



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

Vera B Kaiser, Lana Talmame, Yatendra Kumar, et al.

Genome Res. published online August 20, 2021

Access the most recent version at doi:[10.1101/gr.275407.121](https://doi.org/10.1101/gr.275407.121)

P<P	Published online August 20, 2021 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
Creative Commons License	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution 4.0 International license), as described at http://creativecommons.org/licenses/by/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

CRISPR and RNAi Genetic Screening.
Your new superpower.

LEARN MORE

CELLECTA

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

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

3

4 Vera B. Kaiser¹, Lana Talmane¹, Yatendra Kumar¹, Fiona Semple¹, Marie
5 MacLennan¹, Deciphering Developmental Disorders Study^{1,2}, David R. FitzPatrick¹,
6 Martin S. Taylor^{1*}, Colin A. Semple^{1*}

7 ¹MRC 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 ²The 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 NRF1, are also associated with
30 increased rates of SV breakpoints, apparently independent of recombination.
31 Singleton short insertions (≥ 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 (>10⁶ 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 2015). The method ATAC-seq (Buenrostro et al. 2013) reports chromatin
73 accessibility and provides a snapshot of all active regulatory regions and occupied
74 binding sites in a given tissue. In particular, ATAC-seq footprinting (Sherwood et al.
75 2014; Li et al. 2019), when applied to spermatogonia, has the potential to reveal the
76 binding of hundreds of TFs, as well as PRDM9, in the male germline. In addition,
77 large human genome sequencing projects can be used to reveal patterns of mutation
78 rates, by focussing on extremely rare variants (Messer 2009; Carlson et al. 2018; Li
79 and Luscombe 2020). Making use of such variant datasets as well as novel ATAC-seq
80 data in spermatogonia, we study the mutational landscape at transcription factor
81 binding sites (TFBSs) in accessible human spermatogonial chromatin.

82

83 **Results**

84

85 **Spermatogonial regulatory regions are enriched for rare deletion breakpoints**

86

87 We used ATAC-seq to identify open chromatin sites in FGFR3-positive
88 spermatogonial cells isolated from dissociated human testicular samples. *FGFR3* is
89 most highly expressed in self-renewing spermatogonial stem cells, with low
90 expression also being detected in early differentiating spermatogonia (Guo et al.
91 2018; Sohni et al. 2019); its expression thus overlaps with the onset of *PRDM9*
92 expression in pre-meiotic spermatogonia (Human Protein Atlas:
93 <https://www.proteinatlas.org/ENSG00000164256-PRDM9/celltype/testis> and
94 <https://www.proteinatlas.org/ENSG00000068078-FGFR3/celltype/testis>) (Guo et al.

95 2018). Open chromatin in FGFR3-positive cells was identified using standard peak
96 detection analysis (Methods; Supplemental Datasets 1-3), and multiple metrics
97 (Supplemental Fig. S1A-C) indicated high data quality (Yan et al. 2020). Hierarchical
98 clustering (Ramirez et al. 2016) showed that this novel spermatogonial ATAC-seq
99 dataset displays a genome-wide distribution of peaks consistent with other
100 spermatogonial derived data, and is distinct from ES cell and somatic tissue datasets
101 (Supplemental Fig. S2).

102

103 We assessed the enrichments of different classes of sequence variants at
104 spermatogonial active sites, including singleton SV breakpoint frequencies as a proxy
105 for the mutation rate of such variants. We made use of ultra-rare genomic variants
106 from a variety of human sequencing studies: the Deciphering Developmental
107 Disorders (DDD) study (Deciphering Developmental Disorders Study 2015; McRae
108 et al. 2017) of severe and undiagnosed developmental disorders
109 (<https://www.ddduk.org/>), a large collection of variants from an aggregated database
110 (gnomAD; <http://gnomad.broadinstitute.org/>), and *de novo* variants from trio
111 sequencing studies (<http://denovo-db.gs.washington.edu/>, <https://research.mss.ng/>, An
112 et al. (2018)). Based on the DDD dataset - a combination of high-density arrayCGH
113 and exome sequencing (Deciphering Developmental Disorders Study 2015) - we
114 identified 6,704 singleton deletion variants among 9,625 DDD probands (carrier
115 frequency of ~ 0.002% in the combined dataset) (Supplemental Table S1).

116 Permutation analysis demonstrates that DDD singleton breakpoints are enriched at
117 spermatogonial ATAC-seq sites, their overlap being > 4-times the expected genome-
118 wide rate (Supplemental Table S2), and shifted permutation Z-scores reveal that the
119 enrichment is specific to the ATAC-seq peaks as opposed to wider genomic regions

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

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

138

139 **Locally elevated mutation at spermatogonial TFBSs**

140

141 Compared to larger structural variants, such as those (up to megabase sized) deletions
142 examined above, indels have been shown to occur at a higher rate of about 6 new
143 variants per genome and generation (Besenbacher et al. 2016). Short indels (≤ 4 bp)
144 are thought to arise due to replication slippage (Levinson and Gutman 1987;

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

170 gnomAD singleton variants, with a moderate (~10-60%) increase of mutation rates at
171 TFBSs for all categories of short 1-2bp sequence variants, a but larger increase of
172 ~130% for insertions of 5-20 bp (Figure 3). These results were confirmed using a set
173 of independent positive and negative control sites (Supplemental Fig. S3A, B). We
174 conclude that regulatory sites that are active in spermatogonia show unusual parallel
175 enrichments for both large SV breakpoints and 5-20 bp insertions, consistent with
176 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, and, similarly, 29% (48/167)
189 of motif families were significantly enriched for gnomAD singleton deletion
190 breakpoints (Bonferroni corrected $p = 0.017$); no motif family was found to be
191 depleted for breakpoints in either dataset (Supplemental Tables S3 and S4),
192 suggesting that increased load is a common feature of TFBSs bound by different
193 transcription factors in the germline. Similarly, singleton 5-20bp insertions from the
194 gnomAD database were found to be significantly enriched at 29% (48/167) of

195 families (Bonferroni corrected $p = 0.017$) and, nominally, 84% (140/167) of families
196 showed enrichment for these insertions (Supplemental Table S4). Again, no TFBS
197 family was found to be depleted for these rare variants. Collectively, these results
198 suggest that TFBSs active in spermatogonia incur locally elevated burdens of short
199 insertions and large structural variants across many different binding motifs.

200

201 Certain motif families appear to carry notably higher mutational loads than the
202 general disruption seen across all TFBSs. Based on the insertion fold enrichment
203 (IFE), i.e. the ratio of the observed to expected numbers of insertions (5-20 bp),
204 PRDM9 binding sites are among the most disrupted sites in the genome (IFE = 6.3),
205 and this also holds for PRDM9 sites outside known sites of recombination (IFE = 6.7
206 for 8,139 PRDM9 sites with a distance of at least 500bp from HSs and ssDNA sites,
207 respectively). PRDM9 sites are similarly associated with higher rates of singleton
208 deletion breakpoints (Figure 4A, C), in line with the roles of PRDM9 during
209 recombination, though PRDM9 sites outside known sites of recombination also show
210 this trend (observed overlaps with deletion breakpoints = 9; expected = 1; $p < 10^{-4}$).
211 Two other TFBS families, corresponding to NRF1 (Nuclear Respiratory Factor 1;
212 IFE=7.0) and HINFP (IFE=6.6) exceed the disruption seen at PRDM9 sites, and
213 NRF1 sites are also disrupted at high rates according to DDD and gnomAD
214 breakpoint data (Supplemental Tables S3 and S4). Shifted Z-scores for the enrichment
215 of short insertions (5-20 bp) at both NRF1 and PRDM9 binding sites are in the top
216 four, next to SP/KLF transcription factors (motif families 938 and 992), suggesting
217 strong focal enrichments at these sites (Supplemental Tables S6 and S7). *NRF1* has
218 been shown to be an important testis-expressed gene with meiosis-specific functions
219 (Wang et al. 2017; Palmer et al. 2019), but NRF1 binding sites have, to our

220 knowledge, not been reported to be foci for genomic instability. We find similar
221 enrichments of short insertions (5-20 bp) at TFBSs in SSEA4- and KIT-marked
222 spermatogonial samples produced in previous ATAC-seq studies (Guo et al. 2017;
223 Guo et al. 2018). Reprocessing these previous datasets identically to our own reveals
224 that PRDM9, NRF1 and HINFP sites are again among the top 5 disrupted motif
225 families (Supplemental Tables S8 and S9).

226

227 Although both PRDM9 and NRF1 binding sites are GC-rich, their modest motif
228 similarity suggests that the two factors occupy distinct binding motifs (PWMclus:
229 Pearson's correlation distance $r = 0.35$ for PRDM9 *versus* and NRF1) and should not
230 converge on the same sites. However, in practice, PRDM9 and NRF1 binding sites
231 were often found within the same regulatory regions, such that many (1,199) ATAC-
232 seq peaks contained both the NRF1 and PRDM9 binding motifs. The disruption of
233 motifs within these co-bound peaks was notably higher, with NRF1-motifs being
234 disrupted by short insertions 10.8-fold the expected rate (observed: 108; expected:
235 10), and PRDM9-motifs 11.2-fold the expected rate (observed: 146; expected: 13)
236 when co-occurring with the other factor ($p < 10^{-4}$ in each case). Similarly, 1,311
237 ATAC-seq peaks contained a motif for both CTCF and PRDM9, and CTCF motifs in
238 these peaks were more highly disrupted by short insertions (ratio = 6.3; observed: 69;
239 expected: 11) compared to all CTCF motifs (Supplemental Table S4), as was PRDM9
240 (ratio = 8.2; observed: 115; expected: 14) ($p < 10^{-4}$ in each case).

241 Importantly, the excess of insertions observed at particular motif sites is not a trivial
242 consequence of statistical power (i.e. the number of TFBSs in the genome); for
243 example, the number of binding sites identified for PRDM9 and NRF1 is fewer than
244 many other factors (< 10,000 sites each; Supplemental Tables S3 and S4).

245 In general, mutational loads appear to be dependent on the level of chromatin
246 accessibility (MACS2 peak scores (Zhang et al. 2008)) and the number of factors
247 predicted to bind at ATAC-seq defined regulatory regions, such that regions in the
248 upper quartile of accessibility that are also occupied by more than 4 factors incur the
249 highest indel loads (Supplemental Fig. S4A-D). The significant positive correlation
250 between the rates of binding site disruption via singleton insertions and deletion
251 breakpoints across all motif families (Spearman's $R = 0.52$, $p < 10^{-5}$; Supplemental
252 Fig. S5A-C) suggests that the two types of damage may be mechanistically linked. In
253 support of this idea, singleton short insertions (5-20 bp) and singleton SV deletion
254 breakpoints overlap at the exact nucleotide position more often than expected
255 (genome-wide Z -score = 26.31; $p < 10^{-4}$; see also Supplemental Fig. S6). This overlap
256 is unlikely to be due to erroneous variant calling in the singleton dataset since we
257 observe similar patterns for common variants of the same variant categories (genome-
258 wide Z -score = 62.9, $p < 10^{-4}$).

259

260 **Short insertions generate clustered binding sites within regulatory regions**

261

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

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

284 These unusual patterns of clustered TFBSs at indel breakpoints appear to be specific
285 to spermatogonial ATAC-seq peaks, and do not reflect genome-wide trends. Applying
286 the MEME-ChIP algorithm on 50bp regions flanking singleton insertion and deletion
287 breakpoints, we were able to re-discover the sequence motifs of commonly disrupted
288 binding sites, including the motifs of PRDM9 and NRF1 (Supplemental Table S10).
289 In contrast, genome-wide, the motifs discovered flanking these variants were more
290 likely to be simple repeats and other low complexity sequences that did not match
291 known TFBS motifs, suggesting that processes other than transcription factor binding
292 drive DNA breakage outside of active regulatory sites.

293

294 **Genomic instability at spermatogonial TFBSs impacts enhancers active in neural**
295 **development**

296

297 Since many regulatory regions of the genome are active across a variety of cell types
298 (Andersson et al. 2014), mutation at TFBSs in spermatogonia might disrupt gene
299 regulation in other tissues. The developing brain is of particular interest, given reports
300 of increased SV burdens in neurodevelopmental disorders (Girirajan et al. 2011;
301 Leppa et al. 2016; Collins et al. 2017). We classified developmentally active human
302 brain enhancers (distal regulatory elements) supported by neocortical ATAC-seq data
303 (de la Torre-Ubieta et al. 2018) according to whether they were either active (10,888
304 brain enhancers) or inactive in the male germline (26,162 brain enhancers). We then
305 calculated the odds ratio of a singleton mutation affecting a brain enhancer which is
306 also *active* in spermatogonia, relative to a brain enhancer which is *inactive* in
307 spermatogonia. For DDD singleton deletion breakpoints, the odds ratio was 6.82
308 (95% CI = [5.34,8.71]), and for a singleton gnomAD insertion (5-20 bp), it was 4.69
309 (95% CI = [4.46,4.93]). This suggests that activity in spermatogonia greatly
310 predisposes a brain enhancer to DNA damage, and this damage manifests in
311 enhancers that share activity with the male germline (Figure 7A, B). Brain enhancers
312 that are shared with spermatogonia are, on average, more accessible in the developing
313 brain than those that are inactive in the germline (the median “mean of normalized
314 counts” for the two types of brain enhancers were 104.8 and 54.1, respectively;
315 Wilcoxon test $W = 197340000$, $p\text{-value} < 10^{-15}$), suggesting a link between enhancer
316 activity, the sharing of enhancers across tissues and propensity to mutation. The
317 subset of brain enhancers which overlapped spermatogonial active sites were not
318 enriched for specific motifs, and the number of motif sites for each motif family were

319 highly correlated between brain and spermatogonia (Spearman's $\rho = 0.95$, $p < 10^{-15}$). That is, the propensity to mutation does not appear to be driven by an enrichment
320 of specific motif families in brain enhancers.
321

322

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

324

325 We cannot exclude a small contribution of the TFBS sequence itself on the
326 predisposition to mutation (Kondrashov and Rogozin 2004), but our data suggest that
327 TF binding is a major driver of insertion and deletion mutation in the human
328 germline. This is supported by the fact that we see an increase of disruption of brain
329 enhancers if they are active in spermatogonia (Figure 7) and, more generally, an
330 increase in the mutational load for sites that are active across other somatic tissue if
331 binding also occurs in the germline (Supplemental Table S11). In addition, control
332 motif sites (representing the same TFBS but located outside of ATAC-seq peaks) are
333 subject to lower rates of mutation compared to motifs within spermatogonial ATAC-
334 seq peaks (Figure 6C). Motifs within peaks carry, on average, 73% more mutations
335 than their control counterparts, and for the most highly disrupted motifs, the
336 discrepancy between active and control motifs is even larger. For example, PRDM9
337 motifs are 3.4-fold, HINFP 2.9-fold and NRF1 motifs 2.6-fold more disrupted if they
338 are active in spermatogonia, relative to spermatogonia inactive motifs. We note that
339 this increase in disruption is likely to be a conservative estimate since some control
340 sites may be bound at time points in the germline that our ATAC-seq data cannot
341 ascertain.

342 Since the X Chromosome spends only one third of its time in males - the sex with the
343 higher number of germ cell divisions - a depletion of mutations on the X

344 Chromosome is expected for a male-biased mutational process. We find the X
345 Chromosome to be strongly depleted for short singleton gnomAD insertions (5-20
346 bp), with a ratio of X to autosome variants per uniquely mappable site of 0.78
347 (Supplemental Table S12). However, we note that, despite the overall reduced rate of
348 insertions on the X, ATAC-seq peaks on the X are still subject to increased rates of
349 insertions compared to genome-wide expectations, suggesting that the inferred effects
350 of protein-binding on mutation are larger than the reduction in mutation due to X-
351 linkage (38 observed insertions in X-linked ATAC-seq peaks, whereas 11 were
352 expected; $p < 10^{-4}$).

353

354 To test which candidate genomic feature most reliably predicts DNA damage, we
355 used random forest regression to model the rate of singleton variants within 5 kb
356 genomic windows, based on their overlap with spermatogonial TFBSs, ssDNA sites,
357 LD-based hotspots, average GC content, mappability, gene density, replication time
358 as well as various repeat families (LTRs, SINEs, LINEs and simple repeats). In
359 models of genome-wide short insertion rates or deletion breakpoint rates, measures of
360 replication timing and GC content were important predictors of mutation load as
361 expected (Supplemental Fig. S9). Mappability was an important factor for predicting
362 mutation rates genome-wide, perhaps reflecting the association between segmentally
363 duplicated (low mappability) regions and rapid structural evolution, or perhaps
364 suggesting that a fraction of variants may be erroneously called in the gnomAD
365 dataset. (Only regions with high mappability were included in our more detailed
366 analyses of TFBSs (Figures 3-7 and Supplemental Fig. S7)). However,
367 spermatogonial ATAC-seq derived TFBSs contributed additional predictive power to
368 the models, even at the scale of the entire genome. The same TFBSs appear to be

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

378

379 **Discussion**

380

381 We have demonstrated enrichments of rare and *de novo* SV breakpoints at
382 spermatogonial regulatory sites defined by ATAC-seq, suggesting that these sites
383 suffer high rates of DSBs in the male germline. The same sites show unusual parallel
384 enrichments for short variants, and particularly 5-20bp insertions. These loads appear
385 to be positively correlated with the levels of chromatin accessibility/nucleosome
386 disruption (ATAC-seq peak binding strength) and the number of factors predicted to
387 bind within the region. These results have implications for the evolution of binding
388 site patterns within regulatory regions, and for disrupted regulation in somatic tissues.

389

390 Homotypic clusters of TFBSs are a pervasive feature of both invertebrate and
391 vertebrate genomes, and have long been known to be a common feature of human
392 promoter and enhancer regions (Gotea et al. 2010). Various adaptive hypotheses have
393 been proposed for the presence of such clusters such that they provide functional

394 redundancy within a regulatory region, enable the diffusion of TF binding across a
395 region, and allow cooperative DNA binding of TF molecules (Gotea et al. 2010).
396 More recently it has been suggested that homotypic TFBS clusters may also
397 contribute to phase separation and the compartmentalisation of the nucleus
398 (Kribelbauer et al. 2019). Similarly, the clustered patterns of CTCF sites in the
399 genome have been ascribed critical roles in chromatin architecture and regulation,
400 particularly at regulatory domain boundaries. However, these boundary regions have
401 been shown to exhibit genome instability (Kaiser and Semple 2018) and recurrently
402 acquire new CTCF binding sites in dynamically evolving clusters (Kentepozidou et
403 al. 2020). The data presented here suggest that binding site clusters may arise solely
404 as a selectively neutral consequence of the unusual mutational loads at germline
405 TFBSs, with clusters maintained by recurrent DNA damage and mis-repair.
406
407 We observe significant enrichments of both large SV breakpoints and small insertions
408 together at spermatogonial TFBSs. This parallel enrichment may originate from DNA
409 breakage, followed by misrepair, conceivably via a pathway such as non-
410 allelic homologous recombination (NAHR). It is known that NAHR can create large
411 insertions and deletions (Kim et al. 2016), and PRDM9 activity is implicated in
412 certain developmental disorders arising via NAHR (McVean 2007; Myers et al. 2008;
413 Berg et al. 2010). For example, the locations of PRDM9 binding hotspots coincide
414 with recurrent SV breakpoints causing Charcot-Marie-Tooth disease, and Hunter and
415 Potocki-Lupski/Smith-Magenis syndromes (Pratto et al. 2014). It is possible that the
416 sequence similarity at TFBSs scattered across the genome may make them
417 particularly prone to NAHR. However, the sequence similarity between the low copy
418 repeat units, known to be involved in NAHR, is usually of the size of several kb (Gu

419 et al. 2008), rather than sequences on the scale of TFBSs. The NHEJ pathway can
420 also lead to short insertions after DNA breakage, usually in G0 and G1 phases of the
421 cell cycle. Indeed, NHEJ is the most common repair pathway of DSBs in mammals
422 and it is typically error prone (van Gent et al. 2001; Lieber et al. 2003). During NHEJ,
423 double-strand break ends are resected to form single-stranded overhangs, but when
424 pairing occurs between the tips of the overhangs, sequences near the breakpoints will
425 often be duplicated (Rodgers and McVey 2016). Two previous studies using human–
426 chimpanzee–macaque multiple alignments have shown that high numbers of short
427 insertions have occurred in the human lineage (Kvikstad et al. 2007; Messer and
428 Arndt 2007), and both conclude that these insertions preferentially take place in the
429 male germline, evidenced by decreased mutation rates on the X Chromosome, with
430 similar observations in rodents (Makova et al. 2004).

431

432 The data presented here suggest that different DNA-binding proteins differ widely in
433 their impact on mutation rates. The two proteins with the largest impacts, NRF1 and
434 PRDM9, are both highly expressed in testis, revealing a possible link between the
435 expression level of a gene encoding a DNA-binding protein and the propensity for
436 breakage or inefficient repair at the sites the protein binds. Incidentally, *NRF1* has a
437 pLI score of 0.999, indicating that it is extremely loss-of function intolerant and
438 crucial for the organism’s functioning (Karczewski et al. 2020). A previous study
439 (Montgomery et al. 2013), using 1000 Genomes polymorphism data, failed to find an
440 increase in indels at PRDM9 motifs genome-wide. This highlights the importance of
441 using ATAC-seq data to confine the search for motifs to germline active sites only,
442 combined with singleton variants from large-scale sequencing studies as a more
443 powerful strategy to explore fine scale mutational patterns.

444

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

460

461 **Methods**

462

463 **Identification of spermatogonial binding sites**

464

465 Samples of testicular tissue were obtained from three patients undergoing
466 orchiectomy with total processing completed within ~5-7 hours of explant. Tissue
467 was obtained after informed consent through the Lothian NRS BioResource, and the
468 study was approved by NHS Lothian (Lothian R&D Project Number 2015/0370TB).

469 Tissue samples were disaggregated into cells, and cells were labelled with
470 phycoerythrin (PE)-conjugated antibody against the cell surface marker FGFR3
471 (FAB766P, clone 136334, R&D systems). Spermatogonial cells were isolated using a
472 FACSaria II cell sorter (BD bio- sciences) based on PE fluorescence and cell shape,
473 according to Forward/Side Scatter. Isolated cells were subjected to ATAC-seq using
474 the protocol and reagents described in (Buenrostro et al. 2013), followed by paired-
475 end sequencing on Illumina HiSeq 4000 (75 bp read length). We combined reads
476 from separate sequencing runs into three biological replicates, based on origin and
477 morphological appearance of the FACS sorted cells. Replicate 1: combined
478 sequencing runs H.5.1 and H.5.4; a non-cancer patient; large cells, high side scatter;
479 58,000 and 42,000 cells, respectively. Replicate 2: combined sequencing runs H.5.2
480 and H.5.5; the same non-cancer patient as Replicate 1; large cells; 36,000 and 23,000
481 cells, respectively. Replicate 3: combined sequencing runs H.7.3 and H.10.2; normal
482 tissue from cancer patients; large cells; 69,000 and 24,000 cells, respectively. Raw
483 reads were processed and ATAC-Seq peaks called as described in the Supplemental
484 Methods. For the downstream mutation analyses, ATAC-seq peaks from Replicates 1
485 and 2 (the non-cancer patient) were merged, creating a single peak set. This dataset
486 also formed the basis for the footprinting analysis, which used, as input, the combined
487 short sequencing fragments of Replicates 1 and 2, running “rgt-hint footprinting” with
488 --atac-seq and --bias-correction, followed by “rgt-motifanalysis matching” with the
489 option --remove-strand-duplicates (Li et al. 2019). Input motifs were the 579 position
490 weight matrixes (PWMs) of the JASPAR vertebrate database (Sandelin et al. 2004) as
491 well as the 13-mer PRDM9 motif “CCNCCNTNNCCNC” (Myers et al. 2010) which
492 was also provided as a PWM. The tissue donor for Replicates 1 and 2 was a carrier of
493 the most common (European) alleles of PRDM9, which was confirmed by

494 investigating his allelic state at the SNP (rs6889665) identified by Hinch et al. (2011);
495 this SNP was covered by our ATAC-seq by 10 reads, all of which were “T”.

496 Accordingly, we assume that the donor is a carrier of the A and/or B allele of PRDM9
497 (both of which bind the same DNA motif), and the search for the 13-mer PRDM9
498 motif in this patient’s ATAC-seq data can be used as a proxy for PRDM9 binding in
499 European populations. In addition, Replicate 3 was processed in the same way as the
500 combined Replicates 1 and 2 and served as a positive control to assess the genome-
501 wide enrichment of mutations at spermatogonial accessible sites (Supplemental Fig.
502 S3).

503

504 JASPAR input motifs are often highly similar, resulting in multiple binding proteins
505 being identified by the rgt-hint pipeline to bind at the same ATAC-seq footprint; this
506 is biologically implausible (since only one protein is likely to occupy a given site),
507 and we clustered motifs by similarity, using the default parameters of the PWMclus
508 CCAT package (Jiang and Singh 2014). This resulted in a set of 167 motif families of
509 similar binding motifs (Supplemental Table S7). Using BEDTools (Quinlan and Hall
510 2010), we merged overlapping binding sites that belonged to motifs of the same
511 family (thus calling them only once), and we also merged palindromic binding sites
512 called on both strands. Since PRDM9 is known to leave a characteristic histone
513 methylation mark on bound DNA (Grey et al. 2011; Powers et al. 2016), we
514 intersected the PRDM9 motif sites with testis-derived H3K4me3 marks (called in an
515 PRDM9 A/B heterozygous individuals) from Pratto et al. (2014). This resulted in a
516 stringent set of PRDM9 sites, which were both located in ATAC-seq footprints and
517 also carried the H3K4me3 mark in human testis. ATAC-seq-defined PRDM9 sites
518 showed moderate overlap with DMC1-bound ssDNA sites (Pratto et al. 2014) as well

519 as recombination HSs (Myers et al. 2005), which may reflect the fact that most cells
520 in our experiments are likely to be pre-meiotic: only 10 and 11% of PRDM9 sites
521 were within 500 bp of a ssDNA peak and a recombination HS, respectively, whereas
522 44% of DMC1-bound sites overlap with LD-defined HSs. However, we find that
523 stronger ssDNA peaks are more likely to be near a PRDM9-binding site
524 (Supplemental Fig. S10).

525

526 **Comparisons between ATAC-seq datasets**

527

528 Using the same procedure as described in the Supplemental Methods, we processed
529 raw ATAC-seq reads from previously published datasets in order to call MACS2
530 peaks from short sequencing fragments (Zhang et al. 2008). Datasets included ATAC-
531 seq reads from the germinal zone and cortical plate of the developing brain
532 (SRR6208926, SRR6208927, SRR6208938, SRR6208943) (de la Torre-Ubieta et al.
533 2018), ATAC-seq experiments of KIT+ spermatogonia (sra accessions SRR7905001
534 and SRR7905002) (Guo et al. 2018), SSEA4+ spermatogonia (SRR5099531,
535 SRR5099532, SRR5099533, SRR5099534) (Guo et al. 2017) and ESC cells
536 (SRR5099535 and SRR5099536) (Guo et al. 2017). Adapter sequences within raw
537 sequencing data were identified using `bbmerge.sh` of `bbmap`
538 (<https://sourceforge.net/projects/bbmap/>) and removed using `cutadapt` (Martin 2011),
539 as above. ENCODE ATAC-seq datasets (The ENCODE Project Consortium 2012;
540 Davis et al. 2018) (Liver: ENCFF628MCV, Ovary: ENCFF780JBA, Spleen:
541 ENCFF294ZCT, Testis: ENCFF048IOT, Transverse Colon: ENCFF377DAO) were
542 downloaded as BAM files, converted to BEDPE format, and short fragments were
543 identified for peak calling.

544

545 **Structural Variant Breakpoint data**

546

547 Large SVs, identified by high-density arrayCGH, or a combination of arrayCGH +
548 exome sequencing, were extracted from a cohort of 9,625 DDD patients, using variant
549 calling procedures as described in (Deciphering Developmental Disorders Study
550 2015). We filtered the DDD variants to only keep variants which fulfilled the
551 following criteria: a CNsolidate wscore ≥ 0.468 , a callp < 0.01 and a mean \log_2 ratio
552 of < -0.41 for deletions and 0.36 for duplications; CIPHER “false positives” were
553 removed. Singleton variants were identified as being annotated as “novel” by the
554 DDD release, only seen once among the DDD patients, and not seen in the dgv
555 (MacDonald et al. 2014) and gnomAD V.2 (Collins et al. 2020) structural variant
556 datasets (80% reciprocal overlap criterion). Since there are 9,625 patients in the DDD
557 dataset, the gnomAD V.2 dataset contains SVs from 10,738 genomes and the dgv
558 contains SVs from 29,084 individuals, this puts an upper limit of the frequency of
559 carriers of a singleton variant at $\sim 0.002\%$. Breakpoints were identified as the 5’ and
560 3’ coordinates of SVs, resulting in 13,406 singleton deletion and 3,406 duplication
561 breakpoints; the resolution of the breakpoints was such that the median and mean
562 confidence intervals were 300 bp and 12 kb, respectively. Thus, the DDD dataset has
563 a lower resolution compared to WGS data, but its advantage is that it does not suffer
564 from mapping and variant calling issues associated with the latter (Mahmoud et al.
565 2019).

566 We further identified 11,962 “common” deletion variants in the DDD dataset, which
567 had a minimum variant frequency of 1% in the consensus CNV dataset as described
568 by the DDD study (2015), i.e. pooled CNV datasets of Conrad et al. (2010), The 1000

569 Genomes Project Consortium (2010), the Wellcome Trust Case Control (2010) and
570 the DDD normal controls. We used the 80% reciprocal overlap criterion and grouped
571 common variants using the bedmap options `--echo-map --fraction-both 0.8`, followed
572 by `bedops --merge` (Neph et al. 2012). The breakpoints of common variants are thus
573 the outermost coordinates of all SVs that are collapsed into a given variant. The
574 overlap of such “common” breakpoints with ATAC-seq peaks was assessed
575 independently of SV allele frequencies, i.e. a group of common SVs contributed two
576 breakpoints to the analysis, and this number was further reduced if one breakpoint
577 coordinate was shared between two common SVs, so as to only count each common
578 breakpoint once.

579 We also identified a set of singleton CNVs called with the Manta algorithm (Chen et
580 al. 2016) from the gnomAD V.2 database (Collins et al. 2020) (80% reciprocal
581 overlap criterion with gnomAD V.2, dgv and DDD variants), resulting in a set of
582 73,063 deletion and 15,419 duplication breakpoints seen in ~ 0.002% of individuals
583 but called with a different approach compared to the DDD. Common deletions and
584 duplications ($p \geq 0.05$) were also extracted from the gnomAD V.2 dataset; these
585 variants had also been called with the Manta algorithm and included 5,954 deletion
586 and 1,586 duplication breakpoint sites.

587

588 **Indels and SNP data**

589

590 The recently released gnomAD V.3 variants (indels and SNPs) were downloaded
591 from <https://gnomad.broadinstitute.org/>. Only variants that passed all filters were kept
592 (filtering using `VCFtools --remove-filtered-all`) (Danecek et al. 2011). Multiallelic
593 variants were split using `BCFtools` (Danecek et al. 2021), and `bcftools norm --`

594 IndelGap 2 was applied to indels, to allow only variants to pass that were separated by
595 at least 2 bp. Singleton variants were defined as having an allele count of one, and the
596 allele number was $\geq 100,000$, i.e. the allele frequency of singletons was $p \leq$
597 0.001%.

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

604

605 A total of 854,409 *de novo* SNPs and indels were compiled from three different
606 sources, lifted over to the hg38 assembly using the UCSC liftOver tool as required.
607 First, we downloaded variants from <http://denovo-db.gs.washington.edu/>, including
608 only samples from whole genome sequencing studies (Michaelson et al. 2012; Ramu
609 et al. 2013; The Genome of the Netherlands Consortium 2014; Besenbacher et al.
610 2015; Turner et al. 2016; Yuen et al. 2016; Jonsson et al. 2017; RK et al. 2017; Turner
611 et al. 2017; Werling et al. 2018), which included a total of 404,238 variants from
612 4,560 samples. Additional samples, which were not already included in the denovo-db
613 dataset, were downloaded from the MSSNG database (<https://research.mss.ng/>),
614 version 2019/10/16, which added 2,243 samples and 215,044 *de novo* mutations. A
615 third source of *de novo* variants came from An et al. (2018) - 3,805 samples and
616 255,107 mutations.

617

618 **Circular Permutation**

619

620 To obtain a genome-wide estimate of enrichment of overlap between genomic
621 features (e.g. TFBSs and mutations), we performed circular permutations using the
622 Bioconductor regioneR package (Gel et al. 2016) in R (<https://www.R-project.org>) (R
623 Core Team 2016). We used the permTest() function with parameters ntimes=10000,
624 randomize.function=circularRandomizeRegions, evaluate.function=numOverlaps,
625 genome=hg38_masked, alternative="auto", where hg38_masked =
626 getBSgenome("BSgenome.Hsapiens.UCSC.hg38.masked"). This test evaluates the
627 number of overlaps observed between two sets of genomic features, given their order
628 on the chromosome and the distance between features, i.e. taking their degree of
629 clustering into account; Z-score analysis reveals the degree of local enrichment of
630 overlaps (Supplemental Methods) .

631 For permutations involving SVs, we used the two breakpoints of each SV, and
632 assessed the overlap of breakpoints with another feature of interest (i.e. ATAC-seq
633 sites), treating each breakpoint separately.

634 Circular permutations in regioneR (Gel et al. 2016) were also used assess the mean
635 distance between ATAC-seq peaks and deletion breakpoints, for common and
636 singleton variants separately.

637

638 **Brain enhancer data**

639

640 Active brain enhancers came from de la Torre-Ubieta et al. (2018). Specifically, we
641 used the 37,050 brain enhancers which showed differential accessibility in the
642 germinal zone versus the cortical plate, reflecting activity in the developing brain (de
643 la Torre-Ubieta et al. 2018). Next, we identified brain enhancers that were also active

644 during the male germline formation, i.e. overlapping the spermatogonial ATAC-seq
 645 peaks. To correct for the variable size of the brain active enhancers, we took the
 646 midpoints of each enhancer plus/minus 500 bp on either side, and intersected these
 647 sites with the ATAC-seq peaks using BEDTools' "bedtools intersect" (Quinlan and
 648 Hall 2010), thus classifying brain enhancers as spermatogonial "active" or "inactive".
 649 Next, we intersected these two categories of brain enhancers with the DDD
 650 breakpoint and gnomAD insertion dataset, respectively, to further classify them as
 651 "disrupted" by a singleton variant or "intact". An odds ratio was calculated as

652

$$653 \text{ OR} = (A/(B - A))/(C/(D - C))$$

654 With confidence intervals

$$655 \text{ CI_lower} = \exp(\log(\text{OR}) - 1.96 * \sqrt{1/A + 1/(B-A) + 1/C + 1/(D-C)})$$

$$656 \text{ CI_higher} = \exp(\log(\text{OR}) + 1.96 * \sqrt{1/A + 1/(B-A) + 1/C + 1/(D-C)})$$

657

658 where:

659 A = Disrupted, sperm active

660 B = All sperm active

661 C = Disrupted, sperm inactive

662 D = All sperm inactive

663

664 To analyse the enrichment of short InDels and SNPs around TFBSs and brain
 665 enhancers, we only considered genomic regions with unique mappability in $\geq 95\%$
 666 of the region, using the bedmap option --bases-uniq-f (Neph et al. 2012) and the
 667 mappability file hg38_umap24 (Karimzadeh et al. 2018), converted to bedmap
 668 format.

669

670 Random Forest Regression

671

672 To compare the effects of chromatin state on mutation rates, we performed random
673 forest regression with 200 trees, modelling the outcome variables “singleton
674 breakpoints” and “singleton insertions (5-20 bp)”, from the DDD and gnomAD V.3
675 respectively, within 5-kb wide genomic windows. Predictor variables included
676 “spermatogonial TFBS count”, “ssDNA overlap” (from Pratto et al. (2014)),
677 “recombination HS overlap” (from The 1000 Genomes Project Consortium (2015)),
678 “GC-content”, “Replication timing” (average of Wavelet-smooth signal in 1-kb bins
679 of 15 ENCODE tissues, downloaded from
680 http://hgdownload.soe.ucsc.edu/goldenPath/hg19/encodeDCC/wgEncodeUwRepliSeq
681 /), “Gene density”, “Mappability” (proportion of sites in each window with an
682 umap24 score of 1), and the overlap with “LTRs”, “SINEs”, “LINEs” and “Simple
683 Repeats” (downloaded from the Table Browser at <https://genome.ucsc.edu/>).
684 In a smaller model, we subset the dataset to only include 5-kb bins that also overlap
685 active brain enhancers (de la Torre-Ubieta et al. 2018), then ran the random forest
686 regression model to predict mutation rates within genomic regions that contain active
687 brain enhancers.

688

689 Motif discovery in singleton insertion sites

690

691 In order to find sequence motifs within the 5-20 bp singleton insertion sites from
692 gnomAD V.3, without prior assumptions, we extracted the FASTA sequence for
693 insertions that fell within 10 bp of the top 10 disrupted motif families (motif families

694 992, 193, 796, 907, 579, 825, 984, 171, 991). We ran the MEME 4.11 motif discovery
695 algorithm (Bailey et al. 2009) with “-nmotifs 1” on the inserted sequences. This
696 allowed us to compare the sequence motif of the disrupted TFBSs to any recurrent
697 motif found within the inserted sequences.

698

699 **Control Motif sites**

700

701 Using default search criteria, the FIMO algorithm (Grant et al. 2011) was run on the
702 repeat masked hg38 genome sequence (hg38.fa.masked, downloaded from
703 <https://genome.ucsc.edu/> in March 2020), searching the whole genome for the 579
704 input JASPAR motifs and the 13-mer PRDM9 motif. As with active binding sites,
705 motif matches belonging to the same motif family were merged and reported as a
706 single motif match per family, and only regions with unique umap24 mappabilities for
707 $\geq 95\%$ of sites were kept; motifs that overlapped with spermatogonial ATAC-seq
708 peaks were excluded. Next, these “control” motif sites were down-sampled to 10,000
709 per motif family (using BEDTools’ “bedtools sample” (Quinlan and Hall 2010));
710 circular permutations were performed to compare the observed to expected overlap of
711 the control motif sites (plus/minus 10 bp) with the gnomAD singleton insertions of 5-
712 20 bp.

713

714 The FIMO predicted control sites were also used to assess the degree of “clustering”
715 of motifs at spermatogonia active sites. For this purpose, we intersected the FIMO
716 motifs with a) spermatogonial ATAC-seq sites and b) ENCODE Master regulatory
717 sites downloaded from <https://genome.ucsc.edu/> (DNase I hypersensitivity derived
718 from assays in 95 cell types). For each of the 167 motif families, we calculated the

719 median distance (in basepairs) from a motif located within the active regulatory
720 region to the nearest FIMO motif of the same type. Accordingly, the ratio of the
721 median distance between motif sites (ENCODE/spermatogonia) was larger than one if
722 motifs at spermatogonial sites were, on average, closer to each other than motifs near
723 ENCODE sites, and we used this ratio as a measure of motif clustering. When
724 correlating the IFE with the degree of motif clustering (Figure 6B), we thus largely
725 correct for base compositional biases near active sites (which impact mutation rates –
726 Supplemental Fig. S9) as well as the effects of historical selection on the clustering of
727 motifs near genes, i.e. shorter inter-motif distances in spermatogonia indicate that
728 these sites have specifically high levels of motif density in spermatogonia, beyond the
729 levels expected for binding sites in general.

730

731 **Data access**

732

733 All raw sequencing data generated in this study have been submitted to the European
734 Genome-phenome Archive (EGA; <https://ega-archive.org/>) under accession number
735 EGAS00001005366. The ATAC-seq peak files generated in this study are available
736 as Supplemental Data (Supplemental Datasets 1-3) and at Edinburgh DataShare
737 (<https://doi.org/10.7488/ds/3053>).

738

739 **Competing interest statement**

740

741 The authors have no competing interests to declare.

742

743 **Acknowledgements**

744

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

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

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

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

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

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

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

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

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

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

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

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

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

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

759 grants MC_UU_00007/11, MC_UU_00007/2 and MC_UU_00007/16.

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

761 Bickmore for useful comments to the manuscript.

762

763 **Author Contributions**

764

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

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

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

768 data. V.B.K. performed the analyses. D.D.D. and D.R.F. provided data. D.R.F. helped
769 with the interpretation of the results and provided critical scientific inputs.

770

771 **Figure Legends**

772

773 **Figure 1: Locally elevated structural variation rates at spermatogonial**

774 **regulatory sites.** SV breakpoint count (**A, B**) and circular permutation shifted Z-

775 scores (**C, D**) of deletion breakpoints in the DDD cohort, centred around the

776 midpoints of spermatogonial ATAC-seq peaks. “Singletons” are breakpoints of

777 deletions with a frequency of ~ 0.002% across population samples; “common”

778 variants are seen at a frequency of at least 1% in the DDD consensus dataset (see

779 main text); permutation p-values indicate significant enrichment for both types of

780 variants at ATAC-seq peaks ($p < 10^{-5}$ in each case) (**E, F**) Circular permutation

781 shifted Z-scores of gnomAD deletion breakpoints, centred around spermatogonial

782 ATAC-seq peaks. “Singletons” are breakpoints of deletions with a frequency of ~

783 0.002% across population samples; “common” variants are seen at a frequency of at

784 least 5% in the gnomAD V.2 dataset. Permutation p-values indicate significant

785 enrichment for singleton breakpoints ($p < 10^{-5}$), and a significant depletion for

786 common variants ($p < 0.01$).

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

788 **sites.** Insertion count (**A, B**) and Shifted Z-scores (**C, D**) of gnomAD singleton and

789 common insertions (5-20 bp), centred around spermatogonial TFBSs. Singletons are

790 seen only once in the gnomAD V.3 dataset (allele frequency $\leq 0.001\%$) and are

791 significantly enriched at binding sites ($p < 10^{-4}$); common variants have an allele

792 frequency of at least 5% within gnomAD V.3 and are significantly depleted at binding
793 sites ($p < 10^{-4}$).

794 **Figure 3: Parallel enrichments of short variants and SV breakpoints at**
795 **spermatogonial binding sites.** Circular permutation results are based on 10,000
796 permutations; results for singleton variants and de novo mutation are shown. The Y
797 axis shows the ratio of observed over expected variant counts. Mutation categories
798 with significant enrichment are indicated by asterisks (***) indicating $p < 0.001$. The
799 type of variant tested and the total number of observed variants overlapping
800 spermatogonial TFBSs are indicated below each bar.

801 **Figure 4: Binding factors associated with the highest rates of mutation at**
802 **spermatogonial binding sites.** Plots are centred on the binding sites of a given motif
803 family inside ATAC-seq footprints. **(A)** Singleton and **(B)** common deletion
804 breakpoints in the DDD cohort; singletons are breakpoints of deletions with a
805 frequency of $\sim 0.002\%$ across population samples; common variants are seen at a
806 frequency of at least 1% in the DDD consensus dataset. **(C)** Singleton and **(D)**
807 common insertions (5-20 bp) in the gnomAD dataset. Singletons are seen only once in
808 gnomAD V.3 (allele frequency $\leq 0.001\%$), and common variants have an allele
809 frequency of at least 5% within gnomAD V.3. Only 10 kb regions around TFBSs with
810 $\geq 95\%$ unique mappability (umap24 scores) were included. The top 5 disrupted
811 motifs are shown, listed in order of enrichment of singleton variants in the circular
812 permutations (all enrichments of singletons are associated with p -values $< 10^{-4}$).

813 **Figure 5: Elevated singleton insertion rates at PRDM9 and NRF1 binding sites**
814 **contrast with other short variant classes.** All gnomAD variants have been down-
815 sampled to a total of 650,000 variants per analysis, making the Y axes directly

816 comparable; individual bins are 5 bp in size. Only regions around TFBSs with $\geq 95\%$
817 unique mappability (umap24 scores) were included.

818 **Figure 6: Insertions at spermatogonial TFBSs generate motif clusters in the**
819 **genome. A)** JASPAR database sequence motifs identified in the footprints of
820 spermatogonial ATAC-seq peaks (left) and the motifs identified in the singleton
821 insertions (5-20 bp) (right). The number of insertion sites (N) that were chosen by
822 MEME to construct the motif are shown on the right. **B)** For each motif family, we
823 plot the insertion fold enrichment (IFE) on the X axis and the degree of
824 spermatogonial motif clustering on the Y axis; the least square regression line is
825 indicated in blue. Motif clustering is measured as the distance to the nearest motif at
826 spermatogonial active sites, relative to the distance for motifs at ENCODE active
827 sites. **C)** The insertion fold enrichment (IFE) is contrasted between FIMO control
828 motif sites (left) and spermatogonial active motif sites (right); the Wilcoxon Test was
829 performed to compare the IFE at the two classes of sites.

830 **Figure 7: Neural enhancers with activity in spermatogonia suffer elevated mutation**
831 **rates. A)** Singleton DDD deletion breakpoint and **B)** singleton gnomAD insertion (5-
832 20 bp) count around brain active enhancers. Enhancers were classified as being also
833 active in spermatogonia (red) or inactive in spermatogonia (blue). Plotted is the
834 average number of variants per brain enhancer - in 5 kb windows or 100 bp windows,
835 respectively. In **(B)**, only 10 kb regions around enhancers with $\geq 95\%$ unique
836 mappability (umap24 scores) were included (3,409 brain enhancers that are inactive
837 in spermatogonia and 1,029 that are active).

838

839

840 **References**

841

- 842 Afek A, Shi H, Rangadurai A, Sahay H, Senitzki A, Khani S, Fang M, Salinas R,
843 Mielko Z, Pufall MA et al. 2020. DNA mismatches reveal conformational
844 penalties in protein-DNA recognition. *Nature* **587**: 291-296.
- 845 Altemose N, Noor N, Bitoun E, Tumian A, Imbeault M, Chapman JR, Aricescu AR,
846 Myers SR. 2017. A map of human PRDM9 binding provides evidence for
847 novel behaviors of PRDM9 and other zinc-finger proteins in meiosis. *Elife* **6**.
- 848 An JY, Lin K, Zhu L, Werling DM, Dong S, Brand H, Wang HZ, Zhao X, Schwartz
849 GB, Collins RL et al. 2018. Genome-wide de novo risk score implicates
850 promoter variation in autism spectrum disorder. *Science* **362**.
- 851 Andersson R, Gebhard C, Miguel-Escalada I, Hoof I, Bornholdt J, Boyd M, Chen Y,
852 Zhao X, Schmidl C, Suzuki T et al. 2014. An atlas of active enhancers across
853 human cell types and tissues. *Nature* **507**: 455-461.
- 854 Bailey TL, Boden M, Buske FA, Frith M, Grant CE, Clementi L, Ren J, Li WW,
855 Noble WS. 2009. MEME SUITE: tools for motif discovery and searching.
856 *Nucleic Acids Res* **37**: W202-208.
- 857 Baudat F, Buard J, Grey C, Fledel-Alon A, Ober C, Przeworski M, Coop G, de Massy
858 B. 2010. PRDM9 Is a Major Determinant of Meiotic Recombination Hotspots
859 in Humans and Mice. *Science* **327**: 836-840.
- 860 Berg IL, Neumann R, Lam KW, Sarbajna S, Odenthal-Hesse L, May CA, Jeffreys AJ.
861 2010. PRDM9 variation strongly influences recombination hot-spot activity
862 and meiotic instability in humans. *Nat Genet* **42**: 859-863.
- 863 Besenbacher S, Liu S, Izarzugaza JM, Grove J, Belling K, Bork-Jensen J, Huang S,
864 Als TD, Li S, Yadav R et al. 2015. Novel variation and de novo mutation rates
865 in population-wide de novo assembled Danish trios. *Nat Commun* **6**: 5969.
- 866 Besenbacher S, Sulem P, Helgason A, Helgason H, Kristjansson H, Jonasdottir A,
867 Jonasdottir A, Magnusson OT, Thorsteinsdottir U, Masson G et al. 2016.
868 Multi-nucleotide de novo Mutations in Humans. *PLoS Genet* **12**: e1006315.
- 869 Buenrostro JD, Giresi PG, Zaba LC, Chang HY, Greenleaf WJ. 2013. Transposition
870 of native chromatin for fast and sensitive epigenomic profiling of open
871 chromatin, DNA-binding proteins and nucleosome position. *Nat Methods* **10**:
872 1213-1218.
- 873 Carlson J, Locke AE, Flickinger M, Zawistowski M, Levy S, Myers RM, Boehnke M,
874 Kang HM, Scott LJ, Li JZ et al. 2018. Extremely rare variants reveal patterns
875 of germline mutation rate heterogeneity in humans. *Nat Commun* **9**.
- 876 Chen X, Schulz-Trieglaff O, Shaw R, Barnes B, Schlesinger F, Kallberg M, Cox AJ,
877 Kruglyak S, Saunders CT. 2016. Manta: rapid detection of structural variants
878 and indels for germline and cancer sequencing applications. *Bioinformatics*
879 **32**: 1220-1222.
- 880 Collins RL Brand H Karczewski KJ Zhao X Alföldi J Francioli LC Khera AV
881 Lowther C Gauthier LD Wang H et al. 2020. A structural variation reference
882 for medical and population genetics. *Nature* **581**: 444-451.
- 883 Collins RL, Brand H, Redin CE, Hanscom C, Antolik C, Stone MR, Glessner JT,
884 Mason T, Pregno G, Dorrani N et al. 2017. Defining the diverse spectrum of
885 inversions, complex structural variation, and chromothripsis in the morbid
886 human genome. *Genome Biol* **18**: 36.
- 887 Conrad DF, Pinto D, Redon R, Feuk L, Gokcumen O, Zhang Y, Aerts J, Andrews TD,
888 Barnes C, Campbell P et al. 2010. Origins and functional impact of copy
889 number variation in the human genome. *Nature* **464**: 704-712.

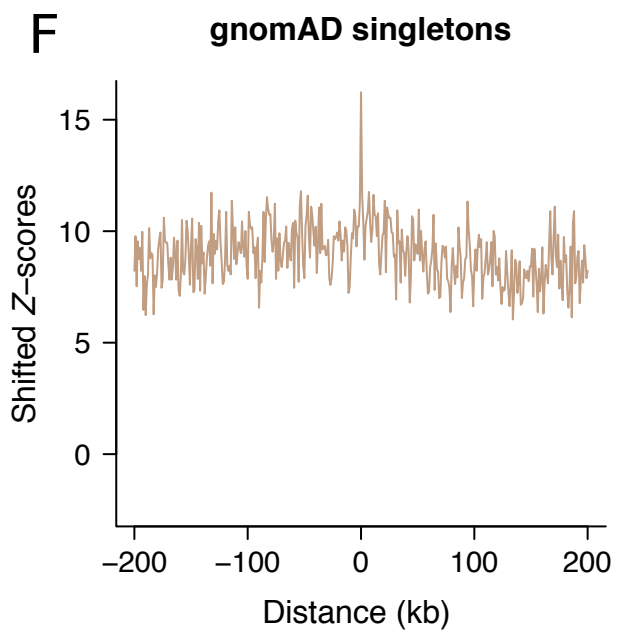
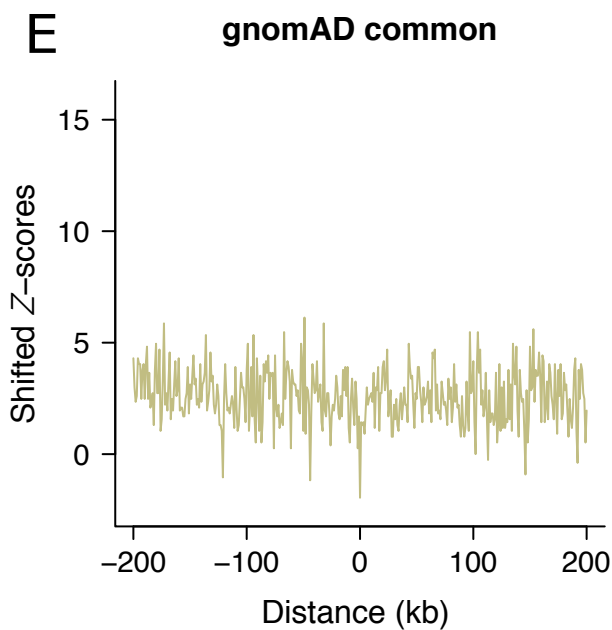
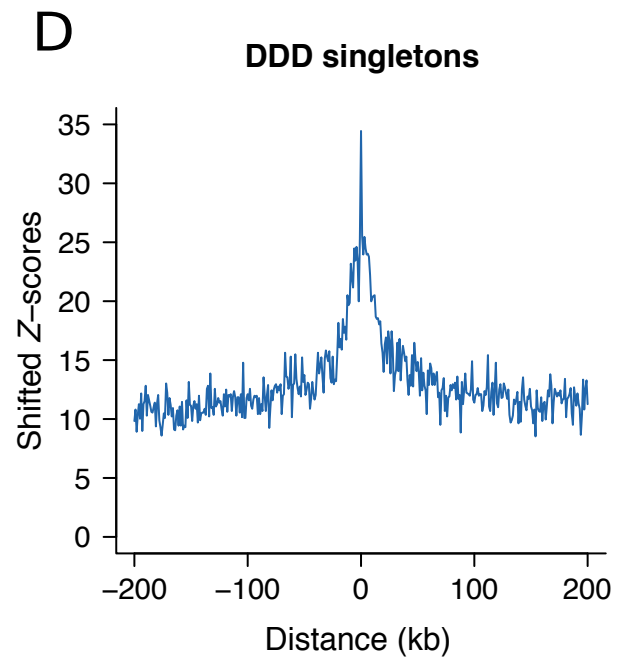
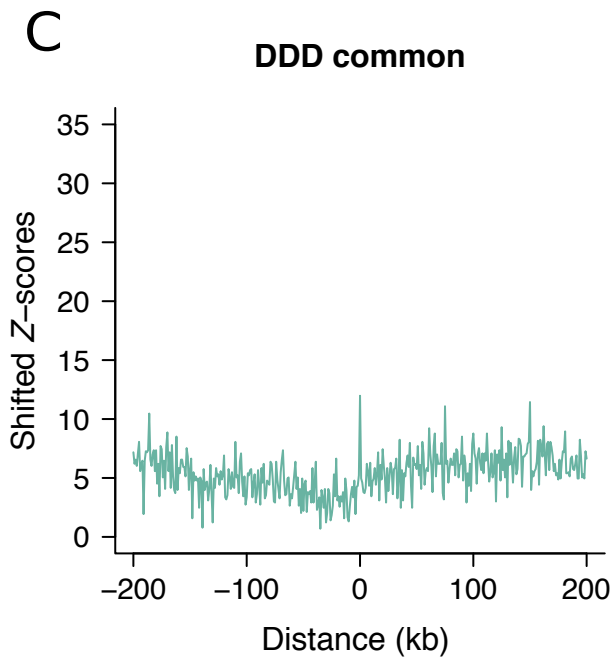
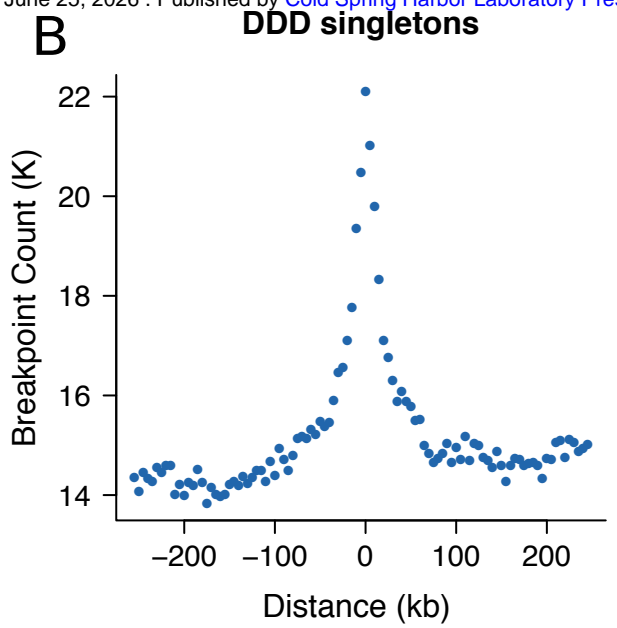
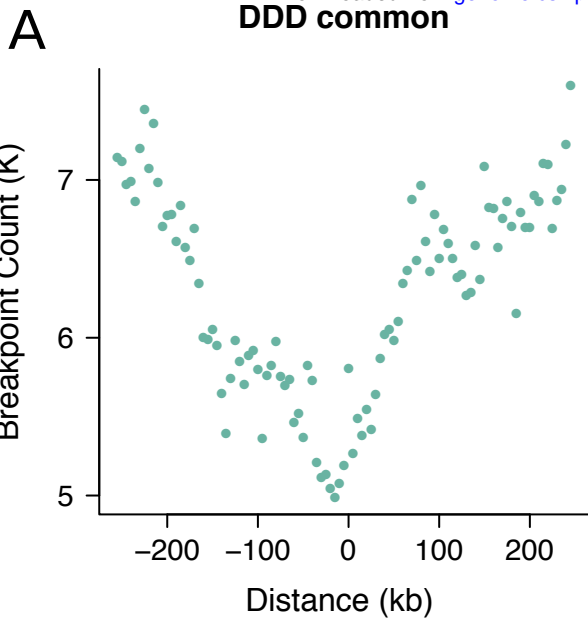
- 890 Danecek P, Auton A, Abecasis G, Albers CA, Banks E, DePristo MA, Handsaker RE,
891 Lunter G, Marth GT, Sherry ST et al. 2011. The variant call format and
892 VCFtools. *Bioinformatics* **27**: 2156-2158.
- 893 Danecek P, Bonfield JK, Liddle J, Marshall J, Ohan V, Pollard MO, Whitwham A,
894 Keane T, McCarthy SA, Davies RM et al. 2021. Twelve years of SAMtools
895 and BCFtools. *Gigascience* **10**.
- 896 Davis CA, Hitz BC, Sloan CA, Chan ET, Davidson JM, Gabdank I, Hilton JA, Jain K,
897 Baymuradov UK, Narayanan AK et al. 2018. The Encyclopedia of DNA
898 elements (ENCODE): data portal update. *Nucleic Acids Res* **46**: D794-D801.
- 899 de la Torre-Ubieta L, Stein JL, Won H, Opland CK, Liang D, Lu D, Geschwind DH.
900 2018. The Dynamic Landscape of Open Chromatin during Human Cortical
901 Neurogenesis. *Cell* **172**: 289-304 e218.
- 902 Deciphering Developmental Disorders Study. 2015. Large-scale discovery of novel
903 genetic causes of developmental disorders. *Nature* **519**: 223-228.
- 904 Gel B, Diez-Villanueva A, Serra E, Buschbeck M, Peinado MA, Malinverni R. 2016.
905 regioneR: an R/Bioconductor package for the association analysis of genomic
906 regions based on permutation tests. *Bioinformatics* **32**: 289-291.
- 907 Gerrard DT, Berry AA, Jennings RE, Birket MJ, Zarrineh P, Garstang MG, Withey
908 SL, Short P, Jimenez-Gancedo S, Firbas PN et al. 2020. Dynamic changes in
909 the epigenomic landscape regulate human organogenesis and link to
910 developmental disorders. *Nat Commun* **11**: 3920.
- 911 Girirajan S, Brkanac Z, Coe BP, Baker C, Vives L, Vu TH, Shafer N, Bernier R,
912 Ferrero GB, Silengo M et al. 2011. Relative burden of large CNVs on a range
913 of neurodevelopmental phenotypes. *PLoS Genet* **7**: e1002334.
- 914 Gonzalez-Perez A, Sabarinathan R, Lopez-Bigas N. 2019. Local Determinants of the
915 Mutational Landscape of the Human Genome. *Cell* **177**: 101-114.
- 916 Gotea V, Visel A, Westlund JM, Nobrega MA, Pennacchio LA, Ovcharenko I. 2010.
917 Homotypic clusters of transcription factor binding sites are a key component
918 of human promoters and enhancers. *Genome Res* **20**: 565-577.
- 919 Grant CE, Bailey TL, Noble WS. 2011. FIMO: scanning for occurrences of a given
920 motif. *Bioinformatics* **27**: 1017-1018.
- 921 Grey C, Barthes P, Chauveau-Le Friec G, Langa F, Baudat F, de Massy B. 2011.
922 Mouse PRDM9 DNA-Binding Specificity Determines Sites of Histone H3
923 Lysine 4 Trimethylation for Initiation of Meiotic Recombination. *Plos Biol* **9**.
- 924 Gu W, Zhang F, Lupski JR. 2008. Mechanisms for human genomic rearrangements.
925 *Pathogenetics* **1**: 4.
- 926 Guo J, Grow EJ, Mlcochova H, Maher GJ, Lindskog C, Nie X, Guo Y, Takei Y, Yun
927 J, Cai L et al. 2018. The adult human testis transcriptional cell atlas. *Cell Res*
928 **28**: 1141-1157.
- 929 Guo J, Grow EJ, Yi C, Mlcochova H, Maher GJ, Lindskog C, Murphy PJ, Wike CL,
930 Carrell DT, Goriely A et al. 2017. Chromatin and Single-Cell RNA-Seq
931 Profiling Reveal Dynamic Signaling and Metabolic Transitions during Human
932 Spermatogonial Stem Cell Development. *Cell Stem Cell* **21**: 533-546 e536.
- 933 Halldorsson BV, Palsson G, Stefansson OA, Jonsson H, Hardarson MT, Eggertsson
934 HP, Gunnarsson B, Oddsson A, Halldorsson GH, Zink F et al. 2019.
935 Characterizing mutagenic effects of recombination through a sequence-level
936 genetic map. *Science* **363**.
- 937 Hinch AG, Tandon A, Patterson N, Song YL, Rohland N, Palmer CD, Chen GK,
938 Wang K, Buxbaum SG, Akyzbekova EL et al. 2011. The landscape of
939 recombination in African Americans. *Nature* **476**: 170-U167.

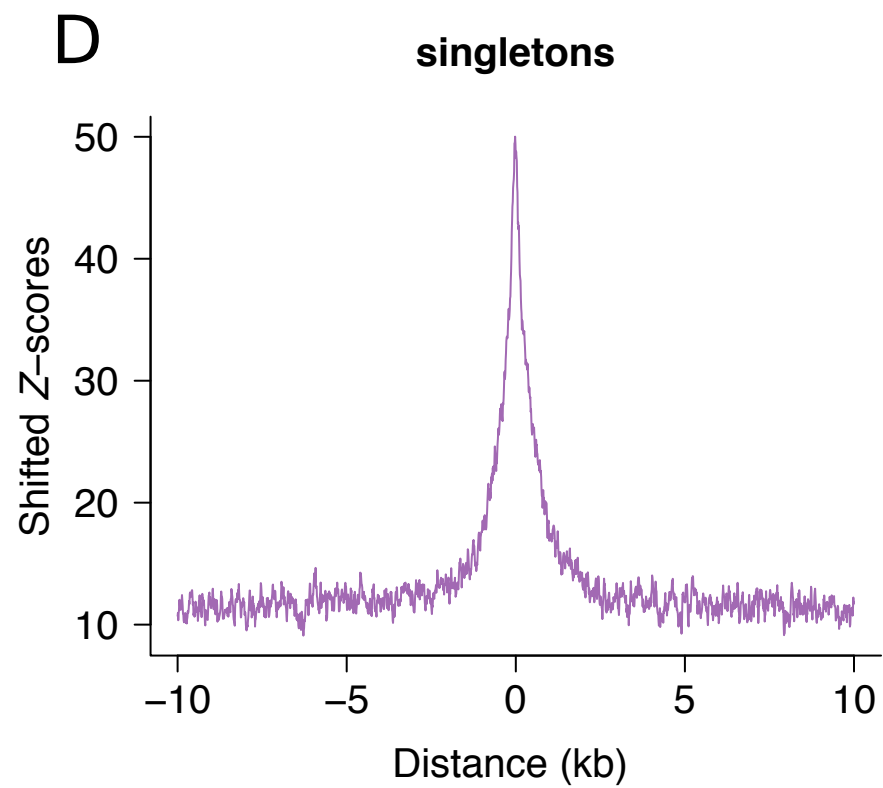
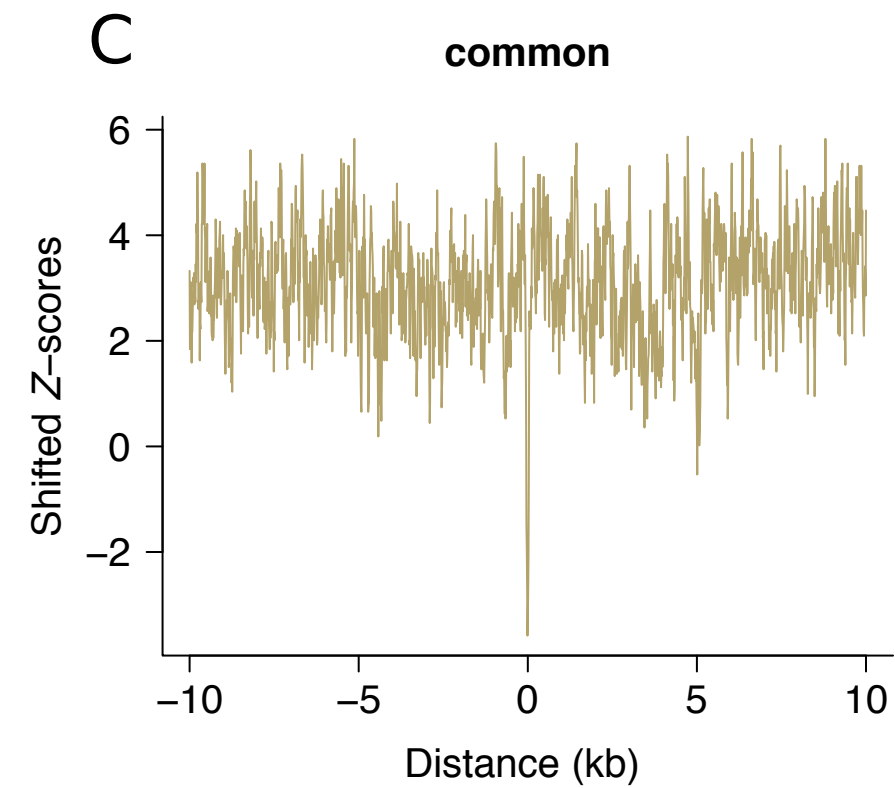
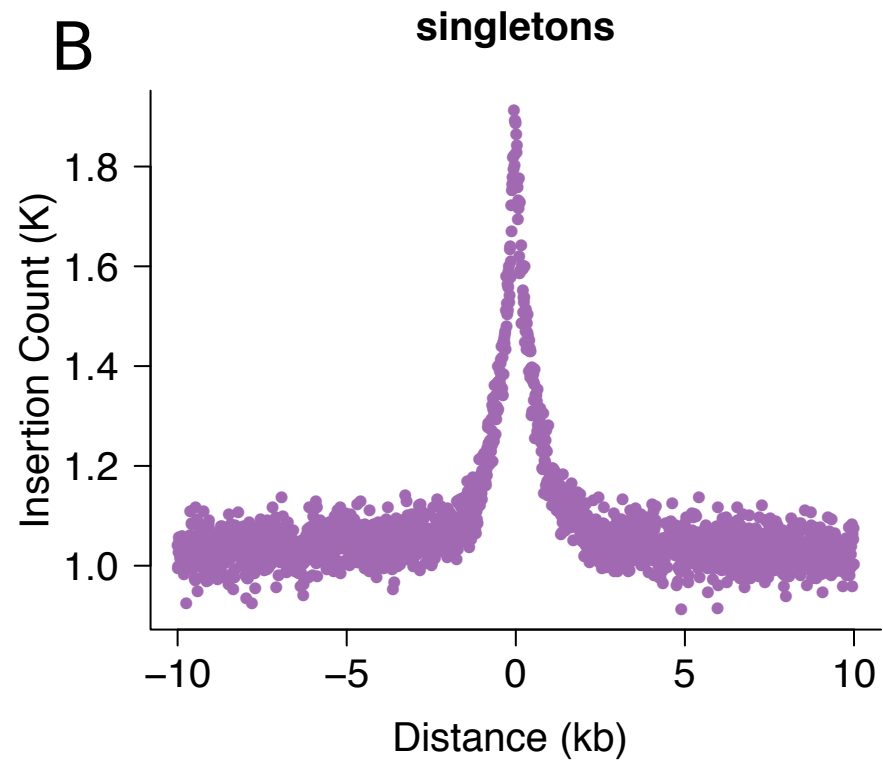
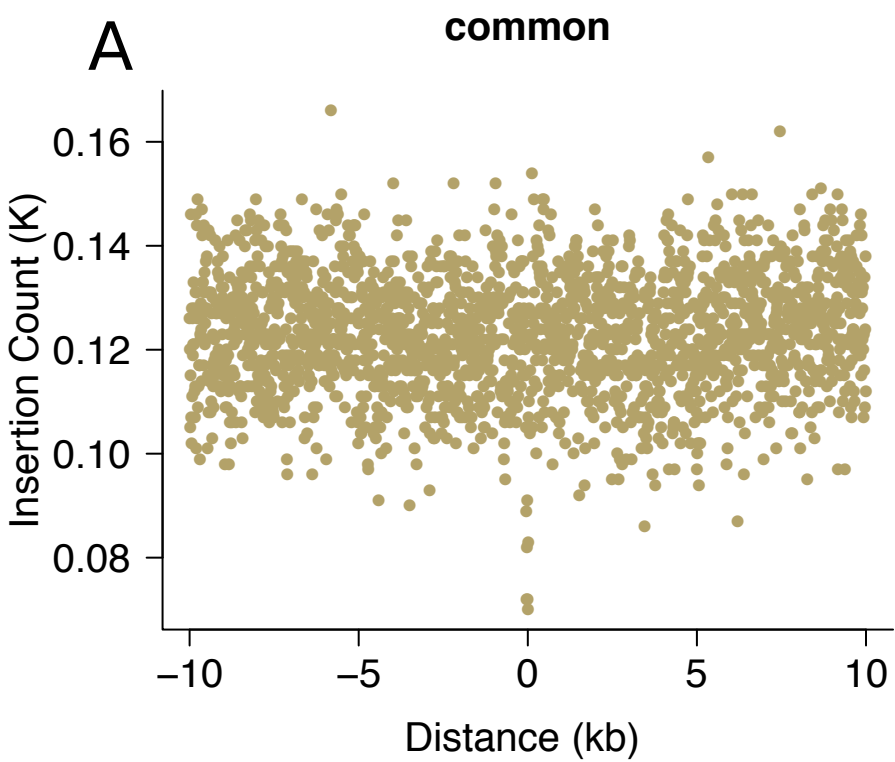
- 940 Hodgkinson A, Eyre-Walker A. 2011. Variation in the mutation rate across
941 mammalian genomes. *Nat Rev Genet* **12**: 756-766.
- 942 Jiang P, Singh M. 2014. CCAT: Combinatorial Code Analysis Tool for transcriptional
943 regulation. *Nucleic Acids Research* **42**: 2833-2847.
- 944 Jonsson H, Sulem P, Kehr B, Kristmundsdottir S, Zink F, Hjartarson E, Hardarson
945 MT, Hjorleifsson KE, Eggertsson HP, Gudjonsson SA et al. 2017. Parental
946 influence on human germline de novo mutations in 1,548 trios from Iceland.
947 *Nature* **549**: 519-522.
- 948 Kaiser VB, Semple CA. 2018. Chromatin loop anchors are associated with genome
949 instability in cancer and recombination hotspots in the germline. *Genome Biol*
950 **19**: 101.
- 951 Kaiser VB, Taylor MS, Semple CA. 2016. Mutational Biases Drive Elevated Rates of
952 Substitution at Regulatory Sites across Cancer Types. *PLoS Genet* **12**:
953 e1006207.
- 954 Karczewski KJ, Francioli LC, Tiao G, Cummings BB, Alfoldi J, Wang Q, Collins RL,
955 Laricchia KM, Ganna A, Birnbaum DP et al. 2020. The mutational constraint
956 spectrum quantified from variation in 141,456 humans. *Nature* **581**: 434-443.
- 957 Karimzadeh M, Ernst C, Kundaje A, Hoffman MM. 2018. Umap and Bimap:
958 quantifying genome and methylome mappability. *Nucleic Acids Res* **46**: e120.
- 959 Kentepozidou E, Aitken SJ, Feig C, Stefflova K, Ibarra-Soria X, Odom DT, Roller M,
960 Flicek P. 2020. Clustered CTCF binding is an evolutionary mechanism to
961 maintain topologically associating domains. *Genome Biol* **21**: 5.
- 962 Kim S, Peterson SE, Jasin M, Keeney S. 2016. Mechanisms of germ line genome
963 instability. *Semin Cell Dev Biol* **54**: 177-187.
- 964 Kondrashov AS, Rogozin IB. 2004. Context of deletions and insertions in human
965 coding sequences. *Hum Mutat* **23**: 177-185.
- 966 Koren A, Polak P, Nemesh J, Michaelson JJ, Sebat J, Sunyaev SR, McCarroll SA.
967 2012. Differential relationship of DNA replication timing to different forms of
968 human mutation and variation. *Am J Hum Genet* **91**: 1033-1040.
- 969 Kribelbauer JF, Rastogi C, Bussemaker HJ, Mann RS. 2019. Low-Affinity Binding
970 Sites and the Transcription Factor Specificity Paradox in Eukaryotes. *Annu*
971 *Rev Cell Dev Biol* **35**: 357-379.
- 972 Kvikstad EM, Tyekucheva S, Chiaromonte F, Makova KD. 2007. A macaque's-eye
973 view of human insertions and deletions: differences in mechanisms. *PLoS*
974 *Comput Biol* **3**: 1772-1782.
- 975 Leppa VM, Kravitz SN, Martin CL, Andrieux J, Le Caignec C, Martin-Coignard D,
976 DyBuncio C, Sanders SJ, Lowe JK, Cantor RM et al. 2016. Rare Inherited and
977 De Novo CNVs Reveal Complex Contributions to ASD Risk in Multiplex
978 Families. *Am J Hum Genet* **99**: 540-554.
- 979 Levinson G, Gutman GA. 1987. High frequencies of short frameshifts in poly-CA/TG
980 tandem repeats borne by bacteriophage M13 in Escherichia coli K-12. *Nucleic*
981 *Acids Res* **15**: 5323-5338.
- 982 Li C, Luscombe NM. 2020. Nucleosome positioning stability is a modulator of
983 germline mutation rate variation across the human genome. *Nat Commun* **11**:
984 1363.
- 985 Li ZJ, Schulz MH, Look T, Begemann M, Zenke M, Costa IG. 2019. Identification of
986 transcription factor binding sites using ATAC-seq. *Genome Biology* **20**.
- 987 Lieber MR, Ma Y, Pannicke U, Schwarz K. 2003. Mechanism and regulation of
988 human non-homologous DNA end-joining. *Nat Rev Mol Cell Biol* **4**: 712-720.

- 989 MacDonalD JR, Ziman R, Yuen RKC, Feuk L, Scherer SW. 2014. The Database of
990 Genomic Variants: a curated collection of structural variation in the human
991 genome. *Nucleic Acids Research* **42**: D986-D992.
- 992 Mahmoud M, Gobet N, Cruz-Davalos DI, Mounier N, Dessimoz C, Sedlazeck FJ.
993 2019. Structural variant calling: the long and the short of it. *Genome Biol* **20**:
994 246.
- 995 Makova KD, Yang S, Chiaromonte F. 2004. Insertions and deletions are male biased
996 too: a whole-genome analysis in rodents. *Genome Res* **14**: 567-573.
- 997 Martin M. 2011. Cutadapt Removes Adapter Sequences From High-Throughput
998 Sequencing Reads. *EMBnetjournal* **17**: 10-12.
- 999 McRae JF Clayton S Fitzgerald TW Kaplanis J Prigmore E Rajan D Sifrim A Aitken
1000 S Akawi N Alvi M et al. 2017. Prevalence and architecture of de novo
1001 mutations in developmental disorders. *Nature* **542**: 433-+.
- 1002 McVean G. 2007. What drives recombination hotspots to repeat DNA in humans?
1003 *Philosophical Transactions of the Royal Society London Series B Biological*
1004 *Sciences* **365**: 1213-1218.
- 1005 Messer PW. 2009. Measuring the Rates of Spontaneous Mutation From Deep and
1006 Large-Scale Polymorphism Data. *Genetics* **182**: 1219-1232.
- 1007 Messer PW, Arndt PF. 2007. The majority of recent short DNA insertions in the
1008 human genome are tandem duplications. *Mol Biol Evol* **24**: 1190-1197.
- 1009 Michaelson JJ, Shi Y, Gujral M, Zheng H, Malhotra D, Jin X, Jian M, Liu G, Greer
1010 D, Bhandari A et al. 2012. Whole-genome sequencing in autism identifies hot
1011 spots for de novo germline mutation. *Cell* **151**: 1431-1442.
- 1012 Montgomery SB, Goode DL, Kvikstad E, Albers CA, Zhang ZD, Mu XJ, Ananda G,
1013 Howie B, Karczewski KJ, Smith KS et al. 2013. The origin, evolution, and
1014 functional impact of short insertion-deletion variants identified in 179 human
1015 genomes. *Genome Res* **23**: 749-761.
- 1016 Myers S, Bottolo L, Freeman C, McVean G, Donnelly P. 2005. A fine-scale map of
1017 recombination rates and hotspots across the human genome. *Science* **310**: 321-
1018 324.
- 1019 Myers S, Bowden R, Tumian A, Bontrop RE, Freeman C, MacFie TS, McVean G,
1020 Donnelly P. 2010. Drive against hotspot motifs in primates implicates the
1021 PRDM9 gene in meiotic recombination. *Science* **327**: 876-879.
- 1022 Myers S, Freeman C, Auton A, Donnelly P, McVean G. 2008. A common sequence
1023 motif associated with recombination hot spots and genome instability in
1024 humans. *Nat Genet* **40**: 1124-1129.
- 1025 Neph S, Kuehn MS, Reynolds AP, Haugen E, Thurman RE, Johnson AK, Rynes E,
1026 Maurano MT, Vierstra J, Thomas S et al. 2012. BEDOPS: high-performance
1027 genomic feature operations. *Bioinformatics* **28**: 1919-1920.
- 1028 Palmer N, Talib SZA, Ratnacaram CK, Low D, Bisteau X, Lee JHS, Pfeifferberger E,
1029 Wollmann H, Tan JHL, Wee S et al. 2019. CDK2 regulates the NRF1/Ehmt1
1030 axis during meiotic prophase I. *J Cell Biol* **218**: 2896-2918.
- 1031 Powers NR, Parvanov ED, Baker CL, Walker M, Petkov PM, Paigen K. 2016. The
1032 Meiotic Recombination Activator PRDM9 Trimethylates Both H3K36 and
1033 H3K4 at Recombination Hotspots In Vivo. *PLoS Genet* **12**: e1006146.
- 1034 Pratto F, Brick K, Khil P, Smagulova F, Petukhova GV, Camerini-Otero RD. 2014.
1035 DNA recombination. Recombination initiation maps of individual human
1036 genomes. *Science* **346**: 1256442.
- 1037 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing
1038 genomic features. *Bioinformatics* **26**: 841-842.

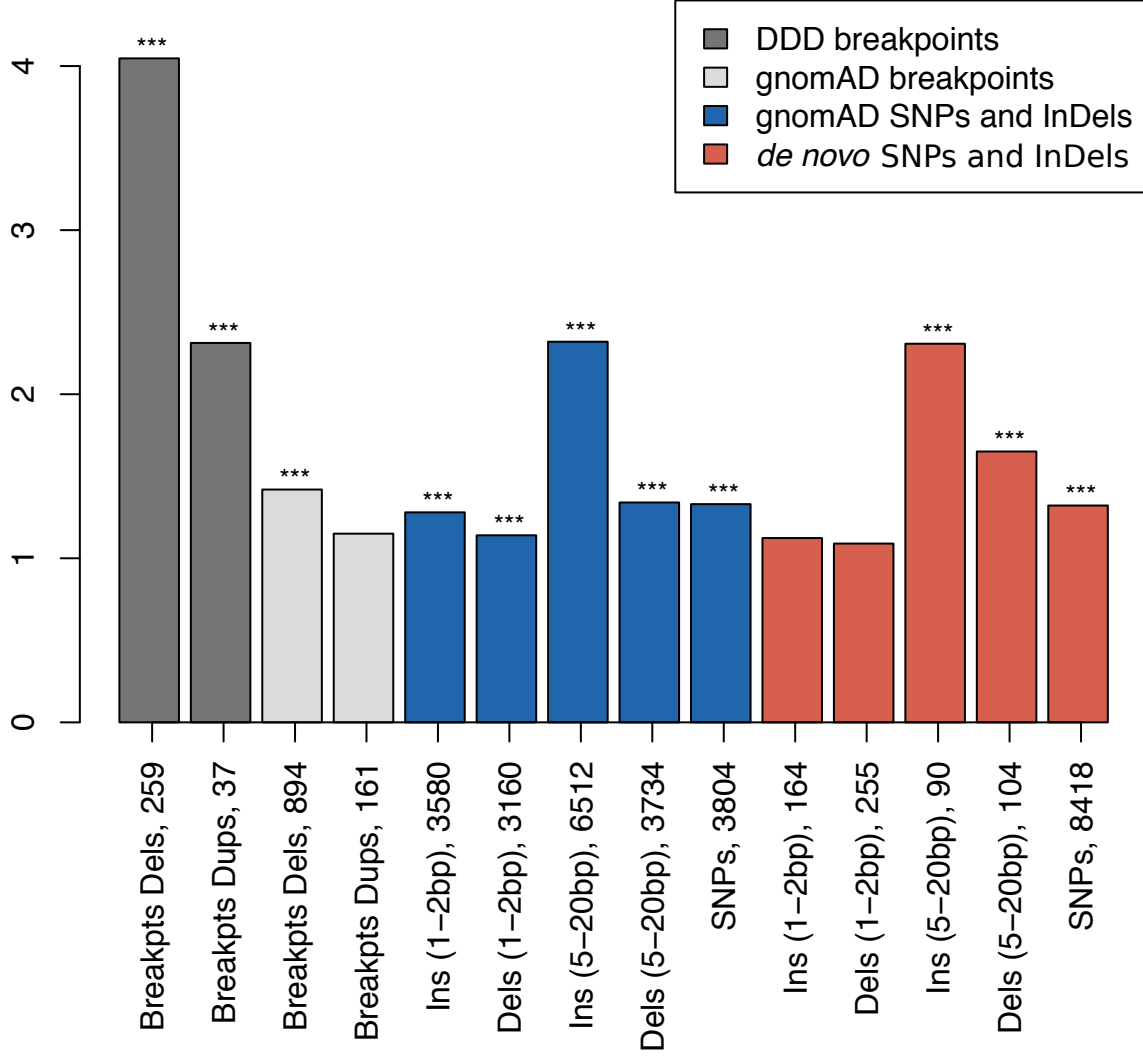
- 1039 R Core Team. 2016. R: A Language and Environment for Statistical Computing.
- 1040 Ramirez F, Ryan DP, Gruning B, Bhardwaj V, Kilpert F, Richter AS, Heyne S,
1041 Dundar F, Manke T. 2016. deepTools2: a next generation web server for deep-
1042 sequencing data analysis. *Nucleic Acids Res* **44**: W160-165.
- 1043 Ramu A, Noordam MJ, Schwartz RS, Wuster A, Hurles ME, Cartwright RA, Conrad
1044 DF. 2013. DeNovoGear: de novo indel and point mutation discovery and
1045 phasing. *Nat Methods* **10**: 985-987.
- 1046 Reijns MAM, Kemp H, Ding J, de Proce SM, Jackson AP, Taylor MS. 2015.
1047 Lagging-strand replication shapes the mutational landscape of the genome.
1048 *Nature* **518**: 502-506.
- 1049 RK CY, Merico D, Bookman M, J LH, Thiruvahindrapuram B, Patel RV, Whitney J,
1050 Deflaux N, Bingham J, Wang Z et al. 2017. Whole genome sequencing
1051 resource identifies 18 new candidate genes for autism spectrum disorder. *Nat*
1052 *Neurosci* **20**: 602-611.
- 1053 Rodgers K, McVey M. 2016. Error-Prone Repair of DNA Double-Strand Breaks. *J*
1054 *Cell Physiol* **231**: 15-24.
- 1055 Sabarinathan R, Mularoni L, Deu-Pons J, Gonzalez-Perez A, Lopez-Bigas N. 2016.
1056 Nucleotide excision repair is impaired by binding of transcription factors to
1057 DNA. *Nature* **532**: 264-267.
- 1058 Sandelin A, Alkema W, Engstrom P, Wasserman WW, Lenhard B. 2004. JASPAR:
1059 an open-access database for eukaryotic transcription factor binding profiles.
1060 *Nucleic Acids Research* **32**: D91-D94.
- 1061 Sherwood RI, Hashimoto T, O'Donnell CW, Lewis S, Barkal AA, van Hoff JP, Karun
1062 V, Jaakkola T, Gifford DK. 2014. Discovery of directional and nondirectional
1063 pioneer transcription factors by modeling DNase profile magnitude and shape.
1064 *Nat Biotechnol* **32**: 171-+.
- 1065 Short PJ, McRae JF, Gallone G, Sifrim A, Won H, Geschwind DH, Wright CF, Firth
1066 HV, FitzPatrick DR, Barrett JC et al. 2018. De novo mutations in regulatory
1067 elements in neurodevelopmental disorders. *Nature* **555**: 611-616.
- 1068 Sohni A, Tan K, Song HW, Burow D, de Rooij DG, Laurent L, Hsieh TC, Rabah R,
1069 Hammoud SS, Vicini E et al. 2019. The Neonatal and Adult Human Testis
1070 Defined at the Single-Cell Level. *Cell Rep* **26**: 1501-1517 e1504.
- 1071 Stamatoyannopoulos JA, Adzhubei I, Thurman RE, Kryukov GV, Mirkin SM,
1072 Sunyaev SR. 2009. Human mutation rate associated with DNA replication
1073 timing. *Nat Genet* **41**: 393-395.
- 1074 The 1000 Genomes Project Consortium. 2010. A map of human genome variation
1075 from population-scale sequencing. *Nature* **467**: 1061-1073.
- 1076 The 1000 Genomes Project Consortium. 2015. A global reference for human genetic
1077 variation. *Nature* **526**: 68-74.
- 1078 The ENCODE Project Consortium. 2012. An integrated encyclopedia of DNA
1079 elements in the human genome. *Nature* **489**: 57-74.
- 1080 The Genome of the Netherlands Consortium. 2014. Whole-genome sequence
1081 variation, population structure and demographic history of the Dutch
1082 population. *Nat Genet* **46**: 818-825.
- 1083 Turner TN, Coe BP, Dickel DE, Hoekzema K, Nelson BJ, Zody MC, Kronenberg ZN,
1084 Hormozdiari F, Raja A, Pennacchio LA et al. 2017. Genomic Patterns of De
1085 Novo Mutation in Simplex Autism. *Cell* **171**: 710-722 e712.
- 1086 Turner TN, Hormozdiari F, Duyzend MH, McClymont SA, Hook PW, Iossifov I,
1087 Raja A, Baker C, Hoekzema K, Stessman HA et al. 2016. Genome Sequencing

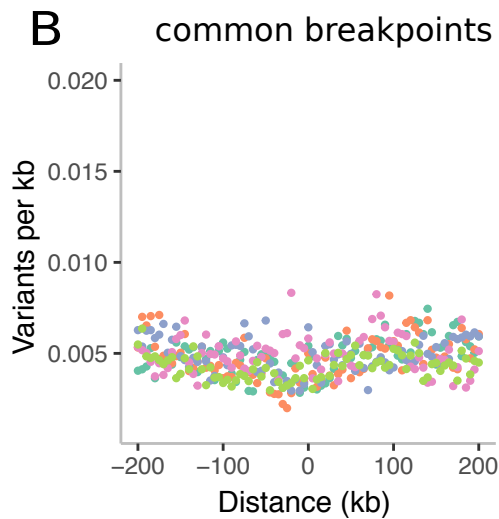
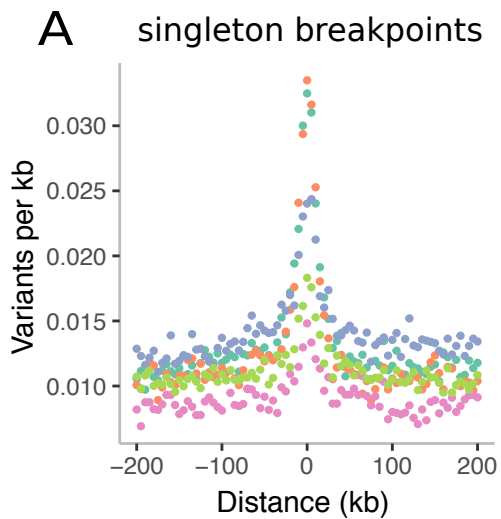
- 1088 of Autism-Affected Families Reveals Disruption of Putative Noncoding
1089 Regulatory DNA. *Am J Hum Genet* **98**: 58-74.
- 1090 van Gent DC, Hoeijmakers JH, Kanaar R. 2001. Chromosomal stability and the DNA
1091 double-stranded break connection. *Nat Rev Genet* **2**: 196-206.
- 1092 Wang J, Tang C, Wang Q, Su J, Ni T, Yang W, Wang Y, Chen W, Liu X, Wang S et
1093 al. 2017. NRF1 coordinates with DNA methylation to regulate
1094 spermatogenesis. *FASEB J* **31**: 4959-4970.
- 1095 Wellcome Trust Case Control Consortium Craddock N Hurles ME Cardin N Pearson
1096 RD Plagnol V Robson S Vukcevic D Barnes C Conrad DF et al. 2010.
1097 Genome-wide association study of CNVs in 16,000 cases of eight common
1098 diseases and 3,000 shared controls. *Nature* **464**: 713-720.
- 1099 Werling DM, Brand H, An JY, Stone MR, Zhu L, Glessner JT, Collins RL, Dong S,
1100 Layer RM, Markenscoff-Papadimitriou E et al. 2018. An analytical framework
1101 for whole-genome sequence association studies and its implications for autism
1102 spectrum disorder. *Nat Genet* **50**: 727-736.
- 1103 Wolfe KH, Sharp PM, Li WH. 1989. Mutation-Rates Differ among Regions of the
1104 Mammalian Genome. *Nature* **337**: 283-285.
- 1105 Yan F, Powell DR, Curtis DJ, Wong NC. 2020. From reads to insight: a hitchhiker's
1106 guide to ATAC-seq data analysis. *Genome Biol* **21**: 22.
- 1107 Yuen RK, Merico D, Cao H, Pellecchia G, Alipanahi B, Thiruvahindrapuram B, Tong
1108 X, Sun Y, Cao D, Zhang T et al. 2016. Genome-wide characteristics of de
1109 novo mutations in autism. *NPJ Genom Med* **1**: 160271-1602710.
- 1110 Zhang Y, Liu T, Meyer CA, Eeckhoute J, Johnson DS, Bernstein BE, Nussbaum C,
1111 Myers RM, Brown M, Li W et al. 2008. Model-based Analysis of ChIP-Seq
1112 (MACS). *Genome Biology* **9**.
- 1113





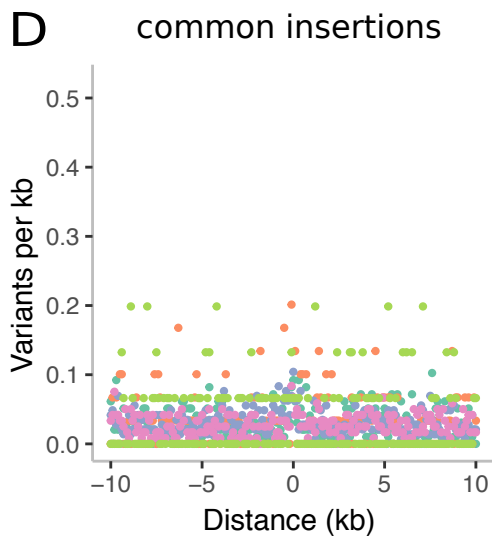
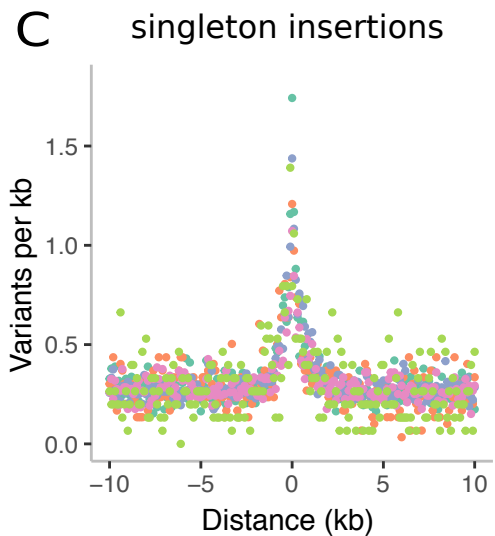
Observed/Expected





TF motifs

- motif_171: NRF1
- motif_579: PRDM9
- motif_672: TFAP2 family
- motif_963: NFYA/NFYB/Dux
- motif_972: ETS family

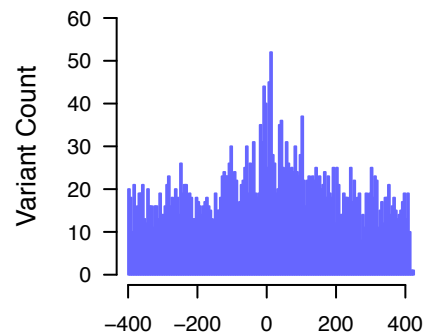


TF motifs

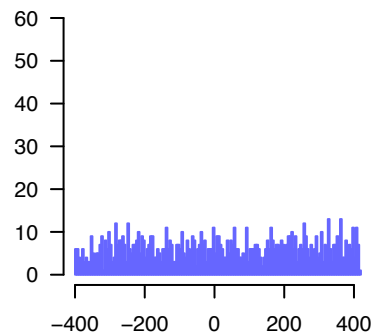
- motif_171: NRF1
- motif_92: HINFP
- motif_579: PRDM9
- motif_825: EGR family
- motif_192: ZBTB33

PRDM9

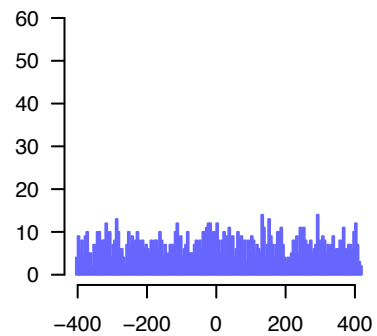
insertions, 5–20 bp



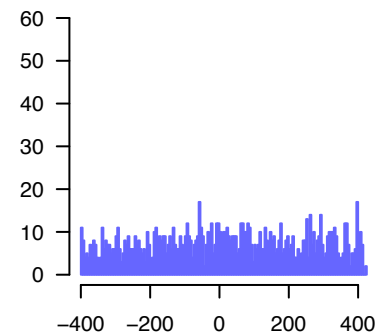
insertions, 1–2 bp



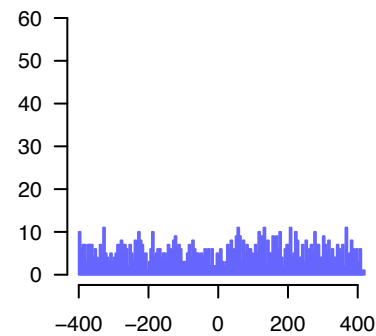
snps



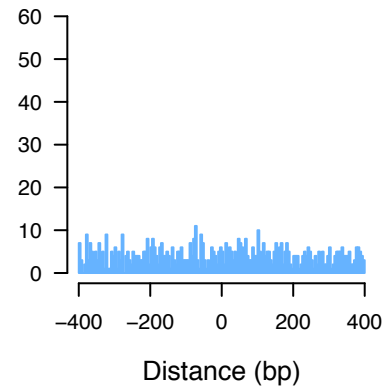
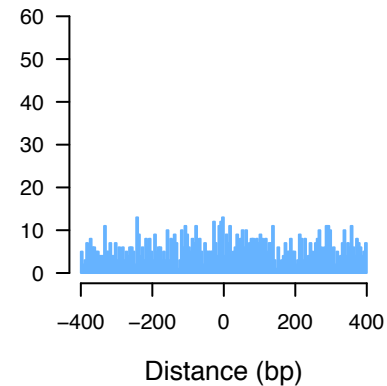
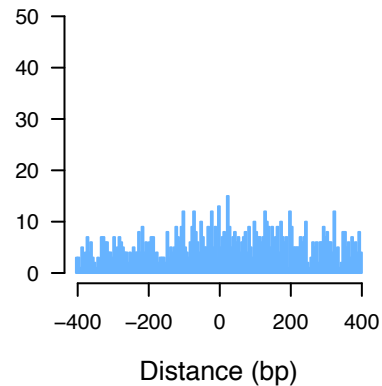
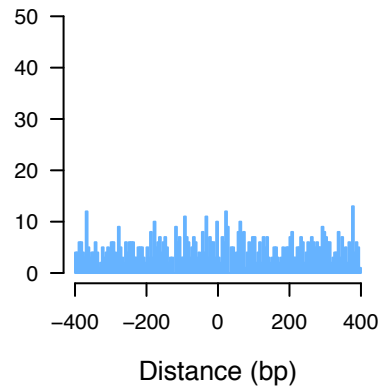
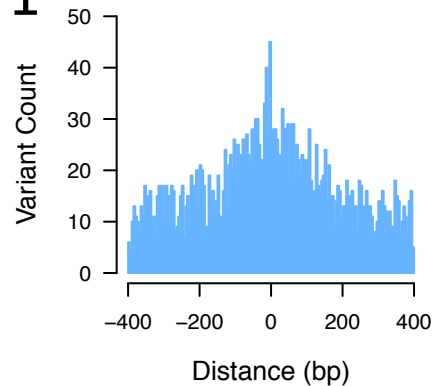
deletions, 5–20 bp



deletions, 1–2 bp



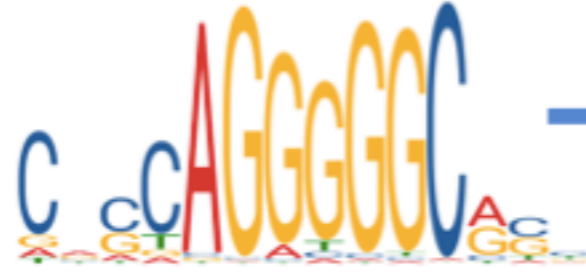
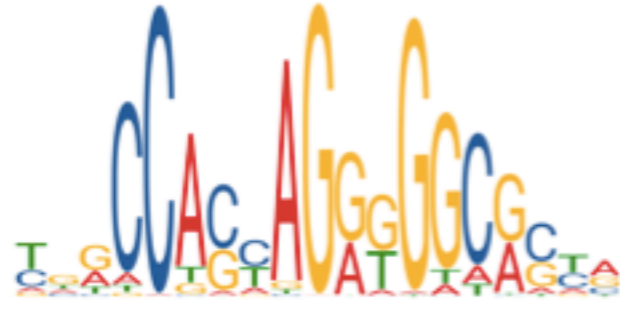
NRF1



A

Motif Family

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



N = 26

Motif 171: NRF1-MA0506.1

Downloaded from genome.cshlp.org on June 25, 2026 . Published by Cold Spring Harbor Laboratory Press



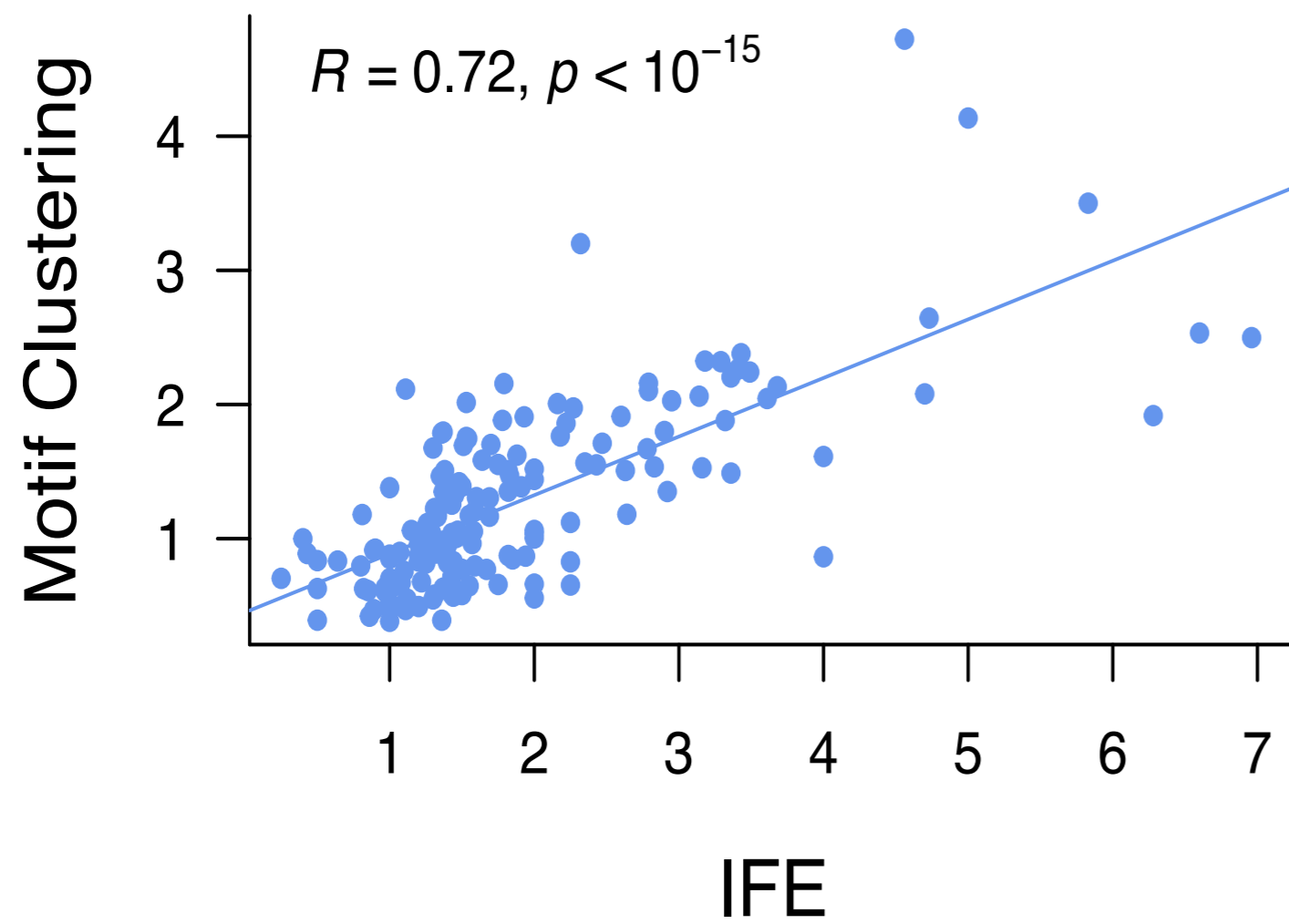
N = 9

Motif 579: PRDM9



N = 91

B



C

