

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

Genome-wide binding analysis of 195 DNA Binding Proteins reveals “reservoir” promoters and human specific SVA-repeat family regulation

Soraya Shehata^{1,2,3¶}, Savannah Spradlin^{4,3¶}, Alison Swearingen^{1,3¶}, Graycen Wheeler^{4,3¶}, Arpan Das^{1,3}, Giulia Corbet^{4,2,3}, Benjamin Nebenfuehr^{1,3}, Daniel Ahrens^{1,3}, Devin Tauber^{4,2,3}, Shelby Lennon^{4,3}, Kevin Choi^{1,3}, Thao Huynh^{4,2,3}, Tom Weiser^{4,3}, Kristen Schneider^{5,2,3}, Michael Bradshaw^{5,2,3}, Maria Lai⁶, Joel Basken⁶, Tim Read⁶, Jon Demasi², Matt Hynes-Grace², Dan Timmons², Michael Smallegan^{1,2*}, and John L. Rinn^{4,2,1*}

1) Molecular, Cellular & Developmental Biology, University of Colorado Boulder, Boulder, Colorado, United States of America

2) BioFrontiers Institute, University of Colorado Boulder, Boulder, Colorado, United States of America

3) Biochemistry 5631 Spring 2020, University of Colorado Boulder, Boulder, Colorado, United States of America

4) Department of Biochemistry, University of Colorado Boulder, Boulder, Colorado, United States of America

5) Department of Computer Science, University of Colorado Boulder, Boulder, Colorado, United States of America

6) Arpeggio Biosciences, Boulder, Colorado, United States of America

* Corresponding author

E-mail: michael.smallegan@colorado.edu (MS)

E-mail: john.rinn@colorado.edu (JR)

¶ These authors contributed equally to this work.

47 **Abstract**

48 A key aspect in defining cell state is the complex choreography of DNA binding events in a
49 given cell type, which in turn establishes a cell-specific gene-expression program. In the past
50 two decades since the sequencing of the human genome there has been a deluge of genome-
51 wide experiments which have measured gene-expression and DNA binding events across
52 numerous cell-types and tissues. Here we re-analyze ENCODE data in a highly reproducible
53 manner by utilizing standardized analysis pipelines, containerization, and literate programming
54 with Rmarkdown. Our approach validated many findings from previous independent studies,
55 underscoring the importance of ENCODE's goals in providing these reproducible data
56 resources. This approach also revealed several new findings: (i) 1,362 promoters, termed
57 'reservoirs,' have up to 111 different DNA binding-proteins localized on one promoter yet do not
58 have any expression of steady-state RNA (ii) The human specific SVA repeat element may
59 have been co-opted for enhancer regulation. Collectively, this study performed by the students
60 of a CU Boulder computational biology class (BCHM 5631 – Spring 2020) demonstrates the
61 value of reproducible findings and how resources like ENCODE that prioritize data standards
62 can foster new findings with existing data in a didactic environment.

63

64

65 **Introduction**

66 In the postgenomic era[1,2] there have been efforts to adapt classical biochemical protocols
67 studying a few DNA regions to genome-wide events. One of the first of these genome-wide
68 assays was Chromatin Immunoprecipitation (ChIP) followed by hybridization of co-precipitate
69 DNA fragments to microarrays (or ChIP-CHIP) representing many thousands of DNA locations

70 (e.g. promoters). This application was first demonstrated in yeast and quickly adapted to many
71 species[3–7]. With the advent of massively parallel sequencing technologies, bound DNA from
72 the biochemical ChIP could be sequenced (ChIP-seq) to unbiasedly detect binding events
73 (reviewed[8,9]). This rapid change in platforms for ChIP analyses resulted in many data sets
74 that differed greatly in their results (ChIP-ChIP versus ChIP-seq)[10,11]. Only three years after
75 sequencing of the human genome it became clear that uniform experimental and data
76 standards were essential to limit a deluge of irreproducible results.

77

78 To this end, the field turned to the publicly available ENCODE consortium as the largest and
79 most standardized repository of ChIP-seq data sets[12–15]. The goal was to develop
80 standardized experimental and computational pipelines. Over the past 17 years since its
81 inception, many thousands of ChIP-seq experiments have been performed. Often these large
82 consortium studies analyze these data sets across cell types and tissues[13,13,16–19]. In
83 contrast, fewer studies have investigated dozens of DNA binding proteins (DBPs) in one cell
84 type.

85

86 Observing how hundreds of DBPs are bound relative to each other in the same cellular context
87 provides a unique perspective. This allows a promoter-centric approach across hundreds of
88 possible DNA binding events. Thus, we can address the underlying properties of combinatorial
89 binding at promoters and, in turn, how this relates to promoter activity. Moreover, this approach
90 allows us to systematically investigate numerous DBPs for possible enrichment in noncoding
91 regions such as repetitive element class and families. Overall, this strategy is limited in cellular
92 diversity, but rich in relative information of binding events at a given promoter.

93

94 By investigating these properties for 195 DBPs in K562 cells, we were able to reproduce known
95 findings from independent data sets. For example, the number of binding events at a promoter

96 correlates with RNA expression output (both nascent and mature transcripts)[17,18]. We also
97 made several new observations. Specifically, we identify 1,362 promoters that do not produce a
98 mature transcript despite having up to 111 DBP binding events. We termed these promoters
99 “reservoirs” because these promoters serve as ‘reservoirs’ for DBPs. Importantly, reservoirs are
100 distinct from super-enhancers and highly overrepresented by long noncoding RNA (lncRNA)
101 promoters. We also observed that the human specific SVA repeat is one of the few repeat
102 families that had specific DBP enrichment, with a total of three DBPs specific to SVA repeats.
103 Looking further we found that SVA repeats reside adjacent to or within enhancers and are often
104 transcribed; suggesting they may have been co-opted in late primates as enhancer elements.
105
106 Overall, we demonstrate the utility of implementing data-science and reproducibility standards to
107 gain new insights combinations of genome-wide DNA binding events. We further note that the
108 design of this study was intended for didactic purposes and carried out by students in a
109 classroom setting.

110

111 **Results**

112 We first set out to survey the encode portal for the largest number of ChIP-seq experiments that
113 satisfied the following criterion: (i) target was considered a DNA binding protein (DBP), the
114 experiment used validated antibodies, sequencing was performed with 100bp paired end reads
115 and were in the same cell setting. We found the maximum number of samples that meet these
116 requirements were performed in K562 cells. Specifically, there are 1,076 FASTQ files comprised
117 of 195 DBPs meeting these criteria in K562. Rather than analyzing the peaks already called by
118 ENCODE for these experiments we chose to re-analyze the raw data using a community-
119 curated pipeline developed by “nf-core”[20]. This approach meets the highest data

120 reproducibility standards by using a container for all software and producing extensive
121 documentation at every stage of analysis within the nf-core/chipseq (v1.1.0) pipeline (Fig 1A).

122

123

124 **Fig 1. Framework of ChIP-seq analyses and peak calling across replicates.** (A) Schematic
125 of data quality requirements (2 or more replicates, 100bp reads, validated antibody) resulting in
126 1,076 FASTQ files representing 195 unique DNA binding proteins. FASTQs were processed
127 using the nf-core/chipseq pipeline (QC and peak calling). All FASTQ files passed nf-core quality
128 control metrics. (B) Browser view of raw data, individual replicate peak calls and our consensus
129 peaks. All scales are from 0 to 1 representing minimum and maximum reads in that window
130 using UCSC auto-scale. Peaks from individual replicates are in gray and consensus peaks
131 called are represented by black boxes.

132

133

134 The nf-core pipeline consists of documented analyses and quality control metrics that results in
135 significant windows or peaks of DNA binding events for each replicate[20]. After the
136 standardized pipeline gave us peak calls, we used this data to support our analysis and
137 exploration of the data. Our approach was to use R and Rmarkdown to document the analyses.
138 Compiling the 11 Rmarkdown files provided in the GitHub repository
139 (https://github.com/boulderrinnlab/CLASS_2020) will reproduce all the results and figures of this
140 study.

141

142 After calling significant peaks (MACS broad peak) for each replicate ChIP experiment for each
143 of the 195 DBPs, we wanted to develop consensus peaks across replicates. Briefly, we filtered
144 to peaks on canonical chromosomes and required that peaks overlap by at least 1nt in all
145 replicates for a given DBP. Peaks that overlap in all replicates are then merged by the union of

146 peak widths (Fig 1B-C). We observed five DBPs that did not have any peaks overlapping across
147 replicates perhaps suggesting that these are promiscuous antibodies, or these proteins have
148 heterogeneous binding across K562 cell populations (MCM2, MCM5, MCM7, NR3C1, TRIM25).

149
150 We next plotted the distribution of the number of consensus peaks for each DBP and found that
151 many DBPs had very few peaks. In order to capture the majority of DBPs and still provide a
152 reasonable number of peaks for analyses (e.g., permutation analyses), we chose a cutoff of 250
153 peaks (15% percentile, Supplemental Fig 1A). This results in 161 proteins to carry forward in the
154 analysis and in losing the following proteins: ARNT BCLAF1 COPS2 CSDE1 DNMT1 eGFP-
155 ETS2 FOXA1 KAT8 KDM4B MCM2 MCM5 MCM7 NCOA1 NCOA2 NCOA4 NR0B1 NR3C1
156 NUFIP1 PYGO2 THRA TRIM25 TRIP13 XRCC3 YBX1 YBX3 ZBTB8A ZC3H8 ZNF318
157 ZNF830.

158

159 **Promoter centric binding properties of 161 DNA Binding**

160 **Proteins**

161 We next plotted the relationship between the number of consensus peaks observed for each
162 DBP and how many promoters overlapped (36,814 lncRNA and mRNA promoters). We observe
163 a linear relationship (slope = 0.31 for mRNA and lncRNA promoters) between the number of
164 peaks and or size of peaks and the number of overlaps with promoter regions (Fig 2A).

165 Somewhat surprising was this trend was even more pronounced when comparing overlaps
166 within gene-bodies rather than promoter regions (Fig 2B). This suggests we could have an
167 observation bias at promoters where promoter binding simply increases with the number of
168 peaks observed for a given DBP and not due to preferential binding at promoters.

169

170

171 **Fig 2. Promoter binding properties of 161 DBPs.** (A) Schematic of promoter overlap strategy.

172 Number of overlapping promoters (y-axis) per number of peaks for each DBP (x-axis). (B) Same

173 as in (A) but for overlapping gene bodies instead of promoters. (C) Binary clustering of 161

174 DBPs based on promoter binding profiles (consensus peaks). Zoom out of specific regions.

175

176

177 To detect preferential binding at promoters, we took a permutation-based approach for each

178 DBP's peak-profile across the genome. Briefly, we took the consensus peaks for each DBP and

179 randomly placed them across the genome, while controlling for (i) number of peaks, (ii) width of

180 peaks and (iii) number of peaks on each chromosome. We then performed a Fisher exact test of

181 the observed binding at promoters versus expected binding in the empirically derived null

182 distributions. We observed that nearly all DBPs exhibit significant overlap with promoters versus

183 the rest of the genome, despite being involved in many different DNA regulatory processes

184 (Supplemental Fig 1B).

185

186 To more closely examine the results of our consensus peak strategy we performed manual

187 inspection of samples with two or more replicates (Fig 1B-C). We find that our peaks are

188 consistent with what would be expected of highly reproducible binding events. We see that most

189 Pol2 and ATF3 peaks show good agreement between replicates. Interestingly in this example

190 ATF3 is not localized to the promoter but in an upstream region that could be a newfound

191 enhancer or upstream regulatory element. Overall, these analyses are consistent with our

192 consensus overlap strategy representing expected and newfound features in peak size profiles.

193

194 **Global analysis of similarities in binding profiles**

195 To determine if there were underlying similarities and differences of 161 DBPs that passed our
196 conservative filtering, we first performed hierarchical clustering (Fig 2C) on binary vectors
197 representing binding events on 36,814 lncRNA and mRNA promoters defined in GENCODE 32
198 where 1 = bound, 0 = not bound for each promoter and DBP. As a quality control check, we
199 looked for clustering of known factors. The binary vector profiles validated that POLR2A,
200 POLR2B, and SUPT5H form a distinct cluster. Known family members, such as ATF3 and ATF2
201 co-cluster together as well, along with the eGFP-ATF3 control. This indicates that these DBPs
202 had similar binding profiles with or without the eGFP tag. However, 11 cases of eGFP-tagged
203 samples clustered together, despite having widely different functions. This may suggest that in
204 some rare cases the tag can alter the binding profiles in a manner that is more consistent with
205 the tag than DBP function.

206

207 As an unbiased approach to find underlying properties in DBP binding profiles, we also
208 performed UMAP[21] dimensionality reduction for the global binding profile of each DBP (Fig
209 3A). Briefly, UMAP uses algebraic topology to reduce the data dimensionality. We further
210 clustered this reduced representation using density-based clustering (HDBscan[22]). We
211 observed a total of seven clusters. Similar to binary clustering, we identify a clear cluster of
212 POLR2A, POLR2B, and SUPT5H and other basal transcriptional associated factors (TAF) as
213 would be expected. This is another example of high reproducibility as POL2 has three different
214 antibodies with 2 replicates each that are all highly concordant with thousands of peaks each.

215

216

217 **Fig 3. Binding properties of DBPs and expression output of promoters.** (A) UMAP
218 dimensionality reduction to identify DBPs with similar promoter binding profiles. (B) Four

219 discrete clusters of binding patterns around promoter TSSs with 3Kb up- and downstream. Line
220 is the average profile of all peaks in each cluster. (C) The number of peaks per DBP versus
221 number of mRNA or lncRNA promoter overlaps. X-axis is the number of DBPs overlapping
222 either lncRNA (red) or mRNA (black) promoters. (D) Chi-squared test for enrichment of DBPs
223 between lncRNAs and mRNAs. The x-axis is the $\log_2(\text{observed over expected})$ and y-axis is the
224 P-value.

225

226

227 Next we compared specific features of the DBP with their position in the reduced space by
228 mapping metadata (e.g., type of DNA binding domain) onto the UMAP points (Fig 2A,
229 Supplemental Fig 2A-D). We found no clear association with (i) type of DNA binding domain, or
230 (ii) annotation as a transcription factor, (iii) RNA-seq expression of bound genes, or other
231 properties. Collectively, these results recapitulate known biological functions of DBPs while
232 including potential new factors across these different promoter regulatory functions.

233

234 **Promoter binding specificity of 161 DNA binding proteins**

235 We next wanted to assess the underlying promoter features associated with each DBP.

236 Specifically, we wanted to determine where each DBP is bound relative to the TSS of 36,814

237 lncRNA and mRNA promoters. To this end, we generated 'binding profile plots' by calculating

238 the read counts across all promoters centered at the TSS with 3kb flanking up- and down-

239 stream (Fig 3B, Supplemental Fig 2E-F). We next clustered the 161 DBPs based on their

240 promoter profile plot. We split the dendrogram into clusters by 'cut-height' ($h = 65$). We

241 observed 4 distinct clusters with at least two DBPs. The first distinction is that about half exhibit

242 a narrow peak profile (71) and half with a broader peak profile (74). In both cases these profiles

243 peak near the TSS. Interestingly, 6 genes (eGFP-PTRF, ZBTB33, SMARCA5, HDGF, eGFP-
244 ZNF512, eGFP-ZNF740) have the inverse pattern: depletion of binding at the TSS with strong
245 enrichment at flanking regions (Supplemental Fig 2E-F).

246
247 Previous studies have identified several differences in binding features at coding (mRNA)
248 compared to noncoding (lncRNA) promoters. Here we wanted to independently test this across
249 161 DBPs to determine if there was an enrichment or depletion at mRNA versus lncRNA
250 promoters. We counted the number of lncRNA and mRNA promoter overlaps separately and
251 observed the same linear trend of more peaks resulting in more binding events for both
252 lncRNAs and mRNAs. However, the slope for mRNA is 0.19 ($R = 0.75$, $P < 1e-10$) and lncRNA
253 is 0.088 ($R = .87$, $P < 1e-10$) suggesting a two-fold reduction on an average lncRNA promoter
254 (Fig 3C).

255
256 We then performed permutation analysis (above) separately for lncRNAs and mRNAs to
257 determine if the observed overlap is greater than expected by chance (Supplemental Fig 3C).
258 Similar to our previous observation, nearly all DBPs were significantly (Fisher-exact $P < 0.05$)
259 enriched at both lncRNA and mRNA promoters yet with a smaller magnitude of enrichment of
260 binding events on lncRNA promoters (similarly to previously reported[17]). We observed two
261 DBPs that were significantly depleted: BRCA1 on mRNA promoters, and ZNF507 on both
262 lncRNA and mRNA promoters. Four DBPs showed neither enrichment or depletion at lncRNA or
263 mRNA promoters. In total 155 of the 161 tested DBPs were enriched at lncRNA and mRNA
264 promoters more than expected by chance (Supplemental Fig 2G).

265
266 Our previous permutation test above demonstrated that most DBPs bind both lncRNA and
267 mRNA promoters more than expected by chance. But this approach does not account for DBPs
268 that may prefer lncRNA or mRNA promoters. Thus, we hypothesized that some DBPs may have

269 a bias in binding for mRNA relative to lncRNA promoters and vice-versa. To test this, we
270 performed a Chi-squared test to compare the number of binding events for each DBP at lncRNA
271 versus mRNA promoters. Interestingly, although most DBPs are enriched on mRNA promoters,
272 there were a few with a relative bias toward lncRNA promoters ($P < 0.05$): BRCA1, eGFP-
273 ZNF507, EWSR1, eGFP-TSC22D4 (Fig 3D). Interestingly, BRCA1 prefers to bind outside of
274 promoters, yet if it does bind a promoter BRCA1 prefers lncRNA over mRNA promoters.
275

276 **Repeat family and class binding preferences for 161 DNA**

277 **binding proteins**

278 In order to determine if DBPs are enriched or depleted in TE classes and families we performed
279 a permutation enrichment analysis. As above, we randomly shuffled peaks around the genome
280 and calculated the number of overlaps with repeat family and classes from RepeatMasker
281 Open-3.0 occurring by chance (Fig 4A).

282

283

284 **Fig 4. Many DBPs are enriched or depleted on repeat families and classes.** (A) Heat map
285 of Z-scores of observed overlaps of each DBP versus the overlap distribution of 1,000 random
286 permutations of each DBPs profile genome wide. Red indicates depletion and blue enrichment
287 (negative versus positive Z-scores respectively). The observed and permuted Z-scores are for
288 overlaps with repeat classes. (B) The same permutation analysis as in (A), but for observed
289 versus permuted overlaps with repeat families. Red indicates depletion and blue enrichment.

290

291

292 We observe that some classes, such as Simple Repeats and tRNAs, were enriched for most
293 DBPs, while others, such as the LINEs and Satellites, were depleted for most DBPs (Fig 4A).
294 The LINE class was depleted of all DBPs with the exception of five DBPs with zinc finger-like
295 motifs (ZNF507, ZNF316, ZNF184, ZNF24 and ZNF512). Additionally, the LTR class was
296 depleted for most DBPs, but enriched for a subset of 23 DBPs (Fig 4B).

297
298 Overall, we found that most small TE families were not significantly enriched or depleted for
299 specific DBPs. However, a subset of 23 DBPs were enriched in the ERV1 family, but depleted in
300 the L1 family. These 23 DBPs are the same that were enriched in the LTR class. This is
301 consistent with ERV1 family TEs being a part of the LTR class. Similarly, the MIR family shows
302 a similar enrichment pattern to the tRNA family (Fig 4B). Thus, using this approach we can
303 provide a map of which DBPs are specifically bound to which repeat family.

304
305 We did observe a subset of 6 DBPs (NUFIP1, ZC3H8, PHF21A, ARHGAP35, NCOA4, PYGO2)
306 enriched in snRNAs, but no other TE family. Each of these DBPs, except NCOA4, contains a
307 zinc-finger-like DNA binding domain, and a few (NUFIP1 and ZC3H8) are known to be a part of
308 the snRNA biogenesis pathway, perhaps suggesting some form of feedback. The L1 family is
309 depleted for almost every DBP, but is highly enriched for ZNF507, an interaction which has
310 been previously described in an undergraduate thesis and confirmed genome-wide here
311 (<https://web.wpi.edu/Pubs/E-project/Available/E-project-042618-111020/unrestricted/MQP.pdf>).

312

313 **The human specific SVA repeat family has enhancer like** 314 **features**

315 Although most families are not enriched for specific DBPs, the human specific SVA repeat
316 family is specifically enriched for three DBPs: ZBTB33, CBFA2T2, and CBFA2T3. Interestingly,
317 all three of these DBPs are known transcriptional repressors. The SVA family is the youngest
318 family of TEs, is enriched in gene-rich areas of the genome[23–25], and can cause human
319 disease[24]. Based on these interesting features we further explored the binding of these factors
320 on the SVA repeat.

321
322 We first retrieved histone modification ChIP data for K562 cells from ENCODE and visualized
323 the coverage centered on the 5,882 SVA repeats with 5kb up- and down-stream. We find that
324 Lysine 4 mono-methylation (K4me1) is the only histone modification enriched on SVA elements
325 – all others were depleted (Supplemental Fig 3A). Moreover, the enrichment of K4me1 is on the
326 5' end of the SVA element suggesting it could be an insulator for enhancers or part of the
327 enhancer element. This pattern is so sharp we were concerned about mappability to the SVA
328 element – despite observing the 5' enrichment of K4me1. We reasoned that ZBTB33, CBFA2T2
329 and CBFA3T3 should be enriched across the SVA element. We performed the same analysis
330 above for the these 3 DBPs and find there is strong mapping to these SVA regions, which
331 suggests a low potential for the histone mark depletion to be an artifact of low mappability
332 (Supplemental Fig 3B). We next looked at the expression level of SVA elements relative to other
333 repeat family members. Interestingly, we observed that SVA elements have more transcription
334 (Supplemental Fig 3C) than LTR family members that are known to function as promoters[26].
335 Together, these results demonstrate that the SVA region has enriched and fully mappable
336 coverage of K4me1, ZBTB33, CBFA2T2, CBFA3T3 and are expressed.

337

338 Of the 5,882 SVA elements genome-wide, 255 SVAs were found to contain consensus peaks
339 for all three enriched DBPs: ZBTB33, CBFA2T2, CBFA3T3. We took the same approach above
340 for this subset of bound SVA elements. We see even stronger enrichment of K4me1 (Fig 5A)
341 and also coverage by ZBTB33, CBFA2T2 and CBFA3T3 (Fig 5B). Interestingly, the shape and
342 position of ZBTB33 is distinctly different than that of CBFA2T2/3T3 (Fig 5B). It suggests that
343 ZBT33 binds on the 5' region near K4me1 and CBFA2T2/3T3 have overlapping positions on the
344 3' end of SVA elements. Closer examination of nascent, steady state RNA-sequencing (see
345 below) and K4me1 CHIP shows a very interesting pattern of the SVA elements being
346 transcribed and or producing bi-directional RNAs in K4me1 enriched (Fig 5C-D). This is very
347 similar to what has been seen for enhancer regions genome wide[27,28]. Thus, the SVA
348 transposon may have evolved (neutrally or positively) to 'co-opt' binding of DBPs adjacent or
349 within enhancer regions.

350

351

352 **Fig 5. Human SVA repeats are enriched for DBPs and enhancer properties.** (A) Heatmap of
353 histone modification reads centered on SVA and 5Kb up- and down-stream for the 255 SVA
354 elements containing ZBTB33 and CBFA2T2/3T3 peaks. Here red indicates enrichment, while
355 blue indicates depletion. Above is the average profile line of enrichment within and outside SVA
356 elements. (B) Same as (A) but coverage of ZBTB33 and CBFA2T2/3T3. The K4me1 plot is
357 same as in (A) for direct comparison. (C) Browser examples in the same format as Fig 1.

358

359

360 **Promoter binding of 161 DNA binding proteins versus** 361 **promoter expression output**

362 Here we set out to investigate how binding events at individual promoters relate to the
363 concomitant expression of the gene-product at that promoter. To this end, we analyzed
364 ENCODE K562 total RNA sequencing data from two replicates. We calculated the average read
365 coverage across replicates and quantified by transcripts per million reads (TPM); while
366 considering the variance between replicates in further analyses. We first asked if the number of
367 binding events at a promoter correlated with expression. We observed a positive correlation (R
368 $= 0.6$, $P < 2.2e-16$) between the number of DBPs bound at a promoter and expression output of
369 the promoter (Supplemental Fig 4A).

370

371 We next wanted to determine if this trend is similar for mRNA and lncRNA promoters
372 separately. Indeed, we see that both lncRNA and mRNA promoters have a positive correlation
373 to binding events and expression output (Fig 6A). We observed that lncRNAs have lower
374 expression in general than mRNA as previously determined[29–32]. Yet despite these
375 expression differences, both exhibit a positive relationship between number of binding events
376 and promoter activity. This is consistent with observations in a previous study using a different
377 yet overlapping subset of 73 DBP ChIP datasets[32].

378

379

380 **Fig 6. Reservoir promoters are comprised of ghosts and zombies.** (A) Number of DBPs
381 bound to a promoter (x-axis) versus \log_{10} (TPM) of transcription as measured by total RNA-seq.
382 (B) Box plot comparing mRNA (black) and lncRNA (red) expression as a function of off, low,
383 medium, and high expressed transcripts. Y-axis is the number of DBPs and X-axis each

384 category. (C) Y-axis is the mean expression level in windows of 5 genes excluding the center
385 gene, with a step (slide) of 1 gene. X-axis is by category of windows containing a reservoir, non-
386 reservoir or super-enhancer. Y-axis is mean expression in each 5 gene window. (D) Density plot
387 of number of DBPs bound at a promoter at expressed (grey) versus non-expressed promoters
388 (red), separated by lncRNA and mRNA promoter types. (E) Nascent TPM expression (y-axis)
389 compared to number of DBPs bound at a promoter. (F) Density plots of DBPs ghost reservoirs
390 (those without nascent expression, PRO-seq TPM < 0.001) vs those with detectable nascent
391 expression (zombie reservoirs).

392

393

394 Although we saw a linear trend with binding events and expression output above, we wanted to
395 refine this analysis to a binned approach. Specifically, we binned lncRNA and mRNA promoters
396 by expression output of: Off: < 0.001 TPM, Low: (0.001,0.137] TPM, Medium: (0.137,3] TPM,
397 and High: >3 TPM. Interestingly, at 'low' and 'off' expressed promoters there is no difference in
398 binding event distributions between lncRNA and mRNA promoters (Fig 6B). Thus, they both
399 have similar numbers of binding events -- and can have dozens of DBPs bound -- despite
400 having little to no expression output. In contrast, mRNA promoters show significant increases in
401 binding events, compared to lncRNA promoters, at medium and high expressed promoters.
402 Thus, in the middle to high ranges of expression is where we begin to see the differences
403 between mRNA and lncRNA promoters. Collectively, these results identify over a thousand
404 promoters that resemble the DBP content of highly-expressed promoters yet do not have any
405 detectable expression by RNA-seq.

406

407 **Promoters with numerous binding events but lack gene-** 408 **expression output**

409 Based on the observation of over a thousand promoters that have numerous DBPs bound, but
410 do not produce a transcript identified by RNA-seq, we wanted to further characterize the global
411 properties of this subset of promoters. First, we made density plots of the number of binding
412 events at promoters. We observed a bimodal distribution of binding events where the cutoff
413 between the distributions is around seven binding events at a promoter (54% percentile). Based
414 on these two distinct distributions, we focus our analysis on those promoters with more than
415 seven binding events (Supplemental 4B) and further required that the RNA-seq output was less
416 than 0.001 TPM. This resulted in 1,362 promoters which had a relatively high number of binding
417 events but lack of RNA-seq output from these promoters. Interestingly, 981 of the 1,362 are
418 comprised of lncRNA promoters (Supplemental Fig 4C). This is a significant over-representation
419 of lncRNAs in these high-binding non-expressed promoters over what would be expected by
420 chance (hypergeometric p-value = 1.1×10^{-88}).

421
422 There are two trivial explanations that could explain these high binding low expression
423 promoters: (i) these are simply super enhancer[33–35] annotations (as they share similar
424 properties of many binding events) and or (ii) the promoter is regulating a neighboring gene.

425
426 Our first concern is that super-enhancers (SE) share the similar property of many binding
427 events, we wanted to determine how many of these regions were super enhancers. For super-
428 enhancer annotations we used the SE-DB[36] that is comprised of 331,601 super-enhancers
429 from 542 tissues and cells, including K562. We first retrieved the SE annotations in K562 with
430 the hg19 reference genome alignments. We then lifted over these annotations from hg19 (732

431 annotated SEs) to hg38. We found 714 annotations have one match to the genome and took a
432 conservative approach of not including the 18 SEs with multi-mapping in the genome (often too
433 many chromosomes). Of these 714 regions, 35 overlapped with the 1,363 reservoirs ($P = .991$
434 Hypergeometric). Thus, reservoirs are distinct from SE annotations and are enriched with
435 repressor complexes unlike SEs.

436
437 Another concern is that these promoters we identified could regulate a neighboring gene; this
438 would be most obvious for bidirectional promoters. Thus, we first defined a set of promoter
439 types: (i) bidirectional, if another promoter on opposite strand overlaps within 1,000 bp upstream
440 of the TSS on the opposite strand (147 / 11%); (ii) multiple nearby promoters, if there is more
441 than one promoter on either strand within 1,000 bp (91 / 7%); (iii) nearby on same strand if there
442 is another promoter upstream within 1,000bp (113 / 8%); (iv) none (1,011 / 74%), if there are no
443 promoters within 1,000 bp (Supplemental Fig 4D). Collectively, very few reservoirs had shared
444 promoters of any type i-iii (26%), thus this cannot likely account for the lack of transcription at
445 the observed or neighboring promoter (since there are so few). Nonetheless we calculated the
446 TPM of promoter(s) neighboring reservoirs. We observed that 68% of these shared promoters
447 did have a neighboring gene expressed (subcategories in Supplemental Fig 4E) for a total 240
448 (15%) of reservoirs that could affect neighboring gene expression. Thus, neighboring promoters
449 of any orientation cannot account for the general lack of expression observed at reservoirs (Chi-
450 squared p -value $2e-22$).

451
452 Although bidirectional expression cannot explain why these promoters seem inert, we wanted to
453 look more globally at the transcription environment of these promoter regions and their 5
454 neighboring genes. Specifically, we used a “sliding-window” approach to calculate the median
455 TPM expression value for windows of 5 genes. Each window is centered on one gene and the
456 mean of the neighboring four genes is calculated excluding the center gene. We first plotted the

457 distribution of windows where the center gene is a reservoir compared to those with non-
458 reservoir center genes. We also removed the 35 reservoirs that were annotated as super-
459 enhancers. We observed that the Wilcoxon test statistic (Fig 6C) between means was
460 significant ($P < 9e-06$), however the means were very similar (mean = 7.2 for reservoir, mean =
461 8.4 for non-reservoir). To be sure this is not an artifact of our permutation analysis we performed
462 the same analysis for windows of genes centered on super-enhancers versus non super-
463 enhancers. Indeed, we see that super-enhancers reside in regions of significantly higher
464 transcriptional activity ($p < 2.5e-12$) with a large fold change (4.5x) in mean expression (mean
465 super-enhancer = 37 TPM, mean = non super-enhancer 8.4 TPM) (Fig 6C).

466
467 Collectively, these results identify a subset of promoters that appear to be a ‘holding place’ for
468 DNA binding events. Thus, we will refer to these promoters as ‘reservoirs’ since they: (i) are
469 distinct from super-enhancer annotations; (ii) are located in more transcriptionally silenced
470 neighborhoods; (iii) share the property of many DNA binding properties as those promoters that
471 are highly expressed and (iv) have no expression output as measured by RNA sequencing.

472

473 **DNA binding properties of reservoir promoters**

474 To understand if reservoir promoters are enriched for certain DBPs, we compared the density of
475 DNA binding events at lncRNA and mRNA reservoir and non-reservoir promoters which had
476 greater than seven binding events. We observed a shift toward fewer binding events for both
477 lncRNA and mRNA reservoirs (Fig 6D). However, it’s notable that there are still reservoirs along
478 the whole range of DBP binding. Although reservoirs have fewer binding events in general, we
479 wanted to determine if there was enrichment of certain DBPs on reservoirs. Using a Chi-
480 squared test to compare the number of bound promoters for reservoirs versus non-reservoirs
481 we observed that 31 DBPs were depleted on reservoirs and only one gene enriched ($P < 0.001$

482 and > 2-fold depletion/enrichment, Supplemental Fig 4F). This is in contrast to the lncRNA and
483 mRNA comparisons above where we saw global depletion of all DBPs on lncRNA promoters.
484 Thus far, reservoirs are deviant from all trends observed for the other ~33,000 promoters tested
485 above.

486
487 We wanted to further globally characterize reservoir promoters using UMAP dimensionality
488 reduction as in Fig 2A. Unlike with all promoters we only observe two distinct clusters across
489 reservoirs (Supplemental Fig 4G). However, gene-ontology analysis revealed that both clusters
490 are strongly enriched for similar processes such as regulation of transcription ($P < 1e-20$).
491 Perhaps as expected, Pol2 and associated transcriptional machinery are some of the most
492 significantly depleted from reservoirs; consistent with their lack of expression. Despite a global
493 depletion of Pol II at reservoirs, we were surprised that over a quarter of reservoirs (417) had
494 Pol II binding events, suggestive of 'paused' transcription. While only one DBP (eGFP-
495 TSC22D4) reached the fold-change threshold, two more were found to be significant ($P <$
496 0.001) with small enrichments. All three are associated with repressive activity. TSC22D4 and
497 CBFA2T2 are both known repressors while EHMT2 facilitates transcription repression through
498 methylation of H3K9. Collectively, these findings show that reservoir promoters are distinct from
499 super enhancers, bound by many DBPs and yet are not transcribed.

500

501

502 **Nascent Expression and chromatin properties of reservoir** 503 **promoter**

504 Since reservoirs don't have mature transcriptional products despite many promoter binding
505 events, we next examined if reservoirs have "nascent" transcription detected via PRO-seq

506 (reviewed³⁴). These approaches are so precise they can identify specific DBP binding sites
507 through PRO-seq nascent RNA read out [37,38]. Thus, we hypothesized that reservoir
508 promoters would exhibit nascent transcription owing to so many DNA binding events. This could
509 also be similar to more well established “paused” promoters as reviewed[39].

510
511 To determine the nascent transcription properties of reservoirs, we obtained two replicate pro-
512 seq data sets that measure the amount of nascent transcription at a promoter. We used
513 “Rsubread”[40] to calculate TPM values of nascent transcription across the same 6 Kb promoter
514 window defined for DBP binding. We first plotted the relationship of nascent sequence at
515 reservoirs versus non-reservoirs (Supplemental Fig 5A). Although statistically different ($P < 3e-$
516 9) the distributions are fairly similar for reservoirs (mean = 0.41) and non-reservoirs (mean =
517 0.51) with a fold change of only 1.25. Thus, consistent with lack of RNA-seq expression,
518 reservoirs also have slightly lower nascent expression than non-reservoirs (Supplemental Fig
519 5A). Next, we compared the relationship between the number of DBPs bound and nascent
520 expression levels (Fig 6C). Similar to what was observed for RNA sequencing and previous
521 studies(17,32) (Fig 3C), nascent transcription also has a significant ($R = .3$, $P < 2e-16$) positive
522 correlation with the number of DBPs bound at that promoter (Fig 6C).

523
524 Interestingly, we observed a subset of reservoirs that have many DNA binding events but do not
525 have nascent transcriptional activity. Specifically, we found 355 (25%) promoters with more than
526 7 and as many as 60 binding events that have neither nascent nor mature expression (PRO-seq
527 $TPM < 0.001$, Fig 6F). We refer to these reservoirs without nascent or mature transcription as
528 ‘ghosts’, as there is no presence of transcriptional activity. We also found 964 promoters with
529 more than seven binding events that had no mature expression but did have nascent
530 expression. These are referred to as ‘zombies,’ as there is some presence of activity.

531

532 We next investigated if the chromatin environment discriminates between ghost and zombie
533 promoters. We therefore retrieved ENCODE ChIP data from K562 for a euchromatic and
534 heterochromatic histone modification; Histone 3 Lysine 27 acetylation (H3K27ac) versus
535 Histone 3 Lysine 27 trimethylation (H3K27me3) respectively. To this end, we downloaded peak
536 files called in two independent replicates for each histone modification from ENCODE analysis
537 pipelines. To validate our re-analysis of these ChIP-seq experiments we first determined if
538 K27ac correlates and K27me3 anticorrelates with global nascent transcription as would be
539 expected. Indeed, we see that those promoters containing K27ac have increased nascent
540 expression ($P < 2e-16$, fold change = 4) (Supplemental Fig 5B). Similarly, we checked the trend
541 for K27me3 status (Supplemental Fig 5C). As expected, we see that promoters containing
542 K27me3 have lower nascent expression ($P < 2e-16$, Fold change = 0.3, Supplemental 5C).

543
544 Having validated that our analysis of PRO-seq faithfully represents known biological processes
545 (e.g., K27ac enriched with higher expression) we wanted to zoom in only on reservoirs. We first
546 compared K27ac status versus nascent transcription levels on reservoirs. As was seen with all
547 promoters we see a significant difference in nascent expression between K27ac containing
548 reservoirs and those without that mark ($P < 0.0006$, fold change = 1.65, Supplemental Fig 5D).
549 Similarly, K27me3 status on reservoirs is negatively associated with nascent expression levels
550 ($P < 0.0002$, fold change = 0.55, Supplemental Fig 5E). However, chromatin environment
551 doesn't fully explain the presence of zombie promoters, as there are promoters with and without
552 nascent expression in each category of chromatin state.

553
554 To understand the difference between ghosts and zombies, we compared DBP binding events,
555 the distribution of nascent transcription, and histone marks. We did not observe a significant
556 difference in distribution of DBPs between ghosts and zombies ($P = 0.064$, fold change = 1.04,
557 Fig 6F, Supplemental Fig 5F). Thus, unlike all other cases tested, the number of DNA binding

558 events cannot account for the difference in those that do and don't have nascent expression.
559 Collectively, these findings demonstrate that more than 60 DBPs bound to the same promoter
560 do not exhibit nascent nor transcript production and are 'ghosted by Pol II'. All properties
561 identified above can be found in S2 Table.
562

563 **Discussion**

564 A fundamental question in biology is to understand when and where DBPs localize on a given
565 promoter and in turn how these combinations affect expression output. Thanks to heroic efforts
566 by ENCODE and other genome consortium efforts we now have standardized DNA binding
567 profiles for hundreds of DBPs[12–15,31]. Moreover, these datasets go through several quality
568 control measures before being released by ENCODE (see ENCODE portal). Thus, these
569 important resources provide two opportunities: one for data-reproducibility standard
570 advancements based on such well documented data; and a second to re-analyze these data-
571 sets to find novel insights into the genome-wide localization of DNA binding proteins.

572
573 This study found a vast majority of ENCODE data to be highly reproducible -- both with known
574 biology and in data quality. However, we do note that it may be recommended to be sure
575 replicates have reproducible peak profiles as we observed a few ChIP-seq experiments that did
576 not have any overlapping replicate peaks. This led us to identify 5 (2%) experiments that did not
577 have any reproducible peaks. However, a majority of the experiments (98%) have peaks that
578 overlap in all replicates as applied in this study. Moreover, taking into account the number of
579 observations (promoters) it is needed to be sure there are sufficient replicable peaks called for
580 each DBP. We found 30 more samples that had fewer than 250 peaks between replicates (14th
581 quartile). Considering the number of observations (promoters) it is also important to be sure

582 there are sufficient peak numbers for permutation analysis and statistical comparisons. Finally,
583 we noted that many of the proteins tagged with “eGFP” had similar binding profiles based on the
584 tag and not DBP function (Fig 2C). We did not see differences in number or sizes of peaks
585 compared to antibody-based ChIP. Yet it is surprising that 15 different DBPs all cluster together
586 based on the “eGFP” tag despite diverse biological roles and all having similar consensus peak
587 profiles.

588

589 These large and standardized data-sets also provide a unique opportunity to search for novel
590 insights into the relationship of DBPs and expression output. Thus, we can compare 161 DBPs
591 from the perspective of a promoter to determine how many bind and how this influences
592 promoter output. Consistent with two recent studies using orthogonal datasets and
593 approaches[17,18] we found that the more DBPs at a given promoter the more it tends to be
594 expressed. This was similar for lncRNA and mRNA promoters alike. This analysis similarly
595 validated these studies finding that mRNA promoters are more enriched in general than
596 lncRNAs for DBPs[17,18].

597

598 Surprisingly, we observed 1,362 promoters had numerous DBPs (more than seven and up to
599 111 DBPs on one promoter) bound yet did not have expression output. In fact, these promoters
600 had similar DBP events as the most highly expressed mRNA promoters. We termed these
601 regions reservoirs as they seem to be a holding spot for DBPs. Notably, reservoirs are highly
602 over-represented for lncRNA promoters relative to mRNA promoters ($p < 2 \times 10^{-12}$). We also
603 determined that reservoirs are not super-enhancers previously defined by having many DBP
604 binding events. Unlike super-enhancers, reservoirs have many different DBPs bound rather
605 than many binding events of cell-specific transcription-factors in a defined region[33–36,41].
606 Another difference from super-enhancers is the lack of Pol II, although we do find that a quarter

607 of reservoirs do have Pol II machinery bound. Perhaps suggesting that they are “paused
608 promoters”[39,42] potentiated with up to 111 DBP binding events.

609
610 Further investigation into reservoirs revealed that almost half produced “nascent” transcription
611 as measured by PRO-seq. This is consistent with the above hypothesis of paused promoters.
612 What is more surprising is that half of the reservoirs also did not produce nascent transcripts
613 within 6Kb of the TSS (ghosts). The distribution of number of DBPs was not different between
614 poised and ghost promoters. Nor could we find enrichment of specific DBPs that separate these
615 categories. Another possibility is that ghosts are positioned in a three-dimensional space with
616 “DBP” hubs[43,44]. Finally, it could be that the large number of binding events at these
617 promoters causes a ‘liquid phase state transition’ owing to so many proteins in a confined
618 space.

619
620 Our permutation-based approach to determine if a DBP prefers a genomic feature allowed us to
621 extend beyond promoters into the noncoding genome. Specifically, we were interested in
622 determining if certain DBPs were specific to repetitive elements, such as transposons, across
623 the genome. Comparing random permutation versus observed overlaps revealed something
624 somewhat surprising: that repeat classes and families such as ‘simple-repeats’ and tRNA
625 repeats were strongly enriched for all DBPs tested. In contrast, Line and Satellite repeats were
626 strongly depleted for all DBPs. Thus, some repeat sequences ‘repel’ DNA binding and some
627 ‘recruit’ DBPs without discretion.

628
629 In some cases, we did observe some interesting biases for DBPs and repeat elements. One
630 example is the human specific repeat family ‘SVA’ as one of the newest evolving repeats in
631 humans compared to primates. Specifically, three genes had a strong bias of binding SVA
632 elements -- all three of which are known transcriptional repressors. Recently studies have

633 identified that primate specific transposons can be co-opted to generate promoters of newly
634 evolving enhancers and even lncRNAs[26,45–47]. Thus, unlike many existing examples of co-
635 option in the case of SVA, it could have selective pressure for binding motifs of the observed
636 repressors and hitherto to unknown repressor motifs – or hitherto unknown promoter regulatory
637 elements.

638
639 Collectively, this exercise in data-science, reproducibility and scale in a singular cellular context
640 has been informative to understand relativistic promoter binding events across 161 DBPs. This
641 has led us to understand new features of the coordination of this binding with respect to
642 promoter expression output. Perhaps most importantly, 15 graduate students learned data-
643 sciences and reproducibility measures that not only provide new insight into reservoir promoters
644 but also a logical framework for future objective teaching exercises of genomic data-science.

645
646 All markdown files needed to reproduce the results and figures of this manuscript can be found
647 here: https://github.com/boulderrinnlab/CLASS_2020.

648

649 **Materials and Methods**

650 **Data, Code and Markdown**

651 Accessions and sample information for the DBPs included in this study can be found in S1
652 table. All data and analyses are publicly available on our GitHub:

653 https://github.com/boulderrinnlab/CLASS_2020.

654 All analyses, code, and compiled markdown are available in S1 File.

655

656

657 **Acknowledgements**

658 BioFrontiers IT team

659 Biochemistry department for allowing the class

660 Michael Snyder for his generosity in visiting our class and helpful feedback.

661 Biophysics training program (T32GM065103)

662 Signaling and cellular regulation training program (T32GM008759)

663

664

665 **References**

666

667 1. Lander ES, Linton LM, Birren B, Nusbaum C, Zody MC, Baldwin J, et al. Initial sequencing
668 and analysis of the human genome. *Nature*. 2001;409(6822):860–921.

669 2. Venter JC, Adams MD, Myers EW, Li PW, Mural RJ, Sutton GG, et al. The Sequence of the
670 Human Genome. *Science*. 2001;291(5507):1304–51.

671 3. Blat Y, Kleckner N. Cohesins Bind to Preferential Sites along Yeast Chromosome III, with
672 Differential Regulation along Arms versus the Centric Region. *Cell*. 1999;98(2):249–59.

673 4. Lieb JD, Liu X, Botstein D, Brown PO. Promoter-specific binding of Rap1 revealed by
674 genome-wide maps of protein–DNA association. *Nat Genet*. 2001;28(4):327–34.

675 5. Ren B, Robert F, Wyrick JJ, Aparicio O, Jennings EG, Simon I, et al. Genome-Wide Location
676 and Function of DNA Binding Proteins. *Science*. 2000;290(5500):2306–9.

677 6. Iyer VR, Horak CE, Scafe CS, Botstein D, Snyder M, Brown PO. Genomic binding sites of the
678 yeast cell-cycle transcription factors SBF and MBF. *Nature*. 2001;409(6819):533–8.

679 7. Weinmann AS, Bartley SM, Zhang T, Zhang MQ, Farnham PJ. Use of Chromatin
680 Immunoprecipitation To Clone Novel E2F Target Promoters. *Mol Cell Biol*. 2001;21(20):6820–
681 32.

682 8. Furey TS. ChIP-seq and beyond: new and improved methodologies to detect and
683 characterize protein-DNA interactions. *Nat Rev Genetics*. 2012;13(12):840–52.

684 9. Nakato R, Sakata T. Methods for ChIP-seq analysis: A practical workflow and advanced
685 applications. *Methods San Diego Calif*. 2020;

686 10. Park PJ. ChIP-seq: advantages and challenges of a maturing technology. *Nat Rev Genetics*.
687 2009;10(10):669–80.

- 688 11. Landt SG, Marinov GK, Kundaje A, Kheradpour P, Pauli F, Batzoglou S, et al. ChIP-seq
689 guidelines and practices of the ENCODE and modENCODE consortia. *Genome Res.*
690 2012;22(9):1813–31.
- 691 12. Consortium EP. A user's guide to the encyclopedia of DNA elements (ENCODE). *Plos Biol.*
692 2011;9(4):e1001046.
- 693 13. Consortium EP. An integrated encyclopedia of DNA elements in the human genome. *Nature*
694 *Publishing Group.* 2012;489(7414):57–74.
- 695 14. Birney E, Stamatoyannopoulos JA, Dutta A, Guigó R, Gingeras TR, Margulies EH, et al.
696 Identification and analysis of functional elements in 1% of the human genome by the ENCODE
697 pilot project. *Nature.* 2007;447(7146):799–816.
- 698 15. Consortium TEP. The ENCODE (ENCyclopedia Of DNA Elements) Project. *Science.*
699 2004;306(5696):636–40.
- 700 16. Ernst J, Kellis M. Chromatin-state discovery and genome annotation with ChromHMM. *Nat*
701 *Protoc.* 2017;12(12):2478–92.
- 702 17. Melé M, Mattioli K, Mallard W, Shechner DM, Gerhardinger C, Rinn JL. Chromatin
703 environment, transcriptional regulation, and splicing distinguish lincRNAs and mRNAs. *Genome*
704 *Research.* 2017;27(1):27–37.
- 705 18. Mattioli K, Volders P-J, Gerhardinger C, Lee JC, Maass PG, Melé M, et al. High-throughput
706 functional analysis of lincRNA core promoters elucidates rules governing tissue specificity.
707 *Genome Research.* 2019;29(3):344–55.
- 708 19. Mahony S, Edwards MD, Mazzoni EO, Sherwood RI, Kakumanu A, Morrison CA, et al. An
709 integrated model of multiple-condition ChIP-Seq data reveals predeterminants of Cdx2 binding.
710 *Plos Comput Biol.* 2014;10(3):e1003501.
- 711 20. Ewels PA, Peltzer A, Fillinger S, Patel H, Alneberg J, Wilm A, et al. The nf-core framework
712 for community-curated bioinformatics pipelines. *Nat Biotechnol.* 2020;38(3):276–8.
- 713 21. McInnes L, Healy J, Saul N, Großberger L. UMAP: Uniform Manifold Approximation and
714 Projection. *J Open Source Softw.* 2018;3(29):861.
- 715 22. McInnes L, Healy J, Astels S. hdbscan: Hierarchical density based clustering. *J Open*
716 *Source Softw.* 2017;2(11):205.
- 717 23. Wang H, Xing J, Grover D, Hedges DJ, Han K, Walker JA, et al. SVA Elements: A Hominid-
718 specific Retroposon Family. *J Mol Biol.* 2005;354(4):994–1007.
- 719 24. Ostertag EM, Goodier JL, Zhang Y, Kazazian HH. SVA Elements Are Nonautonomous
720 Retrotransposons that Cause Disease in Humans. *Am J Hum Genetics.* 2003;73(6):1444–51.
- 721 25. Savage AL, Bubb VJ, Breen G, Quinn JP. Characterisation of the potential function of SVA
722 retrotransposons to modulate gene expression patterns. *Bmc Evol Biol.* 2013;13(1):101.

- 723 26. Kelley D, Rinn J. Transposable elements reveal a stem cell-specific class of long noncoding
724 RNAs. *Genome Biology*. 2012;13(11):R107.
- 725 27. Henriques T, Scruggs BS, Inouye MO, Muse GW, Williams LH, Burkholder AB, et al.
726 Widespread transcriptional pausing and elongation control at enhancers. *Gene Dev*.
727 2018;32(1):26–41.
- 728 28. Calo E, Wysocka J. Modification of enhancer chromatin: what, how, and why? *Mol Cell*.
729 2013;49(5):825–37.
- 730 29. Cabili MN, Trapnell C, Goff L, Koziol M, Tazon-Vega B, Regev A, et al. Integrative
731 annotation of human large intergenic noncoding RNAs reveals global properties and specific
732 subclasses. *Genes & Development*. 2011;25(18):1915–27.
- 733 30. Cabili MN, Dunagin MC, McClanahan PD, Biaesch A, Padovan-Merhar O, Regev A, et al.
734 Localization and abundance analysis of human lincRNAs at single-cell and single-molecule
735 resolution. *Genome Biology*. 2015;16(1):20.
- 736 31. Derrien T, Johnson R, Bussotti G, Tanzer A, Djebali S, Tilgner H, et al. The GENCODE v7
737 catalog of human long noncoding RNAs: analysis of their gene structure, evolution, and
738 expression. *Genome Research*. 2012;22(9):1775–89.
- 739 32. Melé M, Mattioli K, Mallard W, Shechner DM, Gerhardinger C, Rinn J. Chromatin
740 environment, transcriptional regulation, and splicing distinguish lincRNAs and mRNAs. *Biorxiv*.
741 2016;088484.
- 742 33. Hnisz D, Abraham BJ, Lee TI, Lau A, Saint-André V, Sigova AA, et al. Super-enhancers in
743 the control of cell identity and disease. *Cell*. 2013;155(4):934–47.
- 744 34. Parker SCJ, Stitzel ML, Taylor DL, Orozco JM, Erdos MR, Akiyama JA, et al. Chromatin
745 stretch enhancer states drive cell-specific gene regulation and harbor human disease risk
746 variants. *P Natl Acad Sci Usa*. 2013;110(44):17921–6.
- 747 35. Whyte WA, Orlando DA, Hnisz D, Abraham BJ, Lin CY, Kagey MH, et al. Master
748 transcription factors and mediator establish super-enhancers at key cell identity genes. *Cell*.
749 2013;153(2):307–19.
- 750 36. Jiang Y, Qian F, Bai X, Liu Y, Wang Q, Ai B, et al. SEdb: a comprehensive human super-
751 enhancer database. *Nucleic Acids Res*. 2018;47(D1):D235–43.
- 752 37. Cardiello JF, Sanchez GJ, Allen MA, Dowell RD. Lessons from eRNAs: understanding
753 transcriptional regulation through the lens of nascent RNAs. *Biochem Soc Symp*. 2019;1–16.
- 754 38. Azofeifa JG, Allen MA, Hendrix JR, Read T, Rubin JD, Dowell RD. Enhancer RNA profiling
755 predicts transcription factor activity. *Genome Res*. 2018;28(3):334–44.
- 756 39. Core L, Adelman K. Promoter-proximal pausing of RNA polymerase II: a nexus of gene
757 regulation. *Gene Dev*. 2019;33(15–16):960–82.

- 758 40. Liao Y, Smyth GK, Shi W. The R package Rsubread is easier, faster, cheaper and better for
759 alignment and quantification of RNA sequencing reads. *Nucleic Acids Res.* 2019;47(8):e47–
760 e47.
- 761 41. Li Y, Rivera CM, Ishii H, Jin F, Selvaraj S, Lee AY, et al. CRISPR Reveals a Distal Super-
762 Enhancer Required for Sox2 Expression in Mouse Embryonic Stem Cells. *PLoS ONE.*
763 2014;9(12):e114485.
- 764 42. Core LJ, Lis JT. Paused Pol II captures enhancer activity and acts as a potent insulator.
765 *Genes & Development.* 2009;23(14):1606–12.
- 766 43. Rinn J, Guttman M. RNA Function. RNA and dynamic nuclear organization. *Science.*
767 2014;345(6202):1240–1.
- 768 44. Melé M, Rinn JL. “Cat’s Cradling” the 3D Genome by the Act of LncRNA Transcription.
769 *Molecular Cell.* 2016;62(5):657–64.
- 770 45. Kapusta A, Kronenberg Z, Lynch VJ, Zhuo X, Ramsay L, Bourque G, et al. Transposable
771 elements are major contributors to the origin, diversification, and regulation of vertebrate long
772 noncoding RNAs. *Plos Genet.* 2013;9(4):e1003470.
- 773 46. Chuong EB, Elde NC, Feschotte C. Regulatory evolution of innate immunity through co-
774 option of endogenous retroviruses. *Sci New York N Y.* 2016;351(6277):1083–7.
- 775 47. Cosby RL, Chang N-C, Feschotte C. Host–transposon interactions: conflict, cooperation,
776 and cooption. *Gene Dev.* 2019;33(17–18):1098–116.

777
778
779

780 Supporting Information

781
782 **S1 Fig.** (A) Distribution of number of consensus peaks observed for each DBP with cutoff at 15th
783 percentile shown as red line. (B) Permutation analysis of DBP significance of overlapping a
784 promoter versus 1,000 random samplings of the same peak profiles for each DBP genome
785 wide. Showing enrichment and depletion status for DBPs (Fisher Exact $P < 0.01$).

786
787 **S2 Fig.** UMAP dimensionality reduction based on DBP binding profiles and overlaid with: (A)
788 DNA binding domain annotations. (B) enrichment score on reservoir promoters (C) TF
789 annotation status (D) Median RNA-seq expression level of bound promoters. (E) Examples
790 promoter binding profile. Grey line indicates 95% confidence interval and black line is the mean
791 value. (F) Heatmap of each promoter binding profile for individual DBPs centered at TSS. Red
792 indicates degree of binding. Cluster of binding profiles for each DBP. The four clusters are
793 separated by white space. (G) Enrichment for each DBP at lncRNA and mRNA promoters

794 versus 1,000 random samplings of the same profiles for each DBP across the genome. Blue
795 indicates Z-score of observed versus permuted distribution.

796

797 **S3 Fig.** Heatmaps as in Fig 5 for all SVA elements in the human genome. (A) Histone
798 modifications (B) DBPs enriched at SVAs. (C) Expression of SVA elements relative to other LTR
799 containing endogenous retroviruses (ERVs).

800

801 **S4 Fig.** (A) X-axis, number of DBPs bound per promoter for all promoters. Y-axis is the
802 $\log_{10}(\text{TPM})$ expression of resulting transcript as measured by RNA-seq. (B) Cumulative
803 distribution of binding events on promoters. Red line indicates approximately the 50th percentile
804 of binding events occurring at 7 DBPs bound per promoter. (C) Stacked box plots of lncRNA
805 (red) and mRNA (black) promoters in reservoirs versus non-reservoirs. (D) Stacked box plots of
806 promoter types in reservoir (right) versus non reservoir (left) (E) Bar plot of the 25% of reservoir
807 promoters that have other promoters nearby. True equals a neighboring gene promoter is
808 expressed, False is not expressed. (F) X-axis is Chi-squared test value as
809 $\log_2(\text{observed/expected})$, Y-axis is the \log_{10} of Chi-squared P-value. (G) UMAP reduction using
810 only DBP binding to only reservoir promoters.

811

812 **S5 Fig.** (A) Density plot of nascent expression at reservoirs versus non-reservoirs. (B) Box plot
813 of nascent expression without (left) and with (right) H3K27ac modifications. (C) Same as (B) for
814 K27me3. (D-E) Same as (B) for reservoir versus non-reservoir promoters. (F) Boxplot of DBP
815 distribution at ghosts versus non-ghosts.

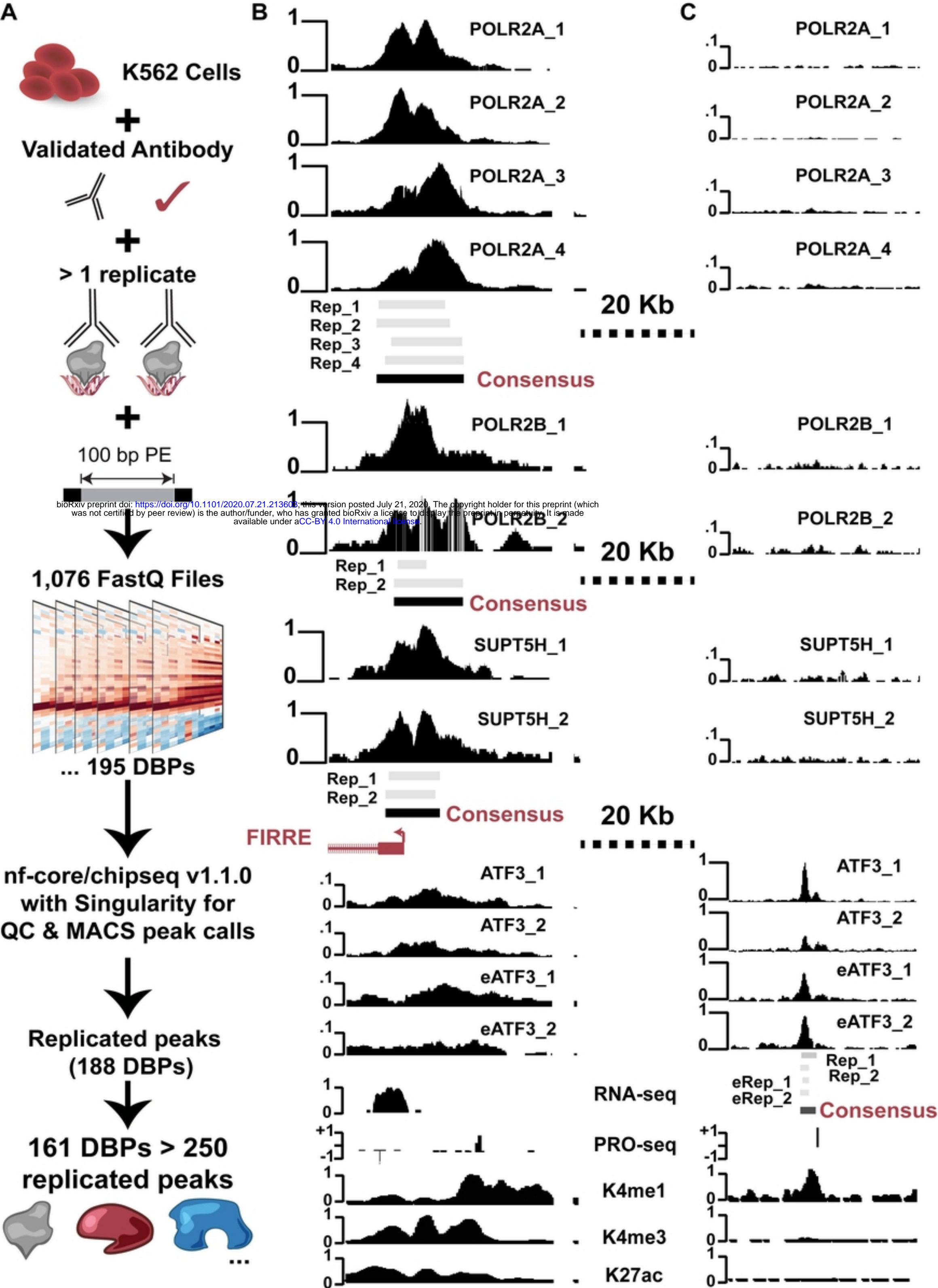
816

817 **S1 Table. Sample information for DNA binding proteins in study.**

818 **S2 Table. Promoter-level summary of DBP properties examined.** Each observation (row) is
819 a promoter and each column a variable investigated in this study.

820 **S1 File. All scripts used to analyze the data and produce figures.**

821



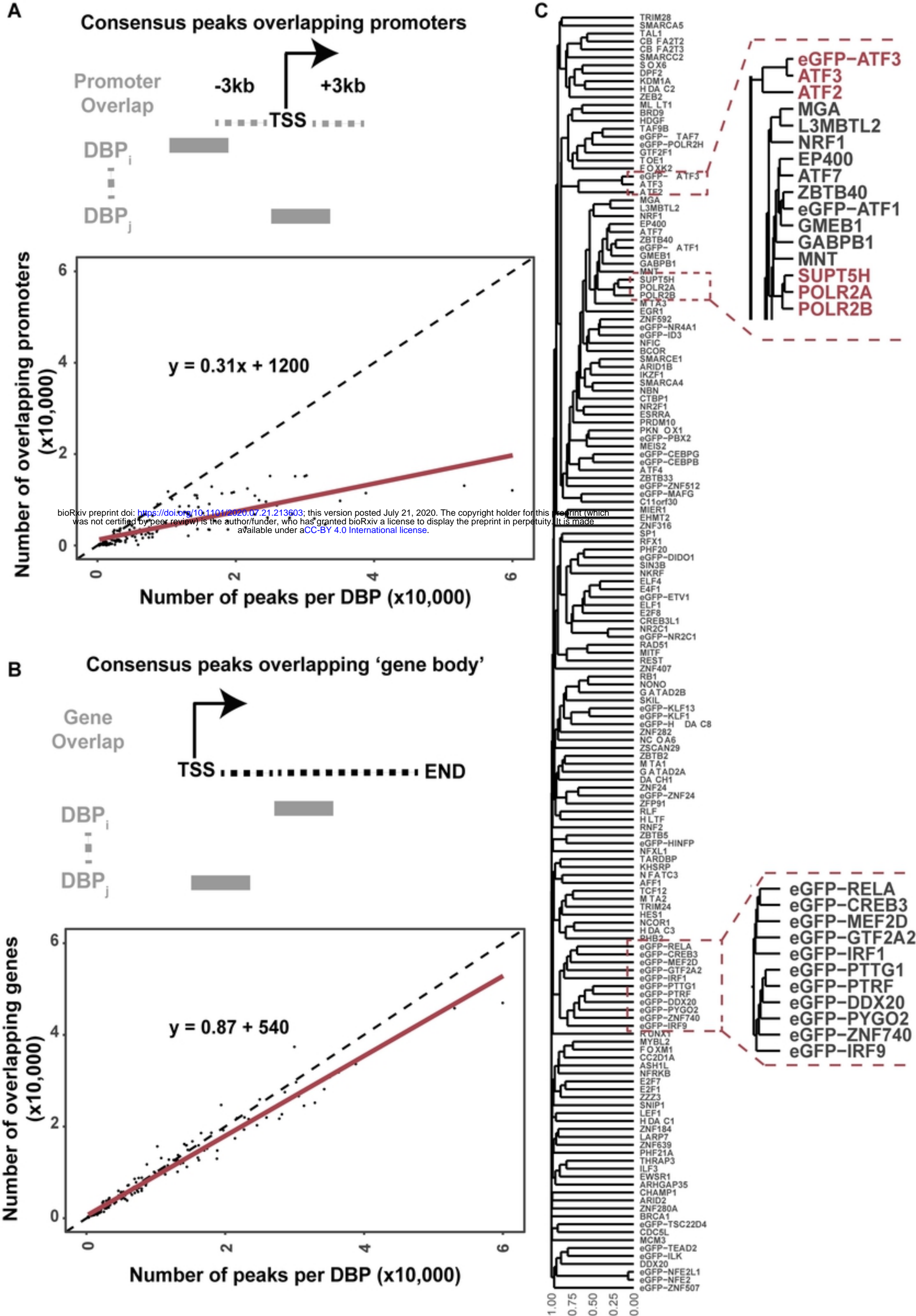
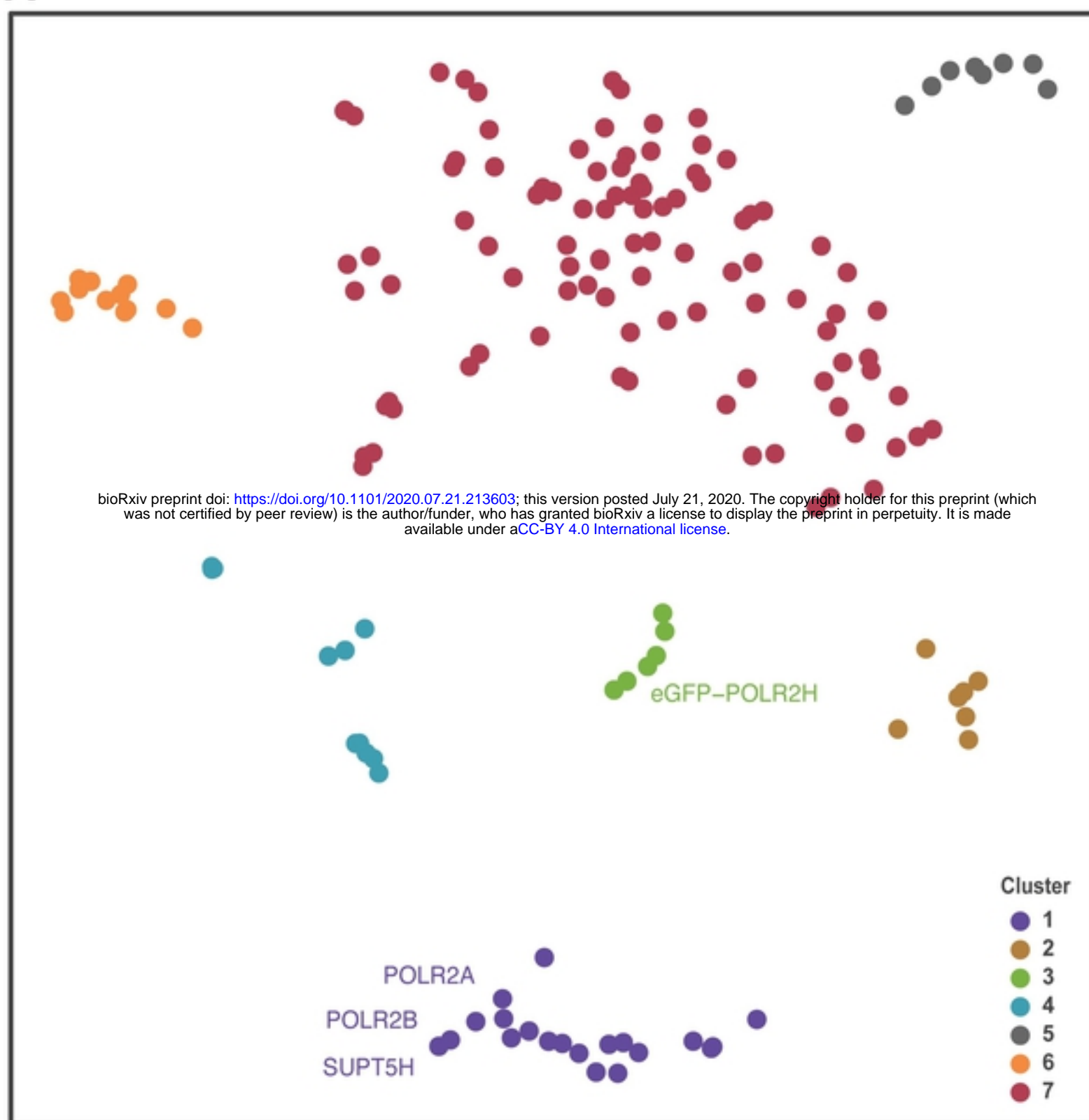
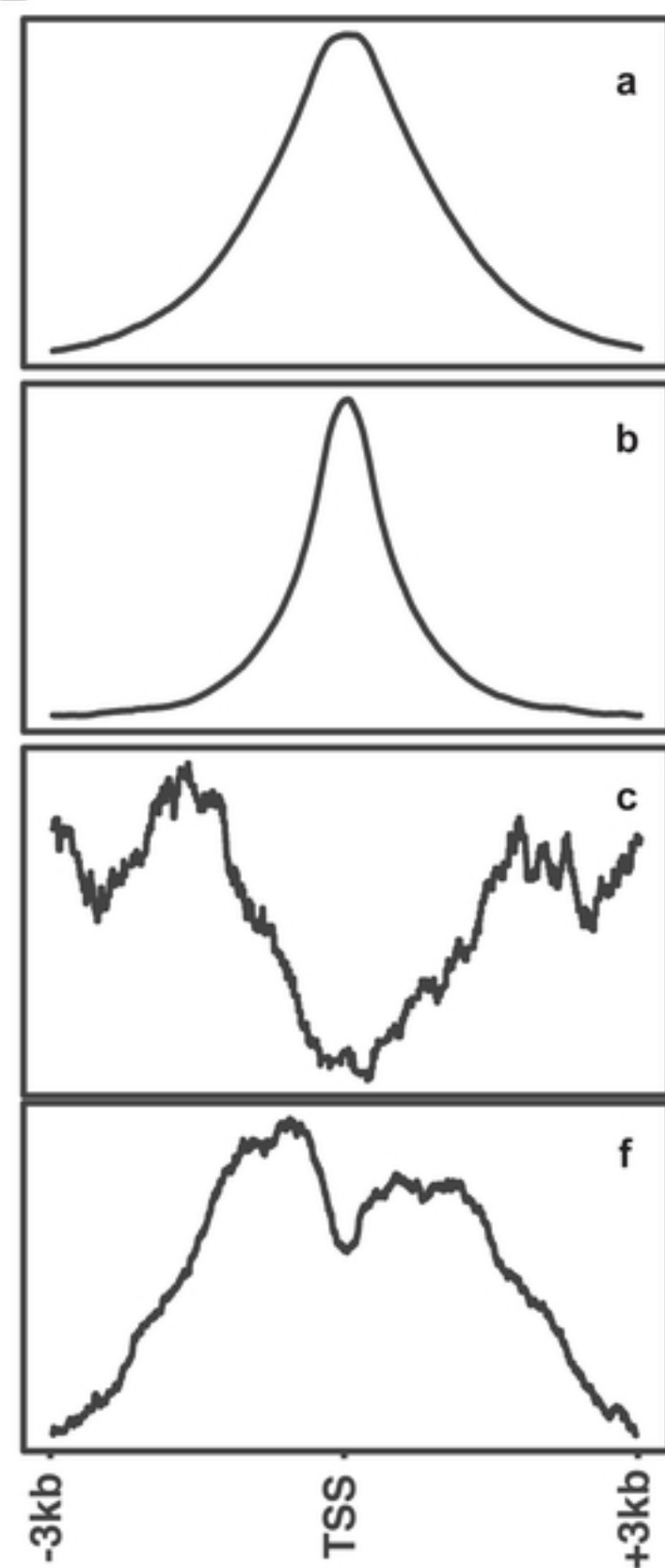


Figure 2

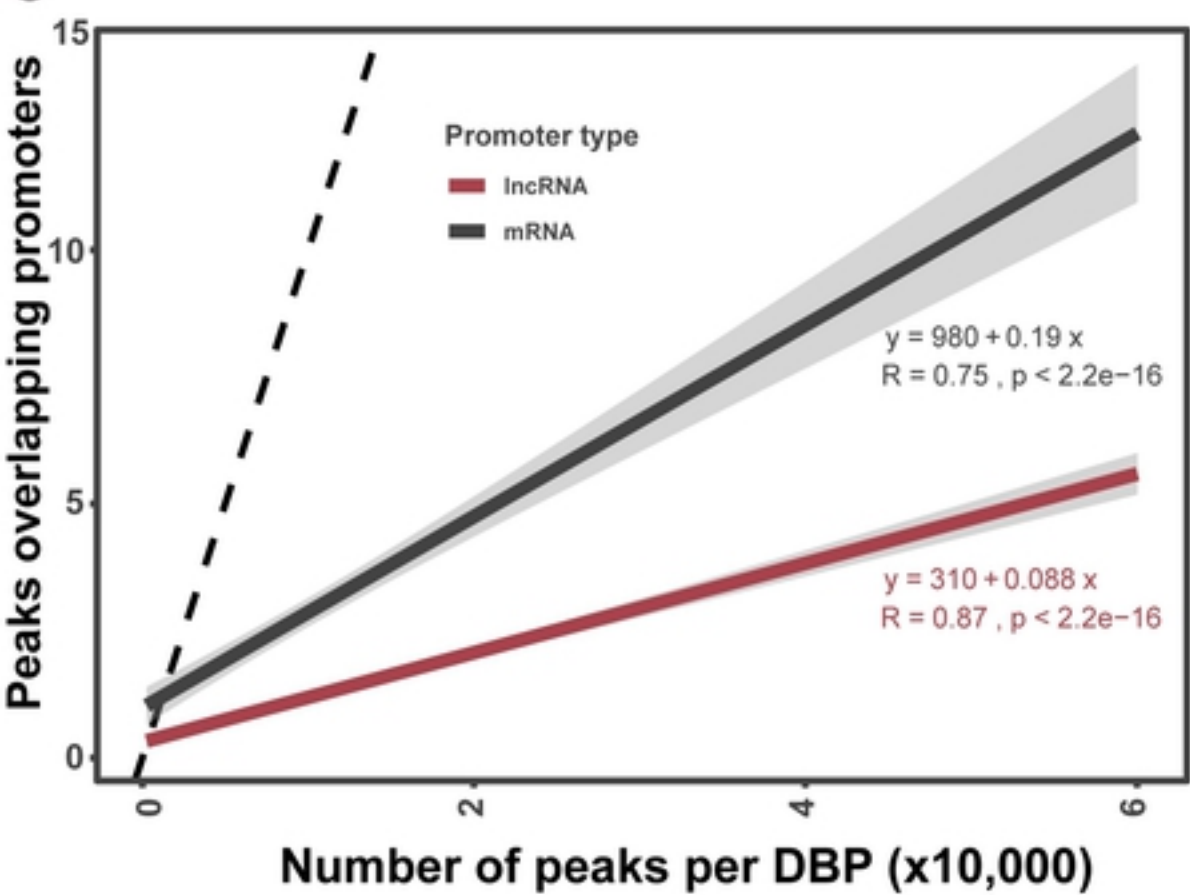
A



B



C



D

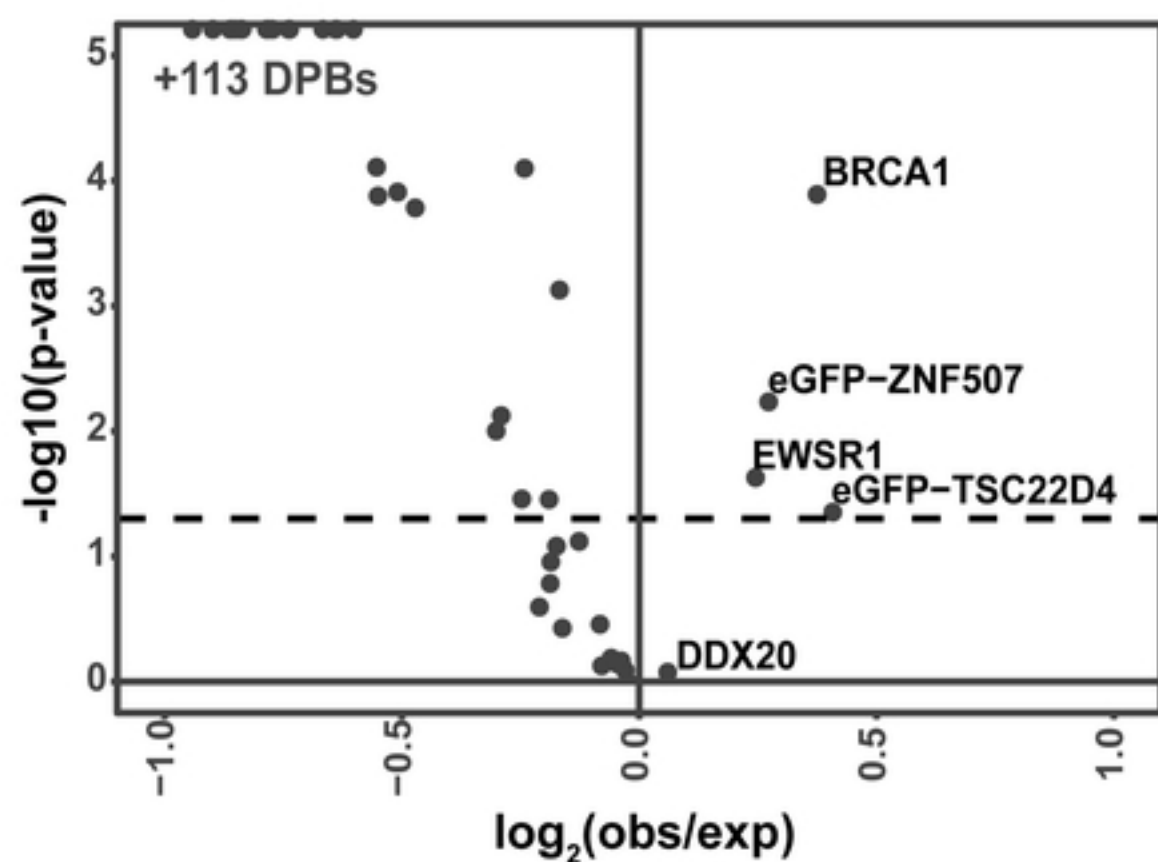


Figure 3

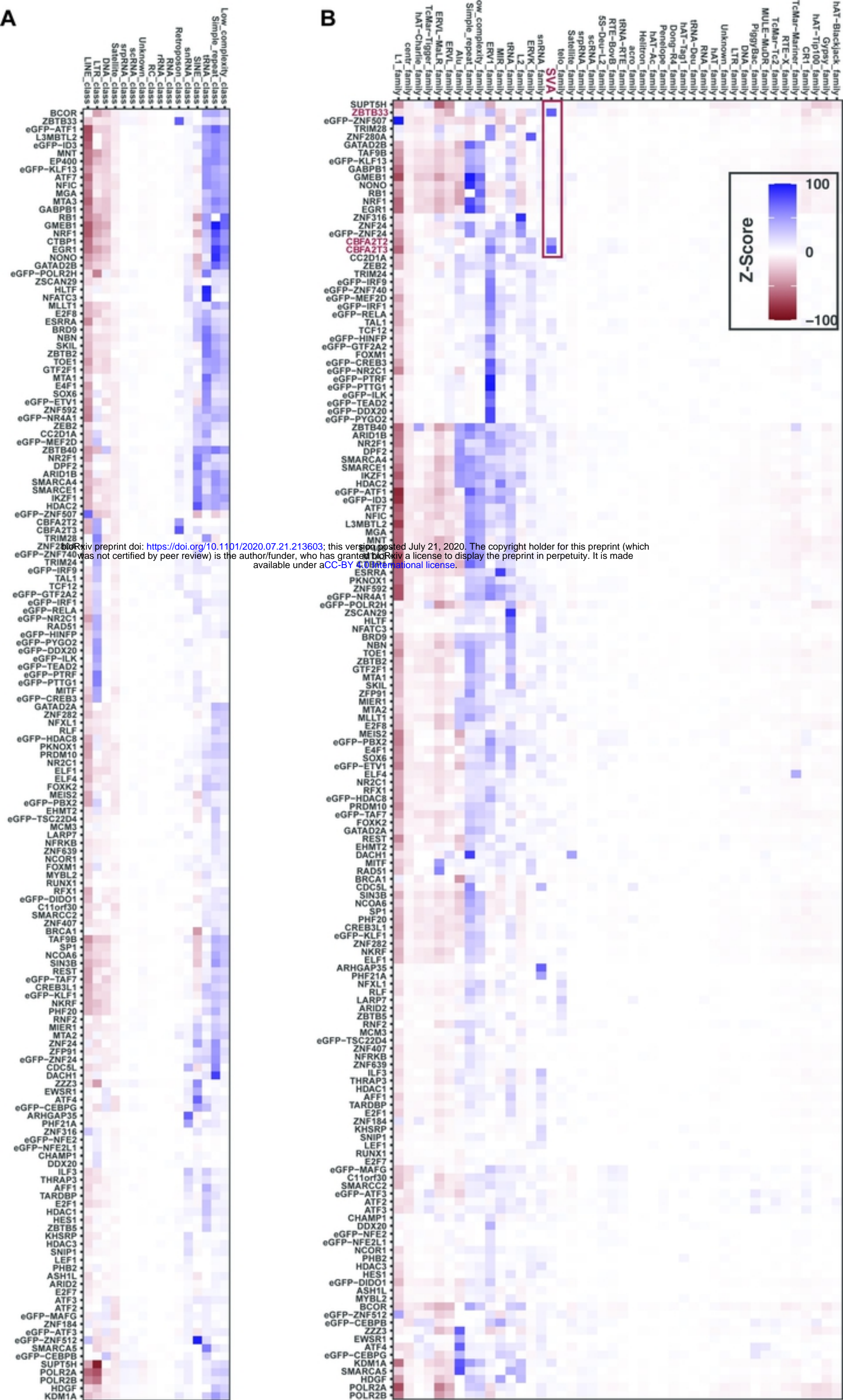


Figure 4

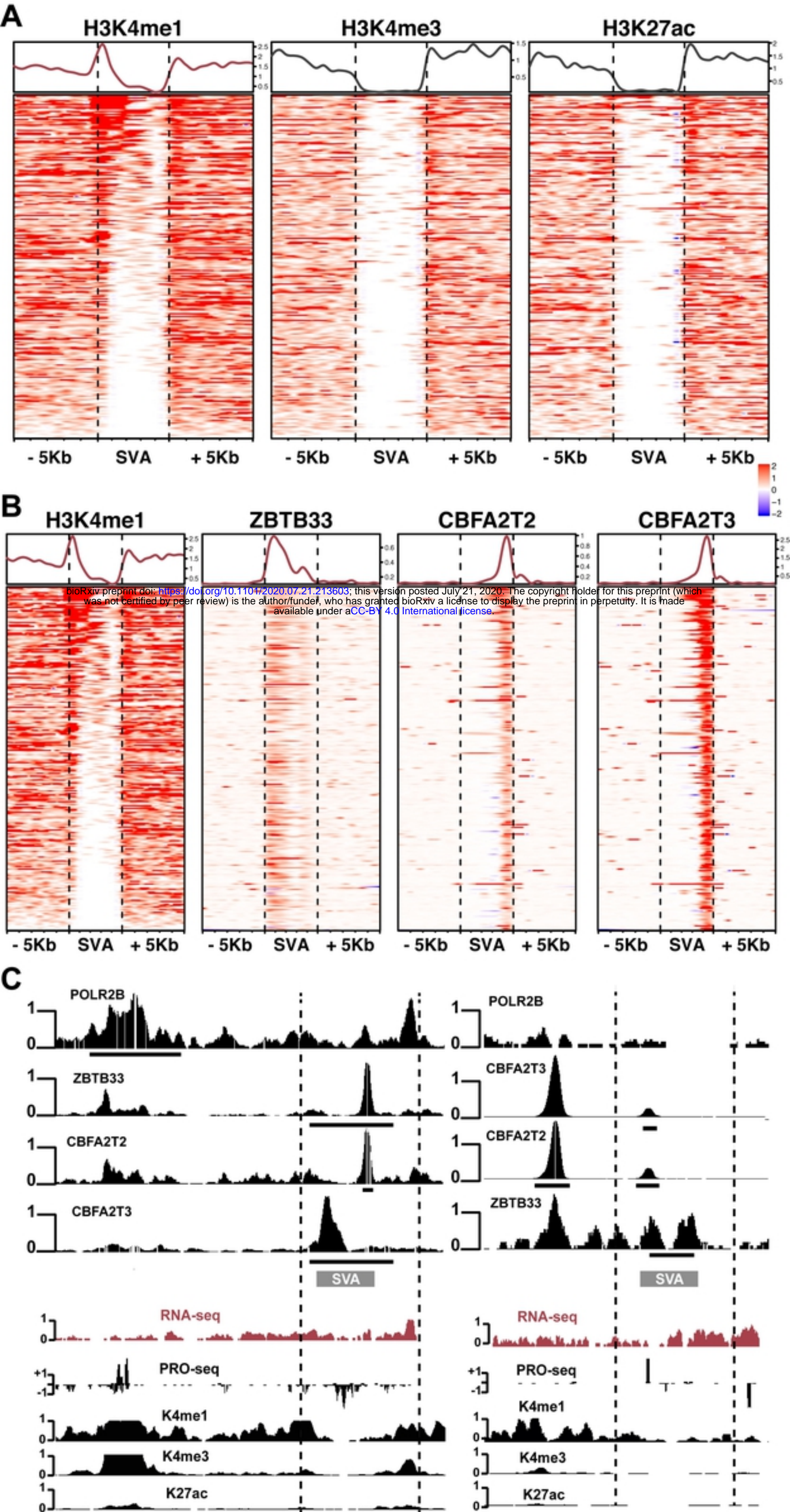


Figure 5

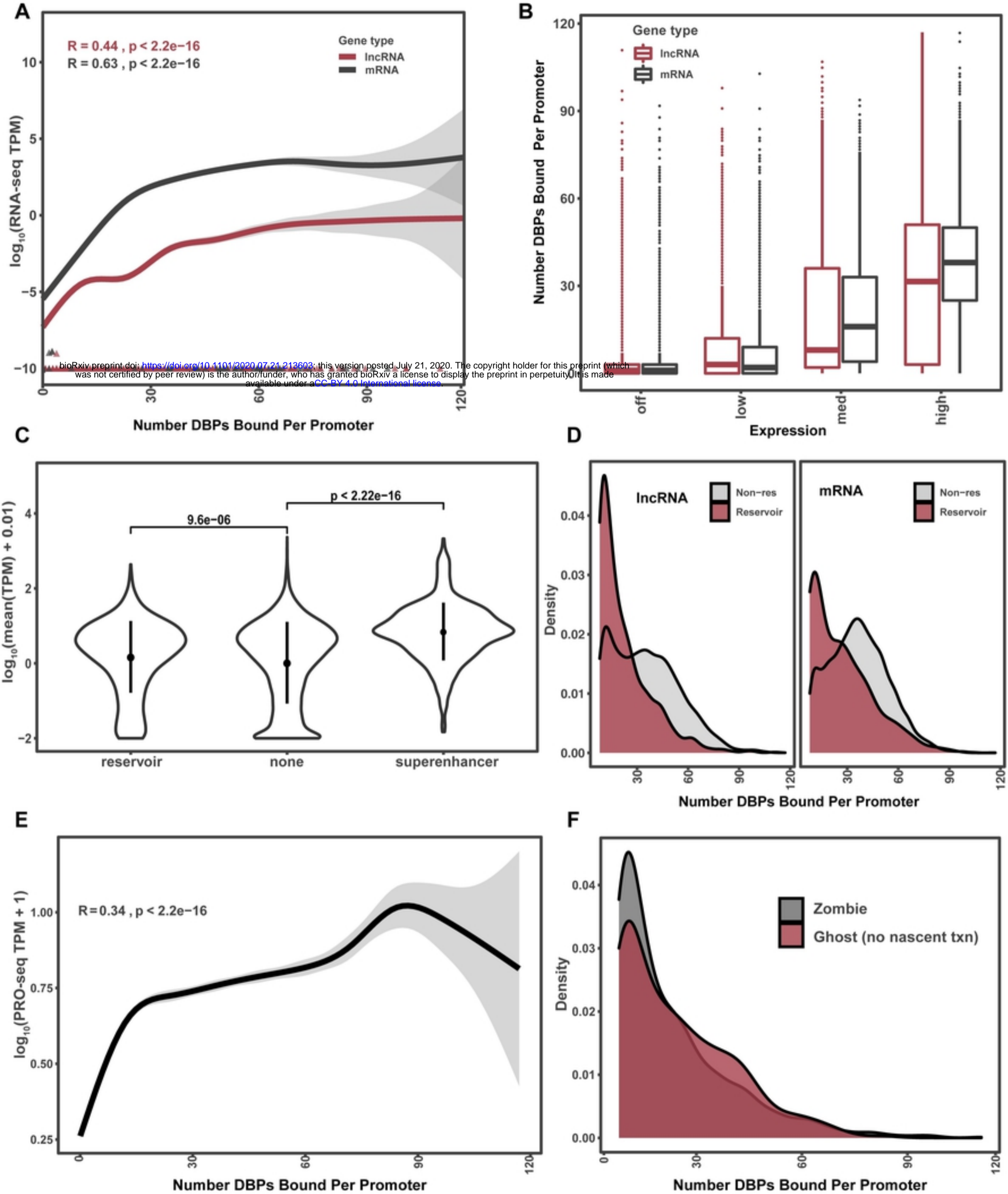


Figure 6