers came from de la Torre-Ubieta et al. (2018). Specifically, we  
698 used the 37,050 brain enhancers which showed differential accessibility in the  
699 germinal zone versus the cortical plate, reflecting activity in the developing brain (de  
700 la Torre-Ubieta et al. 2018). Next, we identified brain enhancers that were also active  
701 during the male germline formation, i.e. overlapping the spermatogonial ATAC-Seq  
702 peaks. To correct for the variable size of the brain active enhancers, we took the  
703 midpoints of each enhancer plus/minus 500 bp on either side, and intersected these  
704 sites with the ATAC-Seq peaks using bedtools intersect (Quinlan and Hall 2010), thus  
705 classifying brain enhancers as spermatogonial “active” or “inactive”. Next, we  
706 intersected these two categories of brain enhancers with the DDD breakpoint and  
707 gnomAD insertion dataset, respectively, to further classify them as “disrupted” by a  
708 singleton variant or “intact”. An odds ratio was calculated as

709

$$710 \text{ OR} = (A/(B - A))/(C/(D - C))$$

711 With confidence intervals

$$712 \text{ CI}_{\text{lower}} = \exp(\log(\text{OR}) - 1.96 * \sqrt{1/A + 1/(B-A) + 1/C + 1/(D-C)})$$

$$713 \text{ CI}_{\text{higher}} = \exp(\log(\text{OR}) + 1.96 * \sqrt{1/A + 1/(B-A) + 1/C + 1/(D-C)})$$

714

715 where:

716 A = Disrupted, sperm active

717 B = All sperm active

718 C = Disrupted, sperm inactive

719 D = All sperm inactive

720

721 To analyse the enrichment of short InDels and SNPs around TFBSs and brain

722 enhancers, we only considered genomic regions with unique mappability in  $\geq 95\%$

723 of the region, using the bedmap option `--bases-uniq-f` (Neph et al. 2012) and the

724 mappability file hg38\_umap24 (Karimzadeh et al. 2018), converted to bedmap

725 format.

726

### 727 **Random Forest Regression**

728

729 To compare the effects of chromatin state on mutation rates, we performed random

730 forest regression with 200 trees, modelling the outcome variables “singleton

731 breakpoints” and “singleton insertions (5-20 bp)”, from the DDD and gnomAD V.3

732 respectively, within 5-kb wide genomic windows. Predictor variables included

733 “spermatogonial TFBS count”, “ssDNA overlap” (from Pratto et al.(2014)),

734 “recombination HS overlap” (from The 1000 Genomes Project et al. (2015)), “GC-

735 content”, “Replication timing” (average of Wavelet-smooth signal in 1-kb bins of 15

736 encode tissues, downloaded from

737 <http://hgdownload.soe.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeUwRepliSeq>

738 /), “Gene density”, “Mappability” (proportion of sites in each window with an

739 umap24 score of 1), and the overlap with “LTRs”, “SINES”, “LINEs” and “Simple

740 Repeats” (downloaded from the Table Browser at <https://genome.ucsc.edu/>).

741 In a smaller model, we subsetted the dataset to only include 5-kb bins that also  
742 overlap active brain enhancers (de la Torre-Ubieta et al. 2018), then ran the random  
743 forest regression model to predict mutation rates within genomic regions that contain  
744 active brain enhancers.

745

#### 746 **Motif discovery in singleton insertion sites**

747

748 In order to find sequence motifs within the 5-20 bp singleton insertion sites from  
749 gnomAD V.3, without prior assumptions, we extracted the fasta sequence for  
750 insertions that fell within 10 bp of the top 10 disrupted motif families (motif families  
751 992, 193, 796, 907, 579, 825, 984, 171, 991). We ran the meme.4.11 motif discovery  
752 algorithm (Bailey et al. 2009) with “-nmotifs 1” on the inserted sequences. This  
753 allowed us to compare the sequence motif of the disrupted TFBSs to any recurrent  
754 motif found within the inserted sequences.

755

#### 756 **Control Motif sites**

757

758 Using default search criteria, the FIMO algorithm (Grant et al. 2011) was run on the  
759 repeat masked hg38 genome sequence (hg38.fa.masked, downloaded from  
760 <https://genome.ucsc.edu/> in March 2020), searching the whole genome for the 579  
761 input Jaspar motifs and the 13-mer PRDM9 motif. As with active binding sites, motif  
762 matches belonging to the same motif family were merged and reported as a single  
763 motif match per family, and only regions with unique umap24 mappabilities for  $\geq$   
764 95% of sites were kept; motifs that overlapped with spermatogonial ATAC-Seq peaks  
765 were excluded. Next, these “control” motif sites were down-sampled to 10,000 per

766 motif family (using bedtools sample (Quinlan and Hall 2010)); circular permutations  
767 were performed to compare the observed to expected overlap of the control motif sites  
768 (plus/minus 10 bp) with the gnomAD singleton insertions of 5-20 bp.

769

770 The FIMO predicted control sites were also used to assess the degree of “clustering”  
771 of motifs at spermatogonia active sites. For this purpose, we intersected the FIMO  
772 motifs with a) spermatogonial ATAC-Seq sites and b) ENCODE Master regulatory  
773 sites downloaded from <https://genome.ucsc.edu/> (DNaseI hypersensitivity derived  
774 from assays in 95 cell types). For each of the 167 motif families, we calculated the  
775 median distance (in basepairs) from a motif located within the active regulatory  
776 region to the nearest FIMO motif of the same type. Accordingly, the ratio of the  
777 median distance between motif sites (ENCODE/spermatogonia) was larger than one if  
778 motifs at spermatogonial sites were, on average, closer to each other than motifs near  
779 ENCODE sites, and we used this ratio as a measure of motif clustering. When  
780 correlating the IFE with the degree of motif clustering (Figure 6), we thus largely  
781 correct for base compositional biases near active sites (which impact mutation rates –  
782 Supplemental Fig. S9) as well as the effects of historical selection on the clustering of  
783 motifs near genes, i.e. shorter inter-motif distances in spermatogonia indicate that  
784 these sites have specifically high levels of motif density in spermatogonia, beyond the  
785 levels expected for binding sites in general.

786

#### 787 **Data access**

788

789 Raw sequencing data generated in this study have been submitted to the European  
790 Genome-phenome Archive (EGA) (accession number EGAS00001005366), and



791 ATAC-Seq peak files are available at Edinburgh DataShare

792 (<https://doi.org/10.7488/ds/3053>).

793

794 **Competing interest statement**

795

796 The authors have no competing interests to declare.

797

798 **Acknowledgements**

799

800 We thank all donors for their participation in genetic research.

801 In particular, we thank the DDD families, study clinicians, research nurses and

802 clinical scientists in the recruiting centres; the Genome Aggregation Database

803 (<http://gnomad.broadinstitute.org/>), MSSNG (<https://www.mss.ng/>) and denovo-db

804 (<http://denovo-db.gs.washington.edu/denovo-db/>) for making their data available. We

805 are grateful to all of the families at the participating Simons Simplex Collection (SSC)

806 sites, as well as the principal investigators (A. Beaudet, R. Bernier, J. Constantino, E.

807 Cook, E. Fombonne, D. Geschwind, R. Goin-Kochel, E. Hanson, D. Grice, A. Klin,

808 D. Ledbetter, C. Lord, C. Martin, D. Martin, R. Maxim, J. Miles, O. Ousley, K.

809 Pelphrey, B. Peterson, J. Piggot, C. Saulnier, M. State, W. Stone, J. Sutcliffe, C.

810 Walsh, Z. Warren, E. Wijsman). We appreciate obtaining access to genetic data on

811 SFARI Base. Approved researchers can obtain the SSC population dataset described

812 in this study by applying at <https://base.sfari.org>.

813 This work was supported by MRC Human Genetics Unit core funding programme

814 grants MC\_UU\_00007/11, MC\_UU\_00007/2 and MC\_UU\_00007/16.

815 We thank Elisabeth Freyer for assistance with the FAC sorting and Wendy A.

816 Bickmore for useful comments to the manuscript.

817

## 818 **Author Contributions**

819

820 V.B.K. and C.A.S. conceived the project, interpreted the results and wrote the

821 manuscript. M.S.T. designed the experiments and managed the acquisition of

822 samples. L.T., Y.K., F.S. and M.M. performed the experiments, L.T. processed raw

823 data. V.B.K. performed the analyses. D.D.D. and D.R.F. provided data. D.R.F. helped

824 with the interpretation of the results and provided critical scientific inputs.

825

## 826 **Figure Legends**

827

### 828 **Figure 1: Locally elevated structural variation rates at spermatogonial**

829 **regulatory sites.** SV breakpoint count (a, b) and circular permutation shifted Z-scores

830 (c, d) of deletion breakpoints in the DDD cohort, centred around the midpoints of

831 spermatogonial ATAC-Seq peaks. “Singletons” are breakpoints of deletions with a

832 frequency of ~ 0.002% across population samples; “common” variants are seen at a

833 frequency of at least 1% in the DDD consensus dataset (see main text); permutation

834 p-values indicate significant enrichment for both types of variants at ATAC-Seq

835 peaks ( $p < 10^{-5}$  in each case) (e, f) Circular permutation shifted Z-scores of gnomAD

836 deletion breakpoints, centred around spermatogonial ATAC-Seq peaks. “Singletons”

837 are breakpoints of deletions with a frequency of ~ 0.002% across population samples;

838 “common” variants are seen at a frequency of at least 5% in the gnomAD V.2 dataset.

839 Permutation p-values indicate significant enrichment for singleton breakpoints ( $p <$   
840  $10^{-5}$ ), and a significant depletion for common variants ( $p < 0.01$ ).

841 **Figure 2: Increased rates of short insertions focussed on spermatogonial binding**

842 **sites.** Insertion count (**a, b**) and Shifted Z-scores (**d, e**) of gnomAD singleton and  
843 common insertions (5-20 bp), centred around spermatogonial TFBSs. Singletons are  
844 seen only once in the gnomAD V.3 dataset (allele frequency  $\leq 0.001\%$ ) and are  
845 significantly enriched at binding sites ( $p < 10^{-4}$ ); common variants have an allele  
846 frequency of at least 5% within gnomAD V.3 and are significantly depleted at binding  
847 sites ( $p < 10^{-4}$ ).

848 **Figure 3: Parallel enrichments of short variants and SV breakpoints at**

849 **spermatogonial binding sites.** Circular permutation results are based on 10,000  
850 permutations; results for singleton variants and de novo mutation are shown. The Y  
851 axis shows the ratio of observed over expected variant counts. Mutation categories  
852 with significant enrichment are indicated by asterisks (\*\*\*) indicating  $p < 0.001$ . The  
853 type of variant tested and the total number of observed variants overlapping  
854 spermatogonial TFBSs are indicated below each bar.

855 **Figure 4: Binding factors associated with the highest rates of mutation at**

856 **spermatogonial binding sites.** Plots are centred on the binding sites of a given motif  
857 family inside ATAC-Seq footprints. **(a)** Singleton and **(b)** common deletion  
858 breakpoints in the DDD cohort; singletons are breakpoints of deletions with a  
859 frequency of  $\sim 0.002$  across population samples; common variants are seen at a  
860 frequency of at least 1% in the DDD consensus dataset. **(c)** Singleton and **(d)** common  
861 insertions (5-20 bp) in the gnomAD dataset. Singletons are seen only once in  
862 gnomAD V.3 (allele frequency  $\leq 0.001\%$ ), and common variants have an allele  
863 frequency of at least 5% within gnomAD V.3. Only 10 kb regions around TFBSs with

864  $\geq 95\%$  unique mappability (umap24 scores) were included. The top 5 disrupted  
865 motifs are shown, listed in order of enrichment of singleton variants in the circular  
866 permutations (all enrichments of singletons are associated with p-values  $< 10^{-4}$ ).

867 **Figure 5: Elevated singleton insertion rates at PRDM9 and NRF1 binding sites**  
868 **contrast with other short variant classes.** All gnomAD variants have been down-  
869 sampled to a total of 650,000 variants per analysis, making the Y axes directly  
870 comparable; individual bins are 5 bp in size. Only regions around TFBSs with  $\geq 95\%$   
871 unique mappability (umap24 scores) were included.

872 **Figure 6: Insertions at spermatogonial TFBSs generate motif clusters in the**  
873 **genome. a)** Jaspar database sequence motifs identified in the footprints of  
874 spermatogonial ATAC-Seq peaks (left) and the motifs identified in the singleton  
875 insertions (5-20 bp) (right). The number of insertion sites (N) that were chosen by  
876 MEME to construct the motif are shown on the right. **b)** For each motif family, we  
877 plot the insertion fold enrichment (IFE) on the X axis and the degree of  
878 spermatogonial motif clustering on the Y axis; the least square regression line is  
879 indicated in blue. Motif clustering is measured as the distance to the nearest motif at  
880 spermatogonial active sites, relative to the distance for motifs at ENCODE active  
881 sites. **c)** The insertion fold enrichment (IFE) is contrasted between FIMO control  
882 motif sites (left) and spermatogonial active motif sites (right); the Wilcoxon Test was  
883 performed to compare the IFE at the two classes of sites.

884 **Figure 7: Neural enhancers with activity in spermatogonia suffer elevated mutation**  
885 **rates. a)** Singleton DDD deletion breakpoint and **b)** singleton gnomAD insertion (5-  
886 20 bp) count around brain active enhancers. Enhancers were classified as being also  
887 active in spermatogonia (red) or inactive in spermatogonia (blue). Plotted is the  
888 average number of variants per brain enhancer - in 5 kb windows or 100 bp windows,

889 respectively. In **b**, only 10 kb regions around enhancers with  $\geq 95\%$  unique  
890 mappability (umap24 scores) were included (3,409 brain enhancers that are inactive  
891 in spermatogonia and 1,029 that are active).

892

893

## 894 **References**

895

- 896 Afek A, Shi H, Rangadurai A, Sahay H, Senitzki A, Xhani S, Fang M, Salinas R,  
897 Mielko Z, Pufall MA et al. 2020. DNA mismatches reveal conformational  
898 penalties in protein-DNA recognition. *Nature* **587**: 291-296.
- 899 Afgan E, Baker D, Batut B, van den Beek M, Bouvier D, Cech M, Chilton J,  
900 Clements D, Coraor N, Gruning BA et al. 2018. The Galaxy platform for  
901 accessible, reproducible and collaborative biomedical analyses: 2018 update.  
902 *Nucleic Acids Res* **46**: W537-W544.
- 903 Altemose N, Noor N, Bitoun E, Tumian A, Imbeault M, Chapman JR, Aricescu AR,  
904 Myers SR. 2017. A map of human PRDM9 binding provides evidence for  
905 novel behaviors of PRDM9 and other zinc-finger proteins in meiosis. *Elife* **6**.
- 906 An JY, Lin K, Zhu L, Werling DM, Dong S, Brand H, Wang HZ, Zhao X, Schwartz  
907 GB, Collins RL et al. 2018. Genome-wide de novo risk score implicates  
908 promoter variation in autism spectrum disorder. *Science* **362**.
- 909 Andersson R, Gebhard C, Miguel-Escalada I, Hoof I, Bornholdt J, Boyd M, Chen Y,  
910 Zhao X, Schmidl C, Suzuki T et al. 2014. An atlas of active enhancers across  
911 human cell types and tissues. *Nature* **507**: 455-461.
- 912 Bailey TL, Boden M, Buske FA, Frith M, Grant CE, Clementi L, Ren J, Li WW,  
913 Noble WS. 2009. MEME SUITE: tools for motif discovery and searching.  
914 *Nucleic Acids Res* **37**: W202-208.
- 915 Baudat F, Buard J, Grey C, Fledel-Alon A, Ober C, Przeworski M, Coop G, de Massy  
916 B. 2010. PRDM9 Is a Major Determinant of Meiotic Recombination Hotspots  
917 in Humans and Mice. *Science* **327**: 836-840.
- 918 Berg IL, Neumann R, Lam KW, Sarbajna S, Odenthal-Hesse L, May CA, Jeffreys AJ.  
919 2010. PRDM9 variation strongly influences recombination hot-spot activity  
920 and meiotic instability in humans. *Nat Genet* **42**: 859-863.
- 921 Besenbacher S, Liu S, Izarzugaza JM, Grove J, Belling K, Bork-Jensen J, Huang S,  
922 Als TD, Li S, Yadav R et al. 2015. Novel variation and de novo mutation rates  
923 in population-wide de novo assembled Danish trios. *Nat Commun* **6**: 5969.
- 924 Besenbacher S, Sulem P, Helgason A, Helgason H, Kristjansson H, Jonasdottir A,  
925 Jonasdottir A, Magnusson OT, Thorsteinsdottir U, Masson G et al. 2016.  
926 Multi-nucleotide de novo Mutations in Humans. *PLoS Genet* **12**: e1006315.
- 927 Buenrostro JD, Giresi PG, Zaba LC, Chang HY, Greenleaf WJ. 2013. Transposition  
928 of native chromatin for fast and sensitive epigenomic profiling of open  
929 chromatin, DNA-binding proteins and nucleosome position. *Nat Methods* **10**:  
930 1213-1218.

- 931 Carlson J, Locke AE, Flickinger M, Zawistowski M, Levy S, Myers RM, Boehnke M,  
932 Kang HM, Scott LJ, Li JZ et al. 2018. Extremely rare variants reveal patterns  
933 of germline mutation rate heterogeneity in humans. *Nat Commun* **9**.
- 934 Chen X, Schulz-Trieglaff O, Shaw R, Barnes B, Schlesinger F, Kallberg M, Cox AJ,  
935 Kruglyak S, Saunders CT. 2016. Manta: rapid detection of structural variants  
936 and indels for germline and cancer sequencing applications. *Bioinformatics*  
937 **32**: 1220-1222.
- 938 Collins RL, Brand H, Karczewski KJ, Zhao X, Alföldi J, Francioli LC, Khera AV,  
939 Lowther C, Gauthier LD, Wang H et al. 2020. A structural variation reference  
940 for medical and population genetics. *Nature* **581**: 444-451.
- 941 Collins RL, Brand H, Redin CE, Hanscom C, Antolik C, Stone MR, Glessner JT,  
942 Mason T, Pregno G, Dorrani N et al. 2017. Defining the diverse spectrum of  
943 inversions, complex structural variation, and chromothripsis in the morbid  
944 human genome. *Genome Biol* **18**: 36.
- 945 Conrad DF, Pinto D, Redon R, Feuk L, Gokcumen O, Zhang Y, Aerts J, Andrews TD,  
946 Barnes C, Campbell P et al. 2010. Origins and functional impact of copy  
947 number variation in the human genome. *Nature* **464**: 704-712.
- 948 Davis CA, Hitz BC, Sloan CA, Chan ET, Davidson JM, Gabdank I, Hilton JA, Jain K,  
949 Baymuradov UK, Narayanan AK et al. 2018. The Encyclopedia of DNA  
950 elements (ENCODE): data portal update. *Nucleic Acids Res* **46**: D794-D801.
- 951 de la Torre-Ubieta L, Stein JL, Won H, Opland CK, Liang D, Lu D, Geschwind DH.  
952 2018. The Dynamic Landscape of Open Chromatin during Human Cortical  
953 Neurogenesis. *Cell* **172**: 289-304 e218.
- 954 Deciphering Developmental Disorders Study. 2015. Large-scale discovery of novel  
955 genetic causes of developmental disorders. *Nature* **519**: 223-228.
- 956 Encode Project Consortium. 2012. An integrated encyclopedia of DNA elements in  
957 the human genome. *Nature* **489**: 57-74.
- 958 Gel B, Diez-Villanueva A, Serra E, Buschbeck M, Peinado MA, Malinverni R. 2016.  
959 regioneR: an R/Bioconductor package for the association analysis of genomic  
960 regions based on permutation tests. *Bioinformatics* **32**: 289-291.
- 961 Genome of the Netherlands C. 2014. Whole-genome sequence variation, population  
962 structure and demographic history of the Dutch population. *Nat Genet* **46**:  
963 818-825.
- 964 Genomes Project Consortium, Abecasis GR, Altshuler D, Auton A, Brooks LD,  
965 Durbin RM, Gibbs RA, Hurles ME, McVean GA. 2010. A map of human  
966 genome variation from population-scale sequencing. *Nature* **467**: 1061-1073.
- 967 Gerrard DT, Berry AA, Jennings RE, Birket MJ, Zarrineh P, Garstang MG, Withey  
968 SL, Short P, Jimenez-Gancedo S, Firbas PN et al. 2020. Dynamic changes in  
969 the epigenomic landscape regulate human organogenesis and link to  
970 developmental disorders. *Nat Commun* **11**: 3920.
- 971 Girirajan S, Brkanac Z, Coe BP, Baker C, Vives L, Vu TH, Shafer N, Bernier R,  
972 Ferrero GB, Silengo M et al. 2011. Relative burden of large CNVs on a range  
973 of neurodevelopmental phenotypes. *PLoS Genet* **7**: e1002334.
- 974 Gonzalez-Perez A, Sabarinathan R, Lopez-Bigas N. 2019. Local Determinants of the  
975 Mutational Landscape of the Human Genome. *Cell* **177**: 101-114.
- 976 Gotea V, Visel A, Westlund JM, Nobrega MA, Pennacchio LA, Ovcharenko I. 2010.  
977 Homotypic clusters of transcription factor binding sites are a key component  
978 of human promoters and enhancers. *Genome Res* **20**: 565-577.
- 979 Grant CE, Bailey TL, Noble WS. 2011. FIMO: scanning for occurrences of a given  
980 motif. *Bioinformatics* **27**: 1017-1018.

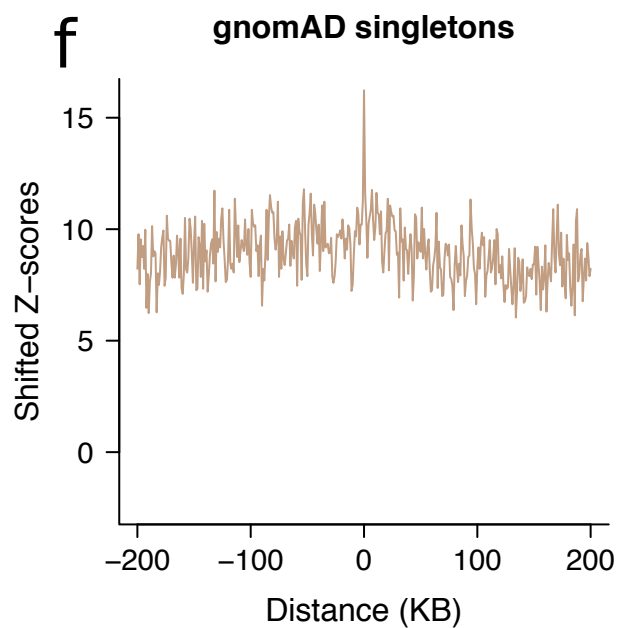
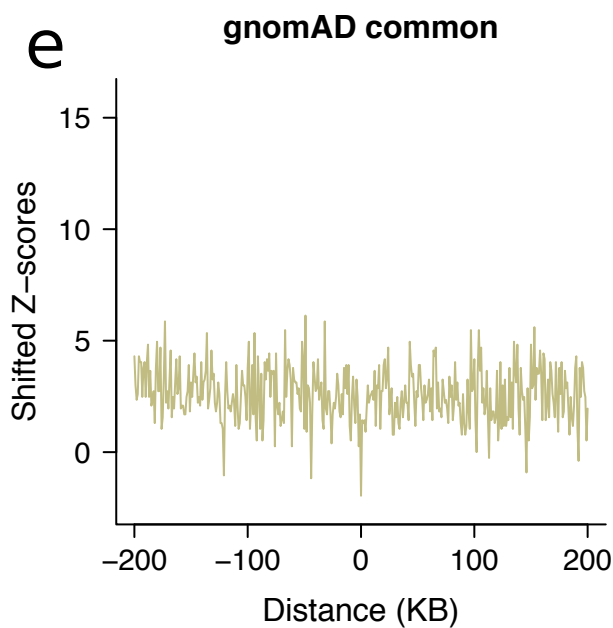
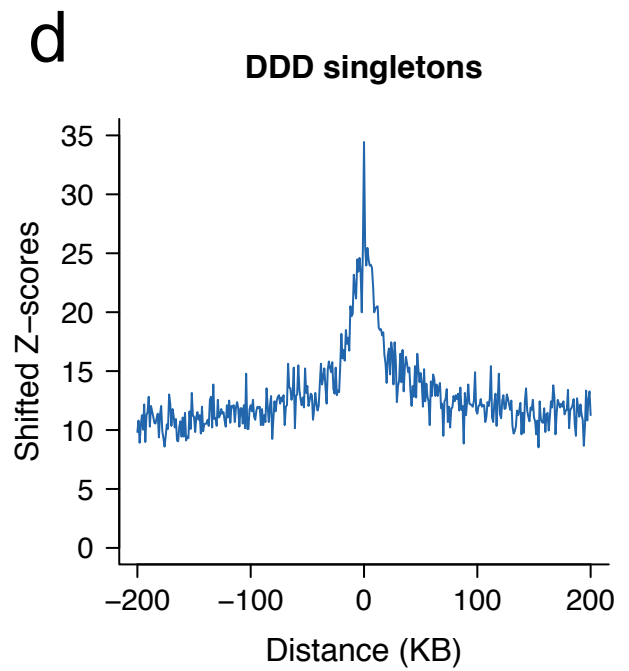
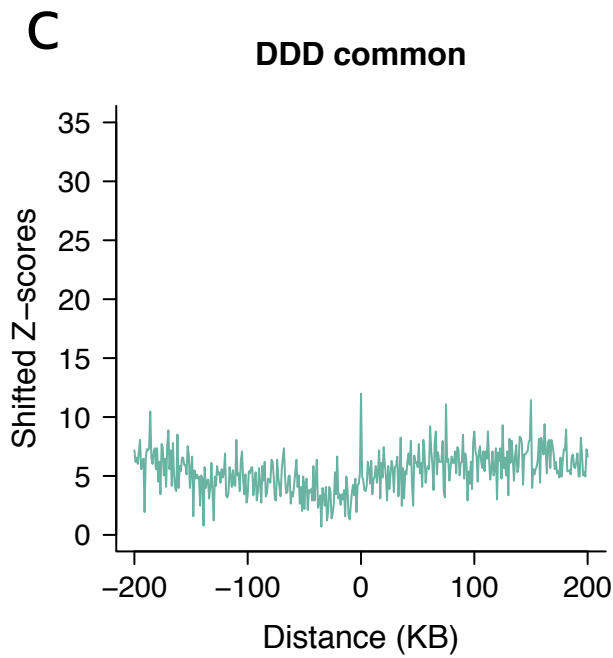
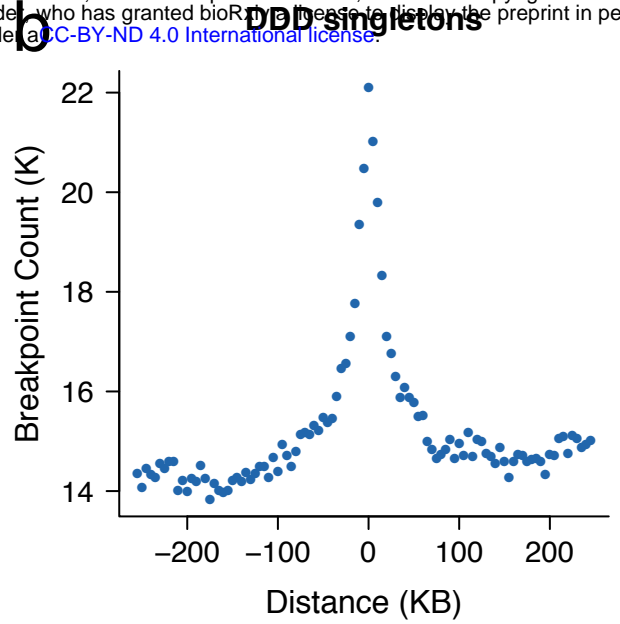
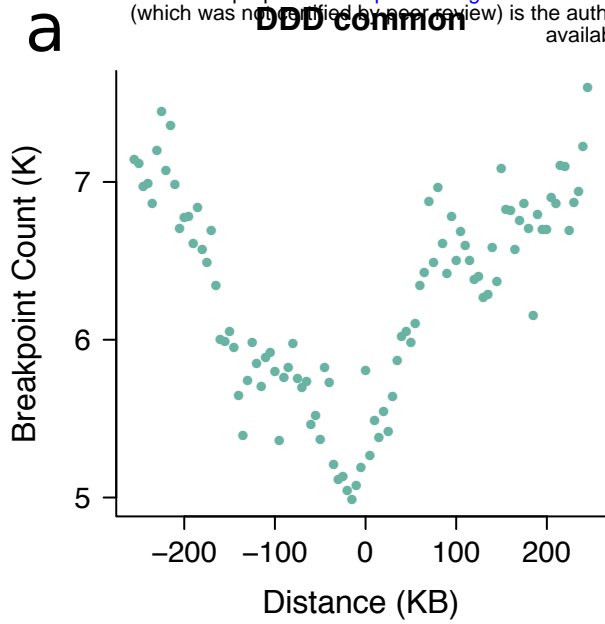
- 981 Grey C, Barthes P, Chauveau-Le Fric G, Langa F, Baudat F, de Massy B. 2011.  
982 Mouse PRDM9 DNA-Binding Specificity Determines Sites of Histone H3  
983 Lysine 4 Trimethylation for Initiation of Meiotic Recombination. *Plos Biol* **9**.  
984 Gu W, Zhang F, Lupski JR. 2008. Mechanisms for human genomic rearrangements.  
985 *Pathogenetics* **1**: 4.  
986 Guo J, Grow EJ, Mlcochova H, Maher GJ, Lindskog C, Nie X, Guo Y, Takei Y, Yun  
987 J, Cai L et al. 2018. The adult human testis transcriptional cell atlas. *Cell Res*  
988 **28**: 1141-1157.  
989 Guo J, Grow EJ, Yi C, Mlcochova H, Maher GJ, Lindskog C, Murphy PJ, Wike CL,  
990 Carrell DT, Goriely A et al. 2017. Chromatin and Single-Cell RNA-Seq  
991 Profiling Reveal Dynamic Signaling and Metabolic Transitions during Human  
992 Spermatogonial Stem Cell Development. *Cell Stem Cell* **21**: 533-546 e536.  
993 Halldorsson BV, Palsson G, Stefansson OA, Jonsson H, Hardarson MT, Eggertsson  
994 HP, Gunnarsson B, Oddsson A, Halldorsson GH, Zink F et al. 2019.  
995 Characterizing mutagenic effects of recombination through a sequence-level  
996 genetic map. *Science* **363**.  
997 Hinch AG, Tandon A, Patterson N, Song YL, Rohland N, Palmer CD, Chen GK,  
998 Wang K, Buxbaum SG, Akyzbekova EL et al. 2011. The landscape of  
999 recombination in African Americans. *Nature* **476**: 170-U167.  
1000 Hodgkinson A, Eyre-Walker A. 2011. Variation in the mutation rate across  
1001 mammalian genomes. *Nat Rev Genet* **12**: 756-766.  
1002 Jiang P, Singh M. 2014. CCAT: Combinatorial Code Analysis Tool for transcriptional  
1003 regulation. *Nucleic Acids Research* **42**: 2833-2847.  
1004 Jonsson H, Sulem P, Kehr B, Kristmundsdottir S, Zink F, Hjartarson E, Hardarson  
1005 MT, Hjorleifsson KE, Eggertsson HP, Gudjonsson SA et al. 2017. Parental  
1006 influence on human germline de novo mutations in 1,548 trios from Iceland.  
1007 *Nature* **549**: 519-522.  
1008 Kaiser VB, Semple CA. 2018. Chromatin loop anchors are associated with genome  
1009 instability in cancer and recombination hotspots in the germline. *Genome Biol*  
1010 **19**: 101.  
1011 Kaiser VB, Taylor MS, Semple CA. 2016. Mutational Biases Drive Elevated Rates of  
1012 Substitution at Regulatory Sites across Cancer Types. *PLoS Genet* **12**:  
1013 e1006207.  
1014 Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alfoldi J, Wang Q, Collins RL,  
1015 Laricchia KM, Ganna A, Birnbaum DP et al. 2020. The mutational constraint  
1016 spectrum quantified from variation in 141,456 humans. *Nature* **581**: 434-443.  
1017 Karimzadeh M, Ernst C, Kundaje A, Hoffman MM. 2018. Umap and Bimap:  
1018 quantifying genome and methylome mappability. *Nucleic Acids Res* **46**: e120.  
1019 Kentepozidou E, Aitken SJ, Feig C, Stefflova K, Ibarra-Soria X, Odom DT, Roller M,  
1020 Flicek P. 2020. Clustered CTCF binding is an evolutionary mechanism to  
1021 maintain topologically associating domains. *Genome Biol* **21**: 5.  
1022 Kim S, Peterson SE, Jasin M, Keeney S. 2016. Mechanisms of germ line genome  
1023 instability. *Semin Cell Dev Biol* **54**: 177-187.  
1024 Kondrashov AS, Rogozin IB. 2004. Context of deletions and insertions in human  
1025 coding sequences. *Hum Mutat* **23**: 177-185.  
1026 Koren A, Polak P, Nemesh J, Michaelson JJ, Sebat J, Sunyaev SR, McCarroll SA.  
1027 2012. Differential relationship of DNA replication timing to different forms of  
1028 human mutation and variation. *Am J Hum Genet* **91**: 1033-1040.

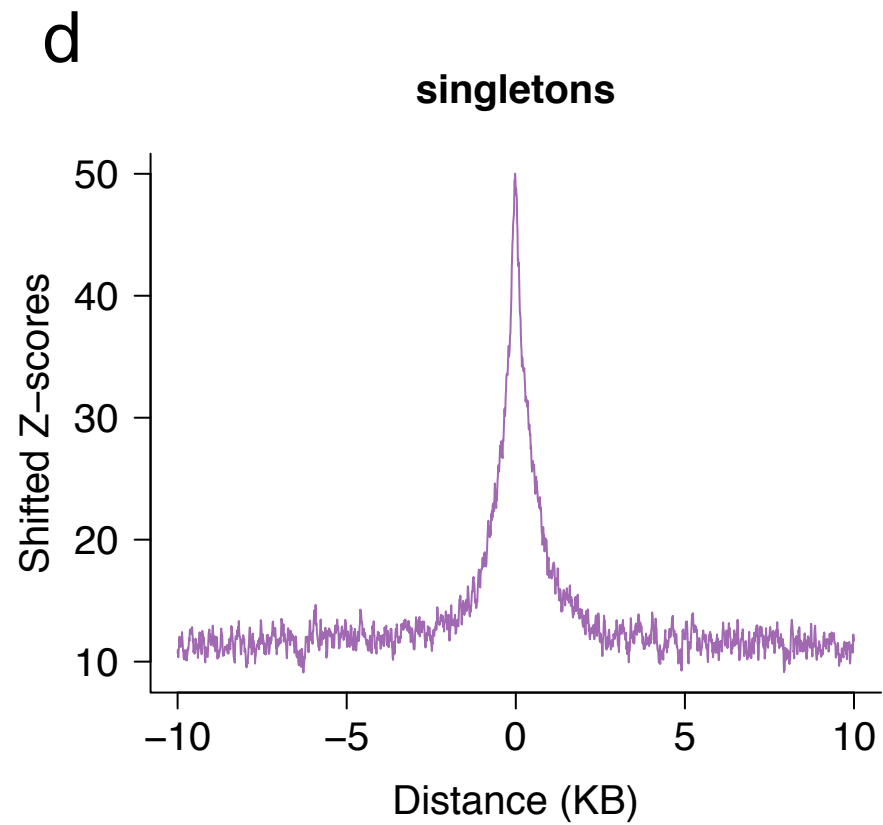
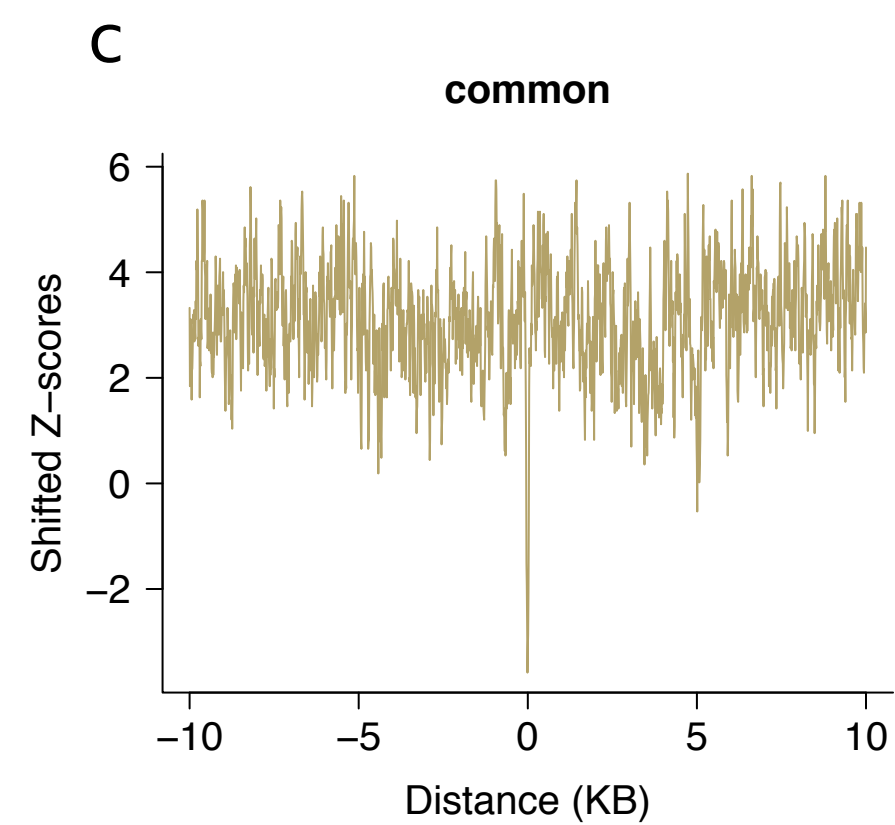
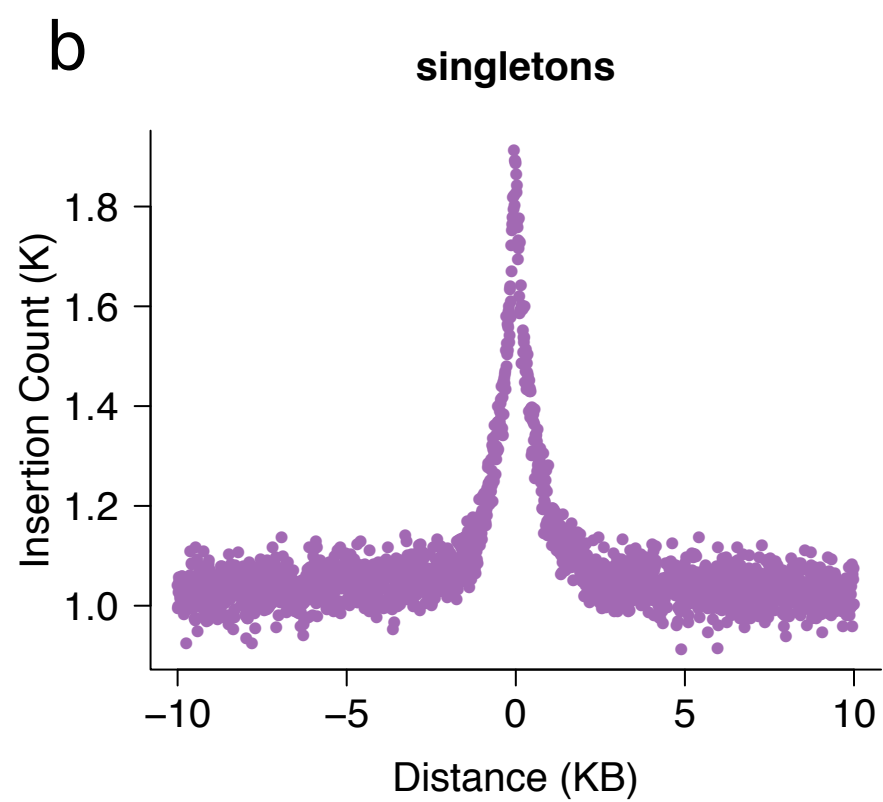
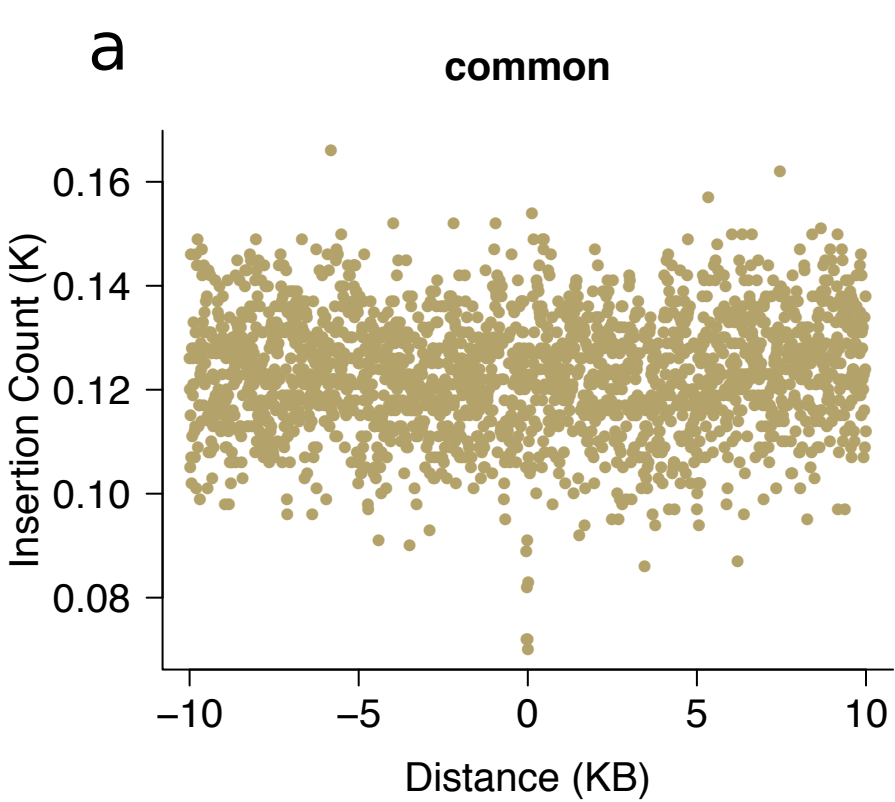
- 1029 Kribelbauer JF, Rastogi C, Bussemaker HJ, Mann RS. 2019. Low-Affinity Binding  
1030 Sites and the Transcription Factor Specificity Paradox in Eukaryotes. *Annu*  
1031 *Rev Cell Dev Biol* **35**: 357-379.
- 1032 Kvikstad EM, Tyekucheva S, Chiaromonte F, Makova KD. 2007. A macaque's-eye  
1033 view of human insertions and deletions: differences in mechanisms. *PLoS*  
1034 *Comput Biol* **3**: 1772-1782.
- 1035 Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat*  
1036 *Methods* **9**: 357-359.
- 1037 Leppa VM, Kravitz SN, Martin CL, Andrieux J, Le Caignec C, Martin-Coignard D,  
1038 DyBuncio C, Sanders SJ, Lowe JK, Cantor RM et al. 2016. Rare Inherited and  
1039 De Novo CNVs Reveal Complex Contributions to ASD Risk in Multiplex  
1040 Families. *Am J Hum Genet* **99**: 540-554.
- 1041 Levinson G, Gutman GA. 1987. High frequencies of short frameshifts in poly-CA/TG  
1042 tandem repeats borne by bacteriophage M13 in Escherichia coli K-12. *Nucleic*  
1043 *Acids Res* **15**: 5323-5338.
- 1044 Li C, Luscombe NM. 2020. Nucleosome positioning stability is a modulator of  
1045 germline mutation rate variation across the human genome. *Nat Commun* **11**:  
1046 1363.
- 1047 Li ZJ, Schulz MH, Look T, Begemann M, Zenke M, Costa IG. 2019. Identification of  
1048 transcription factor binding sites using ATAC-seq. *Genome Biology* **20**.
- 1049 Lieber MR, Ma Y, Pannicke U, Schwarz K. 2003. Mechanism and regulation of  
1050 human non-homologous DNA end-joining. *Nat Rev Mol Cell Biol* **4**: 712-720.
- 1051 MacDonald JR, Ziman R, Yuen RKC, Feuk L, Scherer SW. 2014. The Database of  
1052 Genomic Variants: a curated collection of structural variation in the human  
1053 genome. *Nucleic Acids Research* **42**: D986-D992.
- 1054 Makova KD, Yang S, Chiaromonte F. 2004. Insertions and deletions are male biased  
1055 too: a whole-genome analysis in rodents. *Genome Res* **14**: 567-573.
- 1056 Martin M. 2011. Cutadapt Removes Adapter Sequences From High-Throughput  
1057 Sequencing Reads. *EMBnetjournal* **17**: 10-12.
- 1058 Mcrae JF Clayton S Fitzgerald TW Kaplanis J Prigmore E Rajan D Sifrim A Aitken S  
1059 Akawi N Alvi M et al. 2017. Prevalence and architecture of de novo mutations  
1060 in developmental disorders. *Nature* **542**: 433-+.
- 1061 McVean G. 2007. What drives recombination hotspots to repeat DNA in humans?  
1062 *Philosophical Transactions of the Royal Society London Series B Biological*  
1063 *Sciences* **365**: 1213-1218.
- 1064 Messer PW. 2009. Measuring the Rates of Spontaneous Mutation From Deep and  
1065 Large-Scale Polymorphism Data. *Genetics* **182**: 1219-1232.
- 1066 Messer PW, Arndt PF. 2007. The majority of recent short DNA insertions in the  
1067 human genome are tandem duplications. *Mol Biol Evol* **24**: 1190-1197.
- 1068 Michaelson JJ, Shi Y, Gujral M, Zheng H, Malhotra D, Jin X, Jian M, Liu G, Greer  
1069 D, Bhandari A et al. 2012. Whole-genome sequencing in autism identifies hot  
1070 spots for de novo germline mutation. *Cell* **151**: 1431-1442.
- 1071 Montgomery SB, Goode DL, Kvikstad E, Albers CA, Zhang ZD, Mu XJ, Ananda G,  
1072 Howie B, Karczewski KJ, Smith KS et al. 2013. The origin, evolution, and  
1073 functional impact of short insertion-deletion variants identified in 179 human  
1074 genomes. *Genome Res* **23**: 749-761.
- 1075 Myers S, Bottolo L, Freeman C, McVean G, Donnelly P. 2005. A fine-scale map of  
1076 recombination rates and hotspots across the human genome. *Science* **310**: 321-  
1077 324.



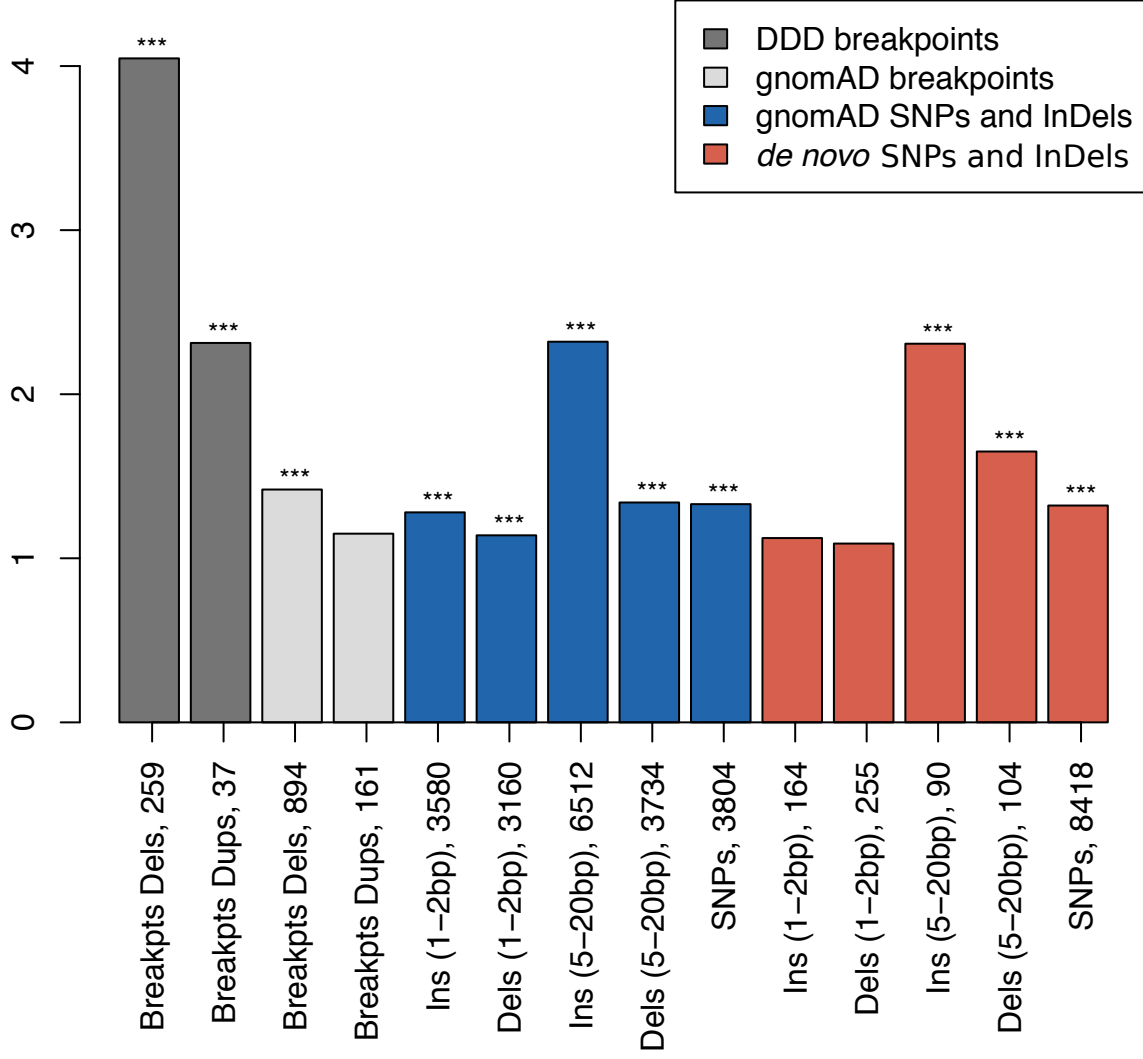
- 1078 Myers S, Bowden R, Tumian A, Bontrop RE, Freeman C, MacFie TS, McVean G,  
1079 Donnelly P. 2010. Drive against hotspot motifs in primates implicates the  
1080 PRDM9 gene in meiotic recombination. *Science* **327**: 876-879.
- 1081 Myers S, Freeman C, Auton A, Donnelly P, McVean G. 2008. A common sequence  
1082 motif associated with recombination hot spots and genome instability in  
1083 humans. *Nat Genet* **40**: 1124-1129.
- 1084 Neph S, Kuehn MS, Reynolds AP, Haugen E, Thurman RE, Johnson AK, Rynes E,  
1085 Maurano MT, Vierstra J, Thomas S et al. 2012. BEDOPS: high-performance  
1086 genomic feature operations. *Bioinformatics* **28**: 1919-1920.
- 1087 Palmer N, Talib SZA, Ratnacaram CK, Low D, Bisteau X, Lee JHS, Pfeiffenberger E,  
1088 Wollmann H, Tan JHL, Wee S et al. 2019. CDK2 regulates the NRF1/Ehmt1  
1089 axis during meiotic prophase I. *J Cell Biol* **218**: 2896-2918.
- 1090 Powers NR, Parvanov ED, Baker CL, Walker M, Petkov PM, Paigen K. 2016. The  
1091 Meiotic Recombination Activator PRDM9 Trimethylates Both H3K36 and  
1092 H3K4 at Recombination Hotspots In Vivo. *PLoS Genet* **12**: e1006146.
- 1093 Pratto F, Brick K, Khil P, Smagulova F, Petukhova GV, Camerini-Otero RD. 2014.  
1094 DNA recombination. Recombination initiation maps of individual human  
1095 genomes. *Science* **346**: 1256442.
- 1096 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing  
1097 genomic features. *Bioinformatics* **26**: 841-842.
- 1098 Ramirez F, Ryan DP, Gruning B, Bhardwaj V, Kilpert F, Richter AS, Heyne S,  
1099 Dundar F, Manke T. 2016. deepTools2: a next generation web server for deep-  
1100 sequencing data analysis. *Nucleic Acids Res* **44**: W160-165.
- 1101 Ramu A, Noordam MJ, Schwartz RS, Wuster A, Hurles ME, Cartwright RA, Conrad  
1102 DF. 2013. DeNovoGear: de novo indel and point mutation discovery and  
1103 phasing. *Nat Methods* **10**: 985-987.
- 1104 Reijns MAM, Kemp H, Ding J, de Proce SM, Jackson AP, Taylor MS. 2015.  
1105 Lagging-strand replication shapes the mutational landscape of the genome.  
1106 *Nature* **518**: 502-506.
- 1107 RK CY, Merico D, Bookman M, J LH, Thiruvahindrapuram B, Patel RV, Whitney J,  
1108 Deflaux N, Bingham J, Wang Z et al. 2017. Whole genome sequencing  
1109 resource identifies 18 new candidate genes for autism spectrum disorder. *Nat*  
1110 *Neurosci* **20**: 602-611.
- 1111 Rodgers K, McVey M. 2016. Error-Prone Repair of DNA Double-Strand Breaks. *J*  
1112 *Cell Physiol* **231**: 15-24.
- 1113 Sabarinathan R, Mularoni L, Deu-Pons J, Gonzalez-Perez A, Lopez-Bigas N. 2016.  
1114 Nucleotide excision repair is impaired by binding of transcription factors to  
1115 DNA. *Nature* **532**: 264-267.
- 1116 Sandelin A, Alkema W, Engstrom P, Wasserman WW, Lenhard B. 2004. JASPAR:  
1117 an open-access database for eukaryotic transcription factor binding profiles.  
1118 *Nucleic Acids Research* **32**: D91-D94.
- 1119 Sherwood RI, Hashimoto T, O'Donnell CW, Lewis S, Barkal AA, van Hoff JP, Karun  
1120 V, Jaakkola T, Gifford DK. 2014. Discovery of directional and nondirectional  
1121 pioneer transcription factors by modeling DNase profile magnitude and shape.  
1122 *Nat Biotechnol* **32**: 171-+.
- 1123 Short PJ, McRae JF, Gallone G, Sifrim A, Won H, Geschwind DH, Wright CF, Firth  
1124 HV, FitzPatrick DR, Barrett JC et al. 2018. De novo mutations in regulatory  
1125 elements in neurodevelopmental disorders. *Nature* **555**: 611-616.

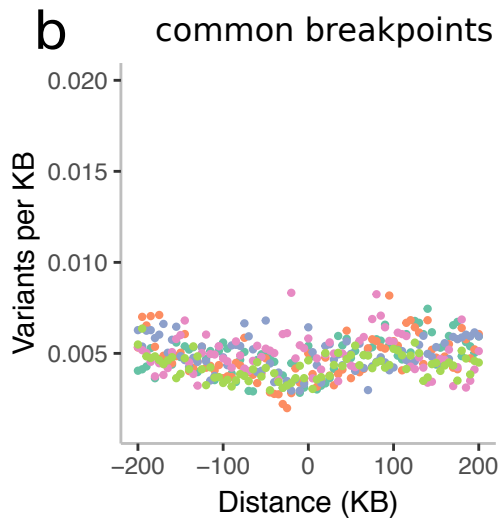
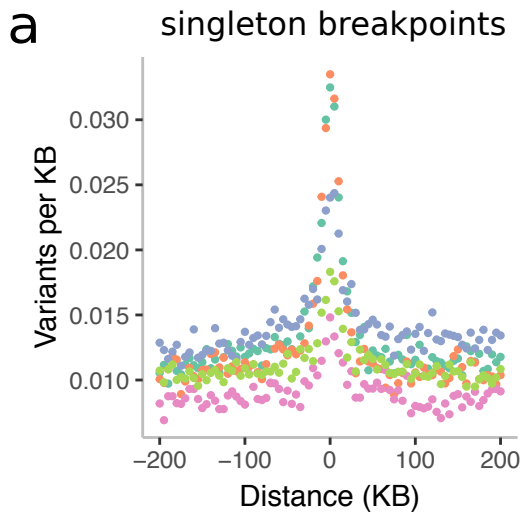
- 1126 Sohni A, Tan K, Song HW, Burow D, de Rooij DG, Laurent L, Hsieh TC, Rabah R,  
1127 Hammoud SS, Vicini E et al. 2019. The Neonatal and Adult Human Testis  
1128 Defined at the Single-Cell Level. *Cell Rep* **26**: 1501-1517 e1504.
- 1129 Stamatoyannopoulos JA, Adzhubei I, Thurman RE, Kryukov GV, Mirkin SM,  
1130 Sunyaev SR. 2009. Human mutation rate associated with DNA replication  
1131 timing. *Nat Genet* **41**: 393-395.
- 1132 The 1000 Genomes Project Consortium, Auton A, Brooks LD, Durbin RM, Garrison  
1133 EP, Kang HM, Korbel JO, Marchini JL, McCarthy S, McVean GA et al. 2015.  
1134 A global reference for human genetic variation. *Nature* **526**: 68-74.
- 1135 Turner TN, Coe BP, Dickel DE, Hoekzema K, Nelson BJ, Zody MC, Kronenberg ZN,  
1136 Hormozdiari F, Raja A, Pennacchio LA et al. 2017. Genomic Patterns of De  
1137 Novo Mutation in Simplex Autism. *Cell* **171**: 710-722 e712.
- 1138 Turner TN, Hormozdiari F, Duyzend MH, McClymont SA, Hook PW, Iossifov I,  
1139 Raja A, Baker C, Hoekzema K, Stessman HA et al. 2016. Genome Sequencing  
1140 of Autism-Affected Families Reveals Disruption of Putative Noncoding  
1141 Regulatory DNA. *Am J Hum Genet* **98**: 58-74.
- 1142 van Gent DC, Hoeijmakers JH, Kanaar R. 2001. Chromosomal stability and the DNA  
1143 double-stranded break connection. *Nat Rev Genet* **2**: 196-206.
- 1144 Wang J, Tang C, Wang Q, Su J, Ni T, Yang W, Wang Y, Chen W, Liu X, Wang S et  
1145 al. 2017. NRF1 coordinates with DNA methylation to regulate  
1146 spermatogenesis. *FASEB J* **31**: 4959-4970.
- 1147 Wellcome Trust Case Control C Craddock N Hurles ME Cardin N Pearson RD  
1148 Plagnol V Robson S Vukcevic D Barnes C Conrad DF et al. 2010. Genome-  
1149 wide association study of CNVs in 16,000 cases of eight common diseases and  
1150 3,000 shared controls. *Nature* **464**: 713-720.
- 1151 Werling DM, Brand H, An JY, Stone MR, Zhu L, Glessner JT, Collins RL, Dong S,  
1152 Layer RM, Markenscoff-Papadimitriou E et al. 2018. An analytical framework  
1153 for whole-genome sequence association studies and its implications for autism  
1154 spectrum disorder. *Nat Genet* **50**: 727-736.
- 1155 Wolfe KH, Sharp PM, Li WH. 1989. Mutation-Rates Differ among Regions of the  
1156 Mammalian Genome. *Nature* **337**: 283-285.
- 1157 Yan F, Powell DR, Curtis DJ, Wong NC. 2020. From reads to insight: a hitchhiker's  
1158 guide to ATAC-seq data analysis. *Genome Biol* **21**: 22.
- 1159 Yu G, Wang LG, He QY. 2015. ChIPseeker: an R/Bioconductor package for ChIP  
1160 peak annotation, comparison and visualization. *Bioinformatics* **31**: 2382-2383.
- 1161 Yuen RK, Merico D, Cao H, Pellecchia G, Alipanahi B, Thiruvahindrapuram B, Tong  
1162 X, Sun Y, Cao D, Zhang T et al. 2016. Genome-wide characteristics of de  
1163 novo mutations in autism. *NPJ Genom Med* **1**: 160271-1602710.
- 1164 Zhang Y, Liu T, Meyer CA, Eeckhoute J, Johnson DS, Bernstein BE, Nussbaum C,  
1165 Myers RM, Brown M, Li W et al. 2008. Model-based Analysis of ChIP-Seq  
1166 (MACS). *Genome Biology* **9**.
- 1167





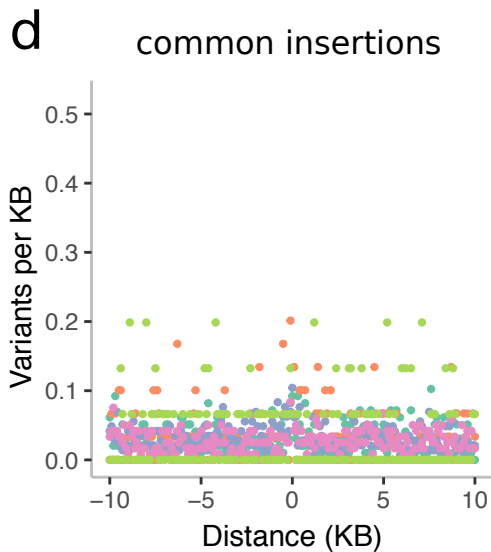
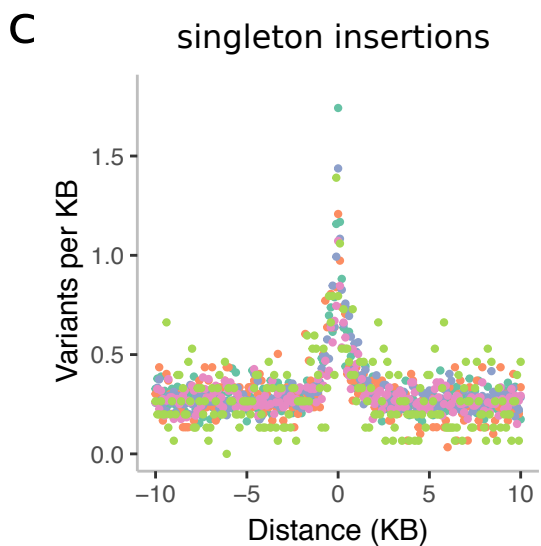
Observed/Expected





TF motifs

- motif\_171: NRF1
- motif\_579: PRDM9
- motif\_672: TFAP2 family
- motif\_963: NFYA/NFYB/Dux
- motif\_972: ETS family

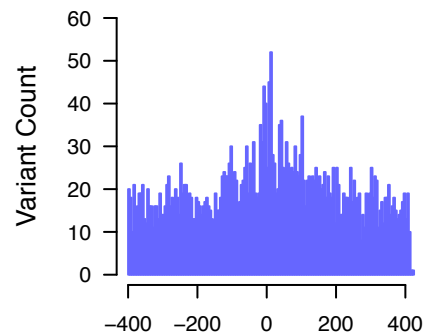


TF motifs

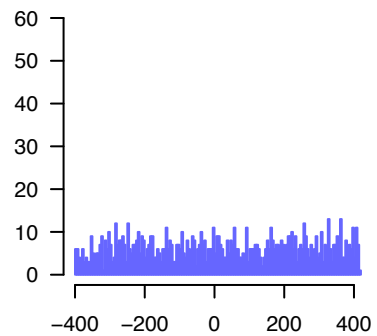
- motif\_171: NRF1
- motif\_92: HINFP
- motif\_579: PRDM9
- motif\_825: EGR family
- motif\_192: ZBTB33

# PRDM9

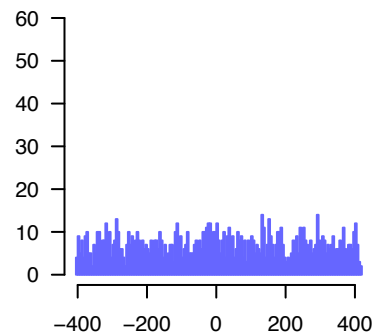
## insertions, 5–20 bp



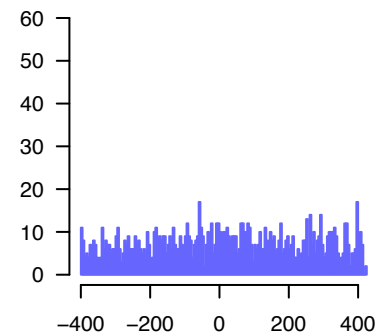
## insertions, 1–2 bp



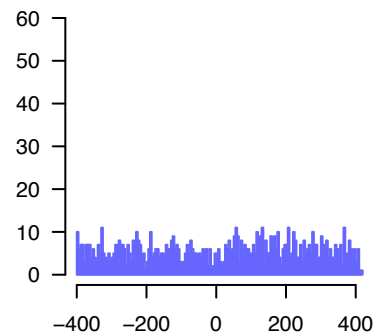
## snps



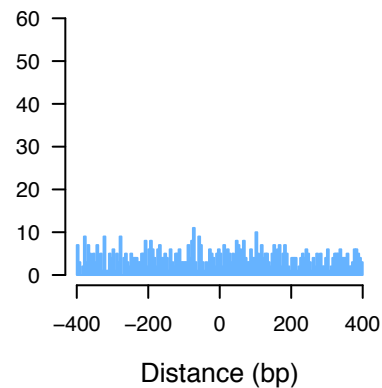
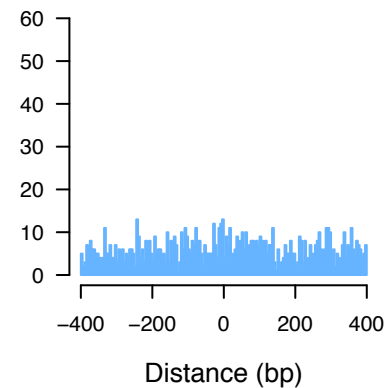
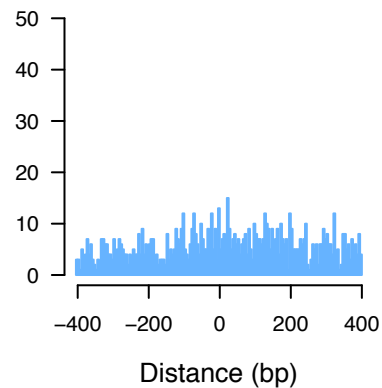
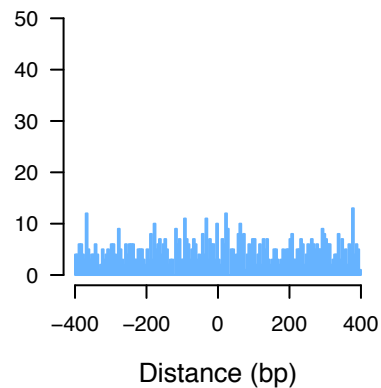
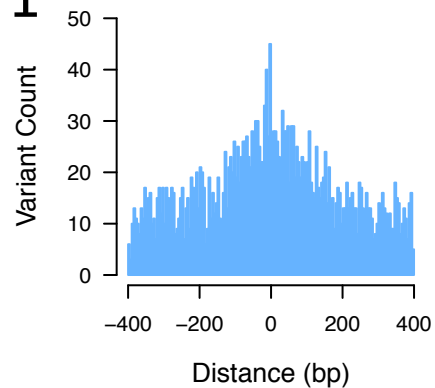
## deletions, 5–20 bp



## deletions, 1–2 bp

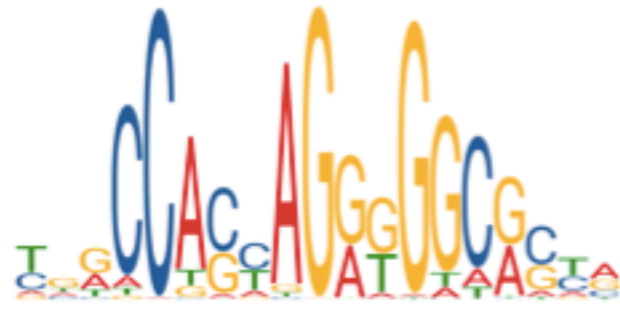


# NRF1



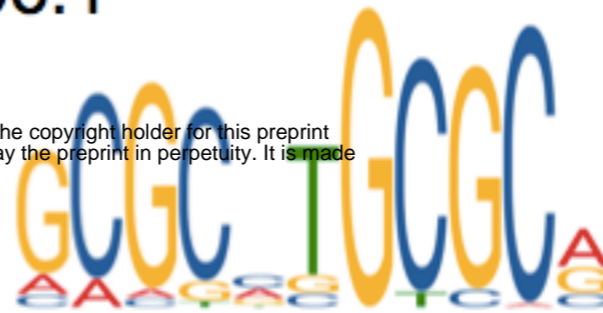
**a****Motif Family****MEME Motif within Insertions**

Motif 984: CTCF-MA0139.1 &amp; CTCFL-MA1102.1



N = 26

Motif 171: NRF1-MA0506.1



N = 9

bioRxiv preprint doi: <https://doi.org/10.1101/2021.06.10.447556>; this version posted June 10, 2021. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-ND 4.0 International license.

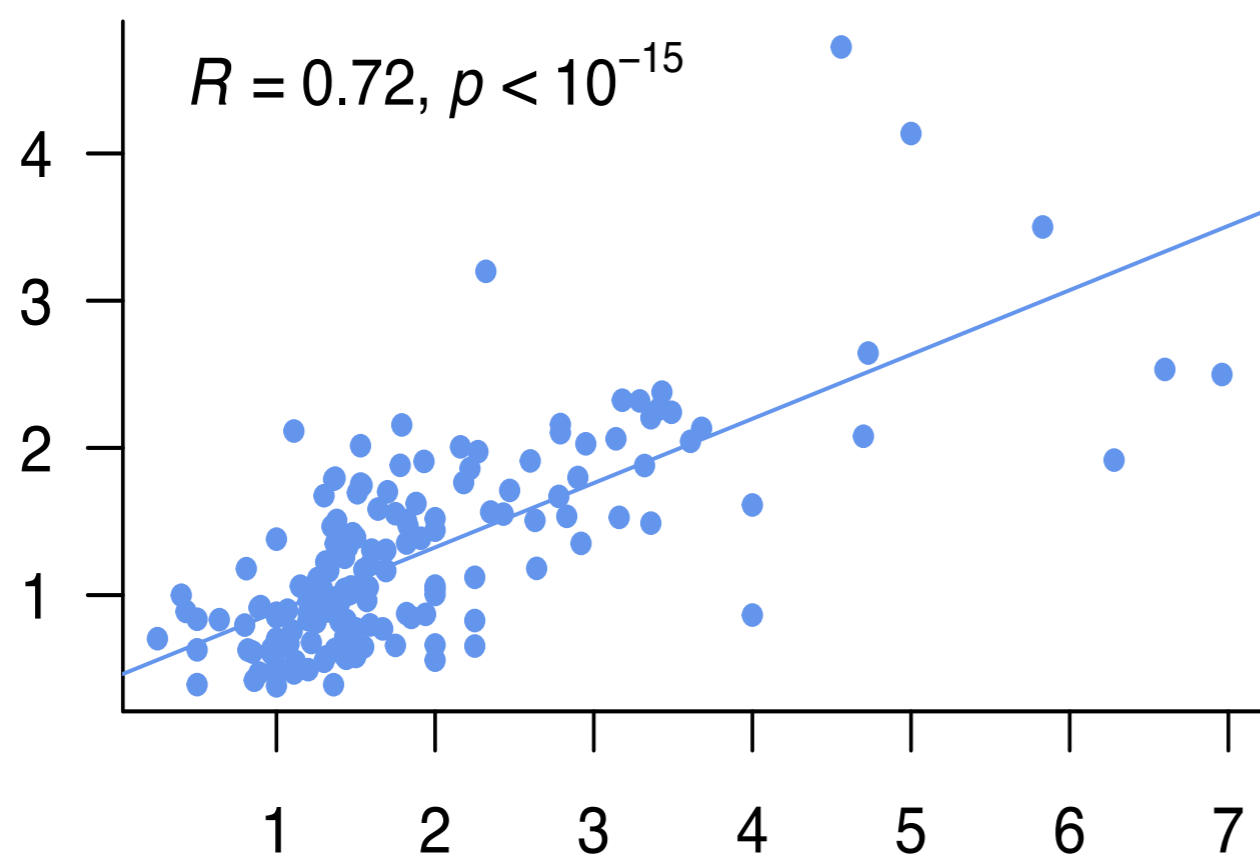
Motif 579: PRDM9



N = 91

**b**

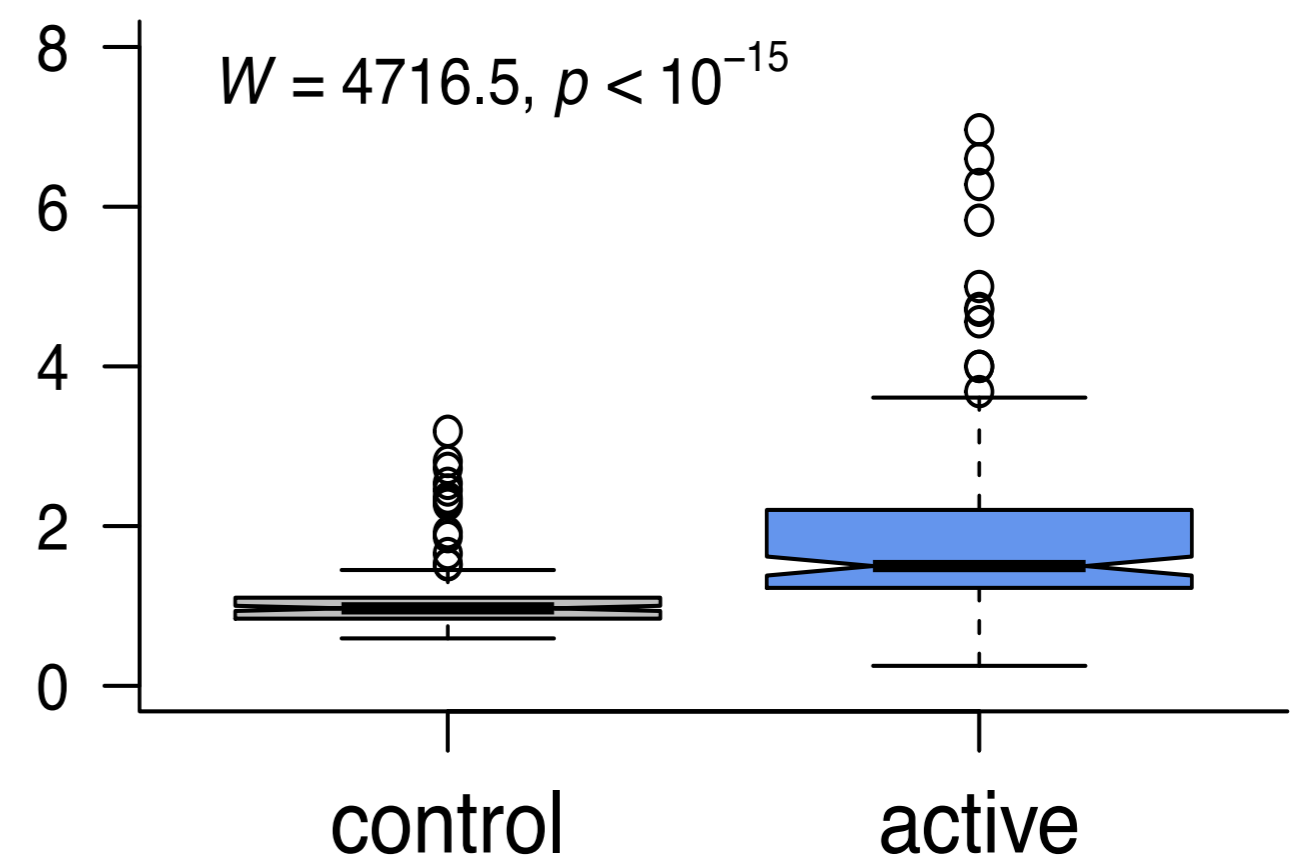
Motif Clustering



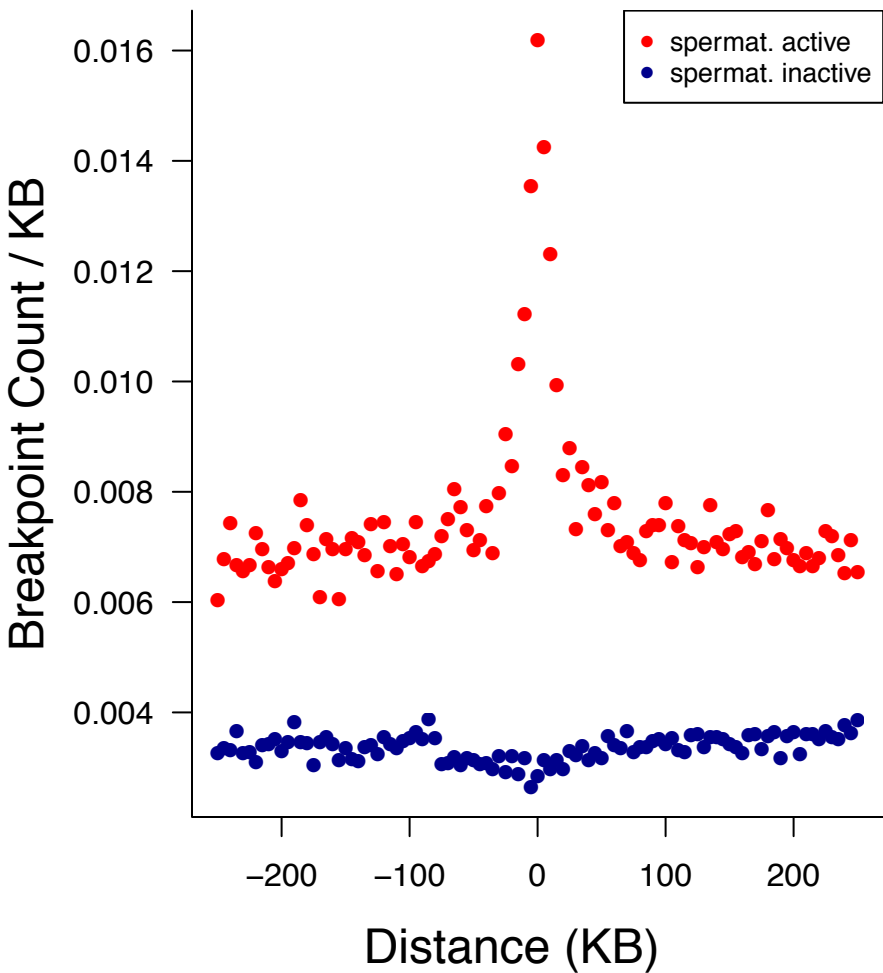
IFE

**c**

IFE





**a****b**