

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

Somatic hypermutation spectra are independent of the local transcriptional and epigenetic landscape

Ursula E. Schoeberl^{1*}, Johanna Fitz^{1*}, Kimon Froussios¹, Renan Valieris², Marina Makharova¹, Iordanis Ourailidis¹, Bernd Bauer¹, Tobias Neumann¹, Eva-Maria Wiedemann¹, Monika Steininger¹, Adriana Cantoran Garcia¹, Marialaura Mastrovito¹, Hugo Mouquet³, Israel Tojal Da Silva² and Rushad Pavri^{1#}

(1) Research Institute of Molecular Pathology (IMP), Campus-Vienna-Biocenter 1, 1030 Vienna, Austria

(2) Laboratory of Bioinformatics and Computational Biology, A. C. Camargo Cancer Center, São Paulo, SP, Brazil

(3) Laboratory of Humoral Immunology, Institut Pasteur, INSERM U1222, Université Paris Cité, 75015 Paris, France

* equal contribution

corresponding author (rushad.pavri@imp.ac.at)

23 **Abstract**

24 Somatic hypermutation (SHM) of immunoglobulin variable regions in B cells modulates antibody-antigen
25 affinity and is indispensable for adaptive immunity. Mutations are introduced by activation-induced
26 cytidine deaminase (AID) in a co-transcriptional manner resulting in discrete mutation spectra. Current
27 models propose that activating epigenetic marks, transcriptional pausing and convergent transcription are
28 necessary for optimal AID recruitment. However, whether these or other transcriptional features can
29 explain the discrete mutation spectra is unknown. To address this, we compared mutation and nascent
30 transcription at single nucleotide resolution. Surprisingly, with this precision, SHM spectra do not correlate
31 with any transcriptional feature at human and mouse variable regions and non-immunoglobulin AID
32 targets. Moreover, SHM is resistant to up to four-fold reduction of both activating epigenetic marks and
33 transcription. We propose that, following AID recruitment to its target genes, the DNA sequence flanking
34 an AID target motif is the key determinant of mutability rather than the local transcriptional and chromatin
35 landscape.

36

37

38

39

40

41

42

43

44

45

46 Introduction

47 Somatic hypermutation (SHM) is the molecular basis for the diversification of antibodies in response to
48 pathogens and vaccines and is hence indispensable for robust long-term immunity (1-3). SHM occurs in
49 activated B cells within microanatomical structures called germinal centers in secondary lymphoid tissue
50 upon engagement with antigens and helper T cells (4). Mutations are generated in the variable regions of
51 the immunoglobulin (IG) heavy (*IGH*) and light chain (*IGK*, *IGL*) genes by activation-induced cytidine
52 deaminase (AID) (5, 6) which acts co-transcriptionally on single-stranded DNA (ssDNA) (7-11). AID also
53 acts at the *IGH* switch regions triggering a reaction cascade leading to class switch recombination (CSR)
54 which yields the various antibody isotypes essential for effector functions following antigen binding (12-
55 15). In addition, SHM occurs, albeit with much lower frequency, at several transcribed non-IG genes (16-
56 20).

57 An important and longstanding question in the field is the relationship of SHM with nascent
58 transcription and activating epigenetic marks. AID associates with the RNA polymerase II (Pol II)
59 machinery via interaction with various elongation (21-25) and RNA processing factors (26, 27). However,
60 since these are generic gene regulatory proteins, such interactions cannot per se explain the specificity of
61 AID targeting. SHM has also been correlated with the epigenetic landscape especially with activating
62 histone marks. In particular, depletion of histone modifying enzymes resulted in decreased SHM and/or
63 CSR accompanied with a decrease in activating histone modifications (24, 28-36). However, the precise
64 role of these marks in SHM is unknown. Moreover, these marks are typical of any transcribed gene hence
65 their presence per se cannot confer specificity. Instead, it appears that specificity is conferred by super-
66 enhancers associated with IG and other AID target loci (37-40). The murine *Igh* super-enhancer is
67 essential for SHM (41) and is itself a target of AID activity (42). Importantly, an AID-targeted enhancer
68 when placed in a genomic domain that lacks AID target genes was able to convert these genes into SHM
69 hotspots, thus providing compelling evidence for the critical role of enhancers in targeting AID (38).

70 The final step of SHM, that is, the actual mutagenesis, is least understood. AID preferentially
71 deaminates cytidine residues within WRCH motifs (where W = A or T, R = A or G and H = A, C or T) (43-
72 45). However, not all WRCH residues are targeted and the mutation frequencies of those that are

73 mutated can differ substantially (46-52). Remarkably, even identical WRCH motifs within a given variable
74 region can differ significantly in mutation frequency (47-52). This differential mutability is a ubiquitous
75 feature of SHM and implies the existence of as-yet-unknown mechanisms that determine which residue is
76 mutated and to what extent. More importantly, differential mutability cannot be explained by the
77 association of AID with Pol II complexes, but rather, suggests that access to ssDNA may be biased
78 towards certain sequence contexts. In this regard, features of nascent transcription that can cause
79 transient ssDNA exposure, such as convergent transcription (53), transcriptional pausing (17, 22, 54-57),
80 initiation (58) and transcription termination (59, 60), have been correlated with SHM and CSR. Moreover,
81 ssDNA patches are detected in variable regions and suggested to arise via negative supercoiling
82 upstream of transcribing Pol II, although these patches do not correlate with SHM patterns (56, 57, 60,
83 61). Indeed, high occupancy of Pol II and SPT5, a Pol II pausing factor implicated in AID targeting (22,
84 56, 57), was able to identify new AID target genes via a machine learning approach (17). Based on these
85 observations, a reasonable hypothesis is that differential mutability of WRCH motifs arises because of
86 local sites of transcription initiation, pausing or convergent transcription which favor AID activity at these
87 sub-regions within the variable region. Since variable regions are very short (<500 bp), answering this
88 question requires high-resolution nascent transcriptional maps of variable regions, which are currently
89 unavailable.

90 In this study, we aimed to address the precise relationship between SHM, nascent transcriptional
91 features and epigenetic marks using a fresh and direct approach. To do so, we performed precision run-
92 on sequencing (PRO-seq and PRO-cap) (62) in Ramos human B cells and murine germinal center B cells
93 to obtain single-nucleotide resolution maps of nascent transcription at multiple human and mouse variable
94 regions as well as at a large collection of non-IG AID target genes (17). We combined this with high-
95 resolution mutational profiling to show that patterns of SHM do not consistently correlate with any nascent
96 transcriptional features. Moreover, via deletion of the *IGH* E μ enhancer in human B cells, we find that
97 SHM is insignificantly affected despite up to four-fold reduction in transcription and active chromatin
98 marks. Importantly, the same variable region sequence, studied in mouse or human regulatory contexts,
99 exhibits a nearly identical SHM pattern. The results strongly suggest that following AID recruitment to the

100 variable gene locale, the SHM spectrum is determined by the DNA sequence flanking the hotspot
101 residues rather than any co-transcriptional or epigenetic feature.

102

103 **Results**

104 **The human and mouse *IGH* variable region locale harbors relatively low levels of histone** 105 **acetylation and H3K4 trimethylation**

106 As a model system of SHM, we chose the Ramos human B cells. Ramos is an IgM-positive Burkitt
107 lymphoma-derived cell line where the single functional, pre-recombined *IGH* variable region consists of
108 the variable gene, VH4-34, the diversity segment, DH3-10, and the joining segment, JH6. Ramos cells
109 constitutively express AID and undergo very low levels of SHM mostly at C:G residues (63).

110 To determine the chromatin landscape of the Ramos *IGH* locus, we performed chromatin
111 immunoprecipitation followed by quantitative PCR (ChIP-qPCR) for histone modifications associated with
112 active promoters and enhancers (histone H3 acetylated on lysine 27, H3K27ac or trimethylated on lysine
113 3, H3K4me3) and transcription elongation in gene bodies (H3 trimethylated on lysine 36, H3K36me3). At
114 active genes, H3K27ac and H3K4me3 are enriched at nucleosomes flanking the transcription start site.
115 Surprisingly, we observed that the variable region is characterized by relatively low levels of all marks
116 that, importantly, was not due to decreased nucleosome occupancy as judged by the levels of histone H3
117 (Fig. S1A). As a result, the first nucleosome from the transcription start site has the lowest levels of
118 H3K27ac and H3K4me3 which is contrary to what is normally observed at highly active genes where the
119 first nucleosome harbors the highest levels of these marks. In contrast, much higher levels of all marks
120 were observed in the intronic regions flanking the E μ enhancer with an expected decrease at E μ itself
121 where nucleosomes are occluded by the presence of transcription factors (Fig. S1A) (64). Importantly, the
122 same profiles were observed at the murine *Igh* in activated, primary B cells expressing the B1-8^{hi} variable
123 region (65) (Fig. S1B). We conclude that variable regions are marked by relatively low levels of activating
124 histone marks and that this landscape is conserved between the human Ramos and mouse B1-8^{hi}
125 variable regions.

126 **The E μ enhancer regulates the levels of chromatin marks in the flanking intronic DNA but not**
127 **within the variable region**

128 The fact that the highest enrichments of chromatin marks were seen on either side of the E μ enhancer
129 suggested that this element may be responsible for the deposition of these marks. This raised the
130 question of whether the loss of E μ would alter the levels of these marks and affect SHM. To minimize
131 effects of ongoing SHM on the variable gene sequence, we first ablated AID in Ramos cells (AID^{-/-}) and
132 subsequently deleted E μ in AID^{-/-} cells (AID^{-/-} E μ ^{-/-}) (Fig. S1C). Specifically, by CRISPR-mediated
133 editing, we deleted a 583 bp region corresponding to the peak of chromatin accessibility measured by
134 assay for transposase-accessible chromatin (ATAC-seq) signal at E μ in AID^{-/-} Ramos cells (AID^{-/-} E μ ^{-/-})
135 (Fig. S1C). Based on PRO-cap analysis of transcription start sites (see below), this deletion eliminates the
136 I μ promoter that encodes the sterile germline *IGHM* transcript and antisense transcription start sites
137 residing in E μ (Fig. S1C). We chose two independent clones, AID^{-/-} E μ ^{-/-} c1 and AID^{-/-} E μ ^{-/-} c2, for
138 subsequent experiments (Fig. S1D).

139 To comprehensively visualize profiles of these histone marks across the variable region and
140 flanking sequences, we performed ChIP followed by deep sequencing (ChIP-seq) for H3K27ac,
141 H3K4me3 and H3K36me3. We generated a custom chromosome spanning the Ramos variable region
142 and promoter. Importantly, the Ramos variable region is highly mappable and hence there is minimal loss
143 of reads due to multimapping with other variable gene families (see later section). The results in AID^{-/-}
144 Ramos cells were consistent with ChIP-qPCR in that the levels of all marks were lowest in the promoter-
145 proximal region and gradually increased into the variable region, peaking in the intronic sequences
146 flanking E μ (Fig. 1A-B). This analysis also revealed the absence of bimodal distributions of H3K27ac and
147 H3K4me3 at the VH4-34 promoter, which is commonly observed at highly active divergently transcribed
148 genes. In AID^{-/-} E μ ^{-/-} cells, all three marks were decreased (Fig. 1A-B), a finding that we confirmed by
149 ChIP-qPCR and which was not due to changes in nucleosome occupancy as determined by histone H3
150 measurements (Fig. 1C). Using antibodies against pan-H3 and pan-H4 acetylation, we observed a
151 significant decrease in histone acetylation across the locus in AID^{-/-} E μ ^{-/-} cells (Fig. 1C). Moreover,

152 another mark of transcription elongation, H3K79me3, was also reduced in AID^{-/-} E μ ^{-/-} cells (Fig. 1C). All
153 marks in AID^{-/-} E μ ^{-/-} cells were significantly reduced up to 4-fold in the intron where they are normally at
154 their peak levels (Fig. 1A-C). However, within the variable region, where these marks are normally at their
155 lowest levels, there was no significant change in AID^{-/-} E μ ^{-/-} cells (Fig. 1A-C). Consequently,
156 AID^{-/-} E μ ^{-/-} cells harbor a chromatin landscape wherein the levels of all marks appear to be comparable
157 between the variable region and the intron (Fig. 1C, compare the measurements from amplicons 2, 3 and
158 5 in AID^{-/-} and AID^{-/-} E μ ^{-/-} cells). We conclude that E μ is responsible for the higher levels of chromatin
159 marks within the intron but has no significant influence on the deposition of these marks within the
160 variable region.

161 These changes in the chromatin landscape in AID^{-/-} E μ ^{-/-} cells were accompanied by a 2-4-fold
162 decrease of nascent transcription in the variable region and intron, as well as a 2-fold decrease of the
163 spliced IgM mRNA (Fig. 2A-B). The decrease in transcription within the variable region and intron implies
164 that E μ regulates optimal transcription initiation from the variable gene promoter. Nevertheless, despite
165 these decreases in transcription, surface IgM levels were unaffected in E μ ^{-/-} cells suggesting a lack of
166 absolute correlation between IgM mRNA and protein levels (Fig. S1E). We conclude that E μ positively
167 and significantly regulates nascent transcription in the variable region but has no significant impact on its
168 epigenetic landscape.

169 **SHM frequency does not correlate with the levels of nascent transcription or activating histone**
170 **marks in the variable region locale**

171 AID^{-/-} E μ ^{-/-} cells provide an ideal system to ask whether changes in the levels of nascent transcription and
172 chromatin marks directly impact on SHM without any contribution from secondary, indirect effects. Hence,
173 we sequenced the variable region as well as the JH6 intron immediately downstream (Fig. 2C, upper
174 panel). This analysis allowed us to compare SHM between a region where chromatin marks were
175 normally low and E μ -independent (variable region) and a region where the marks were normally high and
176 strongly E μ -dependent (JH6 intron). To measure SHM, we infected AID^{-/-} and AID^{-/-} E μ ^{-/-} cells with
177 JP8Bdel, a C-terminal truncation mutant of AID that results in nuclear retention and major increase in

178 mutation rates (66). As a control, we also infected cells with a catalytically inactive AID mutant (E58Q)
179 (66). In both the variable region and JH6 intron, we observed a small (~20%) decrease in SHM frequency
180 in AID^{-/-} E μ ^{-/-} relative to AID^{-/-} JP8Bdel-expressing cells that was not statistically significant (Fig. 2C). This
181 result is consistent with mouse studies where loss of E μ caused only modest decreases in SHM (67, 68).
182 As an alternative readout of SHM activity, we assayed for the loss of surface IgM, which occurs due to
183 AID-induced stop codons or frameshifts (63). The results showed that the magnitude of IgM loss was
184 similar between the control and E μ ^{-/-} cells suggesting that the minor decreases in SHM do not translate
185 into equivalent changes in IgM loss (Fig. 2D-E). We conclude that the E μ enhancer is dispensable for
186 robust SHM in Ramos cells.

187 By comparing the levels of chromatin marks with SHM frequencies in the variable region and JH6
188 intron, we can draw three major conclusions: First, we infer that the levels of SHM do not correlate with
189 those of activating histone marks. Indeed, SHM frequency is slightly lower in the intron (1.6×10^{-3}), which
190 harbors the highest level of these marks, than in the variable region (1.86×10^{-3}), which has the lowest
191 levels of these modifications (Fig. 2C). Second, as exemplified by the JH6 intron, SHM is resistant to
192 major decreases in these marks. Third, we attribute the small reduction in SHM in AID^{-/-} E μ ^{-/-} cells to the
193 decrease in nascent transcription although, here too, the decrease of SHM (~20%; Fig. 2C) is
194 considerably smaller than that of transcription which shows up to 70% loss in the JH6 intron of
195 AID^{-/-} E μ ^{-/-} cells (Fig. 2A, amplicons 2-5). We conclude that although our results do not rule out a role for
196 these activating marks in SHM, they do clearly demonstrate that their presence, profiles and even their
197 significant reduction is not a predictor of SHM.

198 **High-resolution transcriptomic profiling of different human variable regions reveals the absence** 199 **of direct relationships between nascent transcriptional features and SHM spectra**

200 A major, unsolved problem in SHM biology is the differential mutability of WRCH motifs within variable
201 regions which implies the existence of local, sequence-intrinsic mechanisms regulating AID activity post-
202 recruitment (46-48, 50-52). We asked whether SHM spectra are dictated by local features of nascent

203 transcription landscape that could transiently lead to increased ssDNA exposure, such as during
204 transcription initiation, pausing and/or convergent/antisense transcription.

205 Since V genes occur as families with varying degrees of sequence identity, mapping of next-
206 generation sequencing (NGS) reads to recombined variable regions can be problematic because
207 standard alignment workflows eliminate reads that map to more than one location in the genome
208 (multimappers) and only align reads mapping uniquely (typically allowing 2-3 mismatches). In Ramos
209 cells, visual inspection of PRO-seq tracks on the UCSC genome browser revealed robust coverage within
210 the VH4-34 gene with uniquely-mapping reads indicating that this V gene is highly mappable (not shown).
211 Moreover, there were no unique reads mapping to other V gene families suggesting that they are
212 transcriptionally silent in these cells. To address this more systematically, we retrieved the normally
213 discarded, multimapping reads and allowed them to re-align to the recombined Ramos variable region
214 with the important proviso that they only multimap to annotated V, D or J segments and not elsewhere in
215 the genome (Fig. 3A). Separate tracks were created for uniquely mapping reads and the recalled
216 multimapping reads, which were then combined to obtain the total profile (Fig. 3B).

217 We measured nascent transcriptional activity with PRO-seq, which maps the location and
218 orientation of actively engaged RNA polymerase II (Pol II) at single-nucleotide resolution, and the related
219 method, PRO-cap, which maps transcription initiation sites of capped, nascent RNAs (62). Given that
220 variable regions are short (<500 bp), the single-nucleotide resolution of PRO-seq is ideally suited to
221 identify potential sites of pausing, as well as antisense and convergent transcription within variable
222 regions. For simplicity, we use the term “pausing” to refer to accumulation of PRO-seq signal in variable
223 regions although it must be noted this refers to the slower-moving, elongating Pol II molecules and is not
224 equivalent to promoter-proximal pausing which is also reported by PRO-seq. Most notably, promoter-
225 proximal pausing occurs via the binding of NELF which inhibits elongation whereas elongating Pol II lacks
226 NELF and contains additional proteins such as SPT6 and the PAF complex (69-72). Accumulation of Pol
227 II signal in gene bodies may arise due to the slowing down of Pol II as a result of steric barriers such as
228 secondary structures or premature termination signals. However, initiation events within gene bodies

229 could give rise to a paused Pol II state akin to promoter-proximal pausing and can be identified if PRO-
230 seq accumulation is associated with upstream PRO-cap enrichments.

231 Our mapping pipeline revealed that the Ramos variable region is predominantly covered by
232 uniquely-mapping reads with very few multimapping reads (Fig. 3B). Of note, the same pipeline was used
233 for the ChIP-seq analysis of Ramos cells described in Fig. 1 above. Therefore, this approach allowed us,
234 for the first time, to delineate the transcriptional and epigenetic landscapes of IG variable regions. In
235 contrast to the lack of a bimodal peak of activating histone marks (Fig. 1A), the Ramos VH4-34 promoter
236 showed clear bidirectional transcription and promoter-proximal pausing, the latter being noticeable as the
237 PRO-seq enrichments shortly downstream of the promoter-associated PRO-cap initiation signals (Fig.
238 3C). Increased PRO-seq read density was observed towards the 3' end of the variable region suggestive
239 of slowly-elongating or paused Pol II. However, the rest of the variable region harbored relatively low
240 PRO-seq signals (Fig. 3B and magnified view in Fig. 3C). Next, we infected cells with AID(JP8Bdel) and
241 analyzed SHM patterns in the variable region using an optimized version of mutation analysis with paired-
242 end deep sequencing (MutPE-seq) of a PCR-amplified fragment corresponding to the variable region (73)
243 (Fig. 3E). Mutations occurred robustly at C:G residues within WRCH motifs with weak targeting of A:T
244 residues, as expected from Ramos cells. Moreover, the cold spot motif, SYC (where S = C or G and Y =
245 C or T) was poorly mutated (Fig. 3F-H). However, there was no discernible overlap of the mutation profile
246 with that of nascent transcription. In particular, the region of highest PRO-seq read density at the 3' end of
247 the variable region is not preferentially mutated over upstream C:G motifs where PRO-seq signal is much
248 weaker (Fig. 3D-E). Indeed, the most mutated residue is located near the 5' end of the variable region
249 where nascent transcription signals are lower (Fig. 3D). We conclude that there appears to be no
250 discernible feature of the nascent transcriptional landscape that can account for the observed frequency
251 or distribution of mutations.

252 PRO-seq in AID^{-/-} E μ ^{-/-} cells revealed that nascent sense transcription was uniformly decreased
253 across the variable region and intron (Fig. S1F-G), consistent with the RT-qPCR results from these cells
254 (Fig. 2A). Of note, PRO-seq showed that the E μ enhancer harbors the start sites of an antisense
255 transcript, substantially weaker than the sense transcript (note the Y axis scale). This transcript initiates at

256 the 5' boundary of E μ and terminates within the variable region, resulting in convergent transcription
257 within the variable region (Fig. 3C-D). Moreover, the decrease in antisense transcription in AID^{-/-} E μ ^{-/-}
258 cells is stronger than sense transcription resulting in a major decrease in convergent transcription in the
259 variable region and intron (Fig. S1F-G). Given the weak SHM phenotype in AID^{-/-} E μ ^{-/-} cells (Fig. 2C), we
260 conclude that SHM is resistant to major decreases in Pol II levels within the variable region and intron, in
261 both sense and antisense orientations.

262 To determine the transcriptomic profiles of other human variable regions, we developed a
263 strategy to replace the endogenous *IGH* variable region of Ramos cells with any variable region of choice
264 using CRISPR editing (Fig. S2A-B and Methods). We chose two variable regions corresponding to the
265 unmutated precursors of broadly neutralizing antibodies identified in individuals infected with human
266 immunodeficiency virus 1 (HIV1). One variable region consisted of VH4-59*01, DH2-02 and JH6*03
267 (termed VH4-59) (74) and the other of VH3-30*18, DH2-02 and JH6*02 (75). In both cases, the
268 sequences upstream and downstream of the new variable regions, including the VH4-34 promoter, leader
269 peptide and intron are identical to the endogenous Ramos *IGH* sequence. Our alignment pipeline
270 revealed that sub-regions of both VH3-30 and VH4-59 harbored multimapping reads with the most
271 prominent being the framework region 3 preceding CDR3 in VH4-59 which was entirely covered with
272 multimapping reads (Fig. S2C-D).

273 We analyzed the transcriptional and mutational landscape of these two variable regions using
274 PRO-seq, PRO-cap and MutPE-seq, the latter following infection with AID(JP8Bdel) (Fig. 4). At VH4-59,
275 we observed an increase in PRO-seq enrichment near the end of the variable region (Fig. 4A and
276 magnified in Fig. 4B), reminiscent of Pol II pausing and akin to that seen at the endogenous Ramos
277 variable region (Fig. 3B-C). These enrichments were considerably weaker at the VH3-30 variable region
278 (Fig. 4D-E). The most plausible explanation may lie in the fact that VH4-59 and VH4-34 belong to the
279 same variable gene family (VH4) and hence share higher sequence identity than VH3-30 which is
280 phylogenetically more distant from the VH4 family (76). This also suggests that sequence-intrinsic
281 features may be responsible for the observed differences in nascent transcription profiles. Most
282 importantly, and in agreement with the results from the endogenous Ramos variable region, the highest

283 mutation frequencies in either VH4-59 (Fig. 4B) or VH3-30 (Fig. 4E) were often not associated with
284 regions of increased PRO-seq signal.

285 Of note, when comparing the mutation patterns between all three variable regions analyzed here,
286 the strong bias for C:G mutations typical of Ramos cells is seen in all cases (Fig. 3E and Fig. 4D, F).
287 Moreover, there appears to be some bias of mutation load in the CDRs although several major mutation
288 hotspots are also located within the intervening framework regions (Fig. 3E and Fig. 4D, F). The coldspot
289 SYC motifs are weakly targeted or untargeted by AID, as expected (Fig. 3E and Fig. 4D, F). Finally, the
290 differential mutability of WRCH motifs by AID, ranging from highly mutated to unmutated, is seen in all the
291 analyzed variable regions (Fig. 3E and Fig. 4D, F) and is in line with observations in previous studies (47,
292 48).

293 Taken together, we conclude that although pausing, antisense transcription and convergent
294 transcription are all observed to varying degrees in different variable regions, none of these features, or
295 any combination thereof, correlates with the spectra of SHM.

296 **SHM in germinal center B cells at a murine variable region shows no direct relationship with** 297 **nascent transcriptional features**

298 In mice, SHM occurs robustly *in vivo* in germinal center B cells (GCBs) but is inefficient *in vitro* in
299 activated primary B cells (48). Thus, we asked if such differences could be related to differences in the
300 underlying transcriptional landscape. To address this, we made use of the B1-8^{hi} *Igh* knock-in mouse
301 where the murine B1-8^{hi} variable region is knocked-in at the endogenous *Igh* locus (77). In this system,
302 one can compare SHM and transcriptional features of the B1-8^{hi} variable region from both germinal center
303 B cells (GCBs) as well as primary *in vitro* activated B cells. To boost SHM in primary cells, we generated
304 homozygous B1-8^{hi} *Rosa26*^{AIDER} mice where AID fused to the estrogen receptor (AIDER) is expressed
305 constitutively from the *Rosa26* promoter (*Rosa26*^{AIDER}) such that upon addition of 4-hydroxytamoxifen,
306 AIDER is translocated into the nucleus (73).

307 PRO-seq analysis from homozygous B1-8^{hi} mice showed that the Vh1-72 gene used in the B1-8^{hi}
308 recombined variable region shares strong homology with other V genes in the murine *Igh* locus resulting

309 in a very high degree of multimapping (Fig. S2E). Importantly, the CDR3 region was covered with
310 uniquely mapping reads since this region is formed via the junction of the V, D and J segments during
311 VDJ recombination resulting in a unique sequence in the genome (Fig. S2E). PRO-seq and PRO-cap
312 analysis showed that there were no striking differences at the B1-8^{hi} variable region between primary B
313 cells and GCBs (Fig. 5A). Moreover, we observed local zones of Pol II accumulation, internal initiation
314 events and antisense/convergent transcription, but these did not necessarily correlate with zones of
315 mutation (Fig. 5A). As noted previously (48), SHM profiles of B1-8^{hi} GCBs are not identical to primary B
316 cells because GCBs exhibit 5-10 fold higher frequencies of SHM and substantial mutagenesis of A and T
317 residues arising from error-prone DNA repair of AID-induced mismatches at neighboring C:G residues
318 (Fig. 5B-C). The palindromic AGCT hotspot in CDR3 is selectively mutated with high frequency in primary
319 cells, in line with a former report showing that this hotspot is the earliest targeted motif both *in vitro* and *in*
320 *vivo* (48). Importantly, however, closer inspection showed that even in primary B cells, the same WRCH
321 motifs appear to be mutated as in GCBs suggesting that AID targeting specificity is not substantially
322 different between primary B cells and GCBs (Fig. 5B-C), in line with their comparable transcriptional
323 landscapes (Fig. 5A). We suggest that GCBs may have acquired higher rates of mutation, in part, due to
324 many more rounds of cell division within germinal centers compared to primary cultures, although the
325 AGCT hotspot in CDR3 remains enigmatic in this regard. We conclude that the differences in mutation
326 frequency between primary B cells and GCBs cannot be explained by underlying nascent transcriptional
327 features. Consequently, as in the case of human variable regions described above, the nascent
328 transcriptional landscape of the B1-8^{hi} variable region is not predictive of the patterns or frequencies of
329 SHM.

330 **The absence of correlation between SHM and nascent transcription is a global feature of SHM**

331 To extend our analysis to non-IG AID targets, we made use of a previously available SHM dataset from
332 murine GCBs where mutation frequencies at 275 AID target genes were identified by deep sequencing
333 the first 500 bp from the annotated transcription start site (TSS) (17). In this study, mutation analysis was
334 performed in GCBs from mice deficient in base excision and mismatch repair pathways (Ung^{-/-}Msh2^{-/-}).
335 In these mice, the processing of AID-induced U:G mismatches is abolished and, consequently, DNA

336 replication over these lesions leads to C→T and G→A transition mutations that represent the direct
337 targets of AID (78). By comparing these data with PRO-seq and PRO-cap obtained from GCBs, we could
338 compare mutation profiles with those of nascent transcription. Importantly, at many genes, initiation sites
339 or zones defined by PRO-cap are often located at considerable distances from the annotated TSS
340 (defined by the RefSeq database and labeled as TSS in Fig. 5D, Fig. S3 and Fig. S4). Consequently, the
341 500 bp segment sequenced for mutational analysis (17) often begins further upstream or downstream of
342 the initiation sites newly defined by PRO-cap.

343 The large number of genes in this analysis led to new and unexpected observations regarding the
344 nature of AID targeting. At virtually all genes, the differential mutation rates of WRCH motifs are evident
345 with many such motifs being unmutated and many non-WRCH cytidines being as efficiently targeted (Fig.
346 5D and several additional examples in Figs. S3 and S4). As we observed in variable regions, mutation
347 frequencies at individual nucleotides and mutation zones did not necessarily correlate with enrichments of
348 PRO-seq or PRO-cap (Fig. 5D, S3 and S4). For example, at *Bcl6*, the major zone of mutation lies in a
349 region of low PRO-seq signal where many mutations are not in WRCH motifs (Fig. 5D). At *Myc*, mutations
350 are observed upstream of the initiation site defined by PRO-cap indicating that mutagenesis can occur in
351 the antisense orientation (Fig. 5E).

352 In variable and switch regions, mutations initiate ~100-150 bp after the promoter leading to the
353 notion that AID is somehow excluded from initiating or early elongating Pol II (3, 43, 45, 79, 80). Contrary
354 to this, we identified a set of genes where the highest mutation frequencies were observed within 100-150
355 bp of the initiation site defined by PRO-cap (*Irf4* in Fig. 5F and additional examples in Fig. S3A). Indeed,
356 in some of these genes, the most mutated residues were in very close proximity to the initiation site, such
357 as *Mcm7* (Fig. 5G). In some other genes, strongly mutated residues coincided with the strongest peak of
358 initiation (Fig.S3B). These observations indicate that AID can act at the early stages of the transcription
359 cycle at least at some genes, suggesting that AID is not *per se* excluded from initiating or early elongating
360 Pol II.

361 Altogether, we conclude that the lack of predictability of mutation spectra from nascent
362 transcription patterns is observed at both IG variable regions and non-IG AID target genes, indicating that
363 the absence of correlation between SHM and nascent transcriptional patterns is a global feature of SHM.

364 **SHM targeting of the murine B1-8^{hi} variable region is retained in the context of the human *IGH***
365 **locus**

366 To determine whether the gene regulatory context affects the patterns of SHM and transcription, we
367 replaced the endogenous variable region in Ramos cells with the entire murine B1-8^{hi} variable region
368 (Ramos^{B1-8^{hi}}) using the strategy described earlier (Fig. 6A and Fig. S2A-B). In this scenario, transcription
369 and SHM of the murine B1-8^{hi} variable region would be under the control of the Ramos VH4-34 promoter
370 and the human E μ and 3' *IGH* super-enhancers. This allowed us to ask whether the targeting of SHM to
371 the B1-8^{hi} variable region was influenced by differences in the human versus the mouse enhancers and
372 promoters.

373 Importantly, when placed in the context of the human *IGH* locus, the B1-8^{hi} sequence was
374 covered almost exclusively by uniquely mapping reads (Fig. 6B), which is in sharp contrast to the
375 extensive multimapping observed within B1-8^{hi} at the murine *Igh* locus (Fig. S2E). This is because the
376 Vh1-72 gene used in the B1-8^{hi} sequence is much less homologous with the human V genes than the
377 mouse V genes. PRO-seq and PRO-cap revealed that the nascent transcriptional profile of the B1-8^{hi}
378 sequence in Ramos^{B1-8^{hi}} cells (Fig. 6C) shared similarities with the B1-8^{hi} sequence in the murine *Igh*
379 context (Fig. 5A). For example, similarities were noted in CDR3 and flanking sequences, for instance, the
380 presence of the antisense initiation site (PRO-cap track) and shared spikes of nascent transcriptional
381 activity (PRO-seq) (Fig. 6C). More importantly, MutPE-seq following AID(JP8Bdel) infection revealed a
382 very similar SHM spectrum in Ramos^{B1-8^{hi}} cells compared to the murine B1-8^{hi} context, especially with
383 primary murine B cells (Fig. 6D). Of note, the dominant, palindromic AGCT hotspot in CDR3 was the most
384 mutated residue exactly as in B1-8^{hi} murine primary cells and GCBs (Fig. 6D). The other mutated C:G
385 residues in the Ramos^{B1-8^{hi}} sequence were largely the same as those mutated in the murine context albeit
386 the relative mutation rates within each variable region varied (Fig. 6D, compare with Fig. 5B-C).
387 Importantly, mutation frequencies in murine GCBs and AID(JP8Bdel)-expressing Ramos cells were

388 comparable (note the Y axes scales in Fig. 6D) with the major difference being the weak A:T mutagenesis
389 typical of Ramos cells. Of note, the sense PRO-cap peak in CDR3 in the murine B1-8^{hi} (Fig. 5A), located
390 just 3 bp upstream of the strong AGCT hotspot, is absent in Ramos^{B1-8^{hi}} (Fig. 6C) implying that this
391 initiation site is not essential for the high mutation frequency of this hotspot motif. In sum, the SHM
392 spectrum of the B1-8^{hi} variable region in mouse B cells is largely retained in Ramos^{B1-8^{hi}} cells despite
393 dissimilarities in their nascent transcriptional profiles.

394 We conclude that although enhancers are important for recruiting AID to the variable regions and
395 other non-IG genes, the final rules of mutagenesis are apparently “hard-wired” into the DNA sequence
396 rather than specific transcriptional features or chromatin marks.

397

398 Discussion

399 In this study, we aimed to understand the relationship of SHM spectra with local patterns of nascent
400 transcription and epigenetic marks using high-resolution profiling of variable regions. We find that the
401 local SHM patterns and frequencies cannot be predicted from, nor do they appear to be derived from, the
402 underlying transcriptional landscape. Moreover, substantial loss of actively elongating Pol II and nascent
403 transcription in both orientations within variable regions, seen in E μ ^{-/-} cells, does not majorly impact on
404 SHM. These results have important implications for the interpretation of transcription and chromatin
405 modification measurements in the context of SHM or CSR. For example, despite the global correlation of
406 SHM target genes with activating histone marks, our results suggest that they do not play a major role in
407 determining SHM frequency or in defining the SHM spectrum. We suspect that in studies depleting
408 histone modifying enzymes, SHM or CSR defects may result, in part or whole, from secondary effects
409 arising from perturbation of these global gene regulatory factors. In the case of histone modifying
410 enzymes, it is plausible that they function independently of their enzymatic activity or via non-histone
411 substrates, as has been shown for many such enzymes (81-85). Similarly, we infer that although AID
412 target genes are often enriched in paused Pol II and convergent transcription, these features do not
413 appear to play a direct, mechanistic role in mutagenesis. Moreover, the fact that SHM remains robust

414 despite up to 4-fold reduced *IGH* nascent transcription in $E_{\mu}^{-/-}$ cells implies that although SHM requires
415 transcription, it is not directly correlated with the rate of transcription and, by extension, the level of Pol II
416 occupancy in the variable region. This has important implications for interpreting the meaning of
417 correlations between transcription and SHM observed upon gene knockouts or knockdowns. In such
418 cases, our results suggest that the true cause of the SHM defect may not be due to, or may only partly
419 result from, decreased transcription.

420 During transcription, the non-template strand is released as ssDNA to allow RNA synthesis.
421 Hence, this ssDNA bubble has been thought to serve as a source of ssDNA for AID, especially if Pol II is
422 paused or its elongation rate is slowed down, a situation where ssDNA would be available more stably
423 and for longer periods. However, our data show no consistent correspondence of mutations hotspots with
424 sites of Pol II accumulation. From structural studies of paused and elongating Pol II complexes, it is
425 apparent that of the ~11 nt non-template ssDNA bubble, about half lies buried within Pol II (71, 72) (Fig.
426 7A). Importantly, the exposed portion (~5 nt) is entirely covered by SPT5 (72) (Fig. 7A). The PAF complex
427 and SPT6 are located further away from this site (71) (Fig. 7A). Under such structural constraints, it is
428 difficult to conceive how AID could access the ~5 nt exposed ssDNA and catalyze deamination even if Pol
429 II were paused or slowly elongating. Thus, it appears unlikely that the transcription bubble is a substrate
430 for AID which can explain why we observe no correlation between sites of Pol II accumulation and SHM.
431 Moreover, it appears that SHM at recurrent hotspot residues is not affected by distally located mutations
432 elsewhere in the variable region (47). This leads us to infer that sequences directly flanking the target
433 cytidine may be important determinants of mutability.

434 We propose the following four-step model to integrate our data with the literature to date. First,
435 AID is delivered to IG variable regions and non-IG targets via enhancers (38, 40, 42) (Fig. 7B. step 1).
436 This is accomplished by chromatin looping between enhancers and target genes, a process that has
437 been suggested to occur via cohesin-mediated loop extrusion (86). Such interactions can lead to the
438 formation of enhancer-gene hubs via dynamic, multivalent interactions between Pol II, transcription
439 factors and cofactors like the Mediator complex. This results in the stabilization of transcription factors
440 and recruitment of Pol II and cofactors to activate target genes (87-90). Once delivered into the hub, AID

441 associates with Pol II and associated proteins such as SPT5, SPT6 and the PAF complex (Fig. 7B, step
442 2). We propose that the primary purpose of these interactions is to retain AID in the vicinity of Pol II
443 complexes. The ability of AID to engage in multiple, independent dynamic interactions ensures that it is
444 retained within the transcriptional hub, ready to access ssDNA when it is made available. Indeed, the fact
445 that depletion of any of these proteins leads to AID recruitment defects (21-24) suggests that all such
446 interactions are contributing to ensure that AID stays within the transcription hub or that these factors are
447 important for maintaining or forming the hub (90).

448 For successful deamination, ssDNA needs to be exposed and the source of this ssDNA remains
449 enigmatic (Fig. 7B, step 3). As explained above, the transcriptional bubble appears to be inaccessible.
450 Previous work has suggested that DNA upstream of the moving Pol II could provide a source of ssDNA
451 for AID in variable regions since, being negatively supercoiled, it may be more accessible (60, 91).
452 Patches of ssDNA have been detected in the variable region and it has been suggested that they arise
453 from negative supercoiling associated with paused Pol II although their presence did not correlate well
454 with mutation spectra (57, 61, 91). Indeed, AID was shown to target negatively supercoiled DNA, but not
455 relaxed DNA, on both template and non-template strands (92). However, since we find no correlation
456 between mutation and pausing, we infer that these ssDNA patches are likely the result of transcription
457 itself rather than specific sites of pausing.

458 Although, in principle, negative supercoils could create ssDNA patches that allow mutation on
459 both strands, the fact that significant mutational asymmetry is observed upon depletion of RNA Exosome
460 subunits or the DNA:RNA helicase, Senataxin (93-95) argues that the species AID likely targets is one
461 containing an RNA:DNA hybrid. Indeed, the variable region sense transcript is a substrate of the RNA
462 exosome and the latter is required for optimal mutation of the template strand in the B1-8^{hi} variable
463 region, findings that directly implicate the processing of DNA:RNA hybrids in SHM (93). In this regard, and
464 in contrast to IG switch regions, variable regions lack R loops (61, 91). Given the presence of ssDNA
465 patches, the requirement for RNA processing and the lack of R loops, we propose that transcription
466 termination may provide an important source of ssDNA in variable regions. It is plausible that the collapse
467 of the elongation complex could release DNA with the nascent RNA within the transcription bubble still

468 hybridized (Fig. 7B, step 4). The RNA would have to be unwound by DNA:RNA helicases like Senataxin
469 (95) and/or digested with RNaseH followed by degradation from the 5' end by Xrn2 and from the 3' end
470 by the RNA exosome, thereby providing AID access to the template strand (26, 93-96) (Fig. 7B, step 4).
471 In this regard, termination has been correlated with SHM (60, 97) and degradation of Pol II by
472 ubiquitination has been linked to AID targeting (98). We note that if, during multiple rounds of transcription
473 within a cell population, termination occurs stochastically along the variable region and intron, then such
474 events would not be detected by population-based transcriptomic assays like PRO-seq or ChIP-seq.

475 Interestingly, recent genome-wide analyses have shown that termination is frequent in the first 3
476 kb from the promoter (99) which, in the context of the IG genes, would encompass the entire variable
477 region and, depending on which J segment is used, part or whole of the intronic sequence. Indeed, SHM
478 starts 100-150 bp from the V gene promoter and extends for another 1.5-2 kb (43-45). Thus, we suspect
479 that ssDNA patches and DNA:RNA hybrids may arise, at least in part, from premature termination of Pol II
480 within the variable region. A similar process could occur at the 5' ends of non-IG AID target genes or at
481 transcribed enhancers that are also targets of AID. Given that neither we nor others have observed any
482 transcriptional feature correlating with SHM spectra, we favor the idea that deamination occurs on ssDNA
483 made available through co-transcriptional processes that expose ssDNA patches, such as termination,
484 but that it is dissociated from actively elongating Pol II (Fig. 7B, step 4). We note, however, that
485 termination is difficult to measure directly, and how the disassembly of Pol II occurs is poorly understood.
486 Thus, the mechanism and kinetics of resolution of the DNA:RNA hybrid within terminating Pol II are not
487 known.

488 The final step, that is, the actual mutagenesis, determines the observed SHM spectra (step 4,
489 Fig. 7B). AID has a very high affinity for ssDNA and can remain associated with ssDNA for several
490 minutes which may allow sufficient time for deamination, a reaction which is extremely inefficient and can
491 take up to several minutes (100, 101). However, WRCH motifs are not equally mutated, suggesting that
492 the DNA sequence context determines whether ssDNA containing a WRCH motif is mutated and to what
493 extent, a notion which is supported by biochemical studies (101). Indeed, analyses of variable regions
494 has suggested that sequence-intrinsic features may regulate the SHM spectra (47, 50-52). Given that

495 neither ssDNA patches nor nascent transcriptional features correlate with SHM spectra, we favor the
496 notion that sequence context of a WRCH motif determines its mutability and propose that sequence
497 context regulates the residence time of AID on ssDNA and thereby modulates the catalytic efficiency of
498 AID *in situ* (Fig. 7B, step 4).

499 The precise contribution of DNA sequence context towards SHM will be a major focus of
500 investigation for the future with the goal of understanding the rules governing differential mutability.
501 Insights from such studies could be harnessed to further increase the potency and breadth of mature
502 broadly neutralizing antibodies where hotspot saturation can prevent further SHM (47). In this regard, our
503 system of efficiently replacing variable regions in Ramos cells offers an ideal platform for asking how
504 changes in the sequence locale affects the targeting of SHM to a given hotspot.

505

506 **Materials and Methods**

507 **Cell culturing**

508 Ramos were cultured in complete RPMI medium (in-house) supplemented with 10% fetal bovine serum
509 (FBS; Invitrogen), glutamine (Invitrogen), sodium pyruvate (Invitrogen), HEPES (made in-house) and
510 antibiotic/antimycotic (Invitrogen). LentiX packaging cells were cultured in complete DMEM medium (in-
511 house) containing 10% fetal bovine serum (FBS; Invitrogen), glutamine (Invitrogen), sodium pyruvate
512 (Invitrogen), HEPES (made in-house) and antibiotic/antimycotic (Invitrogen).

513 **Mice**

514 The mice were maintained in a C57BL/6 background and housed in the IMBA-IMP animal facility in
515 standard IVC cages with HEPA filtering. All animal experiments were carried out with valid breeding and
516 experimental licenses (GZ: MA58-320337-2019-9, GZ: 925665/2013/20 and GZ: 618046/2018/14)
517 obtained from the Austrian Veterinary Authorities and in compliance with IMP-IMBA animal house
518 regulations.

519

520 **Generation of AID^{-/-} Ramos cells**

521 We generated Ramos cells expressing EcoR, the receptor for the ecotropic envelope protein, Eco-Env,
522 so as to allow efficient infection of Ramos cells with ecotropic lentiviruses. Ramos cells were infected with
523 a pRRL lentiviral vector expressing EcoR and maintained under Puromycin selection. These cells were
524 transfected with a small guide RNA targeting *AICDA*, the gene expressing AID, and recombinant Cas9
525 protein via electroporation using the Neon transfection system (Thermo Fisher Scientific). The following
526 day, single cells were sorted into 96-well plates and allowed to expand. Four weeks later, clones
527 harboring frameshifting indels at the Cas9 target site were identified by genotyping and loss of AID in
528 these clones was confirmed by Western blot analysis. One clone (D3) was used for all subsequent
529 analyses and for generation of new lines.

530 **Generation of E μ ^{-/-} Ramos cells**

531 To delete the 583 bp segment containing E μ , we generated homology repair plasmids having Ef1a
532 promoter-driven floxed GFP and mCherry expression cassettes flanked by homology arms. These
533 plasmids were transfected into AID^{-/-} Ramos cells (clone D3) along with three in vitro synthesized guide
534 RNAs (6 μ g each, designed using CRISPOR) and 7.5 μ g recombinant Cas9 protein (Vienna Biocenter
535 Core Facilities) using the Neon Transfection System (Thermo Fisher Scientific). A week later,
536 GFP/mCherry double-positive single cells were isolated using a BD FACS Aria III sorter (BD
537 Biosciences). Successful knock-in clones were identified by genotyping with PCR and Sanger sequencing
538 of PCR products. Next, 200 μ g recombinant Cre recombinase (Molecular Biology Service, IMP) was
539 added to the culture medium to excise the floxed mCherry/GFP cassettes. GFP/mCherry double-negative
540 clones were isolated one week after electroporation using a BD FACS Aria III sorter and genotyped via
541 PCR and Sanger sequencing of PCR products.

542 **Generation of Ramos cells expressing new, exogenous variable regions**

543 AID^{-/-} Ramos cells (clone D3) were electroporated with Cas9-sgRNA ribonucleoprotein complexes
544 (prepared in-house by the Vienna Biocenter Core Facilities) and homology repair templates containing a
545 pair of unique sgRNA-target sites to excise the entire variable region including the promoter (Fig. S2A).

546 Single, IgM-negative clones were isolated via flow cytometry, expanded and genotyped. One validated
547 IgM-negative clone was selected as the parental clone. This clone was then electroporated with Cas9-
548 sgRNA ribonucleoprotein complexes targeting the unique sgRNA sites and homology repair templates
549 containing new variable regions under the control of the endogenous Ramos VH4-34 promoter. Single,
550 IgM-positive cells were isolated, expanded and genotyped.

551 **Lentiviral infections**

552 Lentiviral pRRL vectors expressing AID(JP8Bdel) coupled with mCherry were transfected along with
553 ecotropic envelope (Eco-env)-expressing helper plasmid into LentiX cells via the standard calcium
554 phosphate methodology. Lentiviral supernatants were used to infect Ramos cells via spinfection (2350
555 rpm for 90 min) in the presence of 8 $\mu\text{g/ml}$ Polybrene (Sigma). Flow cytometry was used to sort mCherry-
556 positive cells for mutation analysis.

557 **Isolation and activation of murine primary, splenic B cells**

558 Mature, naïve B cells were isolated from spleens of 2-4-month-old wild-type C57BL/6, as per established
559 protocols (Pavri et al., 2010). B cells were cultured in complete RPMI medium supplemented with 10%
560 fetal bovine serum (FBS) and antibiotics, Interleukin 4 (IL4; made in-house by the Molecular Biology
561 Service, IMP), 25 $\mu\text{g/ml}$ Lipopolysaccharide (Sigma) and RP105 (made in-house by the Molecular Biology
562 Service, IMP) and harvested after 3 or 4 days. B1-8^{hi} Rosa26^{AID^{ER}} mice were generated by crossing B1-
563 8^{hi} mice (65) with Rosa26^{AID^{ER}} mice (73) and maintained as a homozygous line for both alleles. For SHM
564 assays from activated B1-8^{hi} Rosa26^{AID^{ER}} primary B cells, 2 μM 4-hydroxy tamoxifen (4-HT) (Sigma) was
565 added at the time of activation with IL-4, LPS and RP105.

566 **Isolation of germinal center B cells (GCBs) from immunized B1-8^{hi} mice**

567 Sheep red blood cells (SRBCs) were washed thrice with PBS of which 0.2x10⁹ SRBCs in 100 μl PBS (per
568 mouse) were injected followed by another injection of 1x10⁹ SRBCs in 100 μl PBS five days later. Mice
569 were harvested twelve days after the first immunization. Spleens were harvested and single-cell
570 suspensions were stained with B220 conjugated to fluorescein isothiocyanate (B220-FITC, BD

571 Biosciences, 1:500 dilution), Fas conjugated to PE-cyanine 7 (Fas-PE-Cy7, BD Biosciences, 1:1000
572 dilution) and CD38 conjugated to Allophycocyanin (CD38-APC, ThermoFisher, 1:200 dilution) along with
573 Fc block (BD Biosciences, 1:500 dilution). GCBs (B220⁺ Fas⁺ CD38⁻) were isolated on a BD FACS Aria
574 II sorter (BD Biosciences). Following sorting, nuclei were isolated and used for PRO-seq and PRO-cap
575 (see below). For one experiment, ~10⁷ GCBs from fifteen immunized mice were pooled yielding 3-4x10⁶
576 viable nuclei that were used for run-on. For MutPE-seq, genomic DNA was isolated from 15 x10⁶ GCBs
577 as described previously (102).

578 **PRO-seq and PRO-cap**

579 PRO-seq was performed as described previously (103) with several modifications.

580 To isolate nuclei, murine and *Drosophila* S2 cells were resuspended in cold Buffer IA (160 mM Sucrose, 10
581 mM Tris-Cl pH 8, 3 mM CaCl₂, 2 mM MgAc₂, 0.5% NP-40, 1 mM DTT added fresh), incubated on ice for 3
582 min and centrifuged at 700 g for 5 min. The pellet was resuspended in nuclei resuspension buffer NRB
583 (50 mM Tris-Cl pH 8, 40% Glycerol, 5 mM MgCl₂, 0.1 mM EDTA). For each run-on, 10⁷ nuclei of the
584 sample and 10% *Drosophila* S2 nuclei were combined in a total of 100 μL NRB and incubated at 30°C for
585 3 min with 100 μL 2x NRO buffer including 5μl of each 1mM Bio-11-NTPs. In some PRO-cap
586 experiments, the run-on reaction was performed in the presence of two biotinylated NTPs, Biotin-11-UTP
587 and Biotin-11-CTP (Perkin-Elmer), and unlabeled ATP and GTP. Subsequent steps were performed as
588 described (Mahat et al 2016), except that 3' and 5' ligations were performed at 16°C overnight and
589 CapClip Pyrophosphatase (Biozym Scientific) was used for 5' end decapping. We also used customized
590 adapters (PRO-seq: 3' RNA adapter:

591 5'5Phos/NNNNNNGAUCGUCGGACUGUAGAACUCUGAAC/3InvdT-3' and 5' RNA adapter: 5'-
592 CCUUGGCACCCGAGAAUCCANNNN-3'). RNA was reverse transcribed by SuperScript III RT with
593 RP1 Illumina primer to generate cDNA libraries. Libraries were amplified with barcoding Illumina RPI-x
594 primers and the universal forward primer RP1 using KAPA HiFi Real-Time PCR Library Amplification Kit.
595 PRO-cap: 3' linker (DNA oligo): 5rApp/NNNNAGATCGGAAGAGCACACGTCT/3ddC , 5'linker (RNA):
596 ACACUCUUUCCCUACACGACGCUCUCCGAUCUNNNNNNNNNN, reverse transcription and library
597 amplification as in PRO-seq but using selected TruSeq_IDX_1-48 PCR primers for RT, and the same

598 TruSeq_IDX PCR primers together with TruSeq Universal Adapter forward primer for PCR. For both
599 methods amplified libraries were subjected to gel electrophoresis on 2.5% low melting agarose gel and
600 amplicons from 150 to 350 bp were extracted from the gel (carefully separated from possible linker
601 dimers), multiplexed and sequenced on Illumina platforms with SR50 or longer read mode.

602 **RT-qPCR**

603 RT-qPCR assays with externally spiked-in *Drosophila* S2 cells were described in detail in our previous
604 studies (104). Briefly, *Drosophila* S2 cells were mixed with Ramos cells or B1-8^{hi} B cells at a ratio of 1:4
605 (4×10^5 B cells and 1×10^5 S2 cells) followed by total RNA extraction with TRIzol reagent (Thermo Fisher),
606 DNaseI digestion, and cDNA synthesis with random primers. The $2^{-\Delta\Delta C_t}$ method was used to quantify
607 the data using the *Drosophila Act5c* transcript for normalization.

608 **ChIP-qPCR and ChIPseq**

609 ChIP-seq and ChIP-qPCR were performed as described previously without any modifications (22, 102).
610 Antibodies used in this study are listed in Table S1.

611 **ATAC-seq**

612 ATAC-seq was performed exactly as described in detail in our previous work (102).

613 **Mutational analysis by paired-end deep sequencing (MutPE-seq)**

614 MutPE-seq was performed following the principles described in two previous reports (48, 73) with several
615 adaptations. 80 ng of genomic DNA were amplified by PCR with the Kapa Hifi HS 2x RM (Roche
616 Diagnostics), For the first PCR, we used 20-25 cycles with $0.2 \mu\text{M}$ locus-specific primers fused to a
617 varying number of random nucleotides and to the first part of the Illumina adapter sequences (FW: 5'-
618 CTCTTTCCCTACACGACGCTCTTCCGATCT-(N)_x-gene-specific sequence-3'; RV: 5'-
619 CTGGAGTTCAGACGTGTGCTCTTCCGATCT-gene-specific sequence-3'; see Table S1 for a complete
620 list of primer sequences). Random Ns are used to increase complexity and shift frames of very similar
621 amplicons to improve cluster calling and sample identification. After first PCR, the samples were purified
622 with 0.2x/0.7x SPRI beads (Beckman Coulter), eluted in 10 μl water and amplified for 10 cycles with

623 0.75 μ M primers containing linker sequences and dual barcoding (FW: 5'-
624 AATGATACGGCGACCACCGAGATCTACACXXXXXXXXXACACTCTTCCCTACACGAC-3'; RV: 5'-
625 CAAGCAGAAGACGGCATAACGAGATXXXXXXXXGTGACTGGAGTTCAGACGTGTG-3' where the stretches
626 of X nucleotides serve as barcodes creating a unique dual barcode combination for each sample). PCR
627 products were purified either by extraction from a 2% low-melt agarose gel or with 0.7x SPRI beads and
628 eluted in 20 μ l water. The concentration was determined by Picogreen-based measurements on a
629 Nanodrop machine or by a Fragment Analyzer. Samples were equimolarly pooled for next generation
630 sequencing on an Illumina MiSeq flowcell and sequenced paired-end (PE300).

631 **Primers and sgRNAs**

632 All oligos, primers and adapters used in this study are listed in Table S1.

633

634 **Bioinformatics**

635 **Mapping PRO-seq, PRO-cap, CHIP-seq and ATAC-seq reads to IG variable regions**

636 Reads were trimmed for standard adapters and low quality (Q<30) 3'-end bases and filtered for a
637 remaining length of at least 20nt (excluding the UMI, if applicable) using cutadapt (105). Alignment to the
638 reference genome was done with Bowtie (106). The NCBI GRCm38.6 assembly was used as the mouse
639 reference. To this we added the sequence of the recombined locus containing the variable region as an
640 additional chromosome. Similarly, for human data, the Hg38 assembly was used, with the relevant
641 variable region sequence added as a chromosome. Where applicable, the dmr6 assembly of *Drosophila*
642 *melanogaster* from Flybase ((107) release 6.27) was used as the spike-in reference and was added to the
643 alignment index. During alignment, up to three mismatches were allowed. To accommodate mapping to
644 the repetitive *IGH/Igh* locus, a high degree of multimapping was allowed (194 potential V segments
645 annotated for GRCm38.6). In the case of alignment with spike-ins, only reads mapping exclusively to
646 either genome were considered. For PCR deduplication, UMIs were identified and filtered with UMI-tools
647 ((108) v1.0.0). No sequence differences were allowed in UMIs when collapsing the duplicates.
648 Multimapping reads were then filtered to identify those specific to the variable region sequence, taking

649 into account the repetitiveness of the reference *IGH/Igh* locus. Reads mapping to the fixed variable region
650 sequence were allowed to multimap only to the native *IGH/Igh* locus (mouse chr12:113572929-
651 116009954 or human chr14:105836764-106875071). The native loci were defined as the region between
652 the earliest position belonging to an annotated V segment and the latest position belonging to an
653 annotated J segment. Reads with a mapped location outside these areas were rejected. This filtering was
654 implemented in a custom Python script. The qualifying reads were then classified as reads mapping
655 uniquely to the particular variable region sequence and reads mapping to both the variable region
656 sequence and to the native *IGH/Igh* region. Genome browser tracks were created by quantifying and
657 scaling read coverage of the variable region sequence using bedtools ((109) v2-29.0). Reads were split
658 by strand, and strand labels for PRO-seq were inverted in order to match the strand designations of the
659 respective PRO-cap reads. For both data types, only the first 5' base of the strand-adjusted reads was
660 used for the tracks. Additionally, the 5' end of PRO-seq reads was shifted by 1 nucleotide downstream, to
661 compensate for the fact that the last nucleotide in the RNA was incorporated during the run-on reaction.
662 For normalized tracks, read counts were scaled to RPM (reads per million). Finally, “-” strand coverage
663 tracks were further scaled by -1, for visualization purposes.

664 Code for the workflow and the custom scripts is available on Github at
665 https://github.com/PavriLab/IgH_VDJ_PROcapseq.

666

667 **Analysis of MutPE-seq**

668 Reads were trimmed for standard adapters with cutadapt (105). Poor quality (Q<25) 3' bases were
669 trimmed with trimmomatic (110) by averaging over a sliding window of 5nt. Read pairs were then filtered
670 for minimum remaining length (200nt for read 1, 100nt for read 2) using cutadapt. Read mates were
671 merged down to make combined single-end reads with FLASH (111) allowing 10% mismatch between
672 the mates. Obvious erroneous mergers were removed by selecting combined reads with lengths within
673 ± 30 nt of the amplicon length using cutadapt. The remaining combined reads were aligned with Bowtie2
674 (112), using the “-very-sensitive-local” alignment mode and only the fixed variable region sequence and

675 its immediate vicinity as reference. Samples of V genes VH4-59 and VH4-34 are different only by single
676 nucleotide polymorphisms and were filtered for contamination by the respective second sequence by
677 jointly aligning with the “expected” and the other “contaminating” variable region sequence, discarding all
678 aligned contaminating reads. A pile-up was generated with samtools (113) taking into account only bases
679 with quality of at least 30. The pileups were then quantified with a custom Python script and the resulting
680 mutation counts were processed and visualized with custom scripts in R (v3.5.1), with the help of
681 additional R packages (data.table [<https://cran.r-project.org/web/packages/data.table/index.html>], ggplot2
682 [<https://cran.r-project.org/web/packages/ggplot2/index.html>] , ggrepel [[https://cran.r-
684 project.org/web/packages/ggrepel/index.html](https://cran.r-
683 project.org/web/packages/ggrepel/index.html)], patchwork [[https://cran.r-
686 project.org/web/packages/patchwork/index.html](https://cran.r-
685 project.org/web/packages/patchwork/index.html)]). Background mutation profiles were controlled for by
685 subtracting the corresponding mutation frequencies in control samples from the frequencies in the
686 samples of interest, at each position and for each substitution type. Annotation of hot and cold spots was
687 created by means of regex search for the corresponding patterns in the reference sequence.

688 Code for the workflow and the custom scripts is available on Github at
689 https://github.com/PavriLab/IgH_VDJ_MutPE.

690 **Data availability**

691 All NGS data has been deposited in GEO under accession number GSE202042.

692 **Author contributions**

693 UES performed all PRO-seq, PRO-cap and MutPE-seq assays and analyzed data. JF generated the
694 Ramos lines with new variable regions and performed infections for MutPE-seq. KF and TN developed
695 the bioinformatic pipeline to map reads to variable regions, KF developed the bioinformatic pipeline for
696 MutPE-seq. Marina M performed ChIP-seq, ChIP-qPCR and RT-qPCR in Ramos cells and analyzed data.
697 IO performed the bioinformatic analysis of non-IG AID target genes. BB performed ChIP-qPCR and RT-
698 qPCRs. EMW and RP generated the $E_{\mu}^{-/-}$ lines. MS conducted GCB PRO-seq with UES. ACG performed
699 ATAC-seq. Marialaura M cloned JP8Bdel expression vectors for MutPE-seq. ACG performed ATAC-seq.
700 HM provided the germline-reverted VH4-59 and VH3-30 sequences and provided feedback on the

701 manuscript. ITS performed bioinformatic analysis, provided resources and analyzed data. RP conceived
702 the project, analyzed data and wrote the manuscript with critical inputs from UES and JF.

703

704 **Acknowledgements**

705 We thank Almudena Ramiro and Angel Alvarez-Prado (CNIC, Madrid, Spain) for sharing the mutational
706 datasets from Ung^{-/-}Msh2^{-/-} GCBs. We gratefully acknowledge the Vienna Biocenter Core Facilities
707 (VBCF) for assistance in generating E μ ^{-/-} lines and for next generation sequencing, the IMP/IMBA
708 BioOptics facility for flow cytometry usage, the IMP/IMBA Molecular biology services for Sanger
709 sequencing and reagents, and the IMP/IMBA animal house. We thank Maximilian von der Linde for
710 uploading NGS tracks to GEO. We thank Carrie Bernecky (IST, Vienna) and Clemens Plaschka (IMP) for
711 generating the 3D Pol II elongation complex structural visualizations. This work was funded by Boehringer
712 Ingelheim, The Austrian Industrial Research Promotion Agency (Headquarter Grant FFG-834223), and
713 grants from the Austrian Science Fund to UES (FWF T 795-B30) and RP (FWF P 32043-B).

714

715 **Conflicts of interest**

716 The authors declare no conflicts of interest.

717

718

719

720

721

722

723 **Figure Legends**

724 **Figure 1: The E μ enhancer regulates the chromatin landscape in the intron but not in the variable**
725 **region in Ramos human B cells. (A)** ChIP-seq profiles of H3K27ac, H3K4me3 and H3K36me3 at the
726 *IGH* locus from AID^{-/-} and two independent clones (c1 and c2) of AID^{-/-} E μ ^{-/-} Ramos cells, the latter
727 generated as described in Fig. S1C-D. The sequence annotations indicate the positions of the *IGH*
728 promoter (Prom), the complementary determining regions (CDR1-3), the intronic enhancer (E μ) and the
729 switch μ region (S μ). The bioinformatic approach to include multimapping reads (explained in the text and
730 depicted in Fig. 3A) was applied. **(B)** A magnified view of the locus from A. **(C)** ChIP-qPCR analysis at
731 *IGH* from AID^{-/-} and AID^{-/-} E μ ^{-/-} Ramos cells (c1, c2) to measure the relative levels of H3K27ac,
732 H3K4me3, H3K79me3, pan-acetylation of histones H3 and H4, and histone H3. The amplicons used for
733 PCR are indicated in the schematic diagram at the top. The active *B2M* gene is used as a positive control
734 and a gene desert on chromosome 1 is used as a negative control (Neg.). The data represent three
735 independent experiments. Asterisks indicate $P < 0.05$ using the Student's t-test and ns indicates not
736 significant ($P > 0.05$).

737

738 **Figure 2: Ablation of the E μ enhancer significantly decreases nascent transcription but not SHM.**
739 **(A)** RT-qPCR measurements of nascent transcripts at the Ramos variable region in AID^{-/-} and AID^{-/-}
740 E μ ^{-/-} cells (c1, c2). To account for potential clonal variation, Ramos cells were spiked with *Drosophila* S2
741 cells prior to RNA extraction. The data was normalized to the levels of the *Drosophila* housekeeping
742 gene, *Act5c*. *GAPDH* mRNA was used as a control. Asterisks indicate $P < 0.05$ using the Student's t-test
743 and ns indicates not significant ($P > 0.05$). **(B)** RT-qPCR analysis as in A measuring the spliced IgM
744 mRNA. **(C)** Table of mutation frequencies at the Ramos variable region and the JH6 intron (amplicons
745 shown in the diagram above the table) in AID^{-/-} and AID^{-/-} E μ ^{-/-} cells following infection with
746 AID(JP8Bdel) for 6 days. Statistical analysis was performed with the Student's t-test. **(D)** Flow cytometry
747 analysis of IgM expression in AID^{-/-} and AID^{-/-} E μ ^{-/-} Ramos cells infected with AID(JP8Bdel) for 7 days.

748 **(E)** Bar graph summarizing the flow cytometry analyses from D, showing the percent of IgM loss from
749 three independent experiments.

750

751 **Figure 3: Comparison of transcriptional and mutational landscapes at the endogenous Ramos**
752 **variable region. (A)** Scheme depicting the strategy used to align multimapping reads to the variable
753 region. As described in the text and Methods, since upstream, non-recombined V genes are silent in
754 Ramos cells, a read that maps to the recombined variable region is retained if it is mapping to any V, D or
755 J segment of the human *IGH* locus but nowhere else in the human genome. The same principle is applied
756 when mapping reads to murine variable regions to the mouse genome. **(B)** Integrative Genomics Viewer
757 (IGV) browser snapshot of nascent RNA 3' ends (by PRO-seq) aligned to the Ramos variable region.
758 Multimapping (top track) reads are separated from uniquely mapping reads (middle track), and these two
759 tracks are subsequently combined to generate the total profile (bottom track. **(C)** PRO-cap and PRO-seq
760 5' and 3' end densities, respectively, at the *IGH* locus along with mutation frequencies at the variable
761 region displayed on the Integrative Genomics Viewer (IGV) browser. Mapping of the variable region
762 transcriptome was done via the bioinformatic pipeline outlined in A and exemplified in B (total profiles are
763 shown). The locations of the antigen-binding complementary determining regions (CDR1-3) are
764 highlighted. Mutation analysis via MutPE-seq was performed following infection of AID^{-/-} cells with
765 AID(JP8Bdel)-expressing lentiviruses for 6 days. **(D)** A magnified view of the variable region from C
766 above. **(E)** Details of the mutation spectrum at the Ramos variable region. Mutated cytidines in AID
767 hotspot motifs (WRCH) are displayed as red bars. All other C:G mutation are shown as black bars and
768 A:T mutations as grey bars. The panel under the graph shows the position of both hotspot (WRCH in
769 black, AGCT in red, upper panel) and coldspot (SYC, bottom panel) motifs. **(F)** Waterfall plot with
770 mutations ordered from highest (left) to lowest (right) frequency following the color code described in E.
771 **(G)** Mutation frequency bar plot showing the percentages (indicated within the bars) of the three
772 mutations classes following the color code described in E. **(H)** Bar graph indicating the percentage of the
773 type of mutation indicated on the X axis. The C→T and G→A transition mutations are the signature of AID
774 activity.

775 **Figure 4: Comparison of transcriptional and mutational landscapes at two different human**
776 **variable regions expressed from the Ramos *IGH* locus. (A)** IGV browser snapshots of nascent RNA 5'
777 (PRO-cap) and 3' (PRO-seq) ends and mutation tracks (MutPE-seq) at the VH4-59-DH2-JH6 variable
778 region expressed from the Ramos VH4-34 promoter (CDRs 1-3 highlighted). Mutation analysis was
779 performed following infection of AID^{-/-} cells with AID(JP8Bdel)-expressing lentiviruses for 7 days. **(B)** A
780 magnified view of the VH4-59-DH2-JH6 variable region from A above. **(C)** Detailed mutational analysis of
781 the VH4-59-DH2-JH6 variable region displayed and color-coded as in Fig. 3E. A bar plot, as in Fig. 3G,
782 with the percentage of mutation frequencies is shown on the right. **(D)** Nascent transcriptional analysis of
783 the VH3-30-DH2-JH6 variable region as in A. **(E)** A magnified view of the VH3-30-DH2-JH6 variable
784 region from D above. **(F)** Detailed mutational analysis of the VH3-30-DH2-JH6 variable region (see C
785 above for details).

786

787 **Figure 5: Transcriptional and mutational landscapes of the murine B1-8^{hi} variable region and non-**
788 ***Ig* AID target genes in mice. (A)** Nascent transcriptional analysis at the *Igh* locus in murine B1-8^{hi}
789 primary, splenic B cells stimulated for four days with LPS, IL4 and RP105. IGV browser snapshots show
790 the 5' (PRO-cap) and 3' (PRO-seq) ends of the aligned reads. For MutPE-seq, B1-8^{hi} *Rosa26*^{AIDER}
791 primary B cells (expressing AID fused to the estrogen receptor from the *Rosa26* locus) were activated
792 with LPS, IL4 and RP105 for four days in the presence of 4-hydroxytamoxifen (4-HT) to trigger AIDER
793 nuclear import. PRO-seq, PRO-cap and MutPE-seq were also performed from sorted, splenic germinal
794 center B cells (GCBs) following immunization with sheep red blood cells for 11 days. **(B)** Analysis of SHM
795 spectra of the B1-8^{hi} variable region from B1-8^{hi} *Rosa26*^{AIDER} primary, activated murine splenocytes. The
796 bar graph on the right shows the percentage of each indicated mutation category (see Fig. 3E and 3G for
797 details). **(C)** Analysis of SHM spectra of the B1-8^{hi} variable region from splenic GCBs following
798 immunization with sheep red blood cells for 11 days. **(D)** 5' and 3' ends of nascent RNA (PRO-cap, PRO-
799 seq) and mutation profiles at four selected AID target genes. Mutational data are from Alvarez-Prado et
800 al. (2018) wherein the first 500 bp of the genes were sequenced. The region displayed extends from -100

801 bp from the annotated (RefSeq) TSS up to 50 bp downstream of the sequenced amplicon. The WRCH
802 motifs are indicated as red dots.

803

804 **Figure 6: Analysis of gene regulatory context on the nascent transcriptional landscape and SHM**

805 **of the murine B1-8^{hi} variable region. (A)** Scheme showing the exchange of the endogenous Ramos

806 variable region for the murine B1-8^{hi} variable region to generate the Ramos^{B1-8^{hi}} human *IGH* locus

807 following the approach described in Fig. S2A. **(B)** 3' ends of nascent RNAs by PRO-seq at the Ramos^{B1-}

808 ^{8^{hi}} *IGH* locus showing the distribution of multimapping, unique and total signal. **(C)** Nascent RNA 5' and 3'

809 ends of PRO-cap and PRO-seq respectively and MutPE-seq analysis at the Ramos^{B1-8^{hi}} *IGH* locus.

810 MutPE-seq was performed following infection with AID(JP8Bdel)-expressing lentiviruses for 7 days. **(D)**

811 Details of the mutation spectrum of the B1-8^{hi} variable region expressed from the Ramos *IGH* locus

812 obtained by MutPE-seq. The bar graphs on the right show the percentage of each indicated mutation

813 category (see Fig. 3E and 3G for details). For comparison, the murine B1-8^{hi} mutation spectra from

814 primary murine B cells (middle panel) and murine GCBs (bottom panel), exactly as in Fig. 5B-C, are

815 included.

816

817 **Figure 7: Integrative model for co-transcriptional AID targeting and differential motif mutability. (A)**

818 3D structure of the human RNA Pol II elongation complex visualized with ChimeraX. We superimposed,

819 via their RPB1 subunits, the elongation complex structure (PDB ID 6GMH) onto the transcribing Pol II-

820 DSIF (a complex of SPT5 and SPT4) structure (PDB ID 5OIK). Proteins are shown as surfaces (Pol II,

821 grey; DSIF, salmon; SPT6, brown; PAF complex, yellow) and nucleic acids as cartoons (DNA template

822 strand, blue; DNA non-template strand, cyan; RNA, red). The right panel highlights the trajectory of DNA

823 and RNA buried within Pol II. The exposed non-template strand of the transcription bubble is occluded

824 from interactions with AID due to it being completely covered by SPT5. **(B)** Model for SHM. The cartoon

825 diagram of the Pol II elongation complex reflects the actual structure described in (A) above and in the

826 main text. AID is recruited to the Pol II complex via super-enhancers in the context of a transcriptional hub

827 (step 1) and is then retained in the elongation complex via interactions with elongation factors, SPT5,
828 PAF and SPT6 (step 2). We leave open the possibility that multiple AID molecules can be present in the
829 hub via association with different elongation factors which would serve to increase the local concentration
830 of AID. The next step is the availability of ssDNA (step 3). Importantly, since SPT5 covers the exposed
831 ssDNA, AID has no access to this ssDNA. Sources of ssDNA include upstream negative supercoils (not
832 shown) which may allow AID to target both strands. Although R loops do not appear to form in variable
833 regions, the RNA exosome is necessary to ensure normal distribution of mutations on both strands during
834 SHM which implies that RNA processing is important in SHM. Hence, we propose that transcription
835 termination may be an important source of ssDNA. Upon dissociation of Pol II, the non-template strand is
836 available since the DNA:RNA hybrid would prevent immediate reannealing. The RNA is removed by RNA
837 helicases and RNaseH followed by degradation of the RNA by the RNA exosome, which is known to
838 associate with AID and provide access to the template strand. AID can bind ssDNA in a sequence-
839 independent manner with high affinity and remain bound for several minutes (high on-rate, step 4).
840 Deamination would occur occasionally on exposed WRCH motifs. Importantly, we hypothesize that the
841 probability of deamination is strongly influenced by the sequence context of the WRCH motif and that this
842 may be due to differential off-rates of AID in different sequence contexts. Thus, some WRCH motifs are
843 more frequently deaminated to uridines (U) when they are embedded in a favorable nucleotide context
844 that retains AID on ssDNA for longer periods (step 4, right) than other contexts where deamination is
845 inefficient because of the higher dissociation rate of AID (step 4, left).

846

847 **Figure S1: ChIP-qPCR analysis of the Ramos variable region locale and generation of Ramos**
848 **$E\mu^{-/-}$ cells. (A)** ChIP-qPCR analysis from three independent experiments in $AID^{-/-}$ Ramos cells for the
849 indicated epigenetic marks as well as histone H3. Amplicons (1-7) used are indicated below the locus
850 diagram shown above the graphs. The Neg. amplicon corresponds to a gene desert on chromosome 1
851 and is used as a negative control. **(B)** ChIP-qPCR analysis from three independent experiments in B1-8^{hi}
852 primary splenic cells for the indicated epigenetic marks as well as histone H3. Amplicons (1-6) used are
853 indicated below the locus diagram shown above the graphs. The Neg. (negative region) amplicon

854 corresponds to a gene desert on chromosome 1 and is used as a negative control. **(C)** Strategy to create
855 the $E_{\mu}^{-/-}$ $AID^{-/-}$ Ramos lines. The deleted region (583 bp) corresponds to the peak of accessible
856 chromatin detected by ATAC-seq (top panel). This segment was replaced with a floxed reporter cassette
857 expressing GFP or mCherry using CRISPR. Single clones double-positive for GFP and mCherry
858 expression were isolated followed by excision of the floxed cassette by Cre recombinase. Two clones, c1
859 and c2 were used for all experiments. **(D)** Genotyping PCR analysis to confirm the loss of E_{μ} in $AID^{-/-}$
860 $E_{\mu}^{-/-}$ clones c1 and c2. The location of primers is shown in the diagram above the gel image. **(E)** Surface
861 IgM expression in $E_{\mu}^{-/-}$ $AID^{-/-}$ Ramos clones relative to the parental $AID^{-/-}$ line determined by FACS. **(F)**
862 Analysis of nascent transcription 3' ends (PRO-seq) at the *IGH* locus in $AID^{-/-}$ and $AID^{-/-}$ $E_{\mu}^{-/-}$ Ramos
863 cells (c1, c2). The promoter (Prom.), complementary determining regions (CDR 1-3), intronic enhancer
864 (E_{μ}) and switch μ region (S_{μ}). **(G)** As in F but showing a magnified view of the variable region locale.

865

866 **Figure S2: Generation of Ramos cell lines expressing exogenous variable region.** **(A)** Schematic
867 representation of the workflow (described in Methods) for generating Ramos cells expressing new
868 variable regions using CRISPR-based editing. An IgM-negative line was made by excising the
869 endogenous variable region and replacing it with a unique small guide RNA (sgRNA)-targeting sequence
870 (green). Subsequently, an sgRNA targeting this site is combined with Cas9 and homology repair
871 templates harboring any new variable regions to restore IgM expression. In essence, this system uses the
872 restoration of IgM expression, which can occur only upon correct integration of the new variable regions,
873 as a rapid and sensitive means to identify correctly targeted clones via flow cytometry (Fig. S4B). **(B)**
874 Flow cytometry analysis of Ramos cell clones expressing new human and mouse variable regions used in
875 this study (see scheme in Fig. 2C). Starting from $AID^{-/-}$ cells, the endogenous variable region was deleted
876 ($AID^{-/-}\Delta V$) followed by re-insertion of new variable regions. Shown are three clones of human VH4-59-
877 DH2-JH6 and human VH3-30-DH2-JH6 expressing Ramos cells, and two clones of B1-8^{hi} expressing
878 Ramos cells. **(C)** PRO-seq analysis showing the 3' ends of aligned reads at the human VH4-59-DH2-JH6
879 variable region expressed from the VH4-34 promoter at the human *IGH* in Ramos cells. Tracks of
880 Multimapping, uniquely mapping and total are shown (see Fig. 3A and 3B for detailed description). **(D)**

881 PRO-seq analysis as in C of the human VH3-30-DH2-JH6 variable region expressed from the VH4-34
882 promoter at the human *IGH* in Ramos cells. **(E)** PRO-seq analysis as in C at the murine B1-8^{hi} variable
883 region at the *Igh* locus in mice.

884

885 **Figure S3: Analysis of PRO-cap, PRO-seq and SHM at non-Ig AID target genes in murine GCBs.**

886 **(A)** 5' and 3' ends of nascent RNA (PRO-cap, PRO-seq) and mutation profiles (from Alvarez-Prado et al.
887 (2018) at selected AID target genes as in Fig. 5D-G. Shown are genes where highly mutated residues lie
888 within 150 bp of the transcription initiation site defined by the peak of PRO-cap signals. The WRCH motifs
889 are indicated as red dots. **(B)** As in A but showing genes where highly mutated residues lie near the
890 transcription initiation site defined by the peak of PRO-cap signals.

891

892 **Figure S4: Additional examples of PRO-cap, PRO-seq and SHM at non-Ig AID target genes in**
893 **murine GCBs.** See Fig. 5D for detailed legend.

894

895

896

897

898

899

900

901

902


903 **References**

- 904 1. S. P. Methot, J. M. Di Noia, Molecular Mechanisms of Somatic Hypermutation and Class Switch
905 Recombination. *Adv Immunol* **133**, 37-87 (2017).
- 906 2. J. U. Peled *et al.*, The Biochemistry of Somatic Hypermutation. *Annu Rev Immunol* **26**, 481-511
907 (2008).
- 908 3. V. H. Odegard, D. G. Schatz, Targeting of somatic hypermutation. *Nat Rev Immunol* **6**, 573-583
909 (2006).
- 910 4. G. D. Victora, M. C. Nussenzweig, Germinal centers. *Annu Rev Immunol* **30**, 429-457 (2012).
- 911 5. M. Muramatsu *et al.*, Class switch recombination and hypermutation require activation-induced
912 cytidine deaminase (AID), a potential RNA editing enzyme. *Cell* **102**, 553-563 (2000).
- 913 6. P. Revy *et al.*, Activation-induced cytidine deaminase (AID) deficiency causes the autosomal
914 recessive form of the Hyper-IgM syndrome (HIGM2). *Cell* **102**, 565-575 (2000).
- 915 7. S. K. Petersen-Mahrt, R. S. Harris, M. S. Neuberger, AID mutates E. coli suggesting a DNA
916 deamination mechanism for antibody diversification. *Nature* **418**, 99-103 (2002).
- 917 8. A. R. Ramiro, P. Stavropoulos, M. Jankovic, M. C. Nussenzweig, Transcription enhances AID-
918 mediated cytidine deamination by exposing single-stranded DNA on the nontemplate strand.
919 *Nat Immunol* **4**, 452-456 (2003).
- 920 9. R. Bransteitter, P. Pham, M. D. Scharff, M. F. Goodman, Activation-induced cytidine deaminase
921 deaminates deoxycytidine on single-stranded DNA but requires the action of RNase. *Proc Natl*
922 *Acad Sci U S A* **100**, 4102-4107 (2003).
- 923 10. P. Pham, R. Bransteitter, J. Petruska, M. F. Goodman, Processive AID-catalysed cytosine
924 deamination on single-stranded DNA simulates somatic hypermutation. *Nature* **424**, 103-107
925 (2003).
- 926 11. J. Chaudhuri *et al.*, Transcription-targeted DNA deamination by the AID antibody diversification
927 enzyme. *Nature* **422**, 726-730 (2003).
- 928 12. B. Vaidyanathan, W.-F. Yen, J. Pucella, J. Chaudhuri, AIDing Chromatin and Transcription-
929 Coupled Orchestration of Immunoglobulin Class-Switch Recombination. *Frontiers in Immunology*
930 **5**, (2014).
- 931 13. R. J. Leeman-Neill, J. Lim, U. Basu, The Common Key to Class-Switch Recombination and Somatic
932 Hypermutation: Discovery of AID and Its Role in Antibody Gene Diversification. *The Journal of*
933 *Immunology* **201**, 2527-2529 (2018).
- 934 14. T. Saha, D. Sundaravinayagam, M. Di Virgilio, Charting a DNA Repair Roadmap for
935 Immunoglobulin Class Switch Recombination. *Trends Biochem Sci*, (2020).
- 936 15. J. Chaudhuri, F. W. Alt, Class-switch recombination: interplay of transcription, DNA deamination
937 and DNA repair. *Nat Rev Immunol* **4**, 541-552 (2004).
- 938 16. M. Liu *et al.*, Two levels of protection for the B cell genome during somatic hypermutation.
939 *Nature* **451**, 841-845 (2008).
- 940 17. A. F. Alvarez-Prado *et al.*, A broad atlas of somatic hypermutation allows prediction of
941 activation-induced deaminase targets. *J Exp Med* **215**, 761-771 (2018).
- 942 18. A. Yamane *et al.*, Deep-sequencing identification of the genomic targets of the cytidine
943 deaminase AID and its cofactor RPA in B lymphocytes. *Nat Immunol*, (2010).
- 944 19. H. M. Shen, A. Peters, B. Baron, X. Zhu, U. Storb, Mutation of BCL-6 gene in normal B cells by the
945 process of somatic hypermutation of Ig genes. *Science* **280**, 1750-1752 (1998).
- 946 20. L. Pasqualucci *et al.*, Hypermutation of multiple proto-oncogenes in B-cell diffuse large-cell
947 lymphomas. *Nature* **412**, 341-346 (2001).

- 948 21. S. P. Methot *et al.*, A licensing step links AID to transcription elongation for mutagenesis in B
949 cells. *Nat Commun* **9**, 1248 (2018).
- 950 22. R. Pavri *et al.*, Activation-induced cytidine deaminase targets DNA at sites of RNA polymerase II
951 stalling by interaction with Spt5. *Cell* **143**, 122-133 (2010).
- 952 23. K. L. Willmann *et al.*, A role for the RNA pol II-associated PAF complex in AID-induced immune
953 diversification. *J Exp Med* **209**, 2099-2111 (2012).
- 954 24. N. A. Begum, A. Stanlie, M. Nakata, H. Akiyama, T. Honjo, The histone chaperone Spt6 is
955 required for activation-induced cytidine deaminase target determination through H3K4me3
956 regulation. *J Biol Chem* **287**, 32415-32429 (2012).
- 957 25. Y. Nambu *et al.*, Transcription-coupled events associating with immunoglobulin switch region
958 chromatin. *Science* **302**, 2137-2140 (2003).
- 959 26. U. Basu *et al.*, The RNA exosome targets the AID cytidine deaminase to both strands of
960 transcribed duplex DNA substrates. *Cell* **144**, 353-363 (2011).
- 961 27. U. Nowak, A. J. Matthews, S. Zheng, J. Chaudhuri, The splicing regulator PTBP2 interacts with the
962 cytidine deaminase AID and promotes binding of AID to switch-region DNA. *Nat Immunol* **12**,
963 160-166 (2011).
- 964 28. Z. Duan *et al.*, Role of Dot1L and H3K79 methylation in regulating somatic hypermutation of
965 immunoglobulin genes. *Proc Natl Acad Sci U S A* **118**, (2021).
- 966 29. G. Yu *et al.*, The role of HIRA-dependent H3.3 deposition and its modifications in the somatic
967 hypermutation of immunoglobulin variable regions. *Proc Natl Acad Sci U S A* **118**, (2021).
- 968 30. M. Aida, N. Hamad, A. Stanlie, N. A. Begum, T. Honjo, Accumulation of the FACT complex, as well
969 as histone H3.3, serves as a target marker for somatic hypermutation. *Proc Natl Acad Sci U S A*
970 **110**, 7784-7789 (2013).
- 971 31. B. P. Jeevan-Raj *et al.*, Epigenetic tethering of AID to the donor switch region during
972 immunoglobulin class switch recombination. *J Exp Med* **208**, 1649-1660 (2011).
- 973 32. A. Stanlie, M. Aida, M. Muramatsu, T. Honjo, N. A. Begum, Histone3 lysine4 trimethylation
974 regulated by the facilitates chromatin transcription complex is critical for DNA cleavage in class
975 switch recombination. *Proc Natl Acad Sci U S A* **107**, 22190-22195 (2010).
- 976 33. S. P. Bradley, D. A. Kaminski, A. H. F. M. Peters, T. Jenuwein, J. Stavnezer, The Histone
977 Methyltransferase Suv39h1 Increases Class Switch Recombination Specifically to IgA. *The Journal*
978 *of Immunology* **177**, 1179-1188 (2006).
- 979 34. F. L. Kuang, Z. Luo, M. D. Scharff, H3 trimethyl K9 and H3 acetyl K9 chromatin modifications are
980 associated with class switch recombination. *Proceedings of the National Academy of Sciences*
981 **106**, 5288-5293 (2009).
- 982 35. J. A. Daniel *et al.*, PTIP promotes chromatin changes critical for immunoglobulin class switch
983 recombination. *Science* **329**, 917-923 (2010).
- 984 36. B. Vaidyanathan, J. Chaudhuri, Epigenetic Codes Programing Class Switch Recombination. *Front*
985 *Immunol* **6**, 405 (2015).
- 986 37. J. M. Buerstedde, J. Alinikula, H. Arakawa, J. J. McDonald, D. G. Schatz, Targeting of somatic
987 hypermutation by immunoglobulin enhancer and enhancer-like sequences. *PLoS Biol* **12**,
988 e1001831 (2014).
- 989 38. F. Senigl *et al.*, Topologically Associated Domains Delineate Susceptibility to Somatic
990 Hypermutation. *Cell Rep* **29**, 3902-3915 e3908 (2019).
- 991 39. R. K. Dinesh *et al.*, Transcription factor binding at Ig enhancers is linked to somatic
992 hypermutation targeting. *Eur J Immunol* **50**, 380-395 (2020).
- 993 40. J. Qian *et al.*, B cell super-enhancers and regulatory clusters recruit AID tumorigenic activity. *Cell*
994 **159**, 1524-1537 (2014).

- 995 41. P. Rouaud *et al.*, The IgH 3' regulatory region controls somatic hypermutation in germinal center
996 B cells. *J Exp Med* **210**, 1501-1507 (2013).
- 997 42. S. Peron *et al.*, AID-driven deletion causes immunoglobulin heavy chain locus suicide
998 recombination in B cells. *Science* **336**, 931-934 (2012).
- 999 43. J. M. Di Noia, M. S. Neuberger, Molecular mechanisms of antibody somatic hypermutation.
1000 *Annu Rev Biochem* **76**, 1-22 (2007).
- 1001 44. R. W. Maul, P. J. Gearhart, AID and somatic hypermutation. *Adv Immunol* **105**, 159-191 (2010).
- 1002 45. R. Pavri, M. C. Nussenzweig, AID targeting in antibody diversity. *Adv Immunol* **110**, 1-26 (2011).
- 1003 46. A. G. Betz, C. Rada, R. Pannell, C. Milstein, M. S. Neuberger, Passenger transgenes reveal
1004 intrinsic specificity of the antibody hypermutation mechanism: clustering, polarity, and specific
1005 hot spots. *Proc Natl Acad Sci U S A* **90**, 2385-2388 (1993).
- 1006 47. J. K. Hwang *et al.*, Sequence intrinsic somatic mutation mechanisms contribute to affinity
1007 maturation of VRC01-class HIV-1 broadly neutralizing antibodies. *Proc Natl Acad Sci U S A* **114**,
1008 8614-8619 (2017).
- 1009 48. L. S. Yeap *et al.*, Sequence-Intrinsic Mechanisms that Target AID Mutational Outcomes on
1010 Antibody Genes. *Cell* **163**, 1124-1137 (2015).
- 1011 49. J. Q. Zhou, S. H. Kleinstein, Position-Dependent Differential Targeting of Somatic Hypermutation.
1012 *The Journal of Immunology*, j12000496 (2020).
- 1013 50. N. Spisak, A. M. Walczak, T. Mora, Learning the heterogeneous hypermutation landscape of
1014 immunoglobulins from high-throughput repertoire data. *Nucleic Acids Research* **48**, 10702-
1015 10712 (2020).
- 1016 51. C. Tang, A. Krantsevich, T. MacCarthy, Deep learning model of somatic hypermutation reveals
1017 importance of sequence context beyond hotspot targeting. *iScience* **25**, 103668 (2022).
- 1018 52. J. Q. Zhou, S. H. Kleinstein, Position-Dependent Differential Targeting of Somatic Hypermutation.
1019 *J Immunol* **205**, 3468-3479 (2020).
- 1020 53. F. L. Meng *et al.*, Convergent transcription at intragenic super-enhancers targets AID-initiated
1021 genomic instability. *Cell* **159**, 1538-1548 (2014).
- 1022 54. D. Rajagopal *et al.*, Immunoglobulin switch mu sequence causes RNA polymerase II
1023 accumulation and reduces dA hypermutation. *J Exp Med* **206**, 1237-1244 (2009).
- 1024 55. L. Wang, R. Wuerffel, S. Feldman, A. A. Khamlichi, A. L. Kenter, S region sequence, RNA
1025 polymerase II, and histone modifications create chromatin accessibility during class switch
1026 recombination. *J Exp Med* **206**, 1817-1830 (2009).
- 1027 56. R. W. Maul *et al.*, Spt5 accumulation at variable genes distinguishes somatic hypermutation in
1028 germinal center B cells from ex vivo-activated cells. *J Exp Med* **211**, 2297-2306 (2014).
- 1029 57. A. Tarsalainen *et al.*, Ig Enhancers Increase RNA Polymerase II Stalling at Somatic Hypermutation
1030 Target Sequences. *J Immunol* **208**, 143-154 (2022).
- 1031 58. B. J. Taylor, Y. L. Wu, C. Rada, Active RNAP pre-initiation sites are highly mutated by cytidine
1032 deaminases in yeast, with AID targeting small RNA genes. *Elife* **3**, e03553 (2014).
- 1033 59. X. Wang, M. Fan, S. Kalis, L. Wei, M. D. Scharff, A source of the single-stranded DNA substrate
1034 for activation-induced deaminase during somatic hypermutation. *Nature Communications* **5**,
1035 4137 (2014).
- 1036 60. P. Kodgire, P. Mukkavar, S. Ratnam, T. E. Martin, U. Storb, Changes in RNA polymerase II
1037 progression influence somatic hypermutation of Ig-related genes by AID. *J Exp Med* **210**, 1481-
1038 1492 (2013).
- 1039 61. D. Ronai *et al.*, Detection of chromatin-associated single-stranded DNA in regions targeted for
1040 somatic hypermutation. *J Exp Med* **204**, 181-190 (2007).
- 1041 62. H. Kwak, N. J. Fuda, L. J. Core, J. T. Lis, Precise maps of RNA polymerase reveal how promoters
1042 direct initiation and pausing. *Science* **339**, 950-953 (2013).

- 1043 63. J. E. Sale, M. S. Neuberger, TdT-accessible breaks are scattered over the immunoglobulin V
1044 domain in a constitutively hypermutating B cell line. *Immunity* **9**, 859-869 (1998).
- 1045 64. E. Pinaud *et al.*, The IgH locus 3' regulatory region: pulling the strings from behind. *Adv Immunol*
1046 **110**, 27-70 (2011).
- 1047 65. T. A. Shih, M. Roederer, M. C. Nussenzweig, Role of antigen receptor affinity in T cell-
1048 independent antibody responses in vivo. *Nat Immunol* **3**, 399-406 (2002).
- 1049 66. S. Ito *et al.*, Activation-induced cytidine deaminase shuttles between nucleus and cytoplasm like
1050 apolipoprotein B mRNA editing catalytic polypeptide 1. *Proc Natl Acad Sci U S A* **101**, 1975-1980
1051 (2004).
- 1052 67. T. Perlot, F. W. Alt, C. H. Bassing, H. Suh, E. Pinaud, Elucidation of IgH intronic enhancer
1053 functions via germ-line deletion. *Proc Natl Acad Sci U S A* **102**, 14362-14367 (2005).
- 1054 68. F. Li, Y. Yan, J. Pieretti, D. A. Feldman, L. A. Eckhardt, Comparison of identical and functional Igh
1055 alleles reveals a nonessential role for Emu in somatic hypermutation and class-switch
1056 recombination. *J Immunol* **185**, 6049-6057 (2010).
- 1057 69. Y. Yamaguchi *et al.*, NELF, a multisubunit complex containing RD, cooperates with DSIF to
1058 repress RNA polymerase II elongation. *Cell* **97**, 41-51 (1999).
- 1059 70. S. M. Vos, L. Farnung, H. Urlaub, P. Cramer, Structure of paused transcription complex Pol II-
1060 DSIF-NELF. *Nature* **560**, 601-606 (2018).
- 1061 71. S. M. Vos *et al.*, Structure of activated transcription complex Pol II-DSIF-PAF-SPT6. *Nature* **560**,
1062 607-612 (2018).
- 1063 72. C. Bernecky, J. M. Plitzko, P. Cramer, Structure of a transcribing RNA polymerase II-DSIF complex
1064 reveals a multidentate DNA-RNA clamp. *Nature Structural & Molecular Biology* **24**, 809-815
1065 (2017).
- 1066 73. D. F. Robbiani *et al.*, Plasmodium Infection Promotes Genomic Instability and AID-Dependent B
1067 Cell Lymphoma. *Cell* **162**, 727-737 (2015).
- 1068 74. H. Mouquet *et al.*, Complex-type N-glycan recognition by potent broadly neutralizing HIV
1069 antibodies. *Proc Natl Acad Sci U S A* **109**, E3268-3277 (2012).
- 1070 75. V. Lorin *et al.*, Epitope convergence of broadly HIV-1 neutralizing IgA and IgG antibody lineages
1071 in a viremic controller. *J Exp Med* **219**, (2022).
- 1072 76. F. Matsuda *et al.*, The complete nucleotide sequence of the human immunoglobulin heavy chain
1073 variable region locus. *J Exp Med* **188**, 2151-2162 (1998).
- 1074 77. T. A. Shih, E. Meffre, M. Roederer, M. C. Nussenzweig, Role of BCR affinity in T cell dependent
1075 antibody responses in vivo. *Nat Immunol* **3**, 570-575 (2002).
- 1076 78. C. Rada, J. M. Di Noia, M. S. Neuberger, Mismatch recognition and uracil excision provide
1077 complementary paths to both Ig switching and the A/T-focused phase of somatic mutation. *Mol*
1078 *Cell* **16**, 163-171 (2004).
- 1079 79. S. Longerich, A. Tanaka, G. Bozek, D. Nicolae, U. Storb, The very 5' end and the constant region
1080 of Ig genes are spared from somatic mutation because AID does not access these regions. *J Exp*
1081 *Med* **202**, 1443-1454 (2005).
- 1082 80. K. Xue, C. Rada, M. S. Neuberger, The in vivo pattern of AID targeting to immunoglobulin switch
1083 regions deduced from mutation spectra in msh2^{-/-} ung^{-/-} mice. *J Exp Med* **203**, 2085-2094
1084 (2006).
- 1085 81. K. M. Dorigi *et al.*, Mll3 and Mll4 Facilitate Enhancer RNA Synthesis and Transcription from
1086 Promoters Independently of H3K4 Monomethylation. *Mol Cell* **66**, 568-576 e564 (2017).
- 1087 82. J. Kim *et al.*, Polycomb- and Methylation-Independent Roles of EZH2 as a Transcription
1088 Activator. *Cell Rep* **25**, 2808-2820.e2804 (2018).
- 1089 83. T. Hoshii *et al.*, A Non-catalytic Function of SETD1A Regulates Cyclin K and the DNA Damage
1090 Response. *Cell* **172**, 1007-1021.e1017 (2018).

- 1091 84. C. Wang *et al.*, UTX regulates mesoderm differentiation of embryonic stem cells independent of
1092 H3K27 demethylase activity. *Proc Natl Acad Sci U S A* **109**, 15324-15329 (2012).
- 1093 85. Y. Aubert, S. Egolf, B. C. Capell, The Unexpected Noncatalytic Roles of Histone Modifiers in
1094 Development and Disease. *Trends Genet* **35**, 645-657 (2019).
- 1095 86. X. Zhang *et al.*, Fundamental roles of chromatin loop extrusion in antibody class switching.
1096 *Nature* **575**, 385-389 (2019).
- 1097 87. S. Chong *et al.*, Imaging dynamic and selective low-complexity domain interactions that control
1098 gene transcription. *Science* **361**, (2018).
- 1099 88. X. Darzacq, R. Tjian, Weak multivalent biomolecular interactions: a strength versus numbers tug
1100 of war with implications for phase partitioning. *RNA* **28**, 48-51 (2022).
- 1101 89. J. P. Karr, J. J. Ferrie, R. Tjian, X. Darzacq, The transcription factor activity gradient (TAG) model:
1102 contemplating a contact-independent mechanism for enhancer-promoter communication.
1103 *Genes Dev* **36**, 7-16 (2022).
- 1104 90. B. Lim, M. S. Levine, Enhancer-promoter communication: hubs or loops? *Current Opinion in*
1105 *Genetics & Development* **67**, 5-9 (2021).
- 1106 91. J. Y. Parsa *et al.*, Negative supercoiling creates single-stranded patches of DNA that are
1107 substrates for AID-mediated mutagenesis. *PLoS Genet* **8**, e1002518 (2012).
- 1108 92. H. M. Shen, U. Storb, Activation-induced cytidine deaminase (AID) can target both DNA strands
1109 when the DNA is supercoiled. *Proc Natl Acad Sci U S A* **101**, 12997-13002 (2004).
- 1110 93. B. Laffleur *et al.*, Noncoding RNA processing by DIS3 regulates chromosomal architecture and
1111 somatic hypermutation in B cells. *Nature Genetics* **53**, 230-242 (2021).
- 1112 94. J. Lim *et al.*, Nuclear Proximity of Mtr4 to RNA Exosome Restricts DNA Mutational Asymmetry.
1113 *Cell* **169**, 523-537 e515 (2017).
- 1114 95. D. Kazadi *et al.*, Effects of senataxin and RNA exosome on B-cell chromosomal integrity. *Heliyon*
1115 **6**, e03442 (2020).
- 1116 96. L. Nair, H. Chung, U. Basu, Regulation of long non-coding RNAs and genome dynamics by the
1117 RNA surveillance machinery. *Nat Rev Mol Cell Biol* **21**, 123-136 (2020).
- 1118 97. X. Wang, M. Fan, S. Kalis, L. Wei, M. D. Scharff, A source of the single-stranded DNA substrate
1119 for activation-induced deaminase during somatic hypermutation. *Nat Commun* **5**, 4137 (2014).
- 1120 98. J. Sun *et al.*, E3-ubiquitin ligase Nedd4 determines the fate of AID-associated RNA polymerase II
1121 in B cells. *Genes Dev* **27**, 1821-1833 (2013).
- 1122 99. S. Lykke-Andersen *et al.*, Integrator is a genome-wide attenuator of non-productive
1123 transcription. *Mol Cell* **81**, 514-529.e516 (2021).
- 1124 100. M. Larijani *et al.*, AID associates with single-stranded DNA with high affinity and a long complex
1125 half-life in a sequence-independent manner. *Mol Cell Biol* **27**, 20-30 (2007).
- 1126 101. P. Pham, P. Calabrese, S. J. Park, M. F. Goodman, Analysis of a Single-stranded DNA-scanning
1127 Process in Which Activation-induced Deoxycytidine Deaminase (AID) Deaminates C to U
1128 Haphazardly and Inefficiently to Ensure Mutational Diversity* . *Journal of Biological Chemistry*
1129 **286**, 24931-24942 (2011).
- 1130 102. J. Fitz *et al.*, Spt5-mediated enhancer transcription directly couples enhancer activation with
1131 physical promoter interaction. *Nat Genet* **52**, 505-515 (2020).
- 1132 103. D. B. Mahat *et al.*, Base-pair-resolution genome-wide mapping of active RNA polymerases using
1133 precision nuclear run-on (PRO-seq). *Nat Protoc* **11**, 1455-1476 (2016).
- 1134 104. J. Fitz, T. Neumann, R. Pavri, Regulation of RNA polymerase II processivity by Spt5 is restricted to
1135 a narrow window during elongation. *EMBO J* **37**, (2018).
- 1136 105. M. Martin, Cutadapt removes adapter sequences from high-throughput sequencing reads. *2011*
1137 **17**, 3 (2011).

- 1138 106. B. Langmead, C. Trapnell, M. Pop, S. L. Salzberg, Ultrafast and memory-efficient alignment of
1139 short DNA sequences to the human genome. *Genome Biol* **10**, R25 (2009).
- 1140 107. A. Larkin *et al.*, FlyBase: updates to the *Drosophila melanogaster* knowledge base. *Nucleic Acids*
1141 *Research* **49**, D899-D907 (2020).
- 1142 108. T. Smith, A. Heger, I. Sudbery, UMI-tools: modeling sequencing errors in Unique Molecular
1143 Identifiers to improve quantification accuracy. *Genome Res* **27**, 491-499 (2017).
- 1144 109. A. R. Quinlan, I. M. Hall, BEDTools: a flexible suite of utilities for comparing genomic features.
1145 *Bioinformatics* **26**, 841-842 (2010).
- 1146 110. A. M. Bolger, M. Lohse, B. Usadel, Trimmomatic: a flexible trimmer for Illumina sequence data.
1147 *Bioinformatics* **30**, 2114-2120 (2014).
- 1148 111. T. Magoč, S. L. Salzberg, FLASH: fast length adjustment of short reads to improve genome
1149 assemblies. *Bioinformatics* **27**, 2957-2963 (2011).
- 1150 112. B. Langmead, S. L. Salzberg, Fast gapped-read alignment with Bowtie 2. *Nature Methods* **9**, 357-
1151 359 (2012).
- 1152 113. H. Li *et al.*, The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-2079
1153 (2009).
- 1154

Figure 1

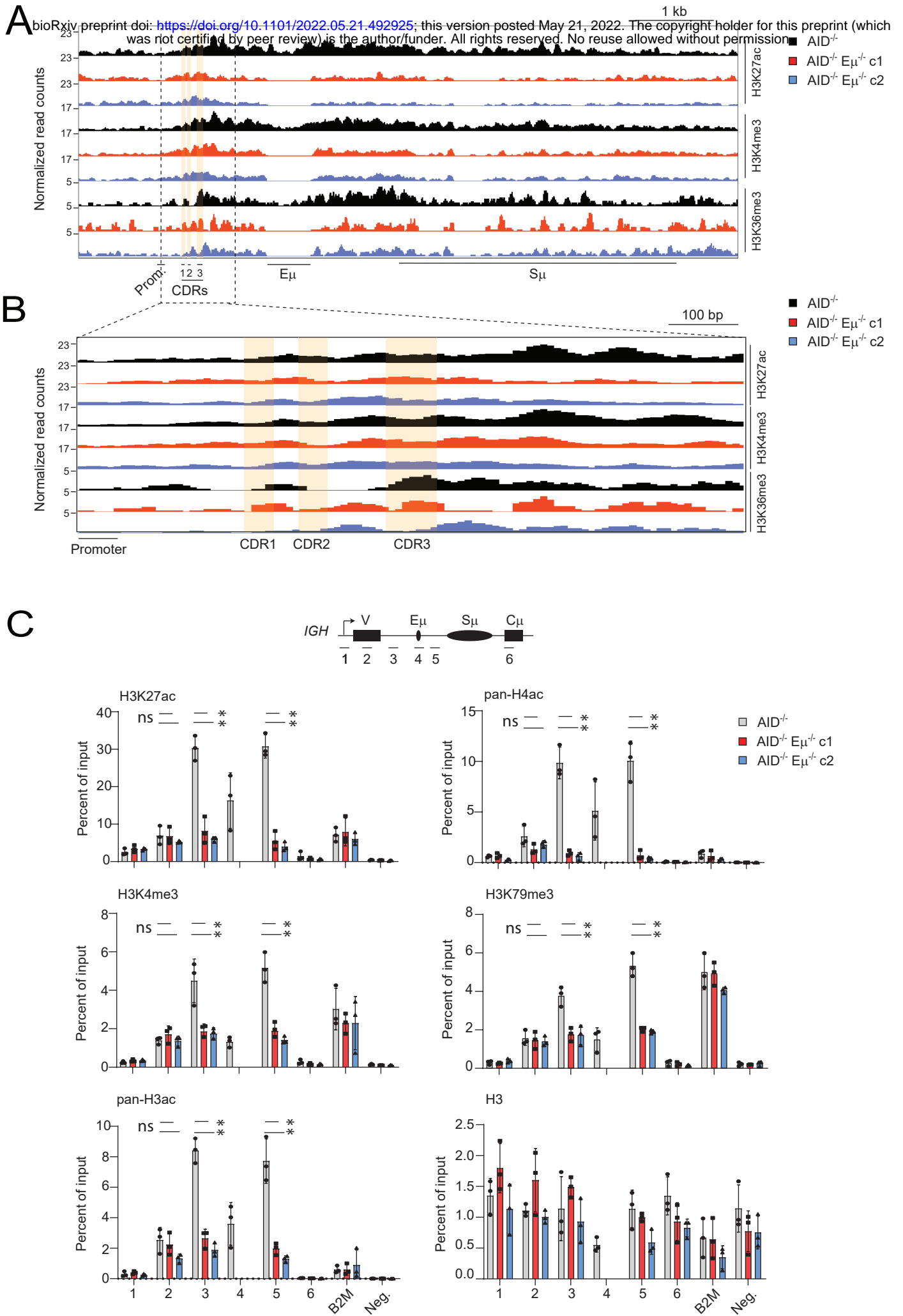


Figure 2

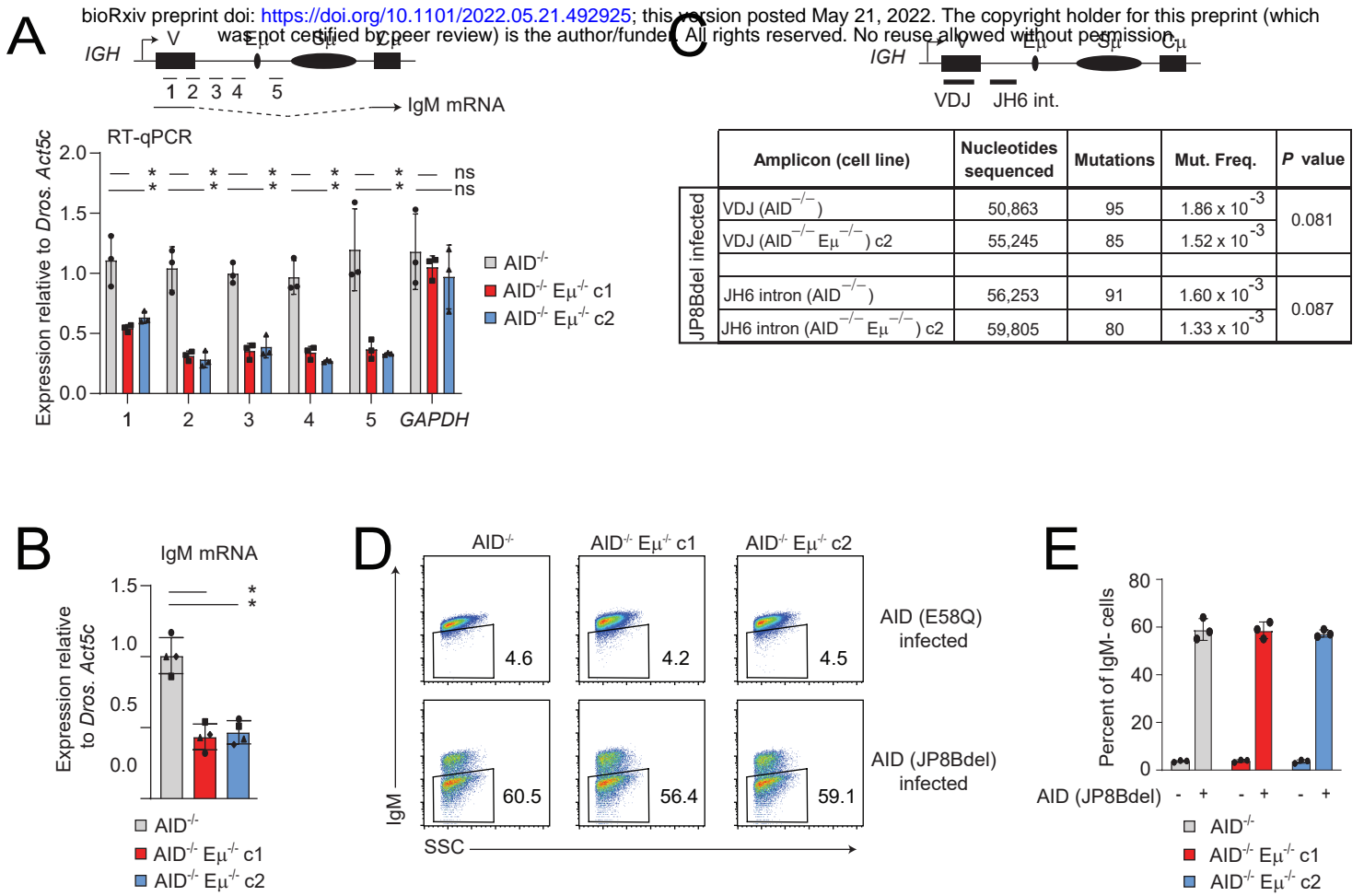


Figure 3

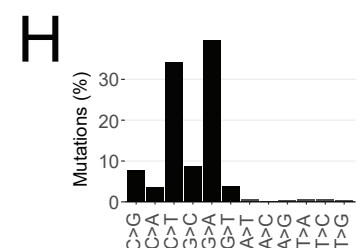
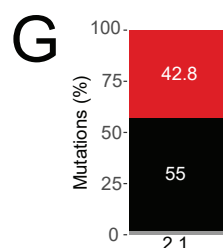
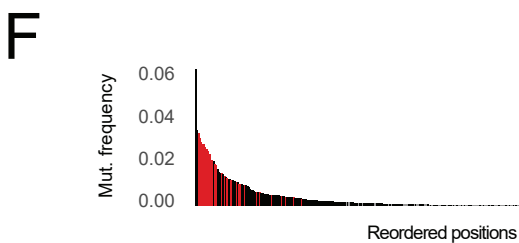
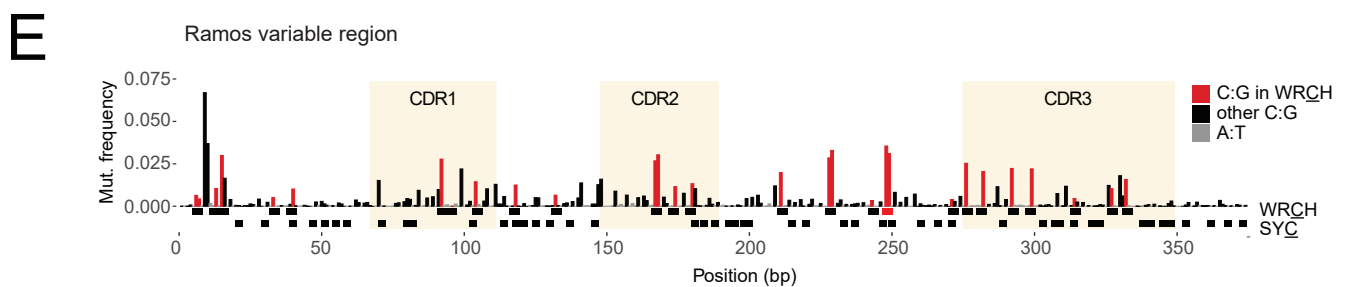
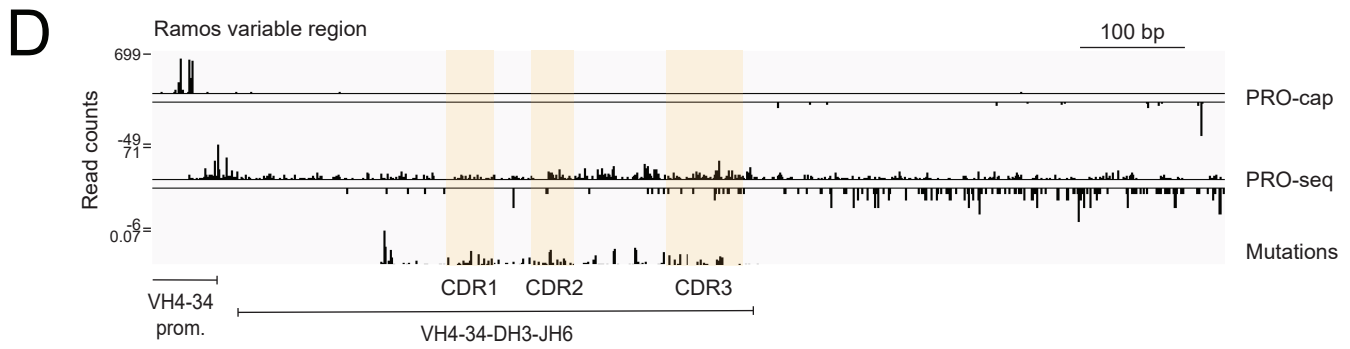
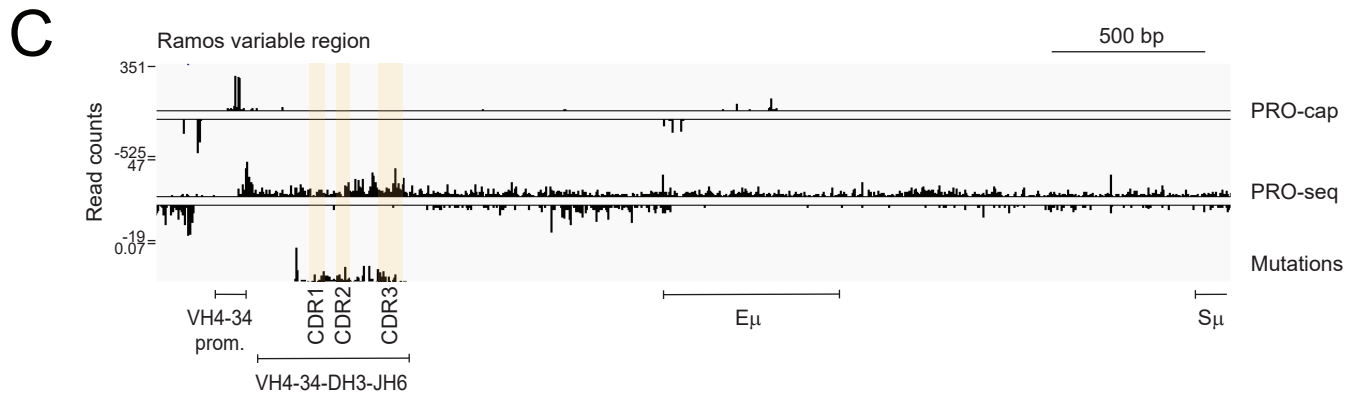
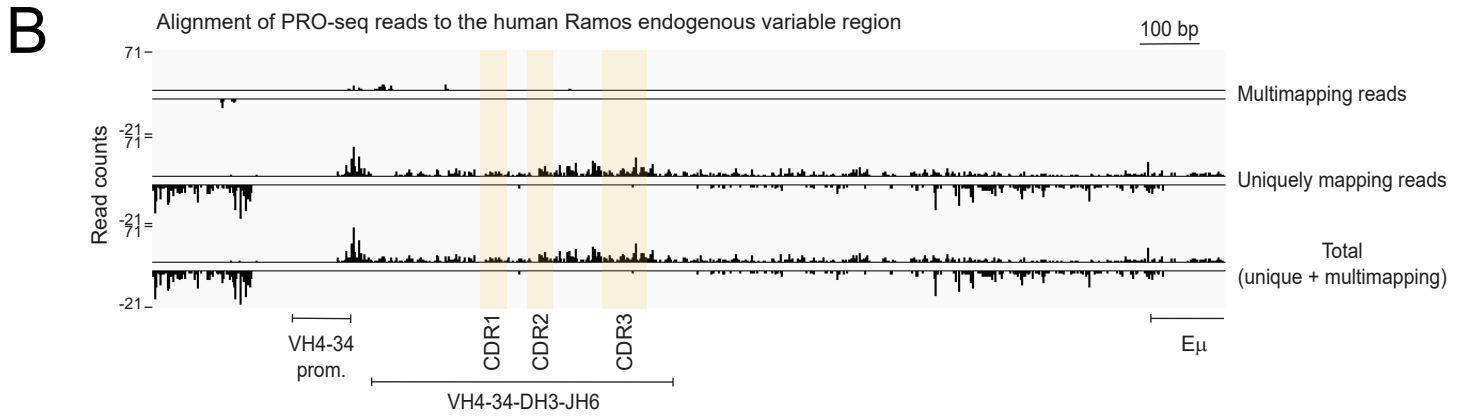
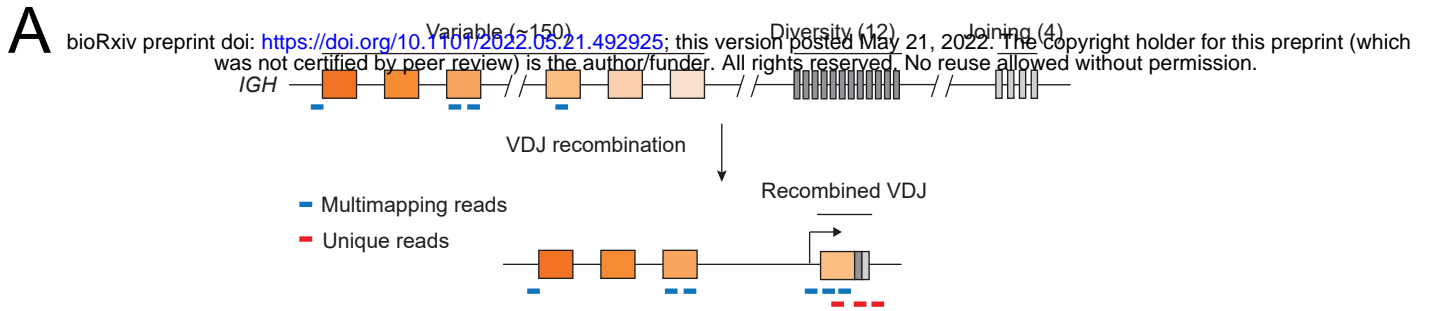
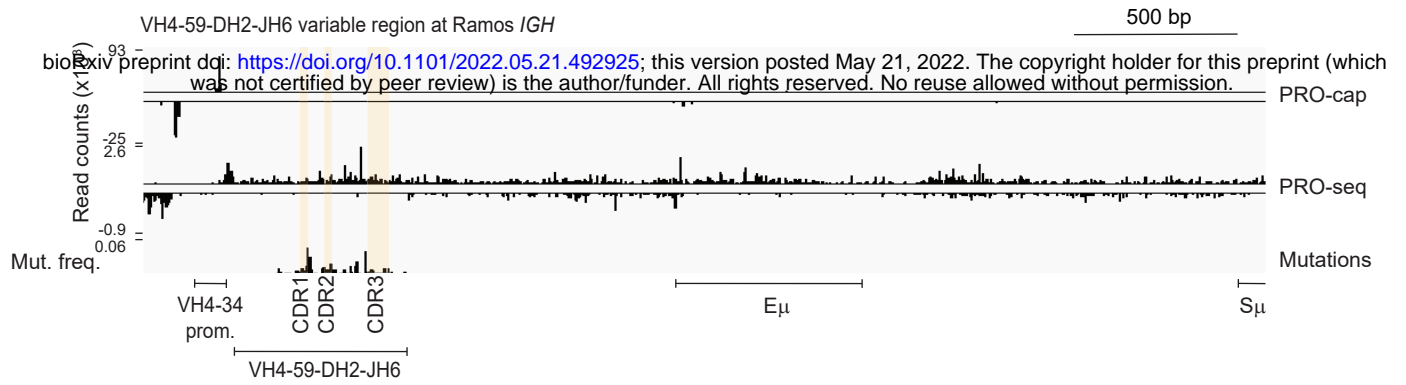
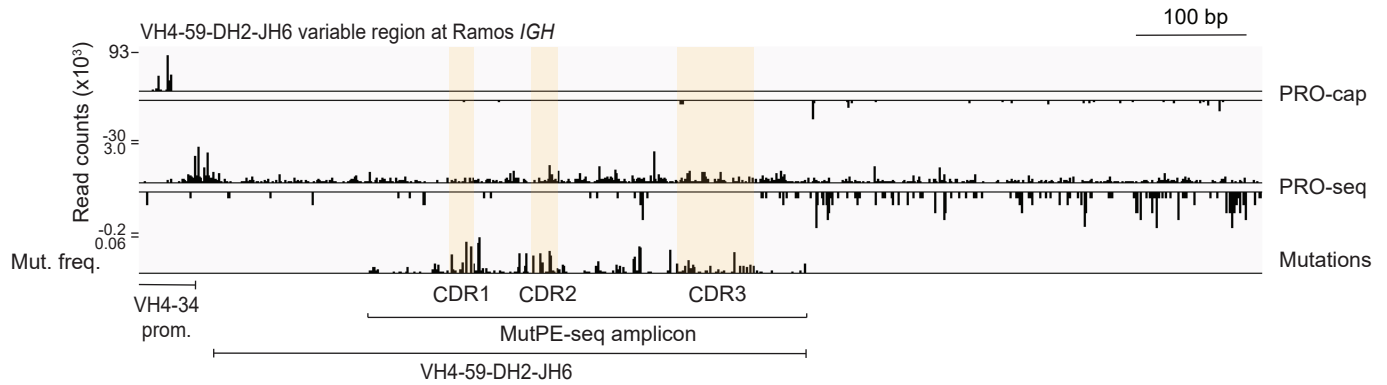


Figure 4

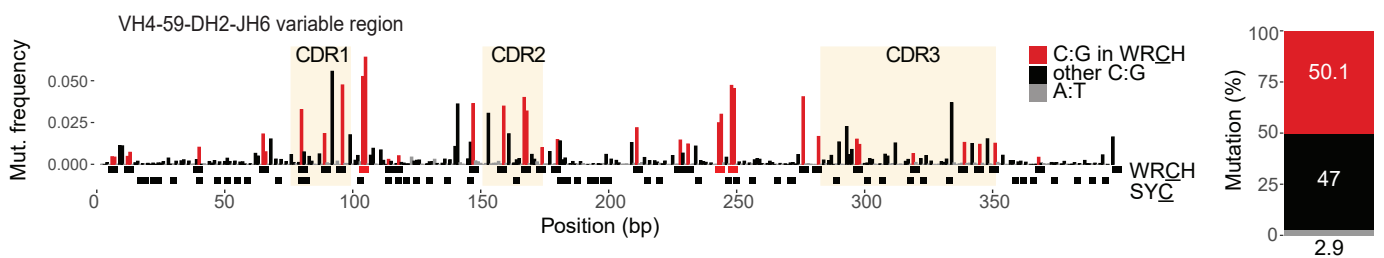
A



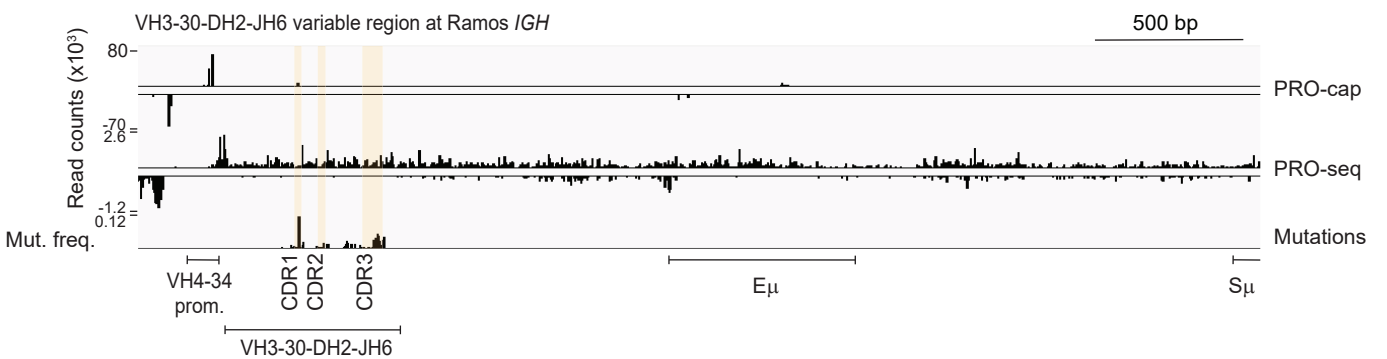
B



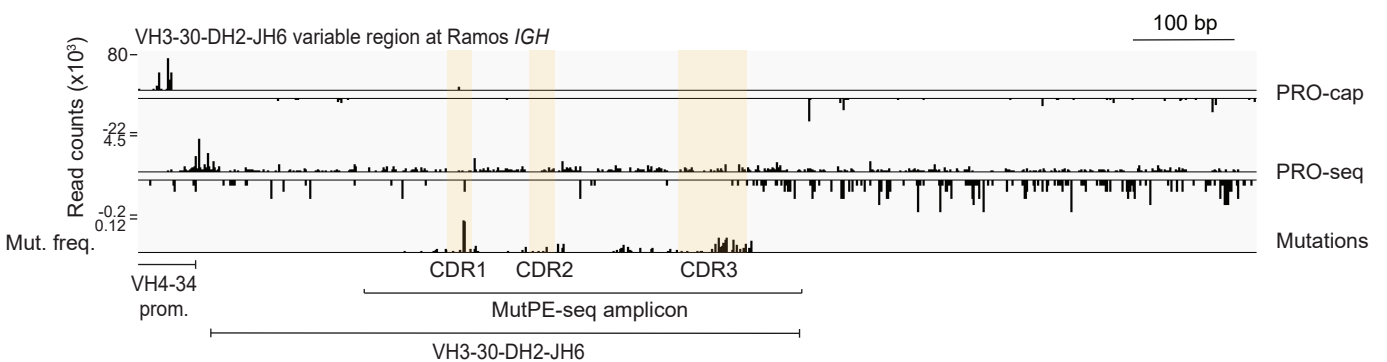
C



D



E



F

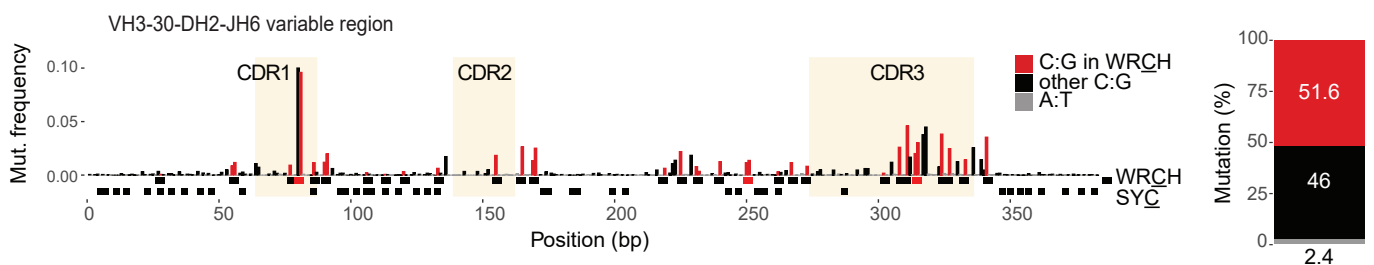


Figure 5

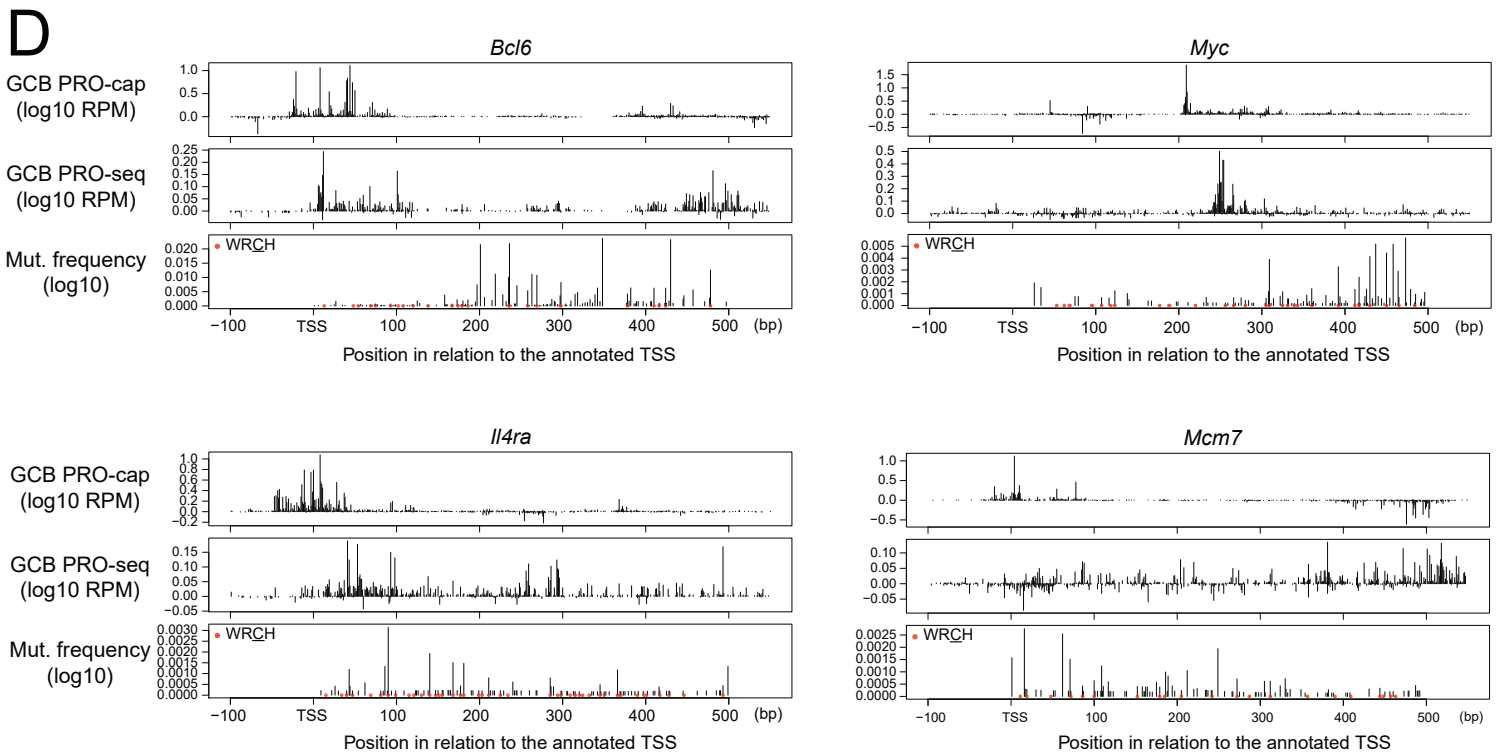
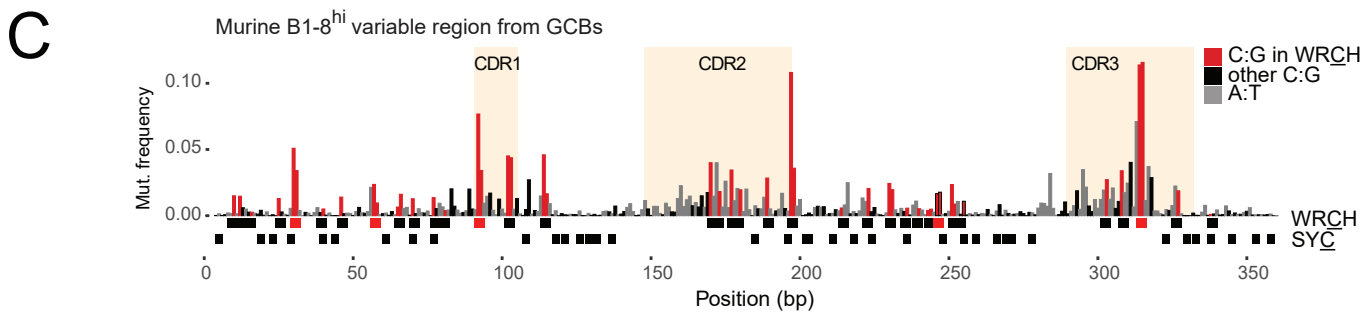
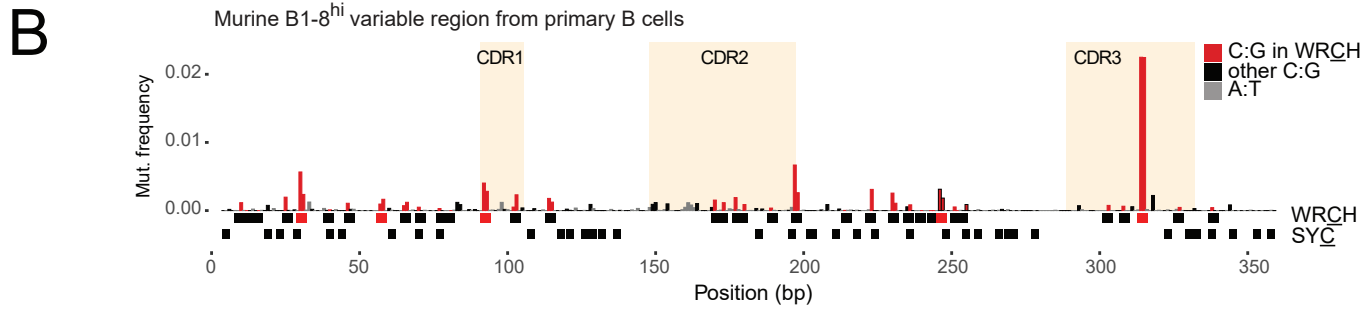
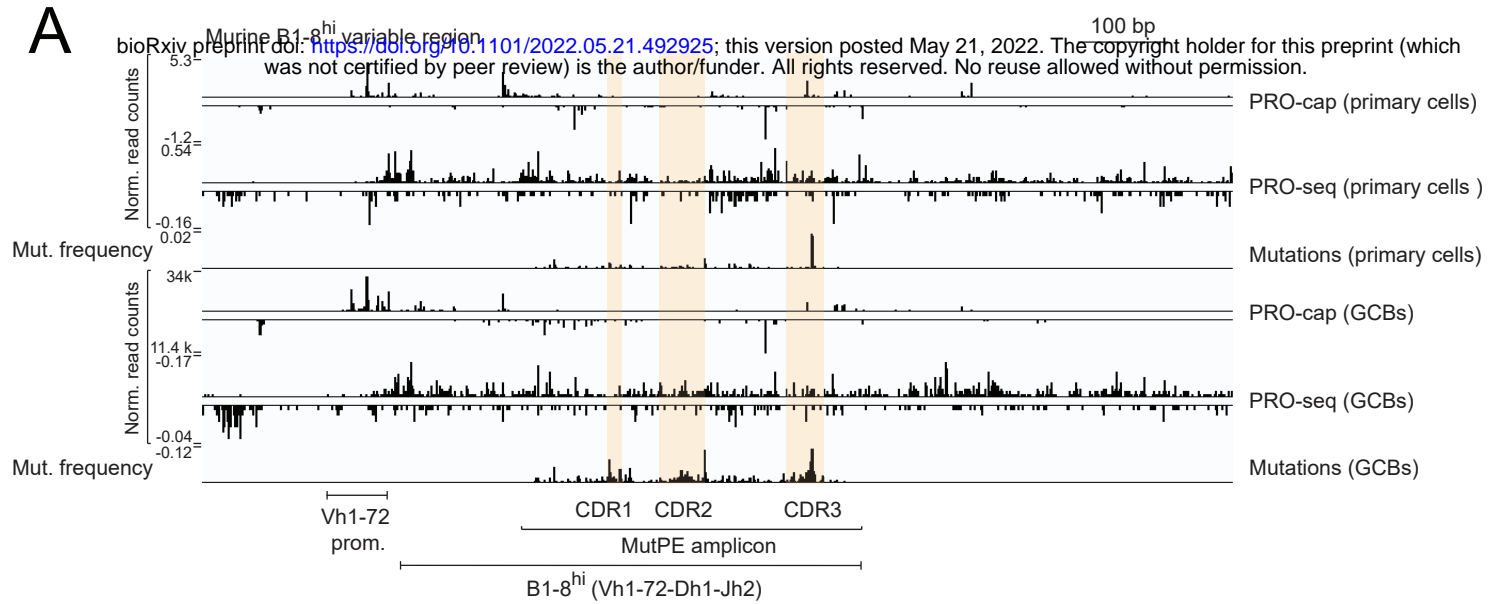


Figure 6

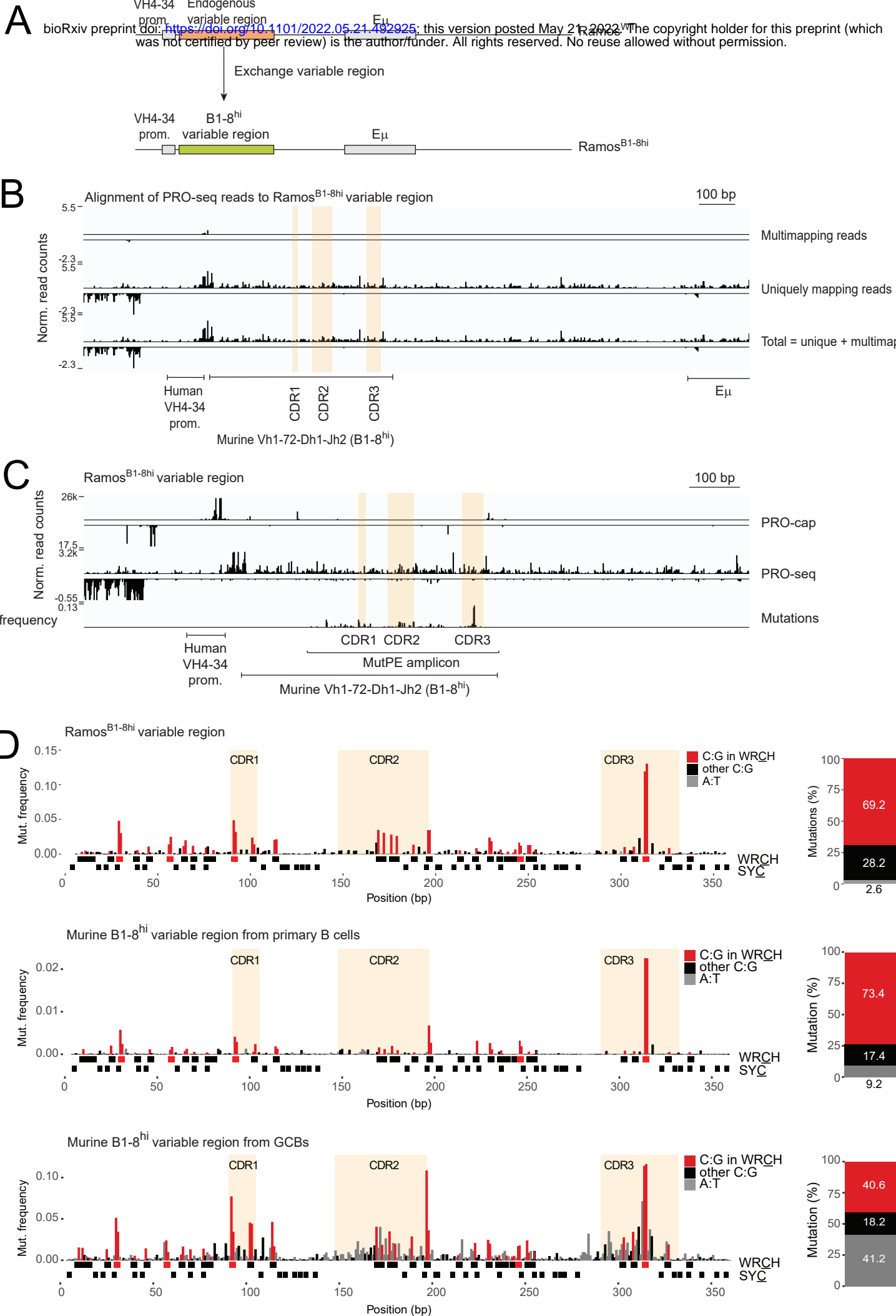
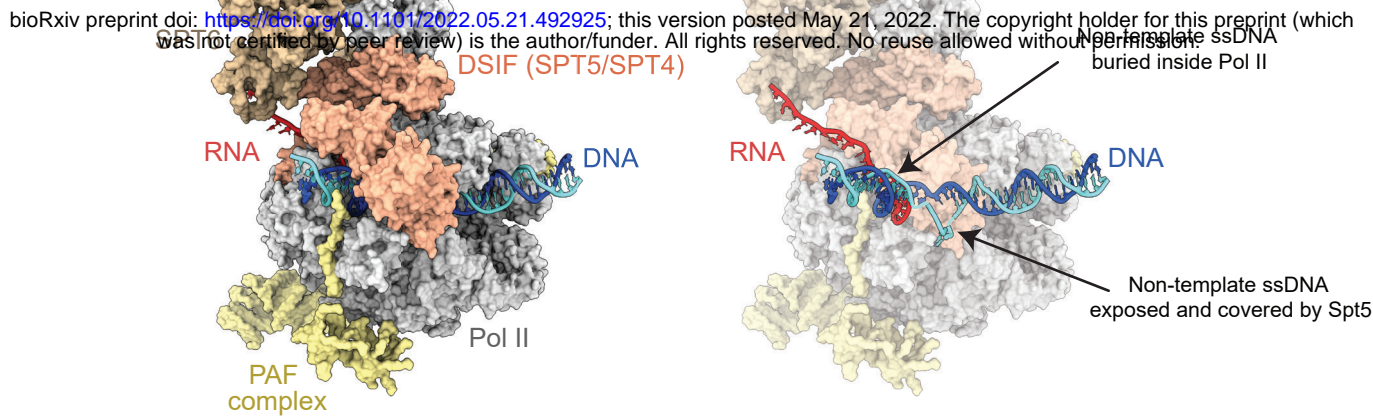


Fig. 7

A



B

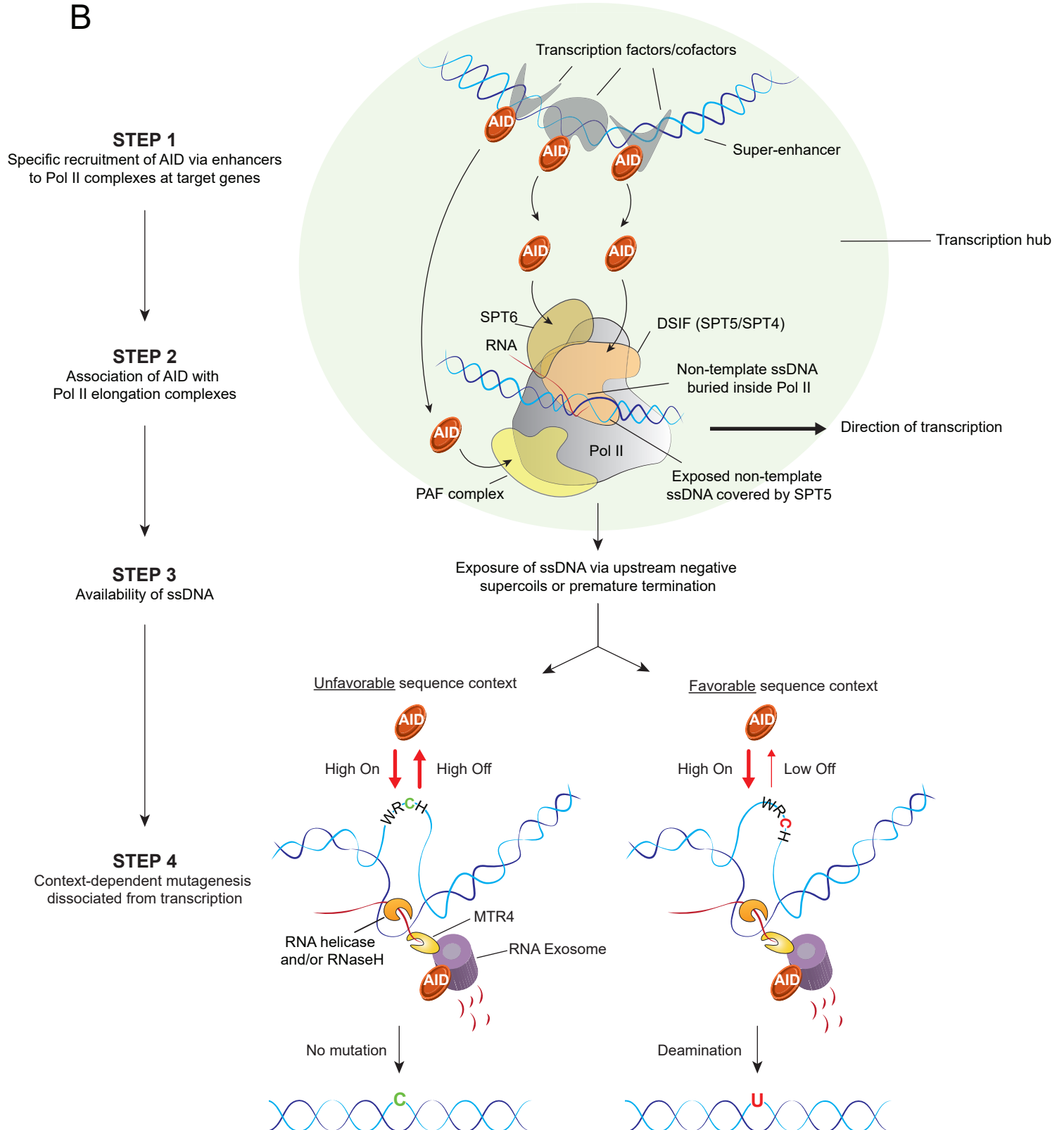


Figure S1

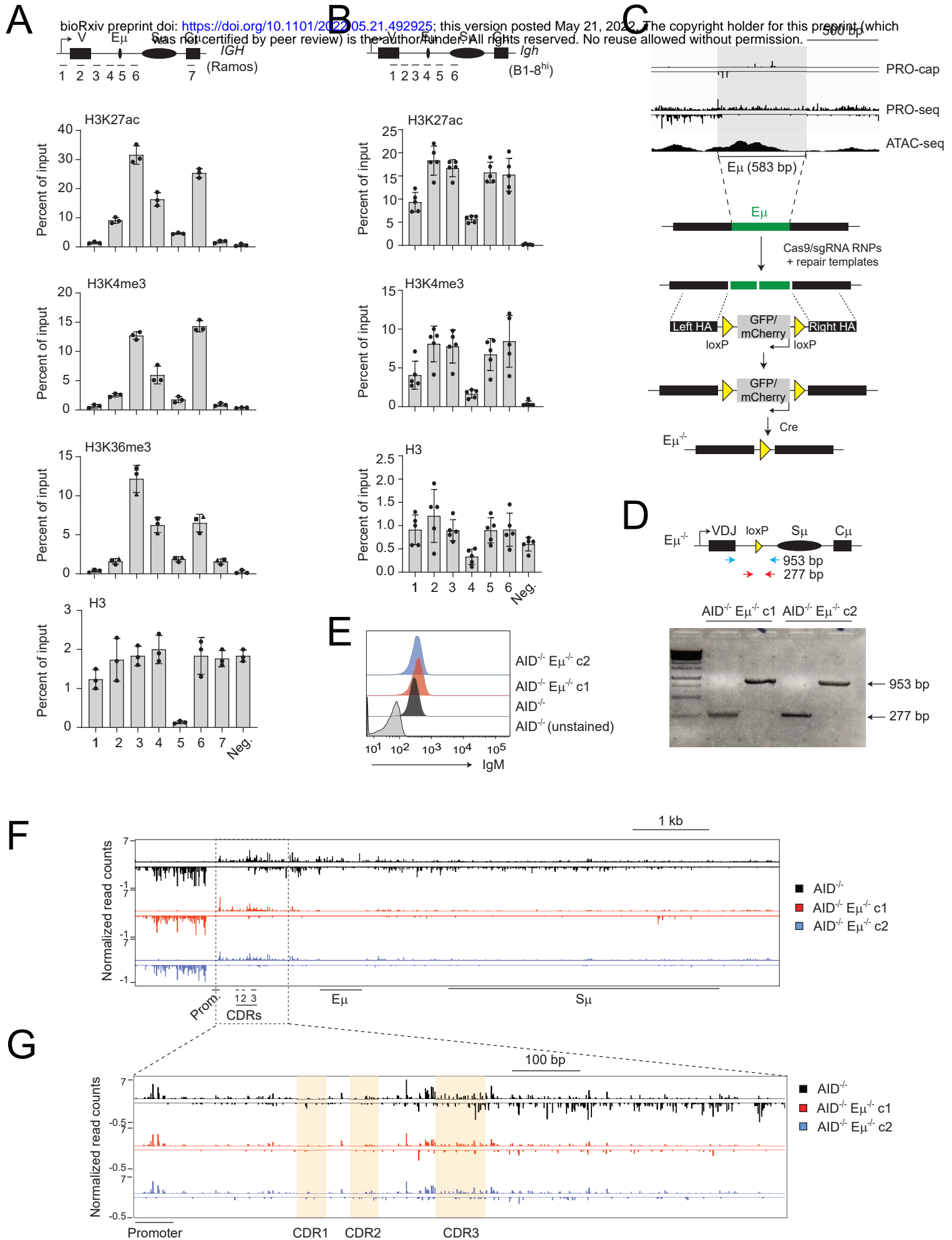
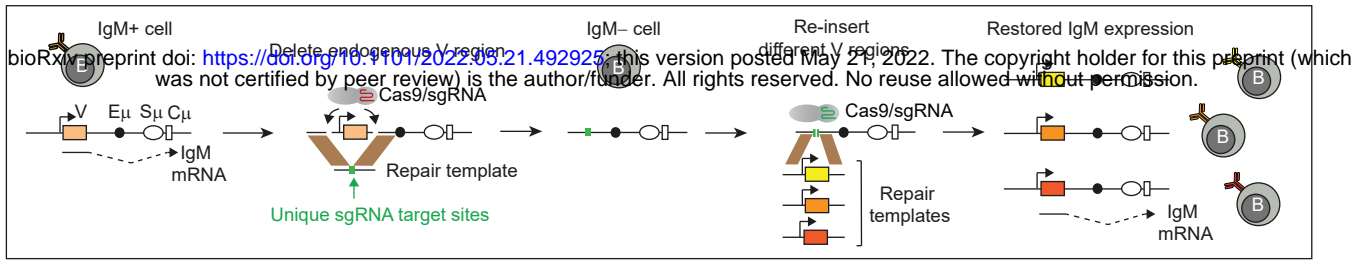
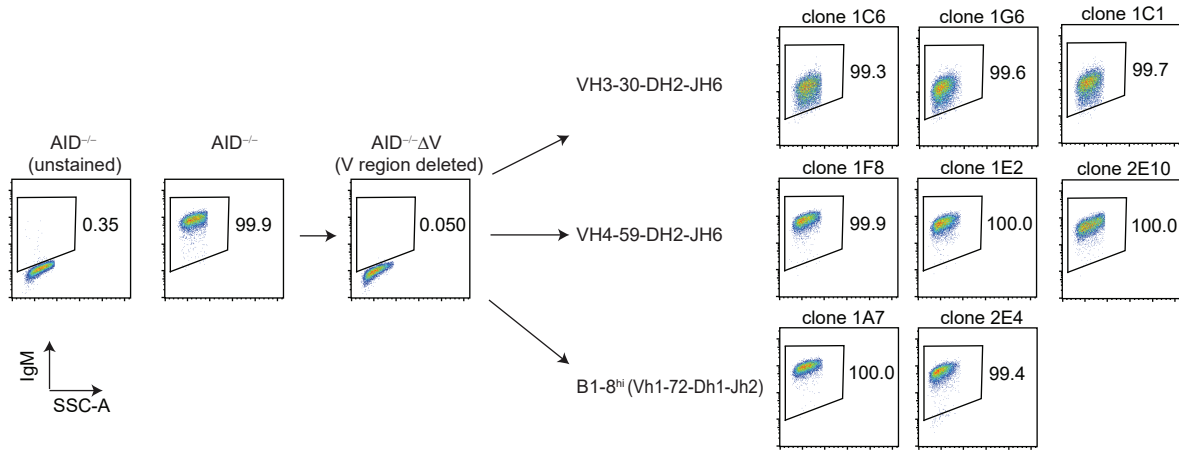


Figure S2

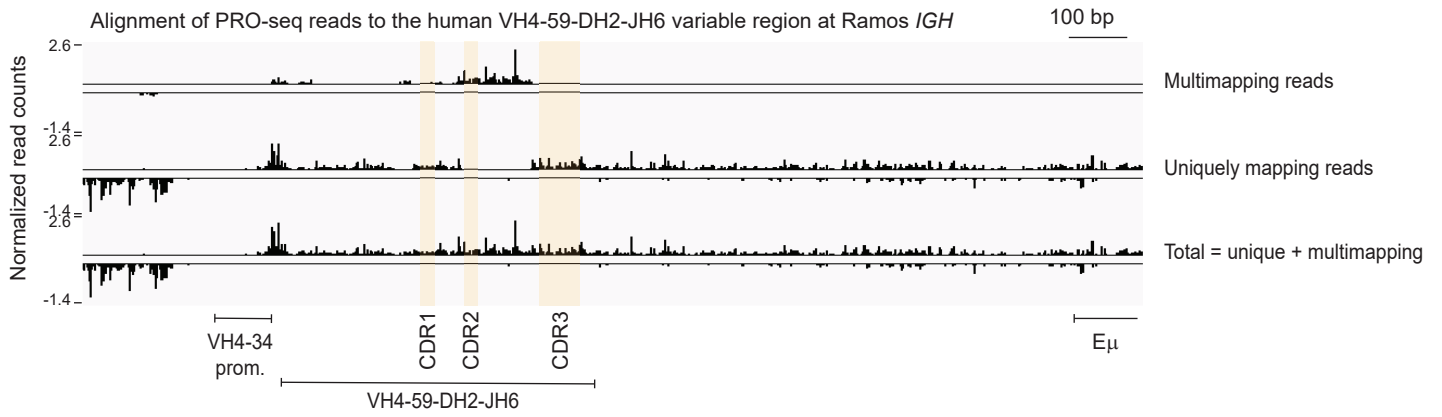
A



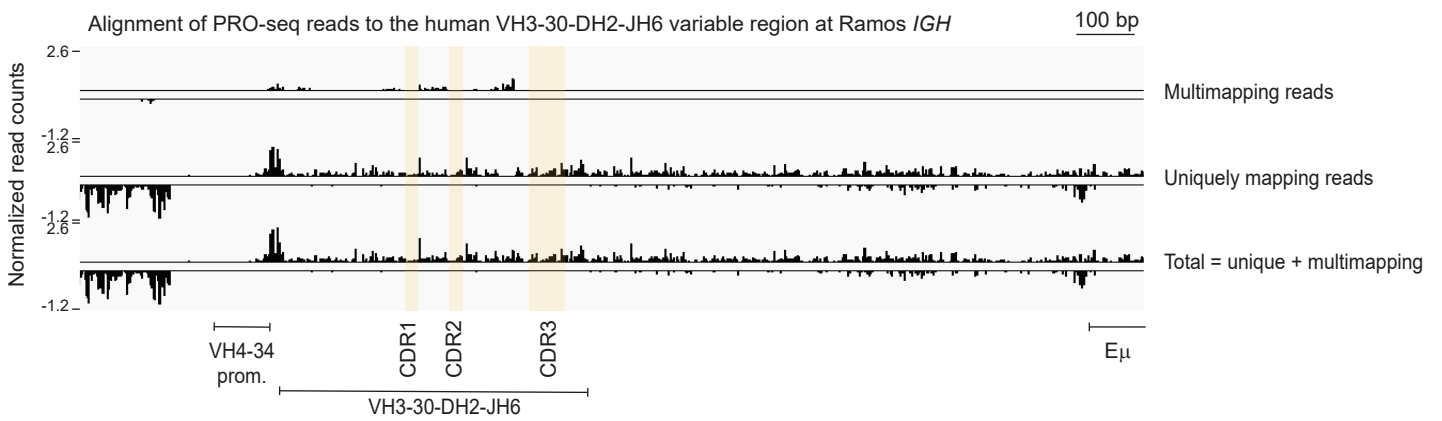
B



C



D



E

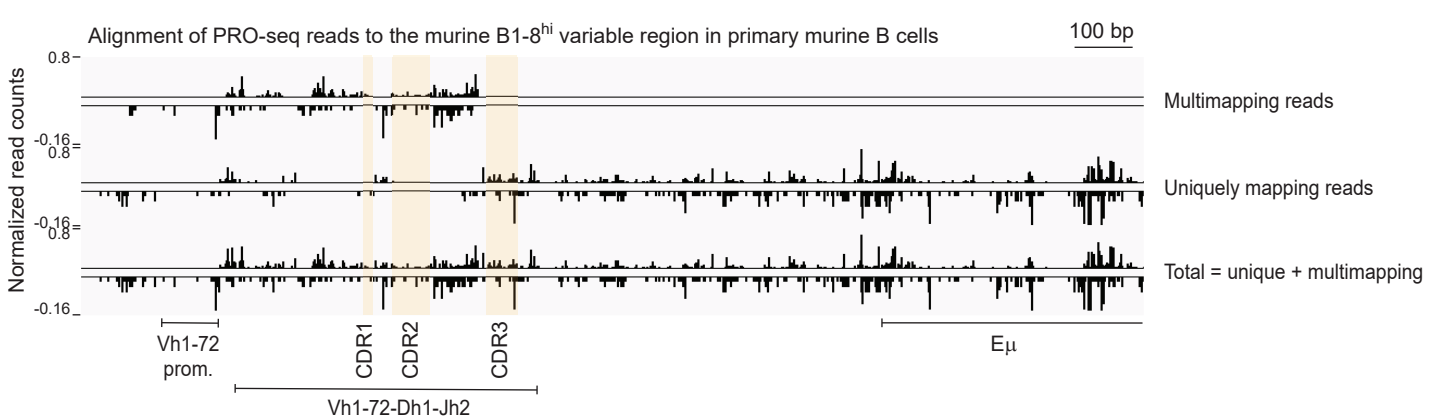
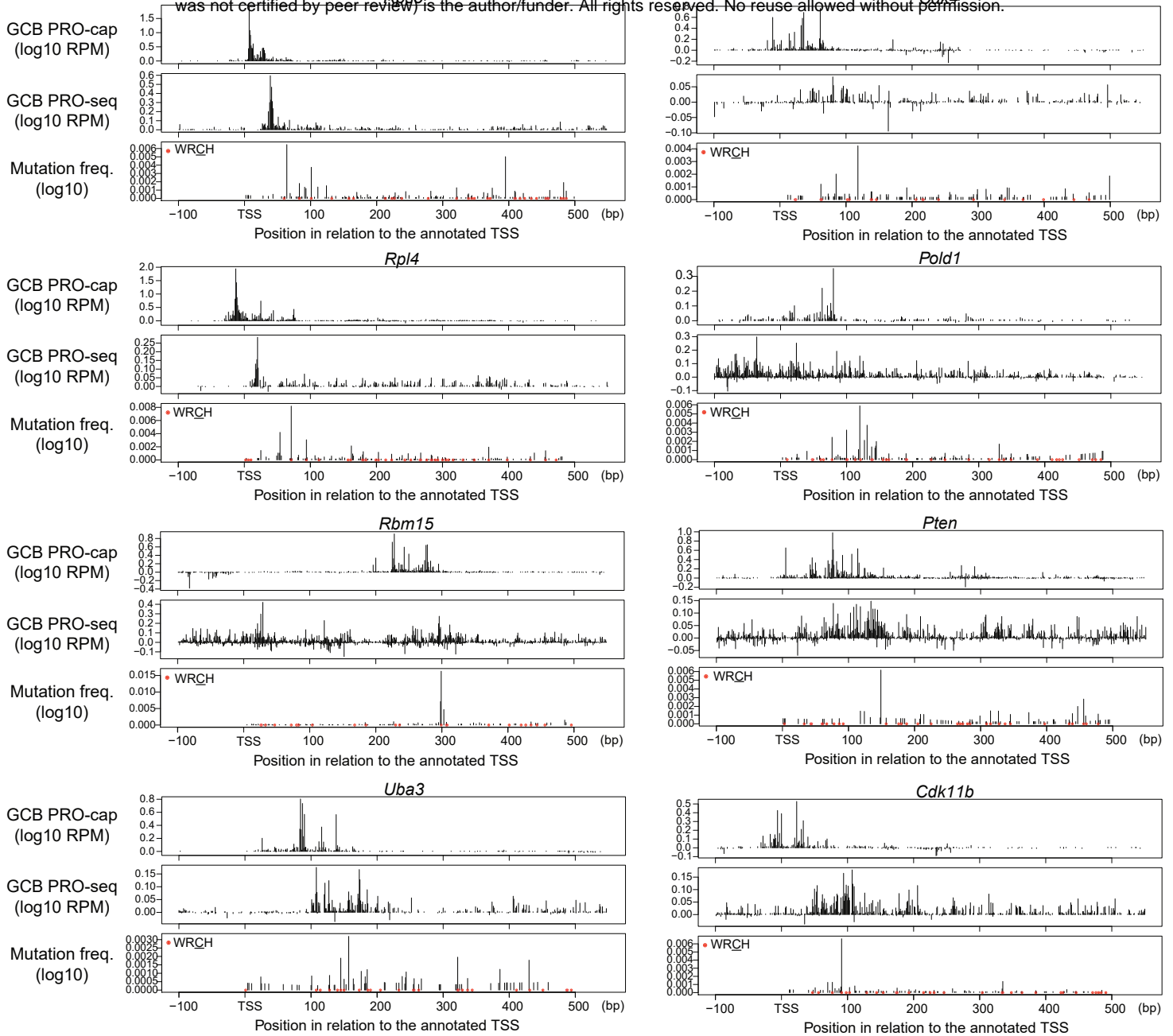


Fig. S3

A

Highly mutated residues within 150 bp of the transcription initiation site (defined by PRO-cap)
 bioRxiv preprint doi: <https://doi.org/10.1101/2022.05.21.492925>; this version posted May 21, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.



B

Highly mutated residues near the transcription initiation site (defined by PRO-cap)

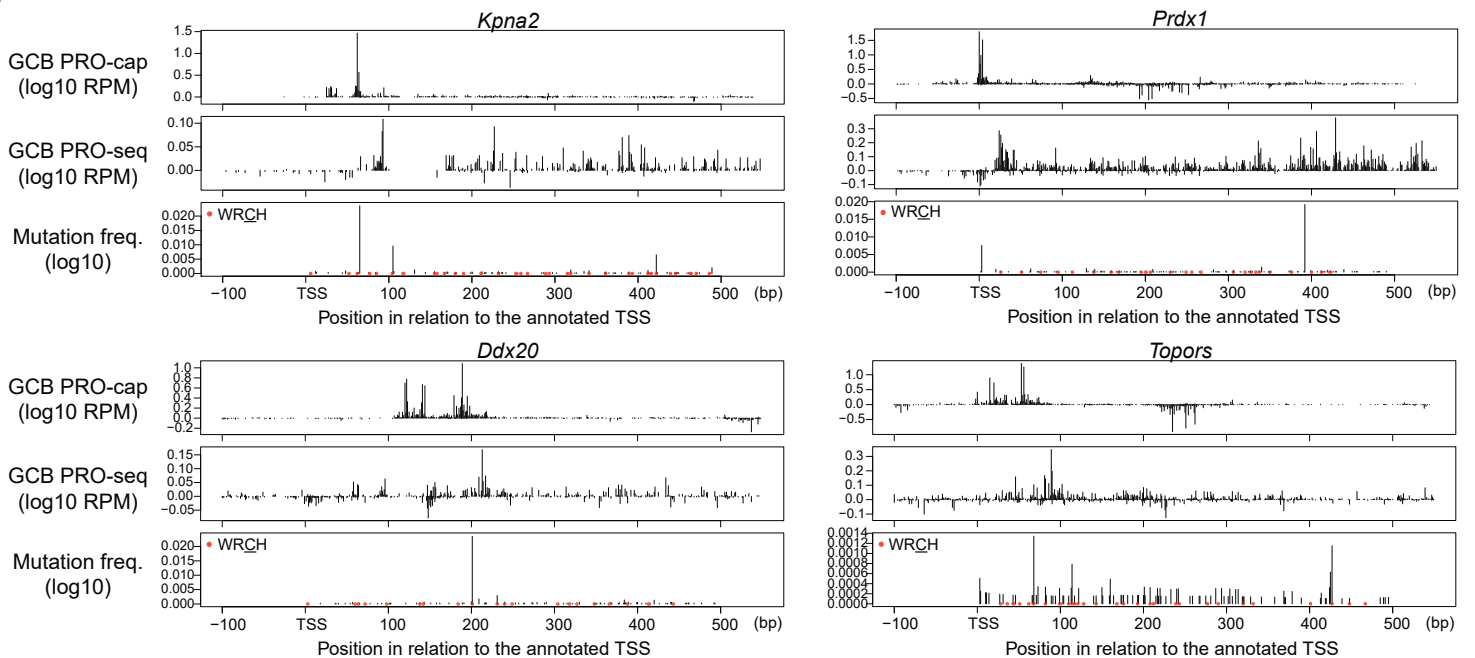


Fig. S4

bioRxiv preprint doi: <https://doi.org/10.1101/2021.05.21.492925>; this version posted May 21, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

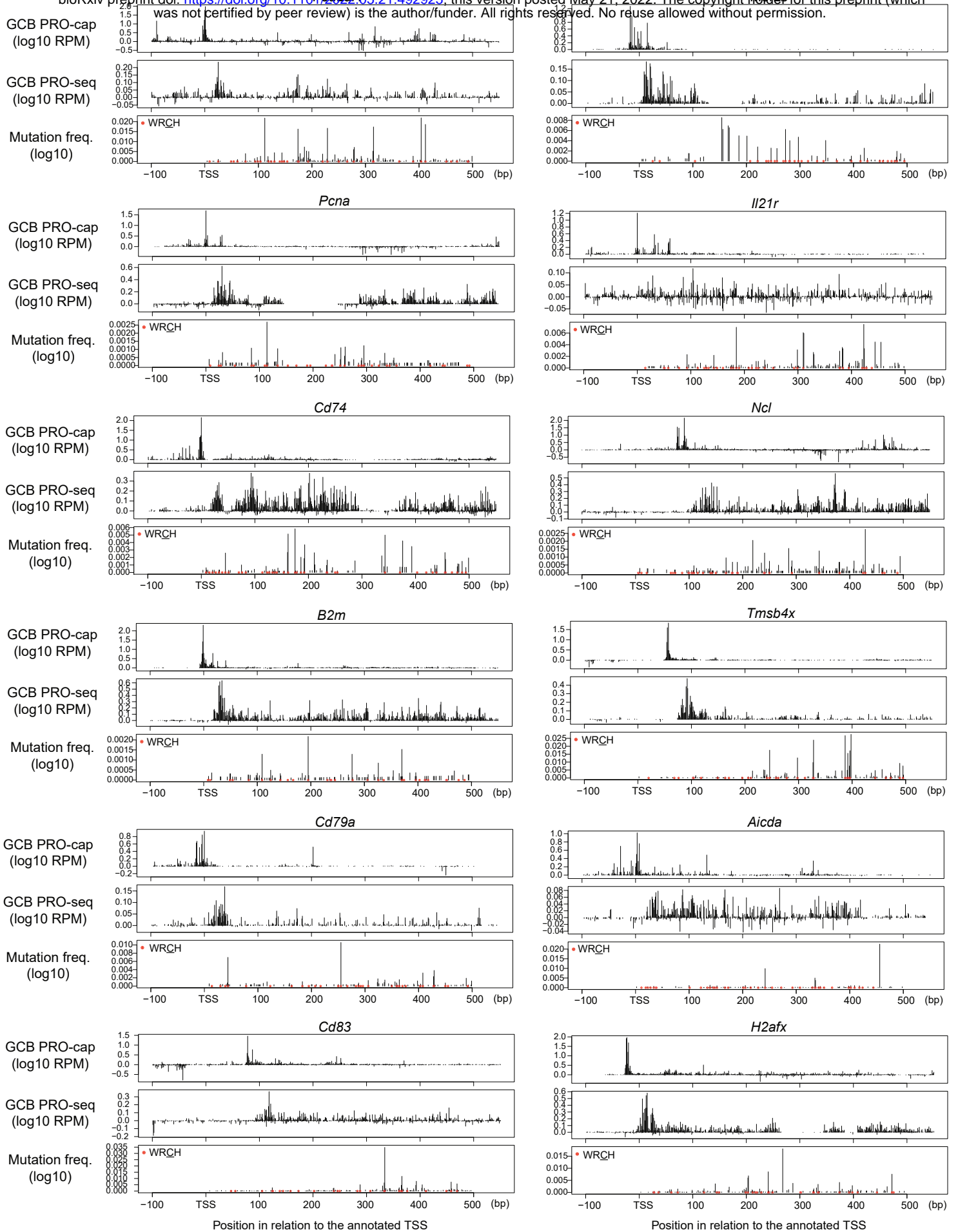


Table S1: List of primers, oligos and sgRNAs

sgRNAs for generating Ramos E μ ^{-/-} cells

sgRNA_E μ 1	CAGTGCTGTCGGCCCCGATG
sgRNA_E μ 2	ACTGTGTTCAACAATCTTTTT
sgRNA_E μ 3	GAGGCTGACCGAAACTGAAA

Genotyping primers for Ramos E μ ^{-/-} cells

	Forward	Reverse	Product sizes
gt_E μ 9+10	ATGTGTCTGGAATTGAGGCCAA	TTTCTTATGCCACCAACTAACA	825 bp WT, 277 bp KO
gt_E μ 11+12	ACAGACGGGAGGTACGGTAT	CCTTAATCACAAAGAGAAAACGAGT	1501 bp WT, 953 bp KO

sgRNAs for generating Ramos variable region knock-in lines

		Reference
sgRNA-Ramos VDJ 5'	CCTTCAGCACATTTCTACC	Voss et al., eLife, 2019
sgRNA-Ramos VDJ 3'	TCCTCGGGCATGTTCCGAG	Voss et al., eLife, 2019
sgRNA-GFP	GAGCTGGACGGCGACGTAAA	This study

RT-qPCR primers (Ramos)

	Forward	Reverse
1. Variable region	CAAGAACATGAAACACCTGTGG	ATCCGCATTCTGAGACTC
2. JH6-intron junction	GAGGTACGGTATGGACGTCTG	AGGACCAACCTGCAATGCTC
3. Intron (5' of E μ)	GGAGCCACATTTGGACGAGA	ACACACAGCGCACCTCATAA
4. Intron (5' of E μ)	ATGCGGGACTGCGTTTTGA	CATCATCTGCTCCAGCTTCG
5. Intron (3' of E μ)	GCGCCCGACATGGTAAGAGA	GACCCAGACAATGGTCACTCAA
IgM mRNA	ACAGACGGGAGGTACGGTAT	CGACGGGAATTCTCACAGG

ChIP-qPCR primers (Ramos)

	Forward	Reverse
1. Promoter	CACAGCCAGCATACACCTCC	CCTGTGGGTGCCTAAGTGAG
2. Variable region	AGGAATGCGGATATGAAGATATGA	CACTGAAGGACCCACCATAGAC
3. Intron (5' of E μ)	GGAGCCACATTTGGACGAGA	ACACACAGCGCACCTCATAA
4. Intron (5' of E μ)	ATGCGGGACTGCGTTTTGA	CATCATCTGCTCCAGCTTCG
5. E μ	GGTCACCGCGAGAGTCTATTT	TTCGGTCAGCCTCGCCTTAT
6. Intron (3' of E μ)	GCGCCCGACATGGTAAGAGA	GACCCAGACAATGGTCACTCAA
7. C μ	CGGTACTTCGCCACAGCA	ACTTGTCCAGGTCTCTCGGTGA
B2M	CTGTGCTCGCGCTACTCTC	AACTTGAGAGGGAAGTCAACG
Neg. (negative control region, chr 1)	GGAAAGGCCCCAGAGTTCAA	CCTCCTGATGTGAAGGCC

ChIP-qPCR primers (B1-8^{hi})

	Forward	Reverse
1. Jh4-intron junction	TGGACGAGGCCTTGAGTGG	AGACTGTGAGAGTGGTGC
2. Intron	CAGTCTCCTCAGGTGAGTCTC	CCCAATGACCCCTTCTGACT
3. Intron	ATGGTGTGGTGGAGTCC	AATCTCCAACACAGCCC
4. E μ	TGGGGCACTTTCTTAGATTTG	GACAGCACTACCCTTTTGAGACC
5. Intron	CTGCAGCAGCTGGCAGG	GGCTGGACAGAGTGTTCAAAACCCAC
6. Intron	GTTGCCTGTTAACCAATAATCATAGACTCATGG	GTATAACTGAAGTAGAGACAGCATCAGTACCTCAAC

ChIP antibodies

	Vendor	Catalog no.
H3K4me3	Abcam	ab8580
H3K27ac	Abcam	ab4729
H3K36me3	Abcam	ab9050
H3 pan-Ac	Active Motif	61937
H4 pan-Ac	Active Motif	39026
H3K79me3	Abcam	ab2621
H3	Abcam	ab1791
IgG	Invitrogen	10500C

FACS antibody

	Vendor	Catalog no.	Dilution
IgM-APC (MHM-88)	BioLegend	314510	1/200
B220-FITC	BD Biosciences	553088	1/500
Fas-PE-Cy7	BD Biosciences	557653	1/1000
CD38-APC	ThermoFisher	17-0389-42	1/200

bioRxiv preprint doi: <https://doi.org/10.1101/2022.05.21.492925>; this version posted May 21, 2022. The copyright holder for this preprint (which was not certified by peer review) is the author/funder. All rights reserved. No reuse allowed without permission.

MutPE-seq primers

Round 1_FWD (B1-8hi)	CTCTTTCCCTACACGACGCTCTCCGATCTNNNNNTCCACAGGTGTCCACTCCCAG	N is A,G,C,T randomly entered during primer synthesis
Round 1_FWD (VH4-34, VH4-59, VH3-30)	CTCTTTCCCTACACGACGCTCTCCGATCTNNNNNTGTTACAGGGGTCCTGTCC	N is A,G,C,T randomly entered during primer synthesis
Round 1_REV	CTGGAGTTCAGACGTGTGCTCTCCGATCTCCTAGAGTGGCCATTCTTACC	
Round 2_FWD	AATGATACGGCGACCACCGAGATCTACACNNNNNNNACACTCTTCCCTACACGAC	N represents 8 positions of individual sample barcode (IdxE501-E5xx)
Round 2_REV	CAAGCAGAAGACGGCATACGAGATNNNNNNGTGAAGTTCAGACGTGTG	N represents 6 positions of individual sample barcode (revIDX1-xx)