

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23

Active RNAP pre-initiation sites are highly mutated by cytidine deaminases in yeast, with AID
targeting small RNAs genes

Benjamin J.M. Taylor, Yee Ling Wu and Cristina Rada

MRC Laboratory of Molecular Biology
Francis Crick Avenue, Cambridge CB2 0QH

Major subject areas: Genes & chromosomes, Immunology

Correspondence: Cristina Rada

email:car@mrc-lmb.cam.ac.uk; phone:+44(0)1223 267098

24

25 **Abstract**

26 Cytidine deaminases are single stranded DNA mutators diversifying antibodies and
27 restricting viral infection. Improper access to the genome leads to translocations and
28 mutations in B cells and contributes to the mutation landscape in cancer, such as kataegis. It
29 remains unclear how deaminases access double stranded genomes and whether off-target
30 mutations favor certain loci, although transcription and opportunistic access during DNA
31 repair are thought to play a role. In yeast, AID and the catalytic domain of APOBEC3G
32 preferentially mutate transcriptionally active genes within narrow regions, 110 base pairs in
33 width, fixed at RNA Polymerases initiation sites. Unlike APOBEC3G, AID shows enhanced
34 mutational preference for small RNA genes (tRNAs, snoRNAs and snRNAs) suggesting a
35 putative role for RNA in its recruitment. We uncover the high affinity of the deaminases for
36 the single stranded DNA exposed by initiating RNA polymerases (a DNA configuration
37 reproduced at stalled polymerases) without a requirement for specific cofactors.

38

39 **Introduction**

40 Cytidine deaminases are a family of polynucleotide mutators that modify cytosines into
41 uracil in viral nucleic acids as part of the innate immune defences (Harris and Liddament,
42 2004). Their success in restricting infection is reflected in the fact that the family has
43 undergone a rapid expansion in primates and humans (Jarmuz et al., 2002). The ancestor
44 founder of the family, activation induced deaminase (AID) functions in the adaptive immune
45 system to mutate antibody genes in B cells as a fast mechanism to promote diversity of the
46 antibody response matching the rapid evolution of pathogens during infection. The
47 evolutionary advantages of these strategies are counterbalanced by the risk of exposing the
48 host genome to active mutagenesis, a frequent cause of oncogenic transformation in
49 leukaemia and lymphomas of B cell origin.

50

51 All members of the AID/APOBEC family are selective in the sequence context of the
52 deaminated cytosine, with the two preceding nucleotides identifying the signature of
53 individual deaminases (Beale et al., 2004). This mutation context signature has identified
54 the human APOBEC3A and 3B proteins as the source of many of the somatic mutations

55 accumulated by cancer genomes (Burns et al., 2013; Nik-Zainal et al., 2012; Roberts et al.,
56 2013). The combined mutational landscape observed in mammalian genomes is
57 complicated by the contribution from multiple cellular processes in addition to enzymatic
58 deamination, such as metabolic oxidation, methyl-CpG deamination and aging, thus
59 elucidating the precise contribution of the APOBECs is far from straightforward (Alexandrov
60 et al., 2013; Lawrence et al., 2013). However the peculiar clustering of same strand
61 mutations at TpC dinucleotides observed in kataegic mutations in breast cancers constitutes
62 a hallmark of the APOBEC3A and 3B deaminases that can be experimentally induced. Repair
63 of double stranded DNA breaks can expose long patches of single stranded DNA with
64 multiple deaminations leading to the mutation clusters observed in association with
65 genomic rearrangements in breast cancer genomes (Nik-Zainal et al., 2012; Roberts et al.,
66 2012; Taylor et al., 2013).

67

68 Physiologically, the activity of such mutators is targeted to specific substrates and restricted
69 from the rest of the genome to limit genomic instability. In the case of AID, expressed upon
70 activation in only a fraction of B cells, by limiting access to the nuclear compartment and
71 preferential recruitment to the immunoglobulin genes; in the case of APOBEC3G, expressed
72 preferentially in lymphoid cells, by its localisation in the cytosol and binding to the viral
73 genome and capsid. The mechanism that preferentially directs AID to the immunoglobulin
74 genes is not fully understood, but active transcription has been repeatedly invoked as a
75 requirement (reviewed in (Storb, 2014)) and many of the proteins found to be associated
76 with AID are also involved in transcription and mRNA processing (Basu et al., 2011; Okazaki
77 et al., 2011; Pavri et al., 2010; Willmann et al., 2012)). Access of AID to off-target loci are
78 documented not only by the anecdotal occurrence of mutations in oncogenes and
79 chromosomal break points bearing the signature of the deaminase (Pasqualucci et al., 2001)
80 (Bcl6, MYC) but also by AID dependent chromosome-break-capture and direct CHIP, where
81 wide spread off-target presence of AID is experimentally detected outside the
82 immunoglobulin locus in mouse B cells (Chiarle et al., 2011; Klein et al., 2011; Yamane et al.,
83 2011).

84

85 In addition to the sporadic off-target mutations induced by AID in B cells, APOBEC3A and 3B
86 are thought to be responsible for many of the non-clustered/non-kataegic mutations at TpC
87 dinucleotides observed not only in breast cancers but in other tumour types where the

88 kataegic signature is not obviously present (Kuong and Loeb, 2013). As with sporadic AID
89 mutations, the circumstances that promote or grant access of the APOBECs to single
90 stranded DNA substrates of the host are not known. We have shown that overexpression of
91 deaminases in yeast faithfully recapitulates the mutation signatures observed in mammalian
92 genomes. Here we have attempted to identify genomic features that promote or are
93 permissive for enzymatic deamination by footprinting mutator activity on multiple genomes.
94 Our results indeed reveal a preferential targeting of the deaminases to defined regions of
95 the genome that is not dependent on cofactors but is rather based on accessibility, with
96 structural features of the DNA at the promoter of actively transcribed genes being the key
97 determinant. We also uncover a potential mechanistic explanation for the targeting and off-
98 target preferences of the antibody diversification mutator AID.

99

100 **Results**

101 **AID and APOBEC3G extensively mutate the yeast genome.**

102 Overexpression of cytidine deaminases in yeast leads to the accumulation of genome wide
103 mutations, which can be monitored by the number of cells resistant to the arginine analogue
104 L-Canavanine through inactivation of the arginine permease CAN1 gene (Figure 1A). We have
105 previously shown that such overexpression leads to an UNG dependent enrichment of
106 kataegic mutations through deamination of cytosines on single stranded DNA intermediates
107 during the repair of double strand breaks (Taylor et al., 2013). To assess the distribution of
108 isolated mutations, we obtained a dataset largely devoid of kataegic mutations by
109 expressing the deaminases in *ungΔ* cells. Overexpression of AID* (an AID hyperactive
110 mutant (Taylor et al., 2013; Wang et al., 2009)) in haploid cells results in highly elevated
111 frequency of Canavanine resistant colonies (164×10^{-6}), but relatively few mutations,
112 averaging 61 SNVs (single nucleotide variations) per genome (Figure 1A and B). Diploid cells
113 can overcome this limit as they avoid the reduction in fitness costs caused by accumulated
114 mutation (Lada et al., 2013; Waters et al., 1973). Our experimental setting confirms this
115 effect; whereas the mutation frequency is reduced almost 40 fold due to the requirement to
116 inactivate both CAN1 alleles, the genome wide SNV increase over 10 fold, averaging 796
117 SNVs per genome for AID* and 592 SNVs for transformants expressing sA3G* (a hyperactive
118 mutant of the catalytic domain of human APOBEC3G (Wang et al., 2009); Figure 1A and B).
119 For comparison, a database of mutations at C•G pairs was generated using the alkylating

120 agent ethyl methane sulfonate (EMS). Alkylation of guanosines promotes base pairing with
121 thymine, thereby causing G>A transitions during replication. Overnight exposure of diploid
122 cells to 0.2% EMS resulted in increased mutation frequency and SNV load per genome
123 similar to that elicited by the deaminases (Figure 1).

124

125 When interrogating the mutations (99.8% of which occur at C:G pairs; A:T mutations were
126 excluded from further analysis; all detected mutations are given in Supplementary file 1),
127 the expected flanking sequence context of WRC was found for AID* and YCC for sA3G*
128 (Figure 1C). In stark contrast, no consensus motif was observed in the EMS data, highlighting
129 the random nature of this mutagenesis. In all three datasets SNVs appeared distributed
130 throughout the genome, with all chromosomes displaying similar overall mutation that is
131 strongly correlated with chromosome length, ruling out major biases in the targeting of
132 mutations (Spearman's correlation coefficient for AID*: $\rho > 0.65$; for sA3G*: $\rho > 0.55$; for
133 EMS: $\rho > 0.68$; Figure 1D).

134

135 **Deaminase induced mutations are highly enriched in a small fraction of the genome.**

136 Whilst mutations are equally distributed amongst chromosomes, they are not uniformly
137 arranged along the chromosome. By combining the SNVs from independent transformants,
138 regions can be observed in AID* and sA3G* genomes which show pronounced mutational
139 peaks (Figure 2A). Only one such region of high mutation density is seen in the EMS treated
140 clones, that of the CAN1 gene. The presence of multiple loci with high mutation density is
141 therefore a deaminase specific process.

142

143 A more detailed look at regions with high density of mutations reveals narrow peaks of
144 accumulated mutation that are in many cases common to both deaminases (Figure 2B), with
145 the most prominent peaks resulting from the proximity of several regions of densely
146 targeted loci. These peaks represent high mutation densities within a bin size of 150 base
147 pairs but surprisingly reflect the accumulation of mutations focussed to very narrow
148 intervals within targeted loci (Figure 2C and D).

149

150 To further delineate mutation favoured loci, we defined regions of high mutation density by
151 identifying overlapping 150 base pair fragments containing higher than expected mutation
152 loads (minimum of six mutations per fragment, originating from three independent

153 transformants). We identify 1227 and 568 such mutation-enriched loci (MELs) in the AID*
154 and sA3G* treated genomes, in contrast to just 1 obtained for EMS treatment (overlapping
155 the body of the CAN1 gene and hence due to canavanine selection). On average 35 such
156 MELs would be expected for simulated datasets of equivalent mutation loads (Figure 2E and
157 Supplementary file 2). MELs span remarkably narrow regions, with a window width
158 averaging 110bp for AID* and 71bp for sA3G* (Figure 2F), and with almost 41% of all AID*
159 and 22% of all sA3G* induced mutations localised to these regions (Table 1 and
160 Supplementary file 2). In total, 25618 of the combined 72196 AID* and sA3G* mutations are
161 occurring in MELs which account for just 1.5% of the genome (Figure 2G).

162

163 Both AID and APOBEC3G target cytosines for deamination within a specific sequence
164 context, leading to the mutation hotspots associated with antibody diversification and the
165 recurrent mutations at CCC trinucleotides observed in HIV-1 genomes during the evolution
166 of viral clades and which accumulate in viral genomes from infected individual (Kijak et al.,
167 2008). We therefore analysed the distribution of AID and APOBEC3G preferred sequence
168 context in the yeast genome and find that the densities of AID and APOBEC3G motifs (WRC
169 and YCC respectively) show no enrichment within the highly targeted regions compared to
170 the remaining genome (Figure 2H). Therefore, the accumulation of mutations in MELs is not
171 a consequence of localised clustering of mutable motifs.

172

173 Reinforcing the notion that MELs are highly favoured targets for mutations, we find these
174 areas are frequency mutated on both alleles: 48% of AID* genomes and 56% of sA3G*
175 genomes have mutations within MELs occurring on both chromosome alleles, compared to
176 just 2-3% predicted for random fragments of equivalent size and mutation loads. MELs also
177 contain most of the homozygous mutations detected (82% of AID* and 78% of sA3G*
178 Targeting of both alleles in MELs suggests they represent highly mutable regions within the
179 genome, with the deaminases returning repeatedly to the same sites (albeit on a second
180 chromosome) to mutate.

181

182 Re-analysis of deaminase mutations we previously reported in haploid yeast (Taylor et al.,
183 2013) identified 39 MELs which overlap with hypermutated MELs in diploid yeast, thus the
184 focusing of mutations to MELs is seemingly unaffected by ploidy, suggesting the skewing of
185 mutations due to selective pressures, such as fitness, is negligible (Figure 2 - figure

186 supplement 1). Equally, we observe no significant strand bias in the hypermutated hotspots
187 associated with AID* MELs suggesting that both strands are targeted in a similar fashion. A
188 broader distribution of sA3G* strand bias more likely reflects the partial skewing in the
189 presence of YCC motifs at MELs (Figure 2 - figure supplement 2).

190

191 In conclusion, deaminases preferentially target narrow focussed regions throughout the
192 genome independent of the sequence density of deaminase targets.

193

194 **MELs exclusively overlap gene promoters**

195 There is a well-recognised relationship between AID induced mutations and transcription
196 both in B cells at immunoglobulin genes, and for off-target loci, with mutations preferentially
197 accumulating towards the promoter proximal region of the transcription unit (Pasqualucci et
198 al., 2001; Rada and Milstein, 2001). The transcription link is interpreted as a mechanism that
199 facilitates access of AID due to the generation of single stranded DNA intermediates
200 (Chaudhuri et al., 2003). We therefore wondered whether AID induced MELs would be
201 found associated with transcription. Contrary to expectation, enrichment analysis reveals
202 that both AID* and sA3G* MELs are depleted within the body of RNAP II transcribed mRNA
203 genes and rather that the deaminase induced mutations are preferentially associated with
204 promoter regions, with over 76% of deaminase targeted hotspots found at promoters,
205 compared to just 24% for simulated fragments (Figure 3A).

206

207 Initiation of replication also transiently generates single stranded DNA at defined locations.
208 However, there is no enrichment of mutated hotspots associated with replication origins
209 (ARS) (Figure 3A). Although this could reflect the relative depletion of mutable motifs within
210 ARS core consensus sequence, we find similar densities of mutable cytosines within the
211 broader sequence context encompassing 200 to 300 base pairs nucleosome depleted
212 regions associated with functional origins (Figure 3 – figure supplement 1), suggesting that
213 single strand availability provided by melting the DNA by the ORC complex might not be
214 sufficient to efficiently target the deaminases.

215

216 Mutation enrichment at promoters is not restricted to hotspots identified within MELs,
217 which exclude 73 % of the total mutations due to the threshold applied in defining enriched
218 loci. Aligning all mutations to mRNA transcriptional starts (TSS) and termination sites (TTS)

219 (Xu et al., 2009), revealed a strong association of deaminase induced mutations with the TSS,
220 with over 57% of AID* and 46% of sA3G* mutations occurring within the promoter region
221 (defined as 500 base pairs upstream and 50 base pairs downstream of the TSS), compared to
222 only 21% of EMS mutations (the expected frequency for randomly distributed mutations).
223 Mutation accumulation is skewed upstream of the TSS (peak at -21bp and -38bp for AID*
224 and sA3G* respectively; Figure 3B), corresponding to the nucleosome free region where the
225 pre-initiation RNAP complex forms before scanning for the TSS (Rhee and Pugh, 2012).
226 Indeed, aligning SNVs to the TATA box/TATA-like element or TSS revealed that not only are
227 the majority of promoter associated mutations occurring between these two features
228 (Figure 3C), there is also a paucity of mutations at the TATA-element suggesting this region is
229 protected by TBP/TFIID binding (this paucity is, at least for AID*, not due to an absence of
230 mutable sequence motifs; Figure 3D). Intriguingly, the peak of AID* and sA3G* induced
231 SNVs centred 30 base pairs from the TATA-element, the region where TBP guides TFIIB to
232 load RNAPII for the formation of the pre-initiation complex (PIC) (Rhee and Pugh, 2012).

233

234 The deaminase mutated hotspots thus identify the position where promoter melting occurs
235 before the scanning polymerase encounters the TSS, suggesting a mechanistic basis for the
236 hypothesis that the deaminases access the promoter coincidentally with the assembly of the
237 pre-initiation complex. Consistent with the notion that initiating polymerases create
238 transient access for the deaminases rather than specifically loading the proteins, we detect
239 robust association of RNAP II with the promoter region of deaminase targeted promoters in
240 yeast but negligible enhancement in the association of either AID or sA3G with mutated
241 promoters compared to unmutated or intergenic regions (Figure 3 - figure supplement 3).
242 Additionally, while there is a correlation between the mutated strand and the direction of
243 transcription (Figure 2 – figure supplement 2), MELs are predominantly composed of
244 mutations occurring in both strands suggesting the PIC makes both strands available during
245 initiation.

246

247 Supporting the idea that the deaminases preferentially mutate promoters due to their ability
248 to recognize the melted DNA associated with the transcription pre-initiation complex, we
249 observe that MELs occur in genes with above average transcriptional activity (García-
250 Martínez et al., 2004) but targeting appears unrelated to any particular transcriptional
251 program (Figure 3 – figure supplement 4). Rather than simply transcription factor binding at

252 the promoter, active initiation by RNAP II is important for MEL development (Wilcox test $p <$
253 0.005 for all groups; Figure 3E). The transition of RNAP II from the pre-initiation complex to
254 the elongation complex is associated with a shift in phosphorylation of the C-terminal
255 domain (CTD) serine 5/7 to serine 2 (Kim et al., 2010). In agreement with the transcription
256 rate analysis, deaminase MELs are associated with both high levels of RNAP II occupancy and
257 CTD-S5P, that parallels the association with the highest transcribed genes (Figure 3F).
258 Indeed, the recurrent association of both AID* and sA3G* MELs with regions enriched for
259 the basal transcription machinery and in particular Spt16 -a chromatin chaperon associated
260 with highly transcribed genes (Formosa, 2013) (Figure 3 – figure supplement 4), reinforces
261 the idea that active transcription and potential pausing (at promoters highly dependent on
262 the FACT/Spt16 complex) determines the deaminases targeting.

263

264 In summary, cytidine deaminases mutate at specific loci through the yeast genome,
265 predominantly within active gene promoter regions.

266

267 **AID targets promoter regions of small RNAP III and RNAP I genes**

268 In B cells, AID is found in association with components of the transcription machinery such
269 as SPT5 and SPT6, and RNAP II itself (Nambu et al., 2003; Okazaki et al., 2011; Pavri et al.,
270 2010), therefore we wondered whether the enrichment of mutations associated with
271 promoters might be a feature restricted to RNAP II dependent genes. Analysis of mutations
272 in highly transcribed non RNAP II dependent transcripts, such as RNAP III dependent tRNA
273 genes, astonishingly reveals an even more pronounced enrichment of targeted hotspots
274 with 78% of the genomic regions corresponding to tRNAs harbouring repeated mutations.
275 While we find that both sA3G* and AID* MELs overlap with tRNAs, AID* MELs are
276 disproportionately overrepresented, with 228 of 275 tRNA genes being highly targeted
277 (Figure 4A). Furthermore, aligning of mutations within 250 base pairs of the TSS of tRNA
278 genes shows that all occur within the tRNA gene body, which is also the site of RNAP III
279 initiation (Figure 4B). As in the case of mRNA promoters, the mutations in tRNAs are highly
280 focussed to narrow hotspots that span the site where loading of the polymerase is thought
281 to occur (Figure 4C).

282

283 The mutation frequency (normalised number of mutations per 550 base pairs) in AID*
284 genomes within tRNA genes is higher than at mRNA gene promoters (p value $< 2e-16$, Wilcox

285 non-parametric test; Figure 4C) and much higher than that observed even in the subset of
286 highly transcribed mRNA promoters. While the differences in mutation frequency between
287 mRNA promoters and tRNAs is still statistically significant for A3G* (p value < 8e-10), this
288 effect is less pronounced. Enhanced mutation is also observed in the promoters of
289 snoRNA/snRNA genes, again particularly in the case of AID* genomes, whereas no
290 statistically significant differences are observed between any of the promoter subsets for
291 mutations driven by EMS. The enhanced mutation attributable to AID* is not likely a feature
292 of RNAP III, since snoRNAs are even more targeted for mutation though all but snR52 are
293 transcribed by RNAP II (Moqtaderi and Struhl, 2004).

294

295 Targeting of tRNA, snRNA and snoRNA genes by the deaminases could be enhanced by the
296 availability of hypermutable motifs, as there is on average one more YCC motif in the tRNA
297 genes (1.5 more in the MEL region itself) targeted by sA3G* than in those tRNA genes not
298 targeted by sA3G*. We see no such difference with AID* target motifs which are present
299 within tRNA, snRNA and snoRNA gene promoters at similar frequency as in other promoters
300 (average 52 to 63 motifs per 550 base pairs promoter window, Figure 4 – figure supplement
301 1). Overall, there is only weak correlation between the number of motifs within the 550
302 base pair promoter window and the number of mutations (Spearman's $\rho = 0.02$ and $\rho = 0.2$
303 for AID* and sA3G* respectively; Figure 4 – figure supplement 2), confirming that motif
304 availability is not the main determinant for targeting.

305

306 Mutations at rRNA genes were poorly mapped due to the repetitive nature of the region on
307 Chr XII (150 to 200 copies of the 9.1 kb unit containing the 35S pre-RNA and the 5S RNA). By
308 including repeatedly mapped reads across the rDNA locus, we could detect several hundred
309 mutations at low allele frequency all within the expected deaminase mutation context,
310 giving confidence in their detection and location (Figure 4 – figure supplement 3A).
311 Mutations were restricted to the well defined ribosomal replication fork barrier (rRFB)
312 located between the 5S and 35S transcriptional units. No enhanced mutation was detected
313 at the promoter regions (which are transcribed in opposite directions by RNAP III and RNAP I
314 respectively). However mutations clustered at the rRFB site for both deaminases (Figure 4 –
315 figure supplement 3B), at a site where induced homologous recombination maintains the
316 size of the ribosomal gene array. Although DNA double-strand breaks (DSB) have been
317 detected at the site, it is likely that *in vivo* persistent breaks are rare in undamaged yeast

318 (Fritsch et al., 2010). Accordingly we did not detect kataegic like clusters in the region, but
319 rather localised mutated hotspots. Thus it is possible that other mechanisms such as cryptic
320 transcription (Houseley et al., 2007) might expose the site to the action of the deaminases,
321 rather than repair of double strand breaks. While AID overexpression in yeast deficient for
322 components of the RNA processing machinery (THO) have enhanced genomic instability,
323 particularly in highly transcribed GC-rich regions prone to R-loop formation (Gómez-
324 González and Aguilera, 2007), in wild type yeast this effect is only mild. Nonetheless we
325 observe positive association of MELs with predicted R-loop potential genes although the
326 paucity of these features across the genomes (between 59 – 78 sites) precludes any
327 predictive dissociation between high density of mutation, R-loop potential and transcription
328 rates (Figure 4 - figure supplement 4).

329

330 **AID but not sA3G binds small RNAs**

331 An alternative explanation for the enhanced targeting of small RNA promoters by AID* is
332 that the RNAs themselves preferentially bind AID, thereby creating co-transcriptional
333 enrichment of AID in the vicinity of their genes. Purified AID binds RNA, with its in vitro
334 deamination activity enhanced by treatment with RNase A (Bransteitter et al., 2003),
335 whereas the non catalytic domain of APOBEC3G is responsible for its ability to bind RNA and
336 form high molecular weight ribonucleic-protein complexes (Bélanger et al., 2013; Huthoff et
337 al., 2009). It is not known whether binding in both cases is specific for any particular RNA
338 species, but based on our current observations we decided to test the ability of human AID
339 and human APOBEC3G to bind in vitro transcribed tRNA as well as polyU RNA. Whereas
340 both Flag-tagged overexpressed human AID and full length human APOBEC3G can be
341 recovered from cell extracts by binding to biotin labelled RNAs, the catalytic domain of
342 APOBEC3G (sA3G) is not (Figure 5A). Furthermore, full length APOBEC3G is efficiently
343 recovered from extracts by the extended linear polyU RNA, a reflection of its ability to
344 oligomerise in an RNA dependent fashion, whereas AID recovery is not enhanced by its
345 binding to linear polyU RNA. Binding of AID to tRNA species was also found for endogenous
346 yeast tRNAs, suggesting that the modifications found in vivo (pseudouridylation and 2'-O-
347 ribose methylation) do not affect the interaction. The single domain APOBEC3A protein
348 shows no RNA binding ability except a limited amount to doubled stranded RNA, despite
349 sharing the preferential targeting to promoters as the rest of the deaminases (Figure 5 –
350 figure supplement 1A). Taken together, this data suggest a degree of specificity in the RNA

351 binding preferences of the deaminases, with AID preference linked to structured rather than
352 linear RNA (Figure 5B). Interestingly the catalytic activity of AID is not required for the
353 binding or the specificity, since similar binding was observed for the inactive mutant AID-
354 E58A (Figure 5A).

355

356 In order to test the RNA binding properties of AID in modulating its targeting preferences we
357 introduced a chimeric snR6 RNA into the RNAP II driven YBR194W gene, which was identified
358 in our dataset as a transcribed but poorly targeted promoter by both deaminases (Figure 5C
359 top panels). Initiation and transcription of the modified locus remained overall unaffected
360 (Figure 5 – figure supplement 2), while comparison of the YBR194W promoter region by
361 Sanger sequencing revealed enhanced mutation focused to the immediate vicinity of the TSS
362 by AID* but not sA3G*. No such focussing of mutations was observed in the unmodified
363 yeast overexpressing AID* (Figure 5C).

364

365 We conclude that the differential preference of AID* for tRNA and snoRNAs in yeast might
366 reflect the ability of AID to preferentially associate with abundant small RNA species, in
367 contrast to the catalytic domain of APOBEC3G (sA3G*) that possesses no RNA binding
368 activity.

369

370 Targeting mutations to initiating promoters is not likely a function of the size of the
371 deaminase, as could be inferred from the results described for both AID and the single
372 domain fragment of APOBEC3G used in our study. Similar promoter associated recurrent
373 mutations can be elicited not only by APOBEC3A (also a single domain deaminase) but also
374 by the double domain APOBEC3B (Figure 5 – figure supplement 1A). It is therefore not
375 entirely unexpected to observe enrichment of mutations at TpC (versus other dinucleotides)
376 in association with promoter regions in a breast cancer genome that has the highest
377 incidence of APOBEC3 kataegic mutations (Figure 5 – figure supplement 1B), suggesting that
378 the deaminases could access dsDNA at initiating or paused RNAPs also in mammalian cells.

379

380 ***Discussion***

381 The involvement of AID and APOBEC3A and 3B in cancer suggests that enzymatic
382 deamination of genomic targets is an infrequent but recurrent consequence of the presence

383 of the deaminases in vertebrates. Despite subcellular compartmentalization, specific
384 targeting and restricted expression limiting AID off-target activity, some genomic regions
385 other than the natural target, the immunoglobulin loci, are predisposed to mutation. BCL6,
386 PIM1 and MYC are recurrent off-targets of AID mutation in B cell malignancies (Pasqualucci
387 et al., 2001); in the case of BCL6 it is estimated that AID induced mutations are also
388 prevalent in non transformed B cells at just 10^3 fold lower frequency than at immunoglobulin
389 genes (Liu et al., 2008) and even in the absence of the mutator phenotypes attributable to
390 malignant transformation, normal B cells frequently show AID induced translocations at the
391 MYC locus (Casellas et al., 2009; Roschke et al., 1997). In cancer genomes, the association of
392 APOBEC mutations with genomic rearrangements suggests that replication stress, persistent
393 DNA lesions and incomplete repair expose single stranded DNA that becomes a substrate for
394 deaminases leading to clustered mutations. It is unclear how APOBEC3A and 3B gain access
395 to single stranded DNA leading to the singlet isolated mutations highly prevalent in mutated
396 cancer genomes that bear the APOBEC signature (Taylor et al., 2013). It is therefore
397 important to understand the genomic context that facilitates off-target activity of the
398 deaminases in the absence of explicit DNA damage.

399
400 Expression of AID and other APOBEC proteins in yeast faithfully recapitulates the signature
401 of mutations observed in mammalian cells in a smaller genome with no background
402 mutations due to unrelated processes, such as DNA repair (Lada et al., 2013; Taylor et al.,
403 2013). In this study we demonstrate the non-random nature of the mutations induced by
404 the deaminases, which is remarkably focussed to just 1.5% of the yeast genome but
405 nonetheless overlaps more than half of the active promoters. AID is known to interact with
406 components of the transcription machinery in mammalian cells (reviewed in (Kenter, 2012)).
407 However, the overlap between highly mutated promoters by both AID and APOBEC3G
408 suggests that rather than conservation of protein-protein interactions of the deaminases
409 with the transcription complex, properties of the promoter itself can determine targeting.

410
411 Enhanced targeting of RNAP III transcribed genes argues against active recruitment of the
412 deaminases by conserved initiation factors, whereas the structural conservation of the DNA
413 template conformation at the core pre-initiation complex of all polymerases (Vannini and
414 Cramer, 2012) supports the idea that the conformation of the DNA template is the common
415 element in the recruitment of the deaminases. Indeed, the site of polymerase loading

416 (within the body of the tRNA genes) rather than the TSS is the preferred target of
417 deamination in the case of the RNAP III transcribed tRNAs in contrast to the 5' region of the
418 RNAP III transcribed *SNR52* snoRNA, where the loading of the RNAP is fixed at the 5'
419 promoter region. Furthermore, the high density of mutations focussed to the small region
420 between the TATA binding protein site (TBP) and the transcription start site (TSS), more
421 precisely identify the pre-initiation complex (PIC) as the target for the deaminases.

422

423 Budding yeast RNAP II promoters show characteristic and highly regulated nucleosome
424 exclusion. This is partly due to sequence composition, with regions enriched for poly dA•dT
425 nucleotides that confer rigidity to the DNA and are therefore thermodynamically less
426 favourable to wrap around nucleosomes (Yuan et al., 2005), and partly due to the regulated
427 and precise positioning of the +1 nucleosome relative to the TSS that includes specific
428 histone variants (H2A.Z and H3.3) that promote chromatin accessibility (reviewed in (Jiang
429 and Pugh, 2009)). Therefore it is highly significant that other nucleosome free regions, such
430 as ARS are not targeted by the deaminases, despite undergoing DNA melting during the
431 initiation of replication. This reinforces our interpretation that intrinsic properties of active
432 promoters, in particular the configuration associated with loading of the polymerase at the
433 pre-initiation complex (open pre-initiation complex) (Grünberg et al., 2012), are sufficient to
434 generate persistent single stranded DNA accessible for deamination. Our data supports the
435 presence of such open PICs in most yeast active promoters.

436

437 Neither the preferential targeting of promoters nor the narrow focus of the MELs is due to
438 the preferential clustering of mutable motifs. Interestingly, the nature of the mutation
439 hotspots within MELs (both at C and G), reveals that both strands of the melted DNA
440 structure associated with active promoters are accessible. Furthermore, protection from
441 mutation is evident at the TBP binding site while the peak of mutations ~30 base pairs
442 downstream identifies the site of RNAP loading and DNA melting mapped by permanganate
443 footprinting (Giardina and Lis, 1993) and high resolution ChIP (Rhee and Pugh, 2012). Our
444 deaminase footprinting data further confirms the persistent open configuration and single
445 stranded nature of this region potentially identifying open pre-initiation promoters.

446

447 Differences in the assembly of the PIC in TATA and TATA-like promoters, do not seem to
448 affect mutation susceptibility, although predictably, TATA box promoters show a more

449 defined distance between the TBP protected footprint and the accessible melted DNA
450 (Figure 3 – figure supplement 2) indicating that it is the structure of the ssDNA rather than
451 the assembly (SAGA or THIID dependent) of the transcription initiation complex itself that
452 determines targeting (Rhee and Pugh, 2012).

453

454 Up to 75 % of human promoters in different cell types are occupied by a pre-initiating form
455 of RNAP II (Guenther et al., 2007), whereas pausing and stalling are much more common in
456 metazoan transcription compared with *S. cerevisiae*. Mammalian promoters are frequently
457 regulated by proximal pausing, with most promoters pausing within 200 base pairs of the
458 TSS (Adelman and Lis, 2012). In the presence of a deaminase, initiating and or paused sites
459 would become accessible for mutation, thus it is intriguing to observe promoter proximal
460 enrichment of mutations at TpC dinucleotides in PD4120a, a breast cancer genome with
461 dramatic accumulation of kataegis that betrays its mutagenesis by APOBEC3B (Nik-Zainal et
462 al., 2012). Our data favours the idea that accessibility of ssDNA at RNAP II stalled sites
463 suffice to recruit APOBECs or indeed AID. This model offers explanation for the association
464 of AID with mammalian SPT5, which functions in modulating the pausing of RNAP II during
465 elongation as transcription stalls, and is consistent with the recurrent targeting by AID of the
466 promoter proximal region of MYC (Duquette et al., 2005) a well characterised promoter-
467 proximal pausing regulated gene (Krumm et al., 1992; Strobl and Eick, 1992).

468

469 The correlation between high transcription rates and enhanced deaminase targeting
470 reinforces the hypothesis that repeated loading of the pre-initiation complex leads to the
471 persistence of a small region of melted DNA that is very efficiently targeted by the
472 deaminases. Indeed the enhanced targeting of tRNA, snoRNAs and snRNA genes could
473 reflect the high transcription rates of these essential RNAs given that RNAP I and III
474 transcripts constitute almost 80% of the total nuclear gene expression in dividing cells
475 (Vannini, 2013). The unexpected finding that tRNAs are disproportionately targeted for
476 mutation by AID compared with APOBEC3G, as are the promoters of other highly structured
477 RNAs (snRNA or snoRNA), and the indication that this difference is not due to motif
478 enrichment at those promoters, brings into focus the potential involvement of the RNA
479 binding properties of the deaminases in promoting targeting. While APOBEC3G has been
480 shown to bind not only HIV RNA, but cellular RNAs, including abundant 7S RNA (Huthoff et
481 al., 2009), this ability is dependent on the N-terminal domain. Mutation targeting of the

482 RNAP initiation complex is not linked to the ability of the deaminases to bind RNA per se, as
483 the catalytic C-terminal domain of APOBEC3G in this study is inert regarding RNA binding.
484 Notably, our results show that AID binds structured RNAs *in vitro* (such as tRNAs), and
485 preferentially targets tRNAs and other small RNA promoters for mutation in yeast,
486 prompting the speculation that binding to abundant RNAs sequesters AID to subnuclear
487 localities such as nucleolar areas, where small RNAs genes also localise during transcription.
488 Indeed nucleolar localisation of overexpressed AID has been reported in mammalian cells,
489 although its significance under physiological levels remains to be tested (Hu et al., 2013).
490 Alternatively preferential recognition of particular RNA structures such as folded tRNAs
491 could determine the recruitment of AID to genomic regions.

492

493 In conclusion, our study uncovers the remarkable preference of mammalian cytidine
494 deaminases to mutate active promoters when expressed in yeast, a preference blind to the
495 type of RNA polymerase (both RNAP II and III genes are targets) and not ascribable to
496 sequence context or targeting by specific cofactors. The precise and narrow location of the
497 recurrent mutations pinpoints the site where the RNAP pre-initiation complex is loaded
498 highlighting the conservation of the TBP (TATA binding protein) site and the formation of the
499 pre-initiation complex, whereas exclusion of mutations from the TBP site confirms the
500 poised nature of active yeast promoters.

501 These results suggest that initiating polymerases create a small but persistent accessible
502 patch of single stranded DNA *in vivo*, which has high affinity for deaminases and where both
503 strands are accessible for mutation. They also strongly support the notion that AID might
504 directly bind to ssDNA at the pre-initiating or stalling RNAP sites without a requirement for
505 specific cofactors and that its targeting is modulated by its ability to interact with structured
506 RNA species.

507

508 **Materials and methods**

509 **Yeast Transformants**

510 Yeast strain BY4743 *ungΔ/ungΔ* was generated by crossing BY4741 *ungΔ* (MATa; *his3Δ1*;
511 *leu2Δ0*; *met15Δ0*; *ura3Δ0*) obtained from Euroscarf deletion collection (Frankfurt,
512 Germany) with the BY4742 *ungΔ* [?]. BY4742 *ungΔ* was generated by removal of the
513 *UNG1* open reading frame by homologous recombination in the parental BY4742 strain,

514 using a PCR generated *URA3* cassette flanked by a 57-bp 5' homology and 51-bp 3' homology
515 arms that include adaptamers for post integration removal of the *URA3* selection cassette
516 (Reid et al., 2002). The YBR194W-snR6 chimeric strain was generated by inserting a *URA3*
517 cassette at the 5' end of the YBR194W gene in BY4741 *ungΔ* cells. Homology arms and the
518 snR6 gene were amplified from genomic DNA using the primers (1) 5'-
519 CCTGCCACTTTCAAAGGCG-3' and 5'-CGAAGGGTACTTCGCGAACTCCTGTCCCTATTACATATT
520 CAACC-3', (2) 5'-GGTTGAATATGTAATAGGGACAGGAGTTCGCGAAGTAACCCTTCG-3' and 5'-
521 GCCAGGCATGCTAATGGCAAACGAAATAAATCTCTTTGTAAAAC-3', (3) 5'-GTTTTACAAAGAGA
522 TTTATTTGTTTTGCCATTAGCATGCCTGGC-3' and 5'-TGGTGGTCATATGCTCGGTG-3'. A PCR
523 fusion of all three fragments with the first and last primer was used to retarget the *URA3*
524 containing locus. 5-Fluoroorotic acid counter-selection was used to isolate targeted colonies
525 that were then mated with BY4742 *ungΔ* to generate the final BY4743
526 *ungΔ/ungΔ* YBR194W-snR6/YBR194W strain. Correct integration of all targeting constructs
527 was confirmed by PCR.

528

529 Yeast transformation and selection, genomic DNA extraction and mutation frequency
530 calculation were performed as described previously (Taylor et al., 2013). Control and AID*
531 expression vectors were as described previously (Taylor et al., 2013). The sA3G* vector was
532 generated by PCR amplification of the C-terminal domain of A3G* fused with a 5' SV40
533 nuclear localization sequence and FLAG tag using primers 5'-
534 GCAAGCTTGCCACCATGCCTAAAAAGAAGCGTAAAGTCGAGATTCTCAGACACTCG-3' and 5'-
535 CCAGAATCAGGAAAACGGAGCAGACTACAAGGACGATGACGACAAGTAGCTCGAGGC-3' and ligating the
536 resultant Hind III-Xho I fragment it into pRS426-GAL1pr-tADHpolyA vector described
537 previously (Taylor et al., 2013).

538

539 Ethyl methanesulfonate (EMS) mutagenesis was performed by culturing BY4743 *ungΔ/ungΔ*
540 yeast overnight in YEPD with 0.2% EMS, after which cells were washed in 5% sodium
541 thiosulfate and plated for viability and canavanine resistance as above.

542

543 **Sample preparation and DNA sequencing**

544 DNA libraries were generated using the multiplexing Nextera DNA Sample Prep Kit (Illumina)
545 according to manufactures instructions. The libraries were sequenced by BGI (BGI, Beijing,
546 China). The de-multiplexed sequence reads were aligned to the reference yeast genome

547 (SacCer_Apr2011/sacCer3) using BWA-MEM (Li and Durbin, 2009). Optical duplicates were
548 removed using Picard (<http://picard.sourceforge.net>) and only uniquely mapped paired
549 reads were retained. On average 43-fold sequence coverage was achieved for each yeast
550 genome.

551

552 **Data analysis**

553 *Mutation calling*

554 An in-house pipeline for mutation calling was used where GATK base quality score
555 recalibration and indel realignment (McKenna et al., 2010) was performed prior to somatic
556 mutation calling by Somatic Sniper (Larson et al., 2012) using the parental BY4743 genome
557 as reference. High confidence single nucleotide variations (SNVs) were filtered using the
558 following criteria: (1) SomaticSniper score > 50, (2) allele frequency ≥ 0.3 , (3) reference or
559 samples read count ≥ 4 , (4) average position as fraction on reads ≥ 0.1 , (5) average distance
560 to 3' end ≥ 0.1 , (6) average base quality ≥ 30 , (7) average read length > 50 bp.

561

562 *Mutation enriched loci (MEL) identification*

563 Within each data set, mutations were pooled with the number of mutations within 150 base
564 pair windows. Based on the assumption of a random distribution of mutation amongst the
565 fragments, a binomial distribution was determined using the following parameters: size
566 equal to the average number of mutations per clone and probability equal to the average
567 number of mutations per clone over the total number of mutable motifs. Mutable motifs
568 were the total number of WRC, YCC, or C bases for AID*, sA3G* and EMS respectively. The
569 99th percentile was used as a threshold to identify significantly mutated windows and
570 adjacent windows merged. To refine the span of each individual mutation enriched loci
571 (MEL), unmutated residues and residues falling in the following categories were removed
572 and the window size adjusted: bases that had a count below the 25th percentile of all the
573 counts in the window; bases which had a mutation count below four standard deviations
574 from the average for the window and all bases with only a single detected mutation (where
575 the median mutation count was above one). A final threshold was applied so that only
576 regions with more than 5 mutations derived from at least four independent transformants
577 were assigned as high confidence MELs. All MELs were manually assessed using a genome
578 browser and are shown in Supplementary file 3.

579

580 The averaged fraction of overlapping regions for simulated MEL dataset were determined by
581 1000 cycles of bootstrap analysis using randomised equivalent number of fragments of
582 identical sizes for each dataset distributed across the genome.

583

584 *Normalised mutation density*

585 The normalised mutation density was calculated by dividing the mutation count for each
586 residue by the total number of mutation for the dataset.

587

588 *RNAP enrichment*

589 ChIP enrichment was determined by taking the sum of the ChIP enrichment scores (Kim et
590 al., 2010) for each promoter fragment (defined as 550 bp upstream and 50 bp downstream
591 from the TSS (Rhee and Pugh, 2012)). Promoters were then grouped according to the
592 transcription rate (García-Martínez et al., 2004) or whether they contained a MEL.

593

594 *Average mutation frequency for mRNA, tRNA, snoRNA and snRNA promoters*

595 Promoter fragments for mRNA genes and transcription rate binning were performed as
596 above. tRNA gene promoter fragments were defined as a 550 bp fragment centred on the
597 middle of the tRNA gene. snoRNA and snRNA promoters were defined as 550 bp upstream
598 and 50 bp downstream from the TSS defined in the Saccharomyces Genome Database
599 (Cherry et al., 2012). Intronic snoRNA genes were assigned the mRNA promoter and
600 polycistronic snoRNA genes were assigned only one promoter. Mutation frequency was
601 calculated by first randomly down-sampling the databases to half the size of the EMS
602 dataset, to allow equivalent numbers of mutations to be compared. The number of
603 mutations occurring on each promoter was then calculated. The process was bootstrapped
604 1000 times to give a directly comparable average number of mutations for each promoter.

605

606 *rDNA mapping*

607 To detect mutations at the repetitive rDNA locus a less stringent algorithm was used. De-
608 multiplexed sequence reads were aligned as before and unmapped reads removed. Reads
609 mapping to the rDNA region (chrXII:434839-508289) were extracted and used for mutation
610 calling by SomaticSniper. Mutations with a SomaticSniper score of above 50, a read depth of
611 10 in both the reference and the sample and no evidence of the mutated base in the
612 reference genome were assigned.

613

614 All analyses were performed using Bioconductor. Scripts are included as Supplementary file
615 4.

616

617 **Immunoprecipitation**

618 *RNA binding*

619 The tl(UAU)D RNA probes were generated by in-vitro transcription (MegaShortScript T7 Kit,
620 Life Technologies) with or without biotin-UTP (Life Technologies), according to manufactures
621 instructions. Free nucleotides were removed using Oligo Clean & Concentrator columns
622 (Zymo). The tl(UAU)D template was generated by annealing the following oligos 5'-
623 AATTTAATACGACTCACTATAGGGCTCGTGTAGCTCAGTGGTTAGAGCTTCGTGCTTATAACG-3' and
624 5'-TGCTCGAGGTGGGGTTTGAACCCACGACGGTCGCGTTATAAGCACGAAGCTCTAACC-3'. The
625 pre-tl(UAU)D template was generated by PCR amplification from yeast genomic DNA using
626 the following primers 5'- AATTTAATACGACTCACTATAGGGCTCGTGTAGCTCAGTGGTTAGAGC-
627 3' and 5'-TGCTCGAGGTGGGGTTTGAACCCACGACGG-3'. Biotinylation of total yeast RNA (Life
628 Technologies), polyuridylic acid, polyadenylic acid-polyuridylic acid (Sigma-Aldrich) and the
629 tl(UAU)D probe were performed using the RNA 3' End Biotinylation Kit (Pierce) according to
630 manufacturers instructions.

631

632 Biotinylated RNA probes (3.6 µg) were refolded by heating to 80° C for 5 minutes in folding
633 buffer (25 mM Tris pH 7.6, 100 mM KCl, 1 mM EDTA), MgCl₂ was then added to a final
634 concentration of 20 mM and the RNA allowed to slowly cool to 10° C before being bound to
635 magnetic beads (Pierce) for 1 hr at 4° C. Unbound probe was removed by washing with RNA
636 buffer (25 mM Tris pH 7.6, 50 mM KCl, 5 mM NaCl, 1.5 mM MgCl₂, 35 mM Glycine, 10 %
637 glycerol) supplemented with 0.5% Triton-X100. The integrity of the RNA was monitored by
638 denaturing gel electrophoresis and staining with toluidine blue.

639

640 Clarified whole cell extracts (in RNA buffer supplemented with 0.3% Triton-X100 and
641 complete protease inhibitors (Roche)) from HEK 293 cells expressing Flag-AID, catalytically
642 inactive AID (E58A mutation), APOBEC3G-Flag and the SV40-NLS tagged catalytic C-terminal
643 domain of APOBEC3G (sA3G)-Flag, were incubated for 1 hr at 4° C in the presence of bead
644 bound biotinylated RNA probes. Unbound proteins were removed by washing the beads
645 four times in RNA buffer supplemented with 0.5% Triton-X100 at 4° C and the bound protein
646 monitored by western using anti Flag antibodies (M2-HRP, Sigma).

647

648 *Chromatin immunoprecipitation*

649 Overnight 60 ml yeast cultures fixed in 1% formaldehyde for 20 minutes and quenched in
650 0.125M glycine (final) were washed twice in cold PBS, resuspended in RIPAl₀ (150 mM NaCl,
651 10 mM Tris-HCl pH 7.5, 1 mM EDTA, 1% Triton X-100, 0.1% SDS, 0.1% Sodium Deoxycholate,
652 1X Complete protease inhibitors) prior to lysis using a MPI TissueLyser (10 cycles of 30
653 seconds on, 5 mins off, 4000 rpm), sonication using a Bioruptor (14 cycles of 30 seconds on,
654 30 seconds off, high intensity) and centrifugation (10 mins 15,000g). Equal amounts of
655 clarified chromatin were incubated overnight at 4° C with 3 µg anti-HA 16B12 (Covance), 2
656 µg anti-H3 ab1791 (Abcam), 2 µg anti-RNAPII S5P ab5131 (Abcam). Purification followed on
657 Protein-G dynabeads for 2 hours with extensive washes [twice in RIPAl₀, twice in RIPAh_i
658 (RIPAl₀ but for 500 mM NaCl), once in RIPA-LiCl (RIPAl₀ but 250 mM LiCl replacing NaCl) and
659 twice in TE] and overnight elution in 25 mM Tris-HCl, 1 mM EDTA, pH to 9.8, 50 µg/ml
660 proteinase K at 65°C. Input DNA was extracted using Genra Puregene (Qiagen) with qPCR
661 performed using QuantiFast SYBR kit (Qiagen) all as per manufactures instructions. Primers
662 used are; YBR019C; 5'-ATCCAGCACCCACCTGTAACC-3' and 5'-AAACTTCTTTGCGTCCATCC-3',
663 YBR020W; 5'-ACCTGAGTTCAATTCTAGCGC-3' and 5'-TCCGGTTTAGCATCATAAGCG-3',
664 YNL067W; 5'-AACCAAACCTAGCCTCCAA-3' and 5'-TGCTGACAGTAACACCTTCTGG-3',
665 YBL003C; 5'-TGTGCACTCTACCAACTGGG-3' and 5'-ATGTCCGGTGGTAAAGGTGG-3', YPL250C;
666 5'-AGAGAGTTGCTCCAGACCCT-3' and 5'-GCATAAAGAAGCGGCTCTGC-3', YEL009C; 5'-
667 GGGGGAGAGTAACCTGTGTT-3' and 5'-TTTCGGCTCGCTGTCTTACC-3', YBR194W; 5'-
668 TCTTCTTGCTCGGGGTCTC-3' and 5'-TGCTGAAGGCCTTTGCAAAG-3', YPL189W; 5'-
669 GCGAAGATTACGGCACTCGA-3' and 5'-ACAGGTACGGGCTATCTGGA-3', YLR183C; 5'-
670 ACATCTGCCACGACACATCA-3' and 5'-TGGTGGAGAGTACGGATCCA-3', YJL105W; 5'-
671 TTTCTTGCTCTTGCGGCTA-3' and 5'-AGTTAGGATCTGAGCCGGGT-3', YPR007C; 5'-
672 ACAGGTTGAGCTTCATGGG-3' and 5'-CGGAATTTTCATCCAGCGGA-3', chrXV;367475-367594
673 5'-ACTTGGCACTTCTTCTCAACA-3' and 5'-TCGCAAAGTTGGCTAACCGT-3', chrX;585916-
674 586020 5'-ATGTCTCCCTGTTACCCGGT-3' and 5'-ACAGGTGCTGTACAAAACA-3', chrIV;76875-
675 76955 5'-GGCAGCACCGAGAATGTTTT-3' and 5'-GCTGTTAGCATATTGGGGGT-3'.

676

677 **Yeast transcript analysis**

678 RNA from 1 ml overnight cultures purified with RNAeasy plus (Qiagen) was used to generate
679 cDNA using oligo-dTs and the GoScript Kit (Promega) followed by qPCR employing

680 QuantiFast SYBR (Qiagen) all as per manufactures instructions. Primers used are TAF10; 5'-
681 ATATTCAGGATCAGGTCTTCCGTAGC-3' and 5'-GTAGTCTTCTCATTCTGTTGATGTTGTTGTTG-3',
682 ACT1; 5'-CTTTCAACGTTCCAGCCTTC-3' and 5'-CCAGCGTAAATTGGAACGAC-3', YBR194W-snR6;
683 5'-CCTGCCACTTTCAAAGGCG-3' and 5'-CAGGGGAAGTCTGATCATCTCTG-3', YBR194W; 5'-
684 GGGTCGTGAAAAGAGAACGG-3' and 5'-ATGTGATGGTGCAGTGCCTC-3'.

685

686 **YBR194W promoter sequencing**

687 The YBR194W promoter region was amplified using the following primers; 5'-
688 ATTGTGGCAGTTCGGCTTTG-3' and 5'-AGGTTTCCCAGTCTGGCTTG-3' and Sanger sequenced
689 using the latter.

690 ***Acknowledgements***

691 We are grateful to David Rueda and Myron Goodman for sharing unpublished results and
692 members of the Rada lab for helpful advice and discussions. The late Michael Neuberger
693 instigated the initial stages of this work and remains in memory an inspiration. This work
694 was supported by the Medical Research Council (MRC reference number MC_U105178806)
695 and through an MRC Centennial Award to BJMT.

696

697 ***References***

698

699 Adelman, K., and Lis, J.T. (2012). Promoter-proximal pausing of RNA polymerase II: emerging
700 roles in metazoans. *Nat. Rev. Genet.* *13*, 720–731.

701 Alexandrov, L.B., Nik-Zainal, S., Wedge, D.C., Aparicio, S.A.J.R., Behjati, S., Biankin, A.V.,
702 Bignell, G.R., Bolli, N., Borg, A., Børresen-Dale, A.-L., et al. (2013). Signatures of mutational
703 processes in human cancer. *Nature* *500*, 415–421.

704 Basu, U., Meng, F.-L., Keim, C., Grinstein, V., Pefanis, E., Eccleston, J., Zhang, T., Myers, D.,
705 Wasserman, C.R., Wesemann, D.R., et al. (2011). The RNA exosome targets the AID cytidine
706 deaminase to both strands of transcribed duplex DNA substrates. *Cell* *144*, 353–363.

707 Beale, R.C.L., Petersen-Mahrt, S.K., Watt, I.N., Harris, R.S., Rada, C., and Neuberger, M.S.
708 (2004). Comparison of the Differential Context-dependence of DNA Deamination by APOBEC
709 Enzymes: Correlation with Mutation Spectra in Vivo. *J. Mol. Biol.* *337*, 585–596.

710 Bélanger, K., Savoie, M., Rosales Gerpe, M.C., Couture, J.-F., and Langlois, M.-A. (2013).
711 Binding of RNA by APOBEC3G controls deamination-independent restriction of retroviruses.
712 *Nucleic Acids Research* *41*, 7438–7452.

713 Bransteitter, R., Pham, P., Scharff, M.D., and Goodman, M.F. (2003). Activation-induced
714 cytidine deaminase deaminates deoxycytidine on single-stranded DNA but requires the
715 action of RNase. *Proc. Natl. Acad. Sci. U.S.a.* *100*, 4102–4107.

716 Burns, M.B., Temiz, N.A., and Harris, R.S. (2013). Evidence for APOBEC3B mutagenesis in
717 multiple human cancers. *Nat Genet* *45*, 977–983.

718 Casellas, R., Yamane, A., Kovalchuk, A.L., and Potter, M. (2009). Restricting activation-
719 induced cytidine deaminase tumorigenic activity in B lymphocytes. *Immunology* *126*, 316–
720 328.

721 Chaudhuri, J., Tian, M., Khuong, C., Chua, K., Pinaud, E., and Alt, F.W. (2003). Transcription-
722 targeted DNA deamination by the AID antibody diversification enzyme. *Nature* *422*, 726–
723 730.

724 Cherry, J.M., Hong, E.L., Amundsen, C., Balakrishnan, R., Binkley, G., Chan, E.T., Christie, K.R.,
725 Costanzo, M.C., Dwight, S.S., Engel, S.R., et al. (2012). *Saccharomyces Genome Database: the*
726 *genomics resource of budding yeast. Nucleic Acids Research* *40*, D700–D705.

727 Chiarle, R., Zhang, Y., Frock, R.L., Lewis, S.M., Molinie, B., Ho, Y.-J., Myers, D.R., Choi, V.W.,
728 Compagno, M., Malkin, D.J., et al. (2011). Genome-wide translocation sequencing reveals
729 mechanisms of chromosome breaks and rearrangements in B cells. *Cell* *147*, 107–119.

730 Duquette, M.L., Pham, P., Goodman, M.F., and Maizels, N. (2005). AID binds to transcription-
731 induced structures in c-MYC that map to regions associated with translocation and
732 hypermutation. *Oncogene* *24*, 5791–5798.

733 Formosa, T. (2013). The role of FACT in making and breaking nucleosomes. *Biochim. Biophys.*
734 *Acta* *1819*, 247–255.

735 Fritsch, O., Burkhalter, M.D., Kais, S., Sogo, J.M., and Schär, P. (2010). DNA ligase 4 stabilizes
736 the ribosomal DNA array upon fork collapse at the replication fork barrier. *DNA Repair*
737 *(Amst.)* *9*, 879–888.

738 García-Martínez, J., Aranda, A., and Pérez-Ortín, J.E. (2004). Genomic run-on evaluates
739 transcription rates for all yeast genes and identifies gene regulatory mechanisms. *Mol. Cell*
740 *15*, 303–313.

741 Giardina, C., and Lis, J.T. (1993). DNA melting on yeast RNA polymerase II promoters. *Science*
742 *261*, 759–762.

743 Ginno, P.A., Lott, P.L., Christensen, H.C., Korf, I., and Chédin, F. (2012). R-loop formation is a
744 distinctive characteristic of unmethylated human CpG island promoters. *Mol. Cell* *45*, 814–
745 825.

746 Gómez-González, B., and Aguilera, A. (2007). Activation-induced cytidine deaminase action is
747 strongly stimulated by mutations of the THO complex. *Proc. Natl. Acad. Sci. U.S.a.* *104*,
748 8409–8414.

749 Grünberg, S., Warfield, L., and Hahn, S. (2012). Architecture of the RNA polymerase II
750 preinitiation complex and mechanism of ATP-dependent promoter opening. *Nat. Struct.*
751 *Mol. Biol.* *19*, 788–796.

- 752 Guenther, M.G., Levine, S.S., Boyer, L.A., Jaenisch, R., and Young, R.A. (2007). A chromatin
753 landmark and transcription initiation at most promoters in human cells. *Cell* *130*, 77–88.
- 754 Harris, R.S., and Liddament, M.T. (2004). Retroviral restriction by APOBEC proteins. *Nat. Rev.*
755 *Immunol.* *4*, 868–877.
- 756 Houseley, J., Kotovic, K., Hage, El, A., and Tollervey, D. (2007). Trf4 targets ncRNAs from
757 telomeric and rDNA spacer regions and functions in rDNA copy number control. *Embo J.* *26*,
758 4996–5006.
- 759 Hu, Y., Ericsson, I., Torseth, K., Methot, S.P., Sundheim, O., Liabakk, N.B., Slupphaug, G., Di
760 Noia, J.M., Krokan, H.E., and Kavli, B. (2013). A combined nuclear and nucleolar localization
761 motif in activation-induced cytidine deaminase (AID) controls immunoglobulin class
762 switching. *J. Mol. Biol.* *425*, 424–443.
- 763 Huthoff, H., Autore, F., Gallois-Montbrun, S., Fraternali, F., and Malim, M.H. (2009). RNA-
764 dependent oligomerization of APOBEC3G is required for restriction of HIV-1. *PLoS Pathog.* *5*,
765 e1000330.
- 766 Jarmuz, A., Chester, A., Bayliss, J., Gisbourne, J., Dunham, I., Scott, J., and Navaratnam, N.
767 (2002). An anthropoid-specific locus of orphan C to U RNA-editing enzymes on chromosome
768 22. *Genomics* *79*, 285–296.
- 769 Jiang, C., and Pugh, B.F. (2009). Nucleosome positioning and gene regulation: advances
770 through genomics. *Nat. Rev. Genet.* *10*, 161–172.
- 771 Kenter, A.L. (2012). AID targeting is dependent on RNA polymerase II pausing. *Seminars in*
772 *Immunology* *24*, 281–286.
- 773 Kijak, G.H., Janini, L.M., Tovanabutra, S., Sanders-Buell, E., Arroyo, M.A., Robb, M.L., Michael,
774 N.L., Birx, D.L., and McCutchan, F.E. (2008). Variable contexts and levels of hypermutation in
775 HIV-1 proviral genomes recovered from primary peripheral blood mononuclear cells.
776 *Virology* *376*, 101–111.
- 777 Kim, H., Erickson, B., Luo, W., Seward, D., Graber, J.H., Pollock, D.D., Megee, P.C., and
778 Bentley, D.L. (2010). Gene-specific RNA polymerase II phosphorylation and the CTD code.
779 *Nat. Struct. Mol. Biol.* *17*, 1279–1286.
- 780 Klein, I.A., Resch, W., Jankovic, M., Oliveira, T., Yamane, A., Nakahashi, H., Di Virgilio, M.,
781 Bothmer, A., Nussenzweig, A., Robbiani, D.F., et al. (2011). Translocation-capture sequencing
782 reveals the extent and nature of chromosomal rearrangements in B lymphocytes. *Cell* *147*,
783 95–106.
- 784 Krumm, A., Meulia, T., Brunvand, M., and Groudine, M. (1992). The block to transcriptional
785 elongation within the human c-myc gene is determined in the promoter-proximal region.
786 *Genes Dev.* *6*, 2201–2213.
- 787 Kuong, K.J., and Loeb, L.A. (2013). APOBEC3B mutagenesis in cancer. *Nat Genet* *45*, 964–965.
- 788 Lada, A.G., Stepchenkova, E.I., Waisertreiger, I.S.-R., Noskov, V.N., Dhar, A., Eudy, J.D.,
789 Boissy, R.J., Hirano, M., Rogozin, I.B., and Pavlov, Y.I. (2013). Genome-wide mutation
790 avalanches induced in diploid yeast cells by a base analog or an APOBEC deaminase. *PLoS*

791 Genet. 9, e1003736.

792 Larson, D.E., Harris, C.C., Chen, K., Koboldt, D.C., Abbott, T.E., Dooling, D.J., Ley, T.J., Mardis,
793 E.R., Wilson, R.K., and Ding, L. (2012). SomaticSniper: identification of somatic point
794 mutations in whole genome sequencing data. *Bioinformatics* 28, 311–317.

795 Lawrence, M.S., Stojanov, P., Polak, P., Kryukov, G.V., Cibulskis, K., Sivachenko, A., Carter,
796 S.L., Stewart, C., Mermel, C.H., Roberts, S.A., et al. (2013). Mutational heterogeneity in
797 cancer and the search for new cancer-associated genes. *Nature* 499, 214–218.

798 Li, H., and Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler
799 transform. *Bioinformatics* 25, 1754–1760.

800 Liu, M., Duke, J.L., Richter, D.J., Vinuesa, C.G., Goodnow, C.C., Kleinstein, S.H., and Schatz,
801 D.G. (2008). Two levels of protection for the B cell genome during somatic hypermutation.
802 *Nature* 451, 841–845.

803 McKenna, A., Hanna, M., Banks, E., Sivachenko, A., Cibulskis, K., Kernytsky, A., Garimella, K.,
804 Altshuler, D., Gabriel, S., Daly, M., et al. (2010). The Genome Analysis Toolkit: a MapReduce
805 framework for analyzing next-generation DNA sequencing data. *Genome Research* 20, 1297–
806 1303.

807 Moqtaderi, Z., and Struhl, K. (2004). Genome-wide occupancy profile of the RNA polymerase
808 III machinery in *Saccharomyces cerevisiae* reveals loci with incomplete transcription
809 complexes. *Mol. Cell. Biol.* 24, 4118–4127.

810 Nambu, Y., Sugai, M., Gonda, H., Lee, C.-G., Katakai, T., Agata, Y., Yokota, Y., and Shimizu, A.
811 (2003). Transcription-coupled events associating with immunoglobulin switch region
812 chromatin. *Science* 302, 2137–2140.

813 Nik-Zainal, S., Alexandrov, L.B., Wedge, D.C., Van Loo, P., Greenman, C.D., Raine, K., Jones,
814 D., Hinton, J., Marshall, J., Stebbings, L.A., et al. (2012). Mutational processes molding the
815 genomes of 21 breast cancers. *Cell* 149, 979–993.

816 Okazaki, I.-M., Okawa, K., Kobayashi, M., Yoshikawa, K., Kawamoto, S., Nagaoka, H.,
817 Shinkura, R., Kitawaki, Y., Taniguchi, H., Natsume, T., et al. (2011). Histone chaperone Spt6 is
818 required for class switch recombination but not somatic hypermutation. *Proc. Natl. Acad.*
819 *Sci. U.S.a.* 108, 7920–7925.

820 Pasqualucci, L., Neumeister, P., Goossens, T., Nanjangud, G., Chaganti, R.S., Küppers, R., and
821 Dalla-Favera, R. (2001). Hypermutation of multiple proto-oncogenes in B-cell diffuse large-
822 cell lymphomas. *Nature* 412, 341–346.

823 Pavri, R., Gazumyan, A., Jankovic, M., Di Virgilio, M., Klein, I., Ansarah-Sobrinho, C., Resch,
824 W., Yamane, A., Reina San-Martin, B., Barreto, V., et al. (2010). Activation-induced cytidine
825 deaminase targets DNA at sites of RNA polymerase II stalling by interaction with Spt5. *Cell*
826 143, 122–133.

827 Rada, C., and Milstein, C. (2001). The intrinsic hypermutability of antibody heavy and light
828 chain genes decays exponentially. *Embo J.* 20, 4570–4576.

829 Reid, R.J.D., Sunjevaric, I., Kedacche, M., and Rothstein, R. (2002). Efficient PCR-based gene

830 disruption in *Saccharomyces* strains using intergenic primers. *Yeast* (Chichester, England) *19*,
831 319–328.

832 Rhee, H.S., and Pugh, B.F. (2012). Genome-wide structure and organization of eukaryotic
833 pre-initiation complexes. *Nature* *483*, 295–301.

834 Roberts, S.A., Lawrence, M.S., Klimczak, L.J., Grimm, S.A., Fargo, D., Stojanov, P., Kiezun, A.,
835 Kryukov, G.V., Carter, S.L., Saksena, G., et al. (2013). An APOBEC cytidine deaminase
836 mutagenesis pattern is widespread in human cancers. *Nat Genet* *45*, 970–976.

837 Roberts, S.A., Sterling, J., Thompson, C., Harris, S., Mav, D., Shah, R., Klimczak, L.J., Kryukov,
838 G.V., Malc, E., Mieczkowski, P.A., et al. (2012). Clustered mutations in yeast and in human
839 cancers can arise from damaged long single-strand DNA regions. *Mol. Cell* *46*, 424–435.

840 Roschke, V., Kopantzev, E., Dertzbaugh, M., and Rudikoff, S. (1997). Chromosomal
841 translocations deregulating c-myc are associated with normal immune responses. *Oncogene*
842 *14*, 3011–3016.

843 Storb, U. (2014). Why does somatic hypermutation by AID require transcription of its target
844 genes? *Adv. Immunol.* *122*, 253–277.

845 Strobl, L.J., and Eick, D. (1992). Hold back of RNA polymerase II at the transcription start site
846 mediates down-regulation of c-myc in vivo. *Embo J.* *11*, 3307–3314.

847 Taylor, B.J., Nik-Zainal, S., Wu, Y.L., Stebbings, L.A., Raine, K., Campbell, P.J., Rada, C.,
848 Stratton, M.R., and Neuberger, M.S. (2013). DNA deaminases induce break-associated
849 mutation showers with implication of APOBEC3B and 3A in breast cancer kataegis. *eLife* *2*,
850 e00534.

851 Vannini, A. (2013). A structural perspective on RNA polymerase I and RNA polymerase III
852 transcription machineries. *Biochim. Biophys. Acta* *1829*, 258–264.

853 Vannini, A., and Cramer, P. (2012). Conservation between the RNA polymerase I, II, and III
854 transcription initiation machineries. *Mol. Cell* *45*, 439–446.

855 Venters, B.J., Wachi, S., Mavrich, T.N., Andersen, B.E., Jena, P., Sinnamon, A.J., Jain, P.,
856 Roller, N.S., Jiang, C., Hemeryck-Walsh, C., et al. (2011). A comprehensive genomic binding
857 map of gene and chromatin regulatory proteins in *Saccharomyces*. *Mol. Cell* *41*, 480–492.

858 Wang, M., Yang, Z., Rada, C., and Neuberger, M.S. (2009). AID upmutants isolated using a
859 high-throughput screen highlight the immunity/cancer balance limiting DNA deaminase
860 activity. *Nat. Struct. Mol. Biol.* *16*, 769–776.

861 Waters, R., Waters, R., Waters, R., Parry, J.M., Parry, J.M., and Parry, J.M. (1973). The
862 response to chemical mutagens of the individual haploid and homoallelic diploid UV-
863 sensitive mutants of the rad 3 locus of *Saccharomyces cerevisiae*. *Molecular & General*
864 *Genetics* : *MGG* *124*, 135–143.

865 Willmann, K.L., Milosevic, S., Pauklin, S., Schmitz, K.-M., Rangam, G., Simon, M.T., Maslen, S.,
866 Skehel, M., Robert, I., Heyer, V., et al. (2012). A role for the RNA pol II-associated PAF
867 complex in AID-induced immune diversification. *Journal of Experimental Medicine* *209*,
868 2099–2111.

869 Wongsurawat, T., Jenjaroenpun, P., Kwoh, C.K., and Kuznetsov, V. (2012). Quantitative
870 model of R-loop forming structures reveals a novel level of RNA-DNA interactome
871 complexity. *Nucleic Acids Research* 40, e16.

872 Xu, Z., Wei, W., Gagneur, J., Perocchi, F., Clauder-Münster, S., Camblong, J., Guffanti, E.,
873 Stutz, F., Huber, W., and Steinmetz, L.M. (2009). Bidirectional promoters generate pervasive
874 transcription in yeast. *Nature* 457, 1033–1037.

875 Yamane, A., Resch, W., Kuo, N., Kuchen, S., Li, Z., Sun, H.-W., Robbiani, D.F., McBride, K.,
876 Nussenzweig, M.C., and Casellas, R. (2011). Deep-sequencing identification of the genomic
877 targets of the cytidine deaminase AID and its cofactor RPA in B lymphocytes. *Nat. Immunol.*
878 12, 62–69.

879 Yuan, G.-C., Liu, Y.-J., Dion, M.F., Slack, M.D., Wu, L.F., Altschuler, S.J., and Rando, O.J.
880 (2005). Genome-scale identification of nucleosome positions in *S. cerevisiae*. *Science* 309,
881 626–630.

882

883 **Figures and Legends**

884 **Figure 1.**

885 Genome wide distribution and signature of unclustered deaminase induced mutations in
886 *ung1Δ* diploid yeast.

887 A) Mutation frequency (expressed as the number of canavanine resistant colonies per 10^6) at
888 the CAN1 locus in *ung1Δ* haploid yeast (data in part from (Taylor et al., 2013)) and
889 *ung1Δ/ung1Δ* diploid yeast transformants expressing AID/APOBEC proteins or upon
890 treatment with 0.2% EMS. Red bars indicate the median mutation frequency (n=12-126
891 colonies).

892 B) Genome wide SNV number in *ung1Δ* haploid and *ung1Δ/Δ* diploid yeast transformants
893 expressing AID/APOBEC proteins or with EMS treatment. Red bars indicate the median
894 mutation per genome (n=25-50 independent clones).

895 C) Sequence context of mutations at G•C pairs in diploid yeast genomes (indicated as
896 mutations at cytosines) exposed to AID*, sA3G* or EMS mutagenesis. The numbers indicate
897 total mutations per dataset, with the height of colour bars proportional to the frequency of
898 each base found in the vicinity of a mutation.

899 D) Distribution of mutations per diploid yeast chromosome expressed as the number of
900 mutations per chromosome in each independent genome against the chromosome length.
901 The bars represent the projected linear trend for mutations at C (in black) or G (in red).

902

903 **Figure 2.**

904 Mutation enriched loci (MELs) identified by focussed deaminase-induced mutation.

905 A) Radial histograms depict the density (Z-score) of pooled mutations for each dataset in 2
906 kb overlapping genomic segments along each chromosome. The CAN1 locus is highlighted in
907 red. The peak highlighted in cyan is further enlarged in panels B, C and D.

908 B) Mutation densities along ChrII in AID* (red), sA3G* (black) and EMS (blue) treated
909 genomes, expressed as the Z-score of mutation density per dataset (y-axis) along
910 chromosome II (x-axis; 200 bp bin size). The region shadowed in cyan is magnified in C.

911 C) Regions of high mutation density identify narrow mutation enriched regions (MELs),
912 shown as green boxes for AID* and purple boxes for sA3G* in the bottom panel. Horizontal
913 lines represent a single genome with each non-clonal mutation at C or G indicated by a dot
914 (black or red respectively). Regions in Chr II and Chr X containing mutation enriched loci
915 shown at the same scale, with the genomic coordinates indicated.

916 D) Mutations in the pronounced MEL on ChrII (highlighted cyan in panels A, B and C) shown
917 in green for AID* and purple for sA3G*. Coordinates are indicated.

918 E) Overlap of detected MELs in AID*, sA3G* and EMS datasets.

919 F) Distribution of MELs width with the median indicated for AID* and sA3G* mutated
920 genomes.

921 G) Fraction of the total deaminase mutations in MELs (black boxes) relative to genomic
922 coverage of MELs.

923 H) Distribution of distances between AID and A3G mutable motifs within MELs versus
924 genome wide mutable motif distances.

925

926 **Figure 2 - figure supplement 1.**

927 Overlap between Haploid and Diploid MELs.

928

929 **Figure 2 - figure supplement 2.**

930 Strand bias in deaminase induced mutations calculated as fraction of mutations at C (+
931 strand) or G (- strand) within each MEL.

932 (A) Strand distribution of mutations within AID* and sA3G* MEL regions. MELs comprising a
933 single base were excluded.

934 (B) Strand distribution of mutations within MEL regions in relation to the direction of
935 transcription of the associated gene.

936 (C) Strand distribution of WRC and YCC deaminase motifs within MEL regions and their
937 flanking 50 base pairs.

938

939 **Figure 3.**

940 Deaminase mutation footprints are focussed to the pre-initiation complex region of active
941 promoters.

942 A) Proportion of promoters, gene bodies, intergenic regions and replication origins (ARS)
943 harbouring a MEL (green) or not (grey) for AID* and sA3G* datasets versus the expected
944 distribution (sim.AID*sA3G*) determined by Monte Carlo simulation of equivalent sized
945 fragments for each MEL dataset distributed randomly across the genome.

946 B) Density of mutations in relation to their distance to the nearest transcription start site
947 (TSS) of mRNA (RNAP II) transcripts compared to the density relative to transcription
948 termination sites (TTS). Data includes all mutations in addition to MELs.

949 C) Deaminase mutations relative to the TATA or TATA-like element for each RNAP II
950 promoters (Rhee and Pugh, 2012) compared to the mutation distance distribution aligned to
951 the transcription start site (TSS).

952 D) Proportion of AID* or sA3G* mutable motifs within RNAP II promoter regions, centred on
953 the TATA-elements (Rhee and Pugh, 2012). Total number of mutations for each dataset is
954 shown at each position (black line).

955 E) Relative transcription rates (see methods) at RNAP II promoters targeted by MELs
956 compared to relative transcription rates for all RNAP II genes in gal induced conditions
957 (García-Martínez et al., 2004).

958 F) Relative enrichment of RNAP II and RNAP II CTD phosphorylation (S2P, S5P and S7P) in
959 promoters containing AID* (red) and sA3G* (black) MELs and all RNAP II promoters (grey)
960 ranked according to transcriptional activity (García-Martínez et al., 2004).

961

962 **Figure 3 - figure supplement 1.**

963 Paucity of deaminase mutations at replication origins is not a consequence of absence of
964 mutable motifs.

965 Proportion of AID* or sA3G* mutable motifs around replication origins (ARS), depicted as in
966 Figure 3D. Total number of mutations for each dataset is shown for at position (black line,
967 scale as in Figure 3D).

968

969 **Figure 3 - figure supplement 2.**

970 Density of mutations (x-axis; AID*, red; sA3G*, black; EMS, blue) in relation to their distance
971 to the nearest TATA box or TATA-like element, grouped according to the TAF1 enrichment
972 status (data from (Rhee and Pugh, 2012)). Data includes all mutations in addition to MELs.

973

974 **Figure 3 - figure supplement 3.**

975 Distribution of the deaminases on chromatin is unrelated to mutation preferences.

976 Enrichment of A) deaminase, B) serine 5 phosphorylated RNAPII and C) Histone H3 at MEL
977 promoters, unmutated promoters and intergenic regions. Enrichment is shown relative to
978 input chromatin (B and C) or further normalised to control cell lines (A). Data from 2-3
979 independent experiments.

980

981 **Figure 3 - figure supplement 4.**

982 Transcription factor binding sites compared to MEL preferences.

983 A) Frequency of each yeast transcription factor at individual promoters as described in
984 (Venters et al., 2011) (blue dots) compared with the frequency that the transcription factor
985 appears in the promoter of genes containing AID* (red dots) and sA3G* (black dots) MELs.
986 Factors are ordered according to number of binding sites in all promoters. Basal
987 transcription factors are the most commonly associated with deaminase targeted promoters
988 (labelled).

989 B) Transcription rates of genes grouped according to Spt16 promoter occupancy and
990 presence of MELs.

991 C) List of transcription factors found to vary in occupancy at MEL targeted promoters versus
992 their overall frequency at all yeast promoters (Venters dataset). Transcription factors
993 showing $\pm 10\%$ variation which are present in at least 25% of MELs are listed.

994

995 **Figure 4.**

996 AID* and sA3G* target both RNAP II and RNAP III promoters

997 A) Number of tRNA genes harbouring (green) an AID* or sA3G* MEL compared with
998 expected number from Monte Carlo simulations.

999 B) Density of mutations in relation to the transcription start site (TSS) of tRNA genes.

1000 Mutations within the 500 base pair interval centred at the TSS are included.

1001 C) Mutation frequency in promoters of mRNA genes (within a window 500 bp upstream and
1002 50 bp downstream of the TSS) compared to the frequency of mutations in the promoters of
1003 tRNA (550 bp window centred on the middle of the tRNA gene), snoRNAs and snRNA genes
1004 (550 bp window as for mRNA genes). mRNA genes are binned according to transcription
1005 rate as in Figure 3. Both RNAP II and III driven snoRNAs are included.

1006 D) Example of MELs in ChrIV and ChrXV corresponding to tRNA tI(UAU)D and tG(CCC)O,
1007 depicted as in Figure 3.

1008

1009 **Figure 4 - figure supplement 1.**

1010 Median number of mutable motifs in promoter regions.

1011

1012 **Figure 4 - figure supplement 2.**

1013 Mutationally enriched loci are not a consequence of increased density of mutable motifs.

1014 The number of deaminase motifs for each MEL versus the number of mutations within each
1015 MEL for AID* and sA3G* datasets.

1016

1017 **Figure 4 - figure supplement 3.**

1018 Mutations in the rDNA locus are restricted to the replication fork block (RFB) site.

1019 A) Sequence context of low allelic frequency mutations detected in the rDNA locus, as
1020 depicted in Figure 1C.

1021 B) Schematic of the rDNA repeat region. Panels show mutations detected in yeast
1022 transformants at the rDNA locus. Each line represents one clone with dots representing
1023 mutations (mutation at C, black; at G, red; at A, green). Clones with no detected mutations
1024 are not depicted.

1025

1026 **Figure 4 - figure supplement 4.**

1027 Deaminase induced mutation distribution in relation to R-loop forming potential. Tables
1028 showing the correlation between R-loops formation predicted by the QmRLFS-finder
1029 (Wongsurawat et al., 2012) or SkewR package (Ginno et al., 2012) and the presence of MELs.

1030

1031 **Figure 5.**

1032 RNA binding by human AID and APOBEC3G.

1033 A) Left panel shows the in vitro transcribed pre-tl(UAU)D tRNA used for affinity purification.
1034 Right panel shows immunoblots for transiently overexpressed AID/APOBEC3G proteins
1035 following RNA-immunoprecipitation with pre-tRNA.

1036 B) Affinity purification with tl(UAU)D probe, total yeast tRNA, homopolymeric single
1037 stranded (polyU) and double stranded (polyA:U) RNA. Left panel shows input proteins, right
1038 panel shows immunoblots for transiently overexpressed AID/APOBEC3 proteins following
1039 RNA-immunoprecipitation. Results representative of at least 3 independent experiments.

1040 C) Deaminase induced mutations in the promoter region of the YBR194W locus. Top panels:
1041 accumulated mutations in the AID*, sA3G* and EMS whole genome datasets. Bottom
1042 panels: mutations detected in Sanger sequenced yeast clones unmodified or harbouring a
1043 chimeric YBR194W-snR6 locus. Each line represents one clone with dots representing
1044 mutations (at C, black; at G, red). Clones with no mutations are indicated.

1045

1046 **Figure 5 – figure supplement 1.**

1047 A) Mutation density relative to the TSS for APOBEC3A and APOBEC3B induced mutations
1048 from *ungΔ* haploid cells (data from (Taylor et al., 2013)). The density at tRNA promoters is
1049 shown separately in red.

1050 B) Mutations in breast cancer genome PD4120a and lung adenocarcinoma LUAD-S01345.
1051 Pie charts show the contribution of mutations at TC over mutations at the remaining
1052 dinucleotides and histograms show mutation density relative to all human TSS (Ensemble
1053 annotation).

1054

1055 **Figure 5 – figure supplement 2.**

1056 Functional comparison of the YBR194W locus in modified yeast clones.

1057 A) Immunoprecipitation of chromatin associated RNAP II or

1058 B) Histone H3 from unmodified or YBR194W-snR6 chimeric yeast.

1059 Black bars show enrichment relative to input in the unmodified strain with the modified
1060 strain in red. An unrelated locus, YJL105W is shown as control. Data from 3 independent
1061 experiments.

1062 C) mRNA levels of YBR194W shown relative to ACF1. Levels at the TAF10 gene are shown as
1063 a control. Data from 3 independent experiments.

1064

1065 **Tables**

1066 **Table 1. Deaminase induced Mutation Enriched Loci (MEL) in yeast genomes**

1067

	Observed			Simulated		
	AID*	sA3G*	EMS	AID*	sA3G*	EMS
MELs	1227	568	1	50	21	3
% MEL mutation	40.7	21.6	0.24	0.75	0.39	0.14

1068

1069 **Supplementary files**

1070 Supplementary file 1: Catalogue of yeast mutations

1071 Supplementary file 2: Coordinates of MELs.

1072 Supplementary file 3: All mutationally enriched regions (MELs). Top panel indicate position
1073 of each non-clonal mutation indicated by a dot (at C, black; at G, red), with horizontal lines
1074 representing a single genome. Middle panel shows MELs (AID*, green; sA3G*, purple; EMS,
1075 grey). Bottom panel displays genomic features (including transcripts, replication origins,
1076 centromers), coloured according to feature type, with arrows indicating the direction of
1077 transcription. The coordinates of the region are indicated. Regions are ranked according to
1078 the number of mutations present.

1079 Supplementary file 4: Scripts used for data analyses.

1080

1081

Figure 1

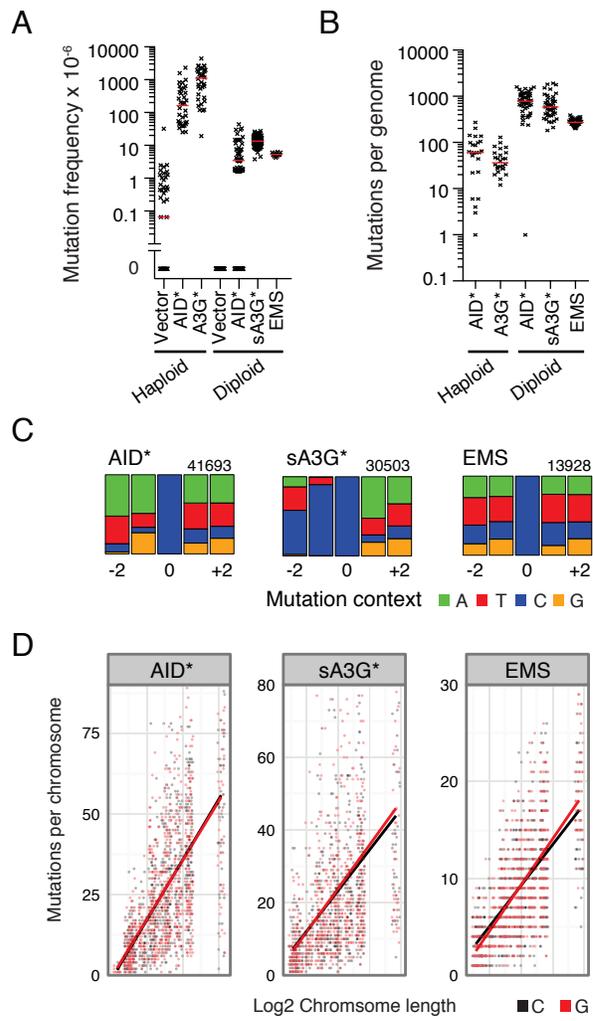
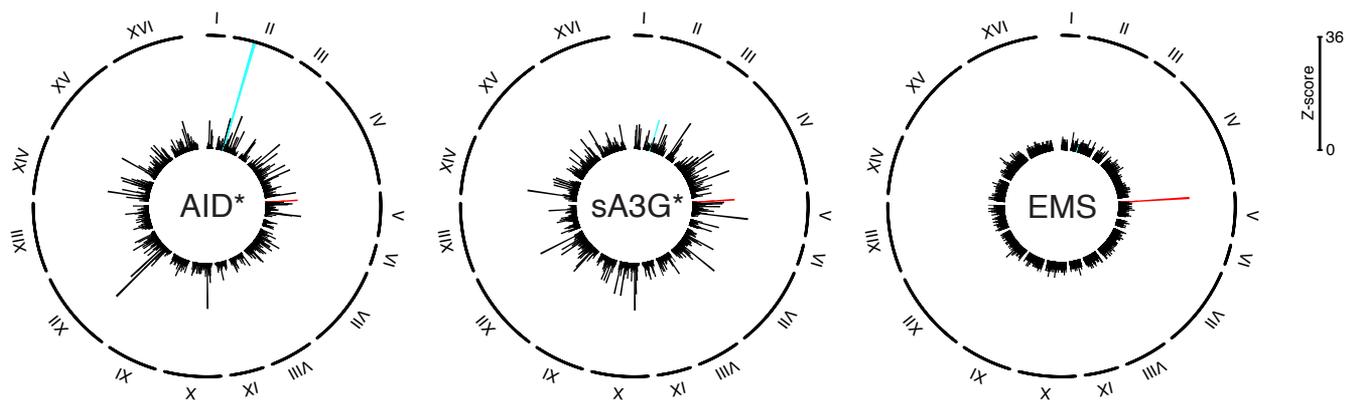
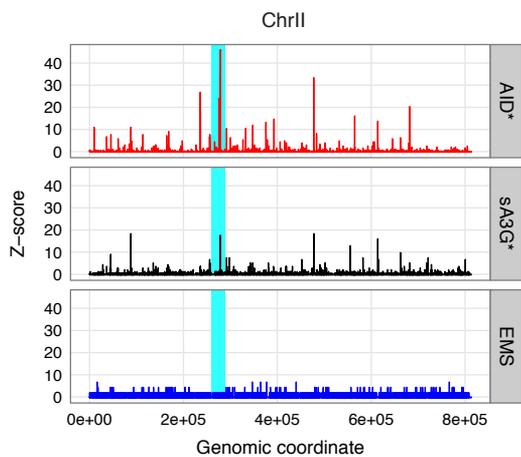


Figure 2

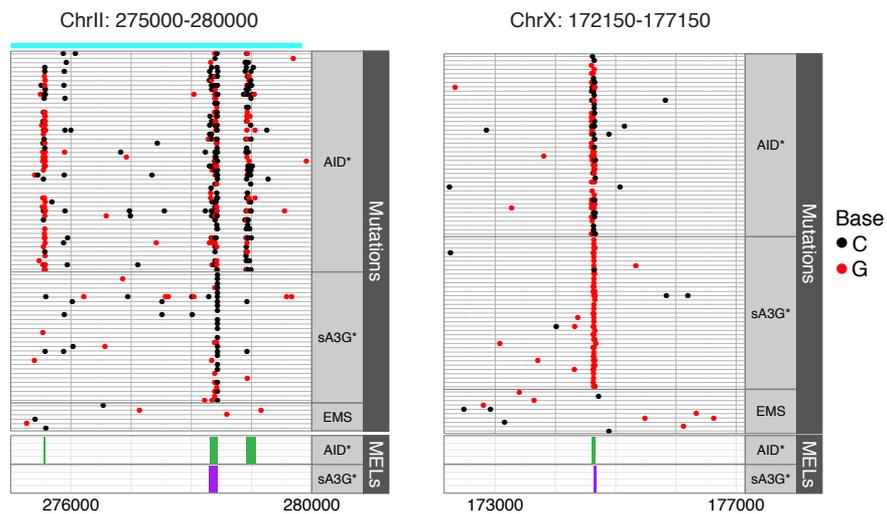
A



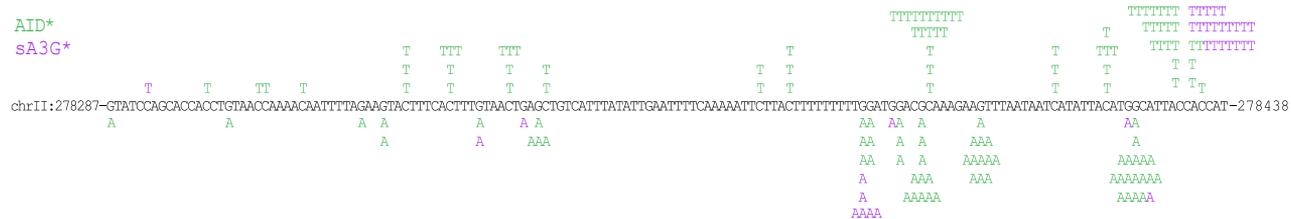
B



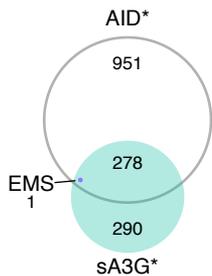
C



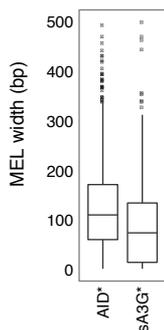
D



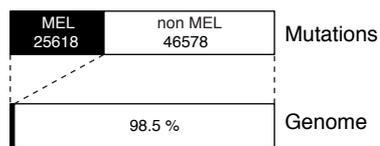
E



F



G



H

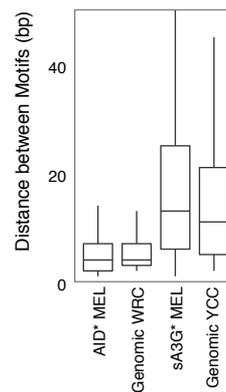
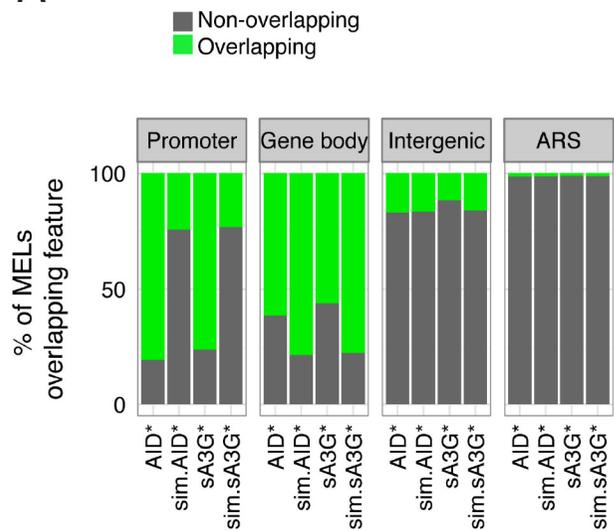
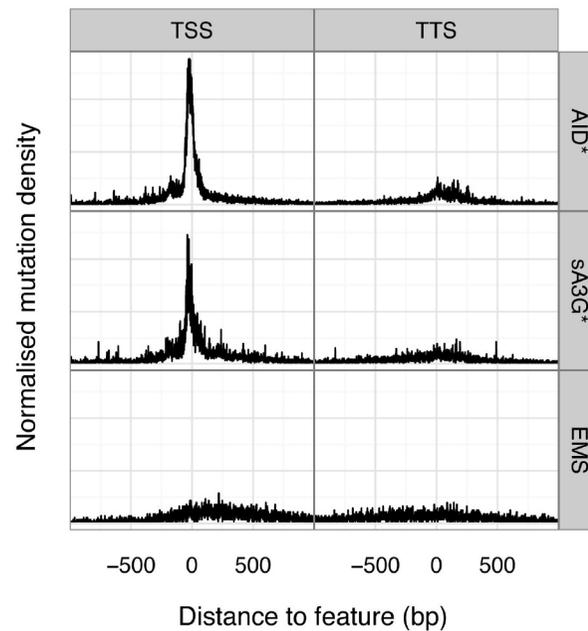


Figure 3

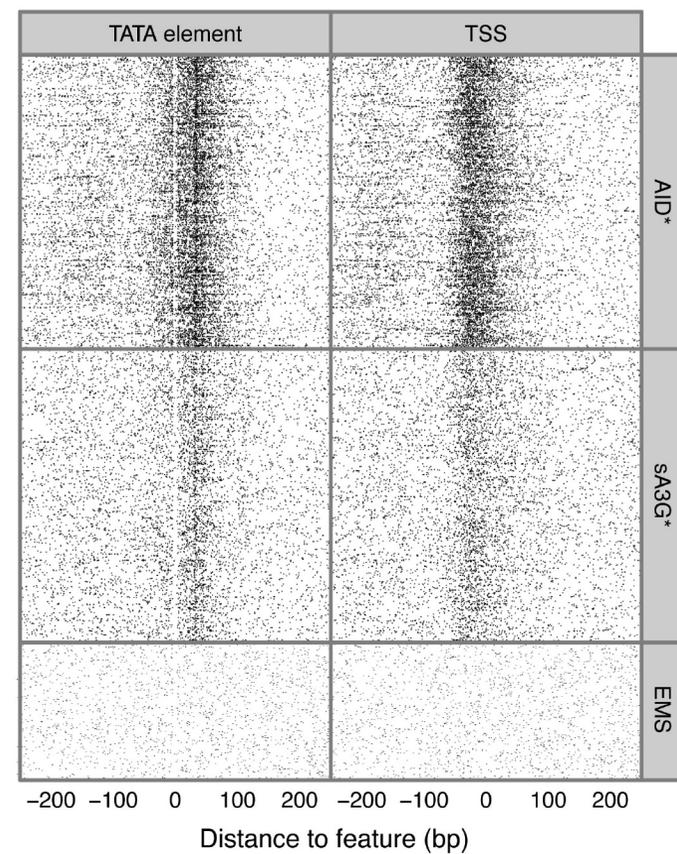
A



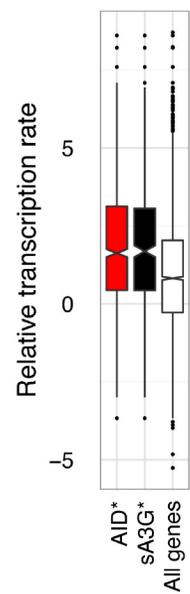
B



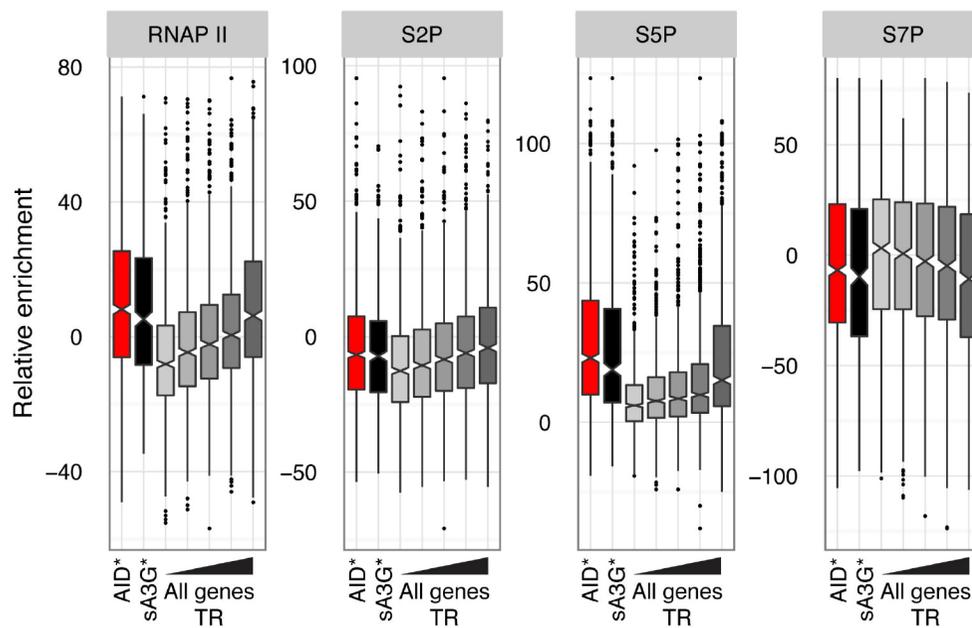
C



E



F



D

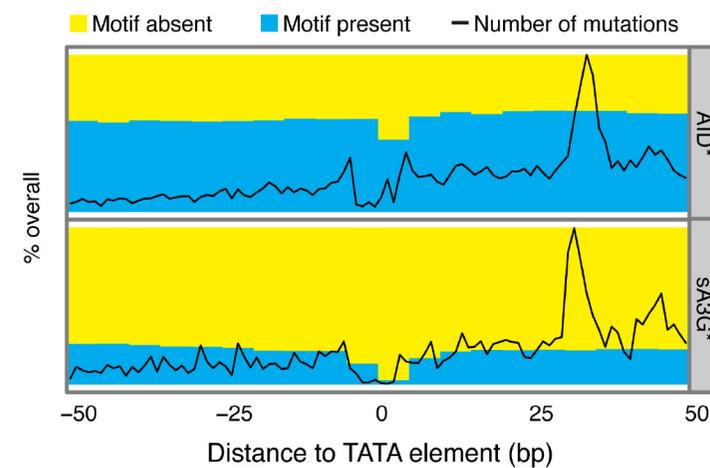


Figure 4

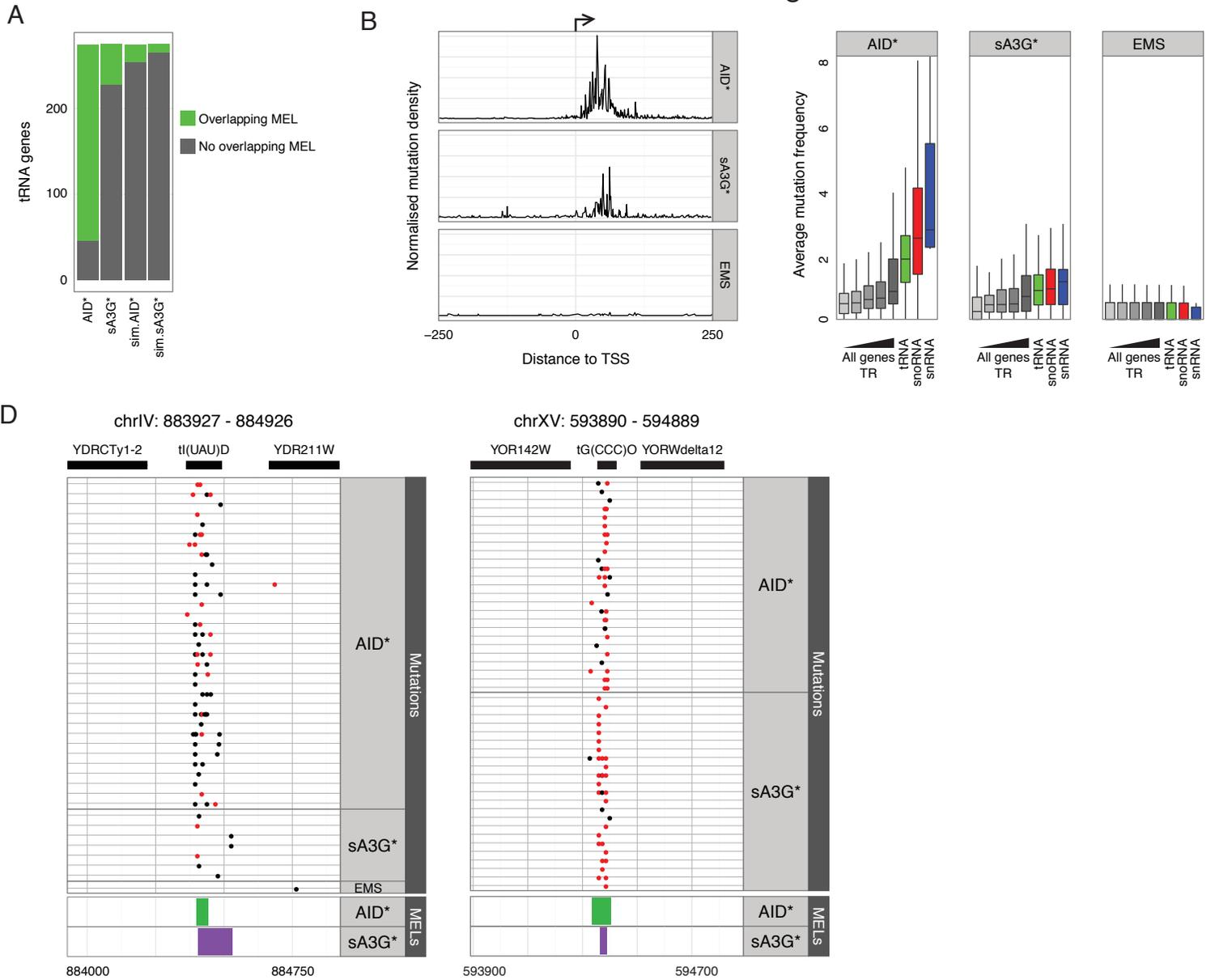
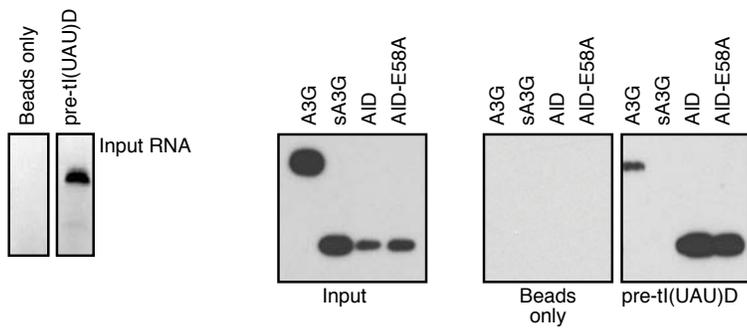
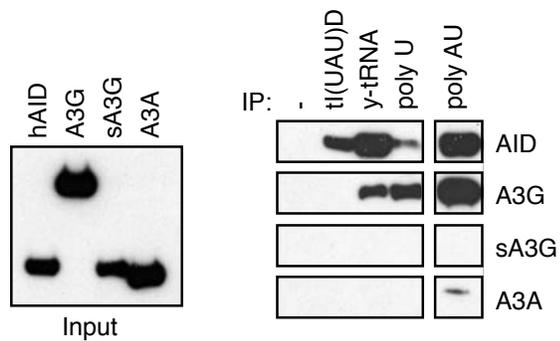


FIGURE 5

A



B



C

