

1 **Epigenetic analyses of the planarian genome reveals conservation of bivalent**  
2 **promoters in animal stem cells.**

3  
4 Damian Kao, Yuliana Mihaylova, Samantha Hughes, Alvina Lai and A. Aziz Aboobaker

5  
6  
7 Department of Zoology, Tinbergen Building, South Parks Road, University of Oxford,  
8 Oxford OX1 3PS

9  
10 Damian Kao, [Damian.Kao@zoo.ox.ac.uk](mailto:Damian.Kao@zoo.ox.ac.uk)  
11 Yuliana Mihaylova, [Yuliana.Mihaylova@zoo.ox.ac.uk](mailto:Yuliana.Mihaylova@zoo.ox.ac.uk)  
12 Samantha Hughes, [Samantha.Hughes@han.nl](mailto:Samantha.Hughes@han.nl)  
13 Alvina Lai, [Alvina.Lai@zoo.ox.ac.uk](mailto:Alvina.Lai@zoo.ox.ac.uk)  
14 Aziz Aboobaker, [Aziz.Aboobaker@zoo.ox.ac.uk](mailto:Aziz.Aboobaker@zoo.ox.ac.uk)

15  
16

17 **Abstract**

18 **Background**

19 Planarian flatworms have an indefinite capacity to regenerate due to a population of  
20 pluripotent adult stem cells (neoblasts). Previous studies have suggested that they have  
21 features in common with both pluripotent mammalian embryonic stem cells and germ  
22 line stem cells. However, little is known about the importance of epigenetic regulation in  
23 these cells, which is likely to be crucial for neoblast biology and regeneration. We set out  
24 to develop analytical and experimental tools for planarians to allow the study of  
25 epigenetic marks in neoblasts and allow direct comparison of this model system with  
26 other animals.

27 **Results**

28 We developed an optimized ChIP-seq protocol for planarian neoblasts that allowed us to  
29 generate genome wide profiles for H3K4me1, H3K4me3 and H3K27me3. These were  
30 found to correlate as expected with genome wide expression profiles from analyses of  
31 planarian RNA-seq data. We found that many genes that are silent in neoblasts and  
32 then switch in post-mitotic progeny during differentiation have both H3K4me3 and  
33 H3K27me3 at promoter regions and are therefore bivalent. Further analysis suggested  
34 that bivalency is present at hundreds of loci in the pluripotent neoblast population.

35 **Conclusions**

36 We confirm that epigenetic regulation is key to neoblast biology and that bivalent  
37 promoters are not confined to vertebrate lineages, but may be a conserved feature of  
38 animal stem cells. Our work further establishes planarian neoblasts as a powerful model  
39 system for understanding the epigenetic regulation of pluripotency and regeneration.

40

## 41 **Background**

42 The potential use of stem cells in regenerative medicine has driven research into  
43 exploring the molecular mechanisms that govern stem cell potency, maintenance and  
44 differentiation. Despite this, we clearly still need to better understand the fundamental  
45 regulatory processes underpinning stem cell function, preferably using *in vivo* model  
46 systems. Highly regenerative planarians can be considered as a relatively simple stem  
47 cell study system that offers a large pool of adult stem cells called neoblasts (NBs).  
48 These cells are the driving force behind the almost limitless capacity to regenerate [1-3].  
49 With a simple bilaterian anatomy, the ability to study gene function using RNAi [4] and a  
50 growing list of well-defined markers [2,5-12], planarians make for a powerful model  
51 system for studying stem cell processes. There is growing evidence for the deep  
52 conservation of molecular regulation in stem cells across metazoans, particularly from  
53 work with planarians [10,13-19]. This suggests research on planarian NBs can also lead  
54 to new insights relevant to mammals.

55 Independent studies from multiple groups have characterized the transcriptome of NBs,  
56 using both bulk approaches [10,13,15,20,21] and single cell sampling and sequencing  
57 approaches [22-24]. The characterization of the transcriptome of NBs has relied on  
58 approaches that facilitate separation of NBs from the rest of the cells in the body or the  
59 targeted removal of NBs. This is followed by differential expression analysis to identify  
60 genes that have enriched NB expression. The most widely used of these approaches  
61 has involved the development of protocols for Fluorescence Activated Cell Sorting  
62 (FACS) that sort dissociated cells on the basis of nuclear to cytoplasmic ratio [25-27].  
63 This approach results in three distinct cell populations, an immediately radiosensitive  
64 population of  $>2N$  cells representing dividing cells in S/-phaseG2/M NBs (called the X1  
65 compartment), a population of  $\leq 2N$  cells that is partially radiosensitive and, over a longer  
66 time frame, containing G1 NBs and also transient undifferentiated post-mitotic progeny  
67 (called the X2 compartment), and a radio-resistant population of post-mitotic  
68 differentiated cells (called the Xins compartment) (see Figure 1A for summary). By  
69 sequencing these cell populations in bulk [9,10,13,21,24,28], performing single cell  
70 multiplex PCR analysis [24] or single cell sequencing [22,23] many groups have  
71 contributed to characterizing the NB transcriptome. These studies have identified novel  
72 NB markers, described gene expression heterogeneity in the NB population and found  
73 other cell type specific markers.

74 Another approach has been to compare expression profiles of whole animals with and  
75 without NBs, where NBs have been removed by ionizing radiation to kill all cycling cells  
76 [20] or by RNAi of a gene [15] to rapidly deplete NBs. These two approaches, combined  
77 with subsequent functional analyses, have allowed the development of more detailed  
78 models of the dynamics of NB proliferation, self-renewal and differentiation during  
79 regeneration and homeostasis [7,9,29-36]. One important feature that these studies  
80 highlight are chromatin-modifying factors involved in epigenetic regulation, the  
81 expression of which is enriched in planarian NBs [10,13,15].

82 A regulatory mechanism of key importance that mammalian embryonic stem cells  
83 (ESCs) and germ line stem cells (GSCs) share is histone modification-based bivalent  
84 control of promoters [37-42]. Bivalent promoters are characterized by the simultaneous  
85 presence of the repressive mark H3K27me3 and the activation mark H3K4me3 around  
86 transcriptional start sites (TSS) [37,38]. This bivalent promoter configuration is  
87 commonly seen on genes 'poised' for activation upon stem cell commitment and  
88 differentiation. In ESCs this is thought to allow pluripotency and the capacity to  
89 sensitively respond to developmental signals to achieve rapid differentiation when  
90 required. Bivalent promoters may achieve this by suppressing the formation of active  
91 RNA polymerase II complexes on one hand (hence 'poised'), and on the other, not allow  
92 other less easily reversible suppressive regulatory mechanisms, like DNA methylation,  
93 to silence genes [40,41]. Bivalent promoters have been described in mammals [37-  
94 39,42] and in zebrafish [43]. However, genome wide examples of promoter bivalency  
95 have so far not been found in invertebrates, for example no evidence for widespread  
96 bivalency has been observed in early *Drosophila* embryos [44]. These data support the  
97 interpretation that bivalency may be vertebrate specific, however more invertebrates  
98 clearly need to be studied and "it remains unclear how universal bivalent domains are  
99 across species" [43]. Given the role bivalency is thought to play in regulating  
100 pluripotency, planarian NBs are a logical place to look for bivalency in invertebrates.

101 To ask whether bivalency is present in NBs, we needed to combine transcriptomic and  
102 epigenetic analyses in the context of the genome. We first identified expressed loci on  
103 the genome and annotated transcribed regions using all available planarian RNA-seq  
104 data. We define the proportion of expression of every locus in X1, X2 and Xins FACS  
105 sorted cell populations allowing us to robustly identify genes silenced or expressed at  
106 very low levels in NB that are then actively transcribed during differentiation. We next

107 developed a robust ChIP-seq protocol for use with X1 NBs and demonstrate clear  
108 correlation of conserved epigenetic marks and gene expression based on the  
109 distribution of Histone 3 Lysine 4 mono-methylation (H3K4me1), tri-methylation  
110 (H3K4me3) and Histone 3 Lysine 27 tri-methylation (H3K27me3) marks. This revealed  
111 many bivalent promoters containing both H3K4me3 and H3K27me3 at similar levels,  
112 particularly at loci that go on to greatly increase their expression in post-mitotic progeny  
113 after asymmetric NB division. This provides strong evidence that this method of  
114 epigenetic regulation may in fact be conserved in animal stem cells. Overall, our work  
115 provides an essential annotation framework to study coding and non-coding loci in the  
116 genome, establishes a robust approach for ChIP-seq in NBs of *S. mediterranea* and  
117 reveals the potential for broad conservation of bivalent promoters in animal stem cells.

118

## 119 **Results**

### 120 **Establishing a genome wide annotation of transcribed loci in *S. mediterranea***

121 A growing number of different transcriptome studies have characterized the RNA-seq  
122 profiles of NBs and other planarian cells by sequencing either FACS-sorted cell  
123 populations (Figure 1A) or animals depleted of stem cells [9,10,13,15,20,21], but these  
124 have not been integrated. We performed a comparison of the results of four published  
125 transcriptomic studies, and found very poor overlap between those genes defined as  
126 having enriched expression in NBs (Figure 1B, Additional File 1).

127 To address this, we used the rapidly growing collection of publicly available *S.*  
128 *mediterranea* transcriptome sequence data to define a set of genome annotations on the  
129 current assembly of the asexual strain to serve as a basis for comparing transcriptional  
130 and epigenetic. This annotation then served as a basis for regulation across different  
131 planarian research projects. Distinct from previous annotations of the genome [45,46],  
132 this annotation includes all transcribed elements present across all available RNA-seq  
133 datasets (Additional file 2). As we have integrated all available RNA-seq data, our  
134 annotation should be particularly useful for describing potential non-coding RNAs and  
135 protein coding genes expressed at low levels, as these may not have been discovered  
136 by individual studies with limited numbers of reads and/or reliant on homology of protein  
137 coding exons.

138 Our annotation is markedly different from the current available annotation [46] of the *S.*  
139 *mediterranea* asexual genome sequence. We annotated more than 11,000 potential new

140 protein coding loci that are expressed at similar overall levels to previously annotated  
141 protein coding genes that were also present in our annotation. These new annotations  
142 were enriched for less well conserved proteins that may not be predicted by homology  
143 based annotation. A total of 6,300 existing annotations were not present in our  
144 expression driven annotation and further analysis of these MAKER specific genes shows  
145 that they generally have no or very low potential expression within the 164 RNA-seq  
146 libraries used for our annotations. (Additional File 2)

147 In summary, our annotation on the current planarian genome assembly shows regions of  
148 active transcription detected by current RNA-seq and transcriptome data, defines many  
149 more protein coding regions than currently available annotations and highlights a large  
150 number of non-coding transcribed loci. Additionally, it facilitates a consistent comparison  
151 specifically between *bona fide* transcriptional activity and the presence of post-  
152 translational histone modifications (ChIP-seq), allowing the relationship between  
153 epigenetic regulation and gene expression to be studied.

154

#### 155 **A genome wide expression profile of FACS-sorted cell populations.**

156 The ability to use FACS is a powerful experimental tool for working with *S. mediterranea*,  
157 providing convenient access to NBs and other cell populations. A growing number of  
158 studies have produced RNA-seq data for the different FACS cell populations that can be  
159 differentiated by nuclear to cytoplasmic ratio [9,10,13,21,24]. Given the discrepancies  
160 we uncovered between different studies that have taken this approach (Figure 1B,  
161 Additional File 1), we decided to reanalyze these datasets and newly available FACS  
162 RNA-seq data in public databases. Hierarchical clustering of the normalized expression  
163 values of each of these libraries revealed a rough congruence between different FACS  
164 cell populations (Additional File 3), and revealed greater heterogeneity among the X1  
165 samples than within the X2 and Xins samples. For example, some X1 samples clustered  
166 with X2 samples (Additional File 3).

167 These inconsistencies are potentially biases introduced by variation in the underlying  
168 biological or technical conditions. To mitigate against technical differences affecting  
169 absolute expression values we transformed absolute expression values into proportional  
170 expression values for each FACs compartment (Figure 1A, Additional File 4). For each  
171 locus, we divided each of the three X1, X2, Xins expression values by the sum of  
172 expression of all compartments for that locus (Additional File 4), obtaining a proportional

173 expression value for a total of 27,206 annotated loci that had at least 10 reads mapped  
174 in at least one FACS RNA-seq library, 18,010 of these are likely to be protein coding.  
175 The advantage of this transformation is that instead of using independent absolute  
176 expression values of the various samples, we can use the relationship among the three  
177 cell populations in each sample. Given similar FACS gating settings, results should be  
178 more consistent between labs despite any technical variations that affect absolute  
179 expression levels. Hierarchical clustering of these proportional values showed a  
180 consistent clustering of FACS samples by cell type, with good separation between  
181 clusters (Additional File 4). We then combined all available FACS RNA-seq data to  
182 reach one set of proportional expression values for each locus in our annotation. This  
183 gave us a new robust expression metric to compare every transcriptional unit in the  
184 genome across FACS cell compartments.

185 In order to achieve a visual representation of the data, we simply used a line to  
186 represent each gene, colored according to the proportion of its total expression in each  
187 of the three cell compartments. This allowed us to create genome wide expression  
188 spectra as an intuitive visualization and analysis tool. For example, we can sort all genes  
189 according to the proportion of their expression in X1 (S/G2/M NBs) (Figure 2A), X2 (G1  
190 NBs and stem cell progeny) (Figure 2B) and Xins (post-mitotic differentiated cells)  
191 (Figure 2C).

192 We used this approach to define which genes were expressed in each FACS  
193 compartment, dividing all genes expressed in FACS RNA-seq data into classes of  
194 enrichment (Figure 2D, Additional File 5). We confirmed our analysis by checking for the  
195 enrichment classes of genes previously described as being expressed in NBs (X1 and  
196 X1/X2 classes), in stem cell progeny (X2 and X2/Xins classes) and in differentiated cells  
197 (Xins class, Figure 3A-I, Additional File 5). In addition, we performed Gene  
198 Ontology(GO) enrichment analyses. We also identified enriched expression for genes  
199 not previously called as differentially expressed due to low levels of absolute expression  
200 in individual studies. An example of this is *Smed-tert*, the gene encoding the protein  
201 subunit of telomerase [47] that is amongst the most enriched X1 genes by proportional  
202 expression but does not appear in previous individual studies because of low absolute  
203 expression (Figure 3A). We also used the ESCAPE database of human pluripotency  
204 factors [48] and found 233 best reciprocal hits to *S. mediterranea*. Looking at the

205 expression of these genes we found them to be enriched in the X1/X2 expression  
206 category (Additional File 5).

207 Taken together our analyses, using all publicly available data, define the transcribed loci  
208 whose expression can be detected in planarian FACS compartments. As well as  
209 defining absolute levels of relative expression, we represent data by proportion of  
210 expression in each FACS compartment. This allows us to generate expression spectra  
211 highlighting loci expressed disproportionately in G2/M stem cells, loci expressed  
212 throughout the cell cycle, loci with most of their expression in transient differentiating  
213 post-mitotic cells and those expressed mainly in post-mitotic differentiated cells. As our  
214 annotations and expression data are in the context of the genome assembly these data  
215 can be integrated with ChIP-seq data.

216  
217 **Expression spectra are supported by RNA-seq of RNAi phenotypes and single cell**  
218 **sequencing analyses**

219 As an independent confirmation of our annotation and expression data we re-examined  
220 previously published RNA-seq after RNAi datasets and single cell RNA-seq data. For a  
221 selection of genes described as being required for stem cell progeny maintenance  
222 (Additional File 6) we visualized the RNA-seq profiles in relation to the defined FACS  
223 expression categories (Figure 2E) and observed that down-regulation of highly enriched  
224 X2 category genes was the main feature of both *Smed-mex3* and *Smed-zfp1* RNAi  
225 datasets (Additional File 6). From this data, it is straightforward to conclude that both  
226 *Smed-mex-3* and *Smed-zfp-1* have a collective effect on many genes that normally  
227 switch on in NB progeny as they differentiate and leave the cell cycle, and this correlates  
228 with the phenotypic effects of RNAi in both cases causing a depletion in stem cell  
229 progeny as stem cell differentiation fails [9,24]. This approach to analyzing RNA-seq  
230 data is useful for identifying patterns in the global effects of RNAi.

231 Recently, two planarian single cell transcriptomic studies have also been used to define  
232 expression profiles of single stem cells and other cell types [22,23]. These have helped  
233 to reveal heterogeneity of expression profiles in planarian stem cells and provide  
234 persuasive evidence for the existence of cycling NB cells that might be committed to  
235 particular lineages [23,24]. We re-mapped available single cell RNA-seq data [22,23] to  
236 identify the the top one thousand genes ranked by expression for each cell type defined  
237 by these two studies. We looked at the position of these genes along expression spectra



238 sorted by X1 proportion (Additional File 6). The single cell data analyzed in this way  
239 follows patterns we would expect and independently validates our proportional  
240 expression spectra. For different NB populations defined by single cell studies (sigma,  
241 gamma, zeta and head X1) we saw enrichment of genes in the X1 and X1/X2 categories  
242 (Additional File 6). Differentiated cell types were enriched for genes in the Xins category.  
243 However, all differentiated cell-types, with the exception of the 'epidermis II' class [22],  
244 have an enrichment of genes in the X2 category as these genes are amongst those with  
245 highest absolute expression in all non-NB cell types, and thus appear amongst the top  
246 1,000 expressed genes in single cell RNA-seq data.  
247 Overall, our annotation and expression analysis is congruent and compatible with  
248 independent data from RNAi coupled RNA-seq experiments and single cell sequencing  
249 data, further validating the success of our approach.

250  
251 **An optimized ChIP-seq protocol reveals H3K4me3 levels at TSSs in cycling cells**  
252 **correlate with gene expression in NBs.**

253 We next wished to combine our new genome annotation with predicted transcriptional  
254 start sites (TSSs) of expressed loci by cell compartment expression with NB derived  
255 epigenetic data. Research into epigenetic mechanisms in planarians is still very much in  
256 its infancy. Previous work characterized loss of function phenotypes of members of the  
257 NuRD complex [32,33,49,50], COMPASS and COMPASS-like families [51,52] and  
258 established a lack of endogenous DNA methylation in the *S. mediterranea* genome [32].  
259 With respect to monitoring epigenetic marks, some of the effects of *mll1/2* and *set1*  
260 RNAi on the H3K4me3 mark of active transcription have been previously reported [52]).  
261 We noted that in this study, the total number of ChIP-seq reads from ~1 million X1  
262 sorted planarian NBs were at relatively low numbers compared to those from *Drosophila*  
263 S2 'carrier' cells, which were added at 10x excess to X1 NBs (Additional File 7). These  
264 data suggested to us that ChIP-seq experiments with FACS sorted NBs might be very  
265 technically challenging.  
266 In order to begin to study epigenetic regulation of NBs, we first developed an optimized  
267 protocol for ChIP-seq on FACS sorted X1 cells for H3K4me3 mark. Relatively high levels  
268 of H3K4me3 have been shown to be broadly characteristic of active promoters [53,54].  
269 We found we were able to generate 13-26 million high quality *S. mediterranea* uniquely  
270 mapped reads using 150,000-200,000 X1 cells per immunoprecipitation, 5-7 times less

271 starting material compared to the a previous planarian ChIP-seq study [52]. We  
272 therefore used *Drosophila* S2 cells to act as a spike-in control for normalization of any  
273 technical replicate differences in immunoprecipitation across samples [55,56].

274 Sequencing X1 sorted cells, we observed high average H3K4me3 peaks around the  
275 TSSs of genes categorized as X1 and X1/X2 enriched, indicative of high expression in  
276 NBs (Figure 4A, Additional File 8). Conversely, we saw much less H3K4me3 at the  
277 promoters of Xins enriched gene. These results validate our planarian ChIP-seq method  
278 and confirm that our annotation is useful for studying global correlations between  
279 epigenetic marks and gene expression in the context of the genome. We saw  
280 intermediate levels of H3K4me3 in the X2 enriched compartment (Figure 4A, Additional  
281 File 8), which includes both NBs and recent post-mitotic progeny. A finer grained look at  
282 the X2 compartment revealed that genes with the highest proportion of X2 expression  
283 had lower levels of H3K4me3 in X1 cells (Figure 4B, Additional File 8), indicative of  
284 enriched expression in post mitotic progeny rather than cycling cells of the X1  
285 compartment (see also Figure 3E). The presence of H3K4me3 at X2 enriched gene  
286 promoters as a whole is, however, higher than that observed in genes enriched for  
287 expression in the differentiated Xins cell compartment (Figure 4B). These observations  
288 are broadly in agreement with previous findings from X1 cells [52] using the previously  
289 available annotations [45,46].

290 A base by base correlation analysis of ChIP-seq signal across the promoter region to  
291 proportional expression in the X1/X2/Xins FACS compartments shows a positive  
292 correlation between X1 proportional expression and levels of H3K4me3 deposition from  
293 near the TSS region to ~1kb downstream (Figure 4C). On the other hand, there is a  
294 negative correlation between H3K4me3 deposition and the proportion of Xins expression  
295 across the same region (Figure 4C). Thus, higher H3K4me3 ChIP-seq signal in X1 cells  
296 tends to reflect higher gene expression in X1 cells and lower H3K4me3 signal reflects  
297 lower expression in X1 cells and higher expression in the Xins compartment. We also  
298 looked at a individual loci of genes known to be expressed in NB and found them all to  
299 have relatively high levels of H3K4me3 and low levels of suppressive marks (Figure 4D).  
300 Overall, the patterns we observe across the genome are consistent with what would be  
301 expected with H3K4me3 being an activating mark. Additionally, it broadly validates our  
302 annotation of transcribed loci, our assignment of proportional expression values for each

303 locus to FACS compartments and our method of ChIP-seq using relatively small  
304 numbers of starting cells.

305  
306 **Levels of the repressive histone marks H3K27me3 and H3K4me1 at TSSs also**  
307 **correlate with gene expression in NBs.**

308 With an optimized ChIP-seq protocol, we decided to investigate two additional key  
309 histone modifications, the repressive mark H3K27me3 important for the assessment of  
310 bivalency [37,57] and H3K4me1 which has also recently been implicated as a repressive  
311 mark at promoter regions, mediated by the MLL3/4 family of histone methyltransferases  
312 [58].

313 We performed ChIP-seq on these two marks in X1 cells and observed ChIP-seq profiles  
314 consistent with these marks being associated with repression of gene expression in  
315 NBs. At loci enriched for X1 expression we observed low levels of H3K27me3 around  
316 the TSS and higher signal for loci with enriched expression in the Xins FACS  
317 compartment. (Figure 5A, Additional File 8). A positive correlation is observed around  
318 the TSS and 1 kb downstream between the levels of H3K27me3 and expression in the  
319 Xins compartment (figure 5B). This fairly broad domain of H3K27me3 is consistent with  
320 previous studies in mammals [59,60] A negative correlation at the TSS is observed  
321 between H3K27me3 signal and genes enriched for X1 expression (Figure 5B). Overall,  
322 this pattern is the opposite to that observed for H3K4me3.

323 ChIP-seq to detect distribution of the H3K4me1 mark revealed a different pattern to that  
324 of either H3K4me3 or H3K27me3. Rather than clear differences in the amount of  
325 H3K4me1 signal between loci with different FACS expression profiles, we observed a  
326 clear shift in the position of signal peaks (Figure 5C). Loci with a high proportion of  
327 expression in the Xins FACS compartment have high levels of H3K4me1 close to the  
328 TSS in X1 cells. Conversely, loci that are expressed in cycling cells (X1 and X1/X2  
329 enriched) have peaks of H3K4me1 signal on average ~1kb downstream of the TSS.  
330 Thus, the peak of H3K4me1 shifts away from the TSS for genes that are actively  
331 expressed and is consistent with observations of a previous study looking at H3K4me1  
332 levels at promoters in mammalian cells [58]. The Spearman correlation of H3K4me1  
333 signal and FACS proportional expression confirms these observations, showing a  
334 positive correlation close to the TSS for Xins enriched loci and a negative correlation for  
335 X1 enriched loci (Figure 5D). The relationship between H3K4me1 and X1 enriched loci

336 is positive further downstream, at which position, we therefore conclude, this  
337 modification does not broadly exert a repressive effect (Figure 5D). We noticed that for  
338 X2 enriched genes H3K4me1 signal had two distinct peaks, one around the TSS and the  
339 other downstream. This suggests two populations of loci, one with raised levels of  
340 H3K4me1 near the TSS suggesting repression, and the another further downstream  
341 suggesting an absence of repression involving H3K4me1 (Figure 5D). One clear  
342 possibility is that the repressive peak near the TSS might be for genes that are off in  
343 NBs and only switch in in post-mitotic progeny, while the other peak represents X2  
344 enriched genes that are expressed in cycling NBs. We also checked individual loci of  
345 genes known to be expressed in differentiated cells, and found they all had relatively  
346 high levels of repressive marks at or near the predicted TSS and low levels of H3K4me3  
347 (Figure 5 E). Conversely genes known to be expressed in NBs had low levels of  
348 repressive marks (Figure 4D).

349 A common method of analyzing ChIP-seq data is to perform a cluster analysis on  
350 coverage profiles to observe whether groups of similar profiles are enriched for a  
351 biological function [61]. While this blind approach to analyzing ChIP-seq profiles can  
352 sometimes yield interesting results when manually checking cluster members, it is often  
353 the case that the broad biological interpretation of clusters is vague due to low-resolution  
354 third party classifications such as gene ontology. Instead of a blind approach, here we  
355 have used proportional expression to categorize loci into distinct groups to observe  
356 broad trends in the ChIP-seq data. Taken together, our work demonstrates that the  
357 dynamics between states of promoter histone methylation are distinct between loci  
358 grouped by expression dynamics, and in agreement with previously studied roles of  
359 these marks described in mammalian cells [53,54,57-60]. The congruity of our  
360 annotation data, expression analysis and ChIP-seq datasets validates our framework for  
361 studying epigenetic regulation in NBs. As well as the genome wide analyses presented  
362 here, it will now be possible to look at the epigenetic regulation of individual planarians  
363 genes or sets of genes of interest in different experimental and environmental conditions  
364 using ChIPseq data.

365  
366 **Evidence for the conservation of bivalent promoter regulation in pluripotent**  
367 **animal stem cells**

368 Having validated our epigenetic analysis and demonstrated conservation of activating  
369 and suppressive marks we next investigated whether promoter of bivalency could be a  
370 regulatory mechanism in NBs. Bivalent promoters were originally observed at genomic  
371 loci for genes that were not expressed or expressed at very low levels in mouse ESCs  
372 [37], and were surprising because they contain both activating H3K4me3 and repressive  
373 H3K27me3 marks. This state is associated with the presence of RNA polymerase in a  
374 poised state and may allow rapid transcriptional responses to incoming signals to  
375 differentiate, at which point histone marks at bivalent promoters resolve so one of the  
376 two marks becomes dominant, resulting in activation or suppression of expression  
377 [40,41]. Bivalent promoters have since been described in various stem cells of different  
378 developmental origin and potency [39,42]. While they have been described outside of  
379 mammals in zebrafish [43], they have not so far been found in any invertebrates,  
380 suggesting they may be a novel epigenetic regulatory feature of vertebrates. Our ChIP-  
381 seq data from FACS sorted cells makes it possible to detect potential bivalent promoters  
382 in NBs if they are present in these cells. We reasoned loci that have relatively low  
383 expression in the X1 fraction and are up-regulated during differentiation and highly  
384 enriched in post-mitotic progeny (high X2 expression) may be good candidates for  
385 regulation by bivalent promoters in NBs.

386 We analysed the ChIP-seq signal as a continuous dataset by transforming the coverage  
387 profile into percent coverage by dividing the coverage at each base by the maximum  
388 coverage in the entire dataset. For each FACS category, we took the top one thousand  
389 most enriched loci and plotted the percent coverage profiles of H3K4me3, H3K4me1,  
390 and H3K27me3 to observe potential bivalency across all these loci. For the top one  
391 thousand X1 enriched loci, we see the expected profiles of a high H3K4me3 peak and  
392 low H3K27me3 peak (Figure 6A). We observe the opposite pattern for the top one  
393 thousand Xins enriched loci (Figure 6B). For the top one thousand X2 enriched loci,  
394 which are enriched for expression in post-mitotic progeny, we see similar percent  
395 coverage peaks for both H3K4me3 and H3K27me3 across these 1000 genes, consistent  
396 with potential bivalency in NBs at many of these promoters (Figure 6C).

397 As an independent source of validation, we also extracted all genes that were  
398 significantly down-regulated more than 2-fold after *Smed-mex3*(RNAi), which blocks the  
399 production of post-mitotic progeny. The ChIP-seq profile of these genes in X1 cells  
400 shows a similar profile to that of the top one thousand X2 enriched loci and is also

401 indicative of potential bivalency (Figure 6D). Genes expressed in the X2 compartment  
402 may stay on as cells differentiate so that they have an X2/Xins expression profile, some  
403 of these genes may also have bivalent promoters. Analysis of this gene sets also  
404 showed suggested some of these loci may be bivalent in X1 NBs (Figure 6E)

405 One caveat of our analysis so far is the possibly that bivalent ChIP signals represent  
406 underlying cell heterogeneity in the sampled X1 cell population [40]. While we know that  
407 cycling NB have some heterogeneity in gene expression that can describe subclasses  
408 with different lineage commitment [24], our focus on promoters of genes that are only  
409 upregulated upon differentiation and not expressed in NBs makes it unlikely the patterns  
410 we observe represent heterogenous epigenetic regulation in the NB population..  
411 Similarly, given that our analysis identified bivalency across large numbers of promoters  
412 it is also possible that our observation is the result of genes that have mostly one or  
413 other mark in NBs leading to an average profile that appears bivalent profile when many  
414 genes are looked at simultaneously. To rule this possibility out we looked at the  
415 correlation (Pearson) between H3K4me3 and H3K27me3 and observed the distribution  
416 of correlations for the top 500 ranked amongst X1, X2 and loci with reduced expression  
417 in *Smed-mex3*(RNAi) (Figure 6F). For X2 and *Smed-mex3* RNAi category loci, we  
418 observe a high density of well correlated H3K4me3 and H3K27me3 profiles indicating  
419 similar paired signals for these marks across the TSS, and indeed closer inspection of  
420 individual genes confirms this to be the case (Figure 7). For X1 enriched loci, we see a  
421 less correlation, with many negatively correlated loci compared to the top X2 enriched  
422 genes and *Smed-mex3* RNAi loci. This analysis suggests that many 100s of loci are in  
423 fact bivalent with respect to H3K4me3 and H3K27me3 in planarian NBs.

424 Overall, our data demonstrate the presence of bivalency at promoters in NBs. This  
425 suggests that this mechanism of gene regulation may be conserved amongst animals  
426 rather than confined to vertebrates [37,62]. It seems likely that the need to have both  
427 embryonic and, where appropriate, adult stem cells, capable of sensitive regulatory  
428 decisions and responses to incoming signals may have arisen very near the origin of  
429 multicellularity. Our work suggests that the evolution of bivalent promoters, arising  
430 earlier than previously thought, may have been an important component of achieving  
431 stem cell flexibility.

432

433 **Conclusion**

434 While there have been successful attempts in the model species *S. mediterranea* to  
435 integrate transcriptome data from different sources to improve overall representation and  
436 annotation [45,46,63-65], different FACS expression datasets from different experiments  
437 and laboratories have not been integrated to improve the quality of gene expression  
438 profiles across these cell compartments. Additionally, many previous approaches  
439 quantifying gene expression have focused on using assembled transcriptomes without  
440 the context of a genome assembly. This means that linking these RNA-seq based  
441 expression datasets directly to epigenetic or transcription factor based regulation using  
442 ChIP-seq is not possible. The goal of our work here was to address these deficits by  
443 combining transcriptome and epigenetic approaches to describe the landscape of  
444 epigenetic regulation at promoter regions in NBs in the context of expression level data.  
445 Our analyses validate our annotation, transcriptome analysis and ChIP-seq protocol and  
446 provide clear demonstration of the existence of bivalent promoters in cycling NBs. Our  
447 analysis is particularly sensitive for detecting genes that switch on after NB  
448 differentiation, due to the structure of the transcriptome and epigenetic datasets  
449 available for analysis. Future work can now use planarians as a model for understanding  
450 how this mode of regulation works, and the similarities and differences with vertebrates.  
451 The discovery that bivalent promoters exist outside of vertebrates adds to the growing  
452 body of evidence that suggests a deeper conservation of stem cell biology amongst  
453 animals than previously appreciated. Previously, endogenous genome stability  
454 mechanisms, splicing and post-transcriptional regulatory mechanisms have all been  
455 shown to be important for NB function [17,28,66]. Additionally, a number of proteins  
456 involved in epigenetically mediated gene regulation have also been shown to be  
457 essential to maintain NB function [10,13,15]. Particularly the previously described cases  
458 of MBNF/CELF mediated splicing regulation [17] and PIWI mediated genome stability  
459 [19], these represent deeply conserved processes that likely mediated stem cell function  
460 in an ancestral animal. Our work suggests that bivalent promoters represent yet another  
461 major conserved mechanism and this regulatory process is not, as previously thought,  
462 vertebrate specific. As well as demonstrating bivalency, our work, through establishment  
463 of an annotation framework and a robust ChIP-seq protocol for NBs, will allow the use of  
464 planarians as a model for epigenetic regulation of stem cell function. For example, the  
465 accessibility of the NB population should allow identification of regulatory targets of

466 chromatin modifying enzymes responsible for pluripotency, self-renewal and  
467 proliferation.

468

## 469 **Materials and methods**

470

### 471 **Data Sources for this study**

472 The NCBI Project accession number for ChIP-seq data produced in this study is  
473 PRJNA338116. All accession numbers for previously published RNA and ChIP-seq data  
474 used in the study are listed in Additional file 9.

475

### 476 **Flow cytometry**

477 A modified version of a planarian FACS protocol [27] was used. The modifications were:  
478 a 35 µm mesh filter was used instead of 100 µm, staining with Hoechst and calcein was  
479 performed simultaneously rather than sequentially and the centrifugation-wash step was  
480 omitted. We used Hoechst 34580 instead of Hoechst 33342. A FACS Aria III machine  
481 equipped with a violet laser was used for cell sorting. BD FACSDiva and FlowJo  
482 software were used for analyses and setting cell population gates.

483

### 484 **ChIP-seq**

485 For each experimental replicate 600,000-700,000 planarian X1 cells (enough for Chip-  
486 seq of 3 histone marks and an input control sample) were FACS-sorted (using 3-day  
487 regenerates) in PBS and pelleted at 4 °C. The pellet was re-suspended in nuclei  
488 extraction buffer (0.5% NP40, 0.25% Triton X-100, 10mM Tris HCl pH 7.5, 3mM CaCl<sub>2</sub>,  
489 0.25mM sucrose, 1mM DTT, 1/10<sup>th</sup> Phosphatase Cocktail Inhibitor 2 (Sigma Aldrich),  
490 1/10<sup>th</sup> Phosphatase Cocktail Inhibitor 3 (Sigma Aldrich)). This was followed by  
491 formaldehyde fixation, that was stopped with 2.5M glycine. The pellet was re-suspended  
492 in SDS lysis buffer (1% SDS, 50mM Tris HCl pH 8, 10mM EDTA) and incubated on ice.  
493 ChIP dilution buffer (0.1% SDS, 1.2mM EDTA, 16.7mM Tris HCl pH 8, 167mM NaCl,  
494 1/1000<sup>th</sup> Phosphatase Cocktail Inhibitor 2, 1/1000<sup>th</sup> Phosphatase Cocktail Inhibitor 3,  
495 1mM DTT) was added in a 2.3:1 ratio to the sample. Samples were sonicated and 1/10<sup>th</sup>  
496 volume 10% Triton X-100 was added. Samples were pelleted at 4 °C and the  
497 supernatant kept for further processing. Test de-crosslinking was performed on 1/8<sup>th</sup>



498 volume of the chromatin solution to verify the DNA fragment range following sonication  
499 was 100-500 bp.

500 Protein A-covered Dynabeads (ThermoFisher) were used for immunoprecipitation (IP).  
501 The amount of reagent used was in a 1:2 ratio to the amount of chromatin used per IP.  
502 The Dynabeads were first pre-blocked with 0.5% BSA/PBS and re-suspended in 0.5%  
503 BSA/PBS (2.5 times their original volume) containing 7 µg of antibody per IP. ChIP  
504 grade antibodies used were anti-H3K4me3 (rabbit polyclonal; Abcam; ab8580), anti-  
505 H3K4me1 (rabbit polyclonal; Abcam; ab8895), anti-H3K27me3 (mouse monoclonal;  
506 Abcam; ab6002).

507 After overnight incubation of the Dynabeads at 4 °C, they were washed 3 times with  
508 0.5% BSA/PBS and re-suspended in 0.5% BSA/PBS, matching their original volume.  
509 1/4<sup>th</sup> of the total chromatin was used for each IP, leaving a final 1/8<sup>th</sup> for input control  
510 libraries. The IP was done on a rotating wheel overnight at 4 °C.

511 Post-IP washes were done 6 times with RIPA buffer (50mM HEPES-KOH pH 8, 500mM  
512 LiCl<sub>2</sub>, 1mM EDTA, 1% NP40, 0.7% Sodium Deoxycholate, cOmplete protease inhibitors  
513 – 1 tablet per 50 ml). Beads were then washed in TE/NaCl (50mM NaCl in TE) and re-  
514 suspended in Elution Buffer (50mM Tris HCl pH 8, 10mM EDTA, 1% SDS). Proteins were  
515 separated from the beads via 15-minute incubation at 65 °C on a shaking heat block  
516 (1400 rpm). Supernatant and input samples underwent overnight heat-based de-  
517 crosslinking at 65 °C. RNaseA (0.2 µg/ml) and Proteinase K (0.2 µg/ml) were used for 1  
518 hour each in order to remove residual RNA and protein. DNA was purified with  
519 phenol:chloroform extraction and ethanol precipitation. DNA was re-suspended in TE  
520 and quantified with Qubit ds DNA HS kit (Thermo Fisher Scientific). The NEBNext Ultra  
521 II (NEB) kit was used for library preparation. Manufacturer's instructions were followed.  
522 Library clean-up was performed with Becton Coulter AMPureXP beads. Libraries were  
523 quantified with Qubit, Agilent Bioanalyzer and using a KAPA Library Quantification  
524 qPCR kit. Libraries were sequenced on an Illumina NextSeq machine.

525

## 526 **Comparison of previous NB transcriptomes**

527 Independently assembled transcriptomes were downloaded from four previous  
528 publications [10,14,20,21]. Transcripts enriched in NBs were extracted based on the  
529 classifications provided in respective publications' supplementary information. A  
530 clustering of these sequences was done by running CAP3 [67] on all transcripts and

531 then extracting transcript groups that assembled. Detailed methods are recorded in an  
532 IPython notebook (Additional File 9).

533

### 534 **Reference assembly and annotations**

535 Transcript sequences from previously assembled transcriptomes (Oxford, Dresden,  
536 SmedGD Asexual, SmedGD Unigenes) and known genes were downloaded from  
537 SmedGD [46], PlanMine [63] and NCBI. These sequences were mapped to the  
538 SmedGD Asexual 1.1 genome with GMAP [68]. PASA [69] was then used to consolidate  
539 the annotations. An independent reference assembly was also performed on 164  
540 available RNA-seq libraries with HISAT2 [70] for mapping and StringTie [71] for  
541 assembly. PASA consolidated annotations and StringTie reference assembly were  
542 merged together with StringTie.

543 To remove redundancy from the annotations we first calculated an intron jaccard  
544 similarity score (intersection of introns / union of introns) for all overlapping transcripts.  
545 Pair-wise jaccard similarity scores of 0.9 or more were kept and used to create a graph  
546 of similar annotations. Maximal-cliques were extracted from this graph as clusters of  
547 redundant annotations. From these cliques, we chose one transcript to be the  
548 representative by prioritizing transcript length, ORF length and BLAST homology. Strand  
549 information was assigned to each transcript by using strand specific RNA-seq libraries,  
550 BLAST homology, and longest ORF length. We ran TransDecoder (utilizing Uniprot and  
551 PFAM for coding evidence) [72] to identify protein coding transcripts. Detailed methods  
552 are recorded in an IPython notebook (Additional File file 9). The genome annotations are  
553 made available here as a gtf file (Additional File 10).

554

### 555 **Proportional expression value generation**

556 Kallisto [73] was used to pseudo-map RNA-seq libraries from four datasets [9,13,24],  
557 accession: PRJNA296017, generating estimated counts and TPM values for each  
558 transcript. Sleuth [74] was used to calculate a normalization factor for each library. For  
559 each locus, the TPM values of member transcripts were summed to generate a loci TPM  
560 value and then normalized accordingly.

561 Not all datasets contained all three X1/X2/Xins populations. The Reddien [24] and  
562 Sanchez (accession: PRJNA296017) datasets only had two of the three populations. In  
563 order to consolidate proportional expression values among all four datasets, pair-wise

564 ratios were first calculated for each dataset (X1:X2, X1:Xins, X2:Xins) using normalized  
565 TPM values. These ratios were then averaged across the datasets.

566 Using two of the three ratios, we can calculate a predicted third ratio (i.e., given X1:X2  
567 and X1:Xins, we can calculate X2:Xins). We then correlate the calculated proportion with  
568 the actual proportion and kept the pairs of actual proportions (in this case, X1 and Xins)  
569 that had the best correlation with the calculated proportion. Detailed methods are  
570 recorded in an IPython notebook (Additional File 9).

571

## 572 **Single cell RNA-seq analysis**

573 Single cell RNA-seq data were downloaded from short-read archive [22,23]. Reads were  
574 pseudo-mapped with Kallisto and the TPM values were used for down-stream analysis.  
575 Cell types of each RNA-seq library were previously defined in the respective publications  
576 by both FACS (X1/X2/Xins) and by gene markers.

577 The top 1,000 expressed loci from each cell type cluster were used for generating the  
578 spectrum density figure and ternary plots.

579

## 580 **ChIP-seq mapping and track generation**

581 Biological triplicate ChIP-seq data from planarian X1 cells for each of three histone  
582 marks considered was analyzed in conjunction with *D. melanogaster* S2 spike-in cells,  
583 used for downstream between IP replicate normalization. The trimmed reads were  
584 mapped to both *S. mediterranea* asexual 1.1 genome (SmedGD, [46]) and *D.*  
585 *melanogaster* r6.10 genome [75] with BWA mem 0.7.12 [76]. Only uniquely mapping  
586 reads were considered further. Paired reads that map to both species were also  
587 removed. Picard tools 1.115 was used to remove duplicate reads. Reads were  
588 separated into the sets that mapped to *D. melanogaster* and *S. mediterranea*  
589 respectively so that numbers of mapped reads could be used for downstream  
590 normalization calculations. For each paired or single mapped read, coordinates  
591 representing the 100bp at the center of the sequenced fragment were parsed and  
592 written to a BED file.

593 To generate coverage tracks in bedgraph format, the bedtools' genomcov function was  
594 used. A normalization factor was calculated using the number of mapped reads  
595 corresponding to the *D. melanogaster* spike-in [55,56]. A scaling factor for the input  
596 ChIP-seq libraries was calculated using the DeepTools [61] python API that utilizes the

597 SES method [77]. Mean normalized coverage was calculated for each sample and input.  
598 The normalized input coverage was then subtracted from the normalized sample  
599 coverage to generate the final coverage track for downstream visualization and analysis.  
600 The normalization process is detailed in an IPython notebook (Additional file 7)  
601 To calculate correlation of ChIP-seq coverage to proportional expression, two vector of  
602 values were used for a group of loci. The first vector is the proportional expression and  
603 the second vector is the coverage at a position in the 5kb region of the loci. A spearman  
604 correlation was performed on both vectors yielding a correlation value for the assayed  
605 position. This correlation value was calculated for every non-overlapping 50 base pair  
606 window in the 5kb region around the TSS.  
607 For bivalency, a percent coverage was used instead of the absolute normalized  
608 coverage. This was generated by calculating the maximum coverage across all 5kb  
609 regions around assay loci. Each absolute coverage value across the loci is then divided  
610 by the maximum coverage resulting in a percent coverage.

611

## 612 **Declarations**

### 613 Ethics approval and consent to participate

614 Not applicable

### 615 Consent for publication

616 Not applicable

### 617 Availability of data and material

618 All data produced in this study is in the form ChIPseq data submitted under The NCBI  
619 Project accession number PRJNA338116.

### 620 Competing interests

621 The authors declare that they have no competing interests

### 622 Funding

623 This work was funded by grants from the Medical Research Council (grant number  
624 MR/M000133/1) and the Biotechnology and Biological Sciences Research Council  
625 (grant number BB/K007564/1) to AA. A.G.L. is funded by a Human Frontier Science  
626 Program fellowship.

### 627 Authors' contributions

628 AA conceived designed the study. DK led and performed all data analysis with help from  
629 YM and AA. YM led the acquisition of all experimental data, with help from AGL. SH and  
630 YM optimized the ChIP-seq protocol. DK, YM and AA wrote the manuscript.

### 631 Acknowledgements

632 We thank members of the AA lab for comments on the manuscript.

633

### 634 **References**

- 635 1. Salo E, Baguna J. Regeneration and pattern formation in planarians. II. and role of  
636 cell movements in blastema formation. *Development*. The Company of Biologists Ltd;  
637 1989;107:69–76.
- 638 2. Wagner DE, Wang IE, Reddien PW. Clonogenic Neoblasts Are Pluripotent Adult Stem  
639 Cells That Underlie Planarian Regeneration. *Science*. 2011;332:811–6.
- 640 3. Newmark PA, Sánchez Alvarado A. Bromodeoxyuridine Specifically Labels the  
641 Regenerative Stem Cells of Planarians. *Developmental Biology*. 2000;220:142–53.
- 642 4. Sánchez Alvarado A, Newmark PA. Double-stranded RNA specifically disrupts gene  
643 expression during planarian regeneration. *Proceedings of the National Academy of  
644 Sciences of the United States of America*. National Academy of Sciences;  
645 1999;96:5049–54.
- 646 5. Eisenhoffer GT, Kang H, Alvarado AS. Molecular Analysis of Stem Cells and Their  
647 Descendants during Cell Turnover and Regeneration in the Planarian *Schmidtea  
648 mediterranea*. *Cell Stem Cell*. 2008;3:327–39.
- 649 6. Guo T, Peters AHFM, Newmark PA. A bruno-like Gene Is Required for Stem Cell  
650 Maintenance in Planarians. *Developmental Cell*. 2006;11:159–69.
- 651 7. Cowles MW, Brown DDR, Nisperos SV, Stanley BN, Pearson BJ, Zayas RM.  
652 Genome-wide analysis of the bHLH gene family in planarians identifies factors required  
653 for adult neurogenesis and neuronal regeneration. *Development*. 2013;140:4691–702.
- 654 8. Lapan SW, Reddien PW. *dlx* and *sp6-9* Control Optic Cup Regeneration in a  
655 Prototypic Eye. Desplan C, editor. *PLoS Genet*. Public Library of Science;  
656 2011;7:e1002226–13.
- 657 9. Zhu SJ, Hallows SE, Currie KW, Xu C, Pearson BJ. A *mex3* homolog is required for  
658 differentiation during planarian stem cell lineage development. *eLife*. eLife Sciences  
659 Publications Limited; 2015;4:304.
- 660 10. Labbé RM, Irimia M, Currie KW, Lin A, Zhu SJ, Brown DDR, et al. A Comparative  
661 Transcriptomic Analysis Reveals Conserved Features of Stem Cell Pluripotency in  
662 Planarians and Mammals. *STEM CELLS*. 2012;30:1734–45.
- 663 11. Scimone ML, Srivastava M, Bell GW, Reddien PW. A regulatory program for

- 664 excretory system regeneration in planarians. *Development*. 2011;138:4387–98.
- 665 12. Rink JC, Vu HTK, Alvarado AS. The maintenance and regeneration of the planarian  
666 excretory system are regulated by EGFR signaling. *Development*. 2011;138:3769–80.
- 667 13. Onal PO, n DGU, Adamidi C, Rybak A, Solana J, Mastrobuoni G, et al. Gene  
668 expression of pluripotency determinants is conserved between mammalian and  
669 planarian stem cells. *The EMBO Journal*. Nature Publishing Group; 2012;31:2755–69.
- 670 14. Adamidi C, Wang Y, Gruen D, Mastrobuoni G, You X, Tolle D, et al. De novo  
671 assembly and validation of planaria transcriptome by massive parallel sequencing and  
672 shotgun proteomics. *Genome Research*. 2011;21:1193–200.
- 673 15. Solana J, Kao D, Mihaylova Y, Jaber-Hijazi F, Malla S, Wilson R, et al. Defining the  
674 molecular profile of planarian pluripotent stem cells using a combinatorial RNAseq, RNA  
675 interference and irradiation approach. *Genome Biol*. BioMed Central Ltd; 2012;13:R19–  
676 23.
- 677 16. Alié A, Hayashi T, Sugimura I, Manuel M, Sugano W, Mano A, et al. The ancestral  
678 gene repertoire of animal stem cells. *Proceedings of the National Academy of Sciences*  
679 *of the United States of America*. 2015;:201514789–8.
- 680 17. Solana J, Irimia M, Ayoub S, Orejuela MR, Zywitza V, Jens M, et al. Conserved  
681 functional antagonism of CELF and MBNL proteins controls stem cell-specific alternative  
682 splicing in planarians. *eLife*. eLife Sciences Publications Limited; 2016;5:1193.
- 683 18. Juliano CE, Swartz SZ, Wessel GM. A conserved germline multipotency program.  
684 *Development*. 2010;137:4113–26.
- 685 19. Shibata N, Kashima M, Ishiko T, Nishimura O, Rouhana L, Misaki K, et al.  
686 Inheritance of a Nuclear PIWI from Pluripotent Stem Cells by Somatic Descendants  
687 Ensures Differentiation by Silencing Transposons in Planarian. *Developmental Cell*.  
688 Elsevier Inc; 2016;37:226–37.
- 689 20. Blythe MJ, Kao D, Malla S, Rowsell J, Wilson R, Evans D, et al. A Dual Platform  
690 Approach to Transcript Discovery for the Planarian *Schmidtea mediterranea* to Establish  
691 RNAseq for Stem Cell and Regeneration Biology. Jaeger J, editor. *PLoS ONE*.  
692 2010;5:e15617–7.
- 693 21. Resch AM, Palakodeti D, Lu Y-C, Horowitz M, Graveley BR. Transcriptome Analysis  
694 Reveals Strain-Specific and Conserved Stemness Genes in *Schmidtea mediterranea*.  
695 Boutros M, editor. *PLoS ONE*. 2012;7:e34447–12.
- 696 22. Wurtzel O, Cote LE, Poirier A, Satija R, Regev A, Reddien PW. A Generic and Cell-  
697 Type-Specific Wound Response Precedes Regeneration in Planarians. *Developmental*  
698 *Cell*. Elsevier Inc; 2015;35:632–45.
- 699 23. Molinaro AM, Pearson BJ. In silico lineage tracing through single cell transcriptomics  
700 identifies a neural stem cell population in planarians. *Genome Biol*. *Genome Biology*;  
701 2016;:1–17.

- 702 24. van Wolfswinkel JC, Wagner DE, Reddien PW. Single-Cell Analysis Reveals  
703 Functionally Distinct Classes within the Planarian Stem Cell Compartment. *Stem Cell*.  
704 Elsevier Inc; 2014;15:326–39.
- 705 25. Hayashi T, Asami M, Higuchi S, Shibata N, Agata K. Isolation of planarian X-ray-  
706 sensitive stem cells by fluorescence-activated cell sorting. *Development, Growth &*  
707 *Differentiation*. 2006;48:371–80.
- 708 26. Higuchi S, Hayashi T, Hori I, Shibata N, Sakamoto H, Agata K. Characterization and  
709 categorization of fluorescence activated cell sorted planarian stem cells by ultrastructural  
710 analysis. *Development, Growth & Differentiation*. 2007;49:571–81.
- 711 27. Romero BT, Evans DJ, Aboobaker AA. FACS Analysis of the Planarian Stem Cell  
712 Compartment as a Tool to Understand Regenerative Mechanisms. In: Orgogozo V,  
713 Rockman MV, editors. *Hox Genes*. Totowa, NJ: Humana Press; 2012. pp. 167–79.
- 714 28. Shibata N, Hayashi T, Fukumura R, Fujii J, Kudome-Takamatsu T, Nishimura O, et  
715 al. Comprehensive gene expression analyses in pluripotent stem cells of a planarian,  
716 *Dugesia japonica*. *Int. J. Dev. Biol.* 2012;56:93–102.
- 717 29. Mangel M. Feedback Control in Planarian Stem Cell Systems. *BMC Systems*  
718 *Biology*. *BMC Systems Biology*; 2016;:1–18.
- 719 30. Salvetti A. DjPum, a homologue of *Drosophila* Pumilio, is essential to planarian stem  
720 cell maintenance. *Development*. 2005;132:1863–74.
- 721 31. Reddien PW. Specialized progenitors and regeneration. *Development*.  
722 2013;140:951–7.
- 723 32. Jaber-Hijazi F, Lo PJKP, Mihaylova Y, Foster JM, Benner JS, Romero BT, et al.  
724 Planarian MBD2/3 is required for adult stem cell pluripotency independently of DNA  
725 methylation. *Developmental Biology*. Elsevier; 2013;384:141–53.
- 726 33. Scimone ML, Meisel J, Reddien PW. The Mi-2-like Smed-CHD4 gene is required for  
727 stem cell differentiation in the planarian *Schmidtea mediterranea*. *Development*.  
728 2010;137:1231–41.
- 729 34. Barberan S, Fraguas S, Cebrià F. The EGFR signaling pathway controls gut  
730 progenitor differentiation during planarian regeneration and homeostasis. *Development*.  
731 2016;143:2089–102.
- 732 35. González-Estévez C, Felix DA, Smith MD, Paps J, Morley SJ, James V, et al. SMG-  
733 1 and mTORC1 Act Antagonistically to Regulate Response to Injury and Growth in  
734 Planarians. Reddien P, editor. *PLoS Genet*. 2012;8:e1002619–17.
- 735 36. Tu KC, Cheng L-C, T K Vu H, Lange JJ, McKinney SA, Seidel CW, et al. Egr-5 is a  
736 post-mitotic regulator of planarian epidermal differentiation. *eLife*. eLife Sciences  
737 Publications Limited; 2015;4:e10501.
- 738 37. Bernstein BE, Mikkelsen TS, Xie X, Kamal M, Huebert DJ, Cuff J, et al. A Bivalent

- 739 Chromatin Structure Marks Key Developmental Genes in Embryonic Stem Cells. *Cell*.  
740 2006;125:315–26.
- 741 38. McDonald D. Chromatin signatures of pluripotent cell lines. 2006;:1–11.
- 742 39. Mikkelsen TS, Ku M, Jaffe DB, Issac B, Lieberman E, Giannoukos G, et al. Genome-  
743 wide maps of chromatin state in pluripotent and lineage-committed cells. *Nature*.  
744 2007;448:553–60.
- 745 40. Voigt P, Tee WW, Reinberg D. A double take on bivalent promoters. *Genes &*  
746 *development*. 2013;27:1318–38.
- 747 41. Lesch BJ, Page DC. Poised chromatin in the mammalian germ line. *Development*.  
748 2014;141:3619–26.
- 749 42. Hammoud SS, Low DHP, Yi C, Carrell DT, Guccione E, Cairns BR. Chromatin and  
750 Transcription Transitions of Mammalian Adult Germline Stem Cells and  
751 Spermatogenesis. *Stem Cell*. Elsevier Inc; 2014;15:239–53.
- 752 43. Vastenhouw NL, Schier AF. Bivalent histone modifications in early embryogenesis.  
753 *Current Opinion in Cell Biology*. 2012;24:374–86.
- 754 44. Schuettengruber B, Ganapathi M, Leblanc B, Portoso M, Jaschek R, Tolhuis B, et al.  
755 Functional Anatomy of Polycomb and Trithorax Chromatin Landscapes in *Drosophila*  
756 Embryos. Kingston R, editor. *PLoS biology*. 2009;7:e1000013–8.
- 757 45. Robb SMC, Ross E, Alvarado AS. SmedGD: the *Schmidtea mediterranea* genome  
758 database. *Nucleic Acids Research*. 2007;36:D599–D606.
- 759 46. Robb SMC, Gotting K, Ross E, Sánchez Alvarado A. SmedGD 2.0: The *Schmidtea*  
760 *mediterranea* genome database. Vize P, Westerfield M, editors. *genesis*. 2015;53:535–  
761 46.
- 762 47. Tan TCJ, Rahman R, Jaber-Hijazi F, Felix DA, Chen C, Louis EJ, et al. Telomere  
763 maintenance and telomerase activity are differentially regulated in asexual and sexual  
764 worms. *Proceedings of the National Academy of Sciences of the United States of*  
765 *America*. *National Acad Sciences*; 2012;109:4209–14.
- 766 48. Xu H, Baroukh C, Dannenfelser R, Chen EY, Tan CM, Kou Y, et al. ESCAPE:  
767 database for integrating high-content published data collected from human and mouse  
768 embryonic stem cells. *Database*. 2013;2013:bat045–5.
- 769 49. Bonuccelli L, Rossi L, Lena A, Scarcelli V, Rainaldi G, Evangelista M, et al. An  
770 RbAp48-like gene regulates adult stem cells in planarians. *Journal of Cell Science*.  
771 2010;123:690–8.
- 772 50. Vásquez-Doorman C, Petersen CP. The NuRD complex component p66 suppresses  
773 photoreceptor neuron regeneration in planarians. *Regeneration*. 2016;3:168–78.
- 774 51. Hubert A, Henderson JM, Ross KG, Cowles MW, Torres J, Zayas RM. Epigenetic



- 775 regulation of planarian stem cells by the SET1/MLL family of histone methyltransferases.  
776 Epigenetics. 2014;8:79–91.
- 777 52. Duncan EM, Chitsazan AD, Seidel CW, Alvarado AS. Set1 and MLL1/2 Target  
778 Distinct Sets of Functionally Different Genomic Loci In&nbsp;Vivo. CellReports. The  
779 Authors; 2015;13:2741–55.
- 780 53. Ruthenburg AJ, Allis CD, Wysocka J. Methylation of Lysine 4 on Histone H3:  
781 Intricacy of Writing and Reading a Single Epigenetic Mark. Molecular Cell. 2007;25:15–  
782 30.
- 783 54. Eissenberg JC, Shilatifard A. Histone H3 lysine 4 (H3K4) methylation in development  
784 and differentiation. Developmental Biology. Elsevier Inc; 2010;339:240–9.
- 785 55. Orlando DA, Chen MW, Brown VE, Solanki S, Choi YJ, Olson ER, et al. Quantitative  
786 ChIP-Seq Normalization Reveals Global Modulation of the Epigenome. CellReports. The  
787 Authors; 2014;9:1163–70.
- 788 56. Hu B, Petela N, Kurze A, Chan K-L, Chapard C, Nasmyth K. Biological  
789 chromodynamics: a general method for measuring protein occupancy across the  
790 genome by calibrating ChIP-seq. Nucleic Acids Research. 2015;:gkv670–20.
- 791 57. Cao R, Wang L, Wang H, Xia L, Erdjument-Bromage H, Tempst P, et al. Role of  
792 Histone H3 Lysine 27 Methylation in Polycomb-Group Silencing. Science.  
793 2002;298:1039–43.
- 794 58. Cheng J, Blum R, Bowman C, Hu D, Shilatifard A, Shen S, et al. A Role for H3K4  
795 Monomethylation in Gene Repression and Partitioning of Chromatin Readers. Molecular  
796 Cell. Elsevier Inc; 2014;53:979–92.
- 797 59. Pauler FM, Sloane MA, Huang R, Regha K, Koerner MV, Tamir I, et al. H3K27me3  
798 forms BLOCs over silent genes and intergenic regions and specifies a histone banding  
799 pattern on a mouse autosomal chromosome. Genome Research. 2008;19:221–33.
- 800 60. Sparmann A, van Lohuizen M. Polycomb silencers control cell fate, development  
801 and cancer. Nat Rev Cancer. 2006;6:846–56.
- 802 61. Ramírez F, Ryan DP, Grüning B, Bhardwaj V, Kilpert F, Richter AS, et al.  
803 deepTools2: a next generation web server for deep-sequencing data analysis. Nucleic  
804 Acids Research. 2016;44:W160–5.
- 805 62. Vastenhouw NL, Zhang Y, Woods IG, Imam F, Regev A, Liu XS, et al. Chromatin  
806 signature of embryonic pluripotency is established during genome activation. Nature.  
807 Nature Publishing Group; 2010;464:922–6.
- 808 63. Brandl H, Moon H, Vila-Farré M, Liu S-Y, Henry I, Rink JC. PlanMine – a mineable  
809 resource of planarian biology and biodiversity. Nucleic Acids Research. 2016;44:D764–  
810 73.
- 811 64. Kao D, Felix D, Aboobaker A. The planarian regeneration transcriptome reveals a

- 812 shared but temporally shifted regulatory program between opposing head and tail  
813 scenarios. *BMC Genomics*. BioMed Central; 2013;14:797.
- 814 65. Rodríguez-Esteban G, González-Sastre A, Rojo-Laguna JI, Saló E, Abril JF. Digital  
815 gene expression approach over multiple RNA-Seq data sets to detect neoblast  
816 transcriptional changes in *Schmidtea mediterranea*. *BMC Genomics*. 2015;16:951–23.
- 817 66. Solana J, Gamberi C, Mihaylova Y, Grosswendt S, Chen C, Lasko P, et al. The  
818 CCR4-NOT Complex Mediates Deadenylation and Degradation of Stem Cell mRNAs  
819 and Promotes Planarian Stem Cell Differentiation. Newmark PA, editor. *PLoS Genet*.  
820 2013;9:e1004003–16.
- 821 67. Huang X, Madan A. CAP3: A DNA sequence assembly program. *Genome*  
822 *Research*. Cold Spring Harbor Laboratory Press; 1999;9:868–77.
- 823 68. Wu TD, Reeder J, Lawrence M, Becker G, Brauer MJ. GMAP and GSNAP for  
824 Genomic Sequence Alignment: Enhancements to Speed, Accuracy, and Functionality.  
825 *Methods Mol. Biol.* New York, NY: Springer New York; 2016;1418:283–334.
- 826 69. Haas BJ, Zeng Q, Pearson MD, Cuomo CA, Wortman JR. Approaches to Fungal  
827 Genome Annotation. *Mycology*. 2011;2:118–41.
- 828 70. Kim D, Langmead B, Salzberg SL. HISAT: a fast spliced aligner with low memory  
829 requirements. *Nat. Methods*. *Nature Research*; 2015;12:357–60.
- 830 71. Perteau M, Perteau GM, Antonescu CM, Chang T-C, Mendell JT, Salzberg SL.  
831 StringTie enables improved reconstruction of a transcriptome from RNA-seq reads.  
832 *Nature Biotechnology*. *Nature Research*; 2015;33:290–5.
- 833 72. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, et al. De  
834 novo transcript sequence reconstruction from RNA-seq using the Trinity platform for  
835 reference generation and analysis. *Nature Protocols*. *Nature Research*; 2013;8:1494–  
836 512.
- 837 73. Bray NL, Pimentel H, Melsted P, Pachter L. Near-optimal probabilistic RNA-seq  
838 quantification. *Nature Biotechnology*. *Nature Research*; 2016;34:525–7.
- 839 74. Pimentel HJ, Bray N, Puente S, Melsted P, Pachter L. Differential analysis of RNA-  
840 Seq incorporating quantification uncertainty. 2016.
- 841 75. Attrill H, Falls K, Goodman JL, Millburn GH, Antonazzo G, Rey AJ, et al. FlyBase:  
842 establishing a Gene Group resource for *Drosophila melanogaster*. *Nucleic Acids*  
843 *Research*. 2016;44:D786–92.
- 844 76. Li H, Durbin R. Fast and accurate long-read alignment with Burrows-Wheeler  
845 transform. *Bioinformatics*. Oxford University Press; 2010;26:589–95.
- 846 77. Diaz A, Park K, Lim DA, Song JS. Normalization, bias correction, and peak calling  
847 for ChIP-seq. *Stat Appl Genet Mol Biol*. 2012;11:Article9.

848

## 849 **Figure legends**

850 Figure 1. Planarian FACS compartments and analysis of currently available neoblast  
851 transcriptome datasets.

852 A) Schematic of FACS cell populations and their relationship to stages in the cell cycle,  
853 stem cell progeny and differentiated cells. B) Venn diagram describing the overlaps from  
854 four independently assembled transcriptomes and genes described as being enriched in  
855 neoblasts, produced by the Aboobaker, Pearson, Rajewsky and Graveley labs  
856 respectively.

857  
858 Figure 2. Proportional transformation of gene expression values in planarian FACS  
859 compartments.

860 Spectrum of genes sorted by X1 (A), X2 (B), and Xins (C) proportion of expression  
861 where each vertical line in the spectrum represents one expressed loci. The proportion  
862 of dark blue, light blue, and orange corresponds proportions of expression in the X1, X2,  
863 and Xins FACS compartments respectively. D) A table presenting colour-coded  
864 classification groups according to proportional expression in different FACS populations  
865 of cells based on the detailed analysis of proportional expression (Additional Files 4 and  
866 5).

867  
868 Figure 3. Gene categories based on proportional expression values.

869 Previously described planarian genes are marked in expression profile following panels  
870 displaying a gene category. A) Genes with 50% or more X1 proportional expression. B)  
871 Genes with 50% or more Xins proportional expression. C) Gene ontology (GO)  
872 enrichment of X1 enriched genes showing terms mainly associated with cell division. D)  
873 GO enrichment of Xins genes showing terms associated with, for example, the  
874 extracellular matrix. E) Genes with 50% or more X2 proportional expression. F) Genes  
875 with the of sum X1 and X2 proportional expression more than or equal to 75% and in  
876 neither falling into X1 nor X2 categories. Gene names in blue and red are not  
877 characterised *S.mediterranea* genes. Blue gene names are genes associated with  
878 methyltransferase activity according to GO. Red names are genes associated with

879 mRNA processing according to GO. G) GO enrichment for X1/X2 genes from F)  
880 showing enrichment of terms involved in RNA and ribosomal processes. H) Genes with  
881 sum X2 and Xins proportional expression more than or equal to 75% and in neither X2  
882 nor Xins categories. I) Genes that are expressed in roughly similar proportions among  
883 X1, X2, and Xins cells.

884  
885

886 Figure 4. ChIP-seq analysis of H3K4me3 in planarian neoblasts.

887 A) ChIP-seq profile of 5kb around predicted transcriptional start sites (TSS) for  
888 H3K4me3 histone marks in X1 cells. B) The amount of H3K4me3 signal decreases with  
889 increasing proportional expression in the X2 compartment, indicative of expression  
890 becoming limited to post-mitotic progeny rather than NBs. C) The correlation for each 50  
891 bp window across the 5kb region around TSS to the X1 (dark blue), X2 (light blue), and  
892 Xins (orange) proportional expression value. A positive correlation value means that the  
893 higher the ChIP-seq signal, the higher the proportional expression value. A negative  
894 correlation means that the lower a ChIP-seq signal, the higher the proportional  
895 expression value. D) Example ChIP-seq profiles of individual planarian neoblast genes.

896

897 Figure 5. ChIP-seq analysis of repressive marks in planarian neoblasts.

898 A) ChIP-seq profile of 5kb around transcriptional start sites (TSS) for H3K27me3 in X1  
899 cells B) The correlation between H3K27me3 ChIP-seq signal with the X1 (dark blue), X2  
900 (light blue), and Xins (orange) proportional expression value for each 50 bp window  
901 across the 5kb region around TSS. A positive correlation value means that the higher  
902 the ChIP-seq signal, the higher the proportional expression value. A negative correlation  
903 means that the lower a ChIP-seq signal, the higher the proportional expression value. C)  
904 ChIP-seq profile of 5kb around transcriptional start sites (TSS) for H3K4me1 in X1 cells  
905 D) The correlation between H3K4me1 ChIP-seq signal with the X1 (dark blue), X2 (light  
906 blue), and Xins (orange) proportional expression value for each 50 bp window across  
907 the 5kb region around TSS. E) Example ChIP-seq profiles of individual planarian genes  
908 expressed in differentiated cells.

909

910 Figure 6. Bivalency of activation and repressive histone marks signal and shifting of  
911 H3K4me1 signal. For figure A-D, the mean ChIP-seq profiles shown were transformed

912 into percent coverage (y axis) by dividing each coverage value by the max coverage  
913 among all loci. The percent coverage of genomic region around TSS (x axis) was  
914 plotted. A) H3K4me3 and H3K27me3 profile of the top 1000 ranked X1 enriched genes.  
915 B) H3K4me3 and H3K27me3 profile of the top 1000 ranked Xins enriched genes. C)  
916 H3K4me3 and H3K27me3 profile of the top 1000 ranked X2 enriched genes. D)  
917 H3K4me3 H3K27me3 profile of genes down-regulated after *Smed-mex3* RNAi. E)  
918 H3K4me3 and H3K27me3 profile of X2/X2ins enriched genes with less than 10% X1  
919 proportional expression. F) The distribution of correlations between H3K4me3 signal  
920 5KB around TSS and H3K27me3 signal 5KB around TSS for X1, X2, and *Smed-mex3*  
921 RNAi loci.

922  
923 Figure 7. ChIP-seq profiles of high ranked X2 genes focusing on annotated transcription  
924 factors as examples. H3K4me3 and H3K27me3 profiles are displayed for the 5kb region  
925 surrounding the transcription start site of each loci.

926  
927

## 928 **Additional files**

929 Additional File 1. PDF Format.

930 Bar chart showing the number of neoblast transcripts from each of the dataset (blue), as  
931 well as the number of neoblast transcripts that have at least one other match in another  
932 dataset (red).

933

934 Additional file 2. PDF Format.

935 Asexual genome annotation workflow and metrics. A) A schematic of the workflow used  
936 to annotate the asexual genome. This process involved utilizing previously *de novo*  
937 assembled transcripts and available RNA-seq datasets. PASA (Program to Assemble  
938 Spliced Alignments) was used to create a merged reference assembly. HISAT2 was  
939 used for mapping RNA-seq data to the genome and StringTie for defining transcripts.  
940 The two sets of reference assemblies (from consolidated transcriptomes and from RNA-  
941 seq data) were merged with StringTie and filtered for redundancies, resulting in a final  
942 annotation set. B) Proportion of the annotations from the final annotation set that are

943 likely coding (with TransDecoder evidence) and non-coding (no TransDecoder  
944 evidence). Proportion of loci without TransDecoder evidence but with a BLAST hit to the  
945 non-redundant (NR) protein database (e-value  $\leq 1e-5$ ) is also shown. C) A comparison  
946 of the new annotation in this study with MAKER annotations available on *Schmidtea*  
947 *mediterranea* Genome Database (SmedGD) showing cumulative percentages of  
948 annotations at a range of expression value thresholds for SmedGD MAKER annotations  
949 as a whole, for the 'Oxford' annotations as a whole, for SmedGD MAKER exclusive  
950 annotations, and 'Oxford' exclusive annotations, and 'Oxford' coding annotations.

951  
952 Additional File 3. PDF Format.  
953 Mapping available FACS RNA-seq libraries to new annotations. A) FACS libraries from  
954 the Rajewsky, Reddien, Pearson, and Sanchez labs were downloaded and mapped to  
955 the annotations with Kallisto and normalized using Sleuth. Normalization was performed  
956 for datasets within each lab. B) A hierarchical clustering of the FACS samples using  
957 normalized transcripts per million (TPM). Wherever possible, a Short Read Archive Run  
958 (SRR) ID is provided for the original dataset

959  
960 Additional File 4. PDF Format.  
961 Proportional transformation of expression values. Expression values were converted to  
962 proportional expression values resulting in consistent clustering of samples.

963  
964 Additional File 5. PDF Format.  
965 The ternary plots within this PDF file describe subsets of expressed loci with each of the  
966 three axis representing X1, X2, and Xins proportional expression. Dots represent loci,  
967 which are categorized according to proportional expression. X1 is dark blue, X2 is light  
968 blue, Xins is orange, X1/X2 is green, X2/Xins is red, X1/Xins is purple, and non-enriched  
969 is grey.

970 The second page shows the best reciprocal hits to the ESCAPE database. Human  
971 pluripotency factors were obtained from the ESCAPE database. Best reciprocal hits  
972 were found between human and *S.mediterranea* genes. A) Ternary plot showing the  
973 distribution of the *S.mediterranea* best reciprocal hits of human pluripotency factors. B)  
974 Shows the same set of data but extracted from the proportional spectrum. C) Shows a  
975 pie chart indicating the percentage of the 233 genes that belong to each category.

976 The third, fourth, fifth pages show ternary plots of loci with X1/X2, Xins, and X1 enriched  
977 GO terms.

978  
979 Additional File 6. PDF Format.  
980 RNA-seq profiles of selected RNAi datasets. The first page of the PDF file shows RNA-  
981 seq profiles divided into four segments representing genes enriched in X1, X2, Xins, and  
982 X1/X2. The proportional values of each category are plotted as dark blue (X1  
983 proportion), light blue (X2 proportion), and orange (Xins proportion) on the right of each  
984 profile. The RNA-seq profile is displayed as red (up-regulated) and blue (down-  
985 regulated) lines that are significantly differentially expressed with a p-value less than  
986 0.05 and fold-change value  $< -1.5$  ( $\log_2$  fold-change  $-0.58$ ) or  $> 1.5$  ( $\log_2$  fold-change  
987  $0.58$ ). The length of each line represents the  $\log_2$  fold-change. A) RNA-seq profile of  
988 *Smed-mex3* RNAi performed on whole worms. B) RNA-seq profile of *Smed-zfp-1*  
989 performed on X2 and X1 cells. C) RNA-seq time-course data for *Smed-CHD4* RNAi  
990 performed on whole worms. D) RNA-seq time-course data for *Smed-p53* RNAi  
991 performed on whole worms. E) RNA-seq data for *Smed-coe* RNAi performed on whole  
992 worms.

993 The second page of the PDF file shows single-cell RNA-seq data. A) Transcripts per  
994 million (TPM) values for RNA-seq libraries for each cell type were averaged and the top  
995 one thousand genes extracted. Each row represents a cell type and the intensity of the  
996 color represents the density of genes at the position on the proportional expression  
997 spectra. B) Ternary plots of the cell types where the three axes represent X1, X2, and  
998 Xins proportional expression.

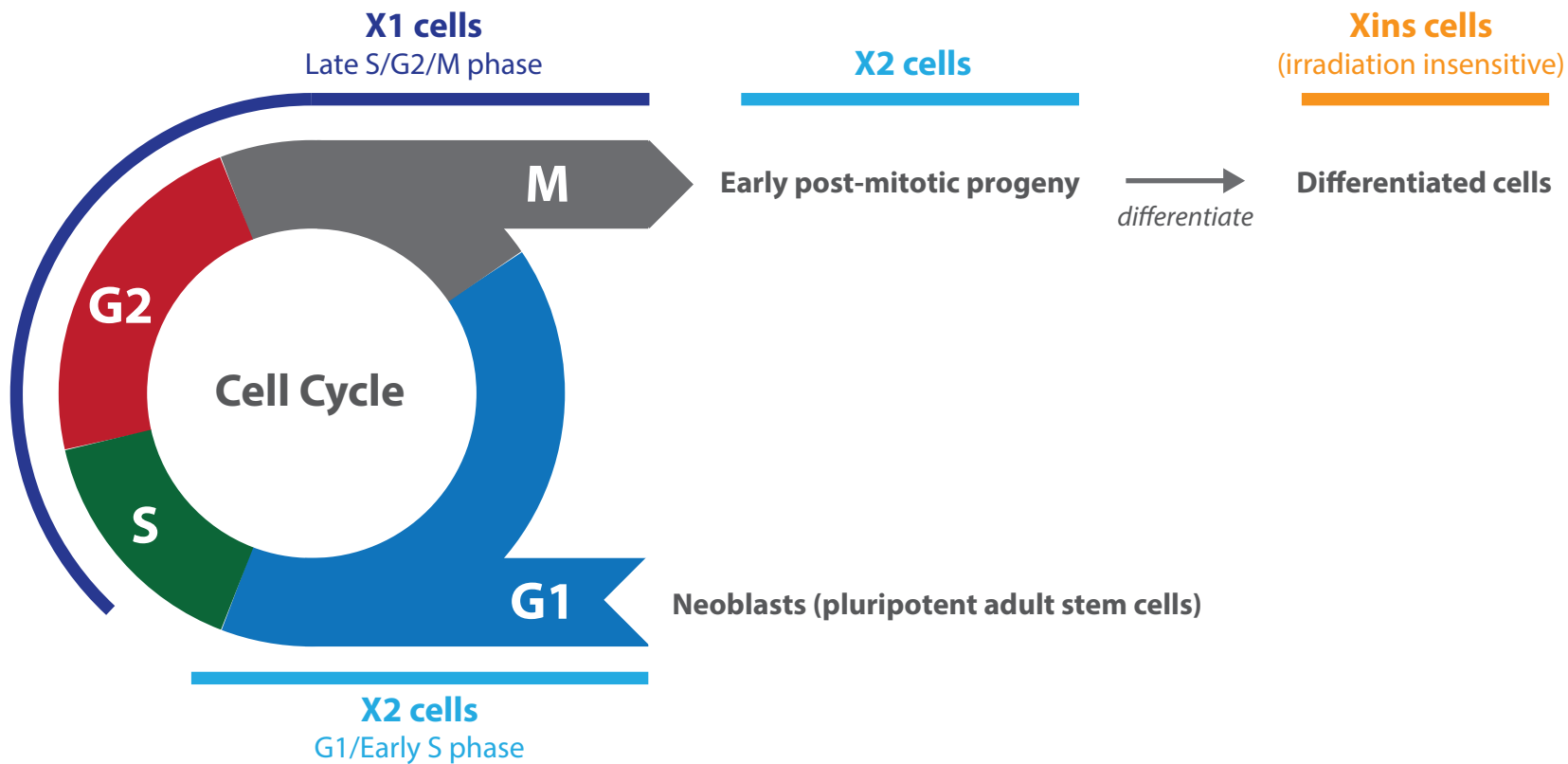
999  
1000 Additional File 7. PDF format.  
1001 Summary of ChIP-seq mapping data from all available planarian ChIP seq data,  
1002 demonstrating the improved data yield from the methods developed in the current study.

1003  
1004 Additional File 8. PDF format.  
1005 Genome wide ChIP signal presented for all genes in each proportional gene expression  
1006 category, average profiles are presented above each genome wide plot.

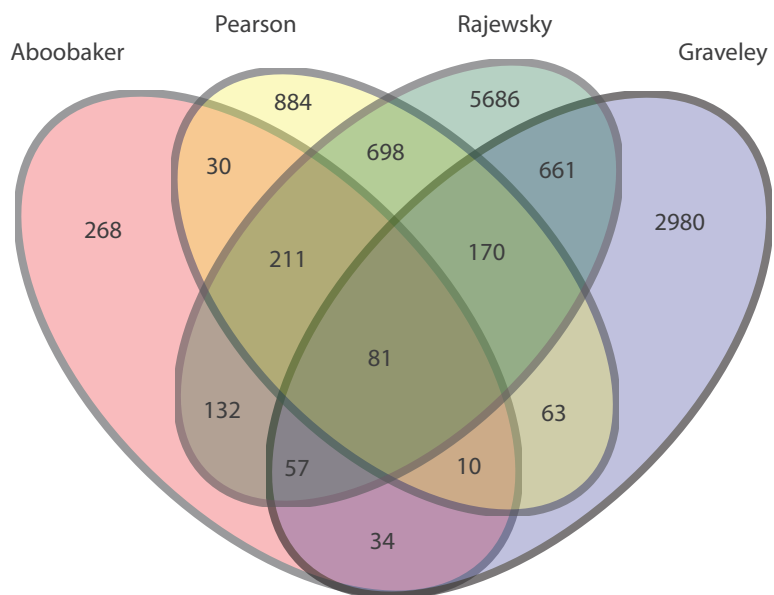
1007  
1008 Additional File 9. HTML format.

- 1009 Jupyter notebook of all analysis performed.
- 1010
- 1011 Addition File 10. ZIP format.
- 1012 GTF annotation file of the asexual genome



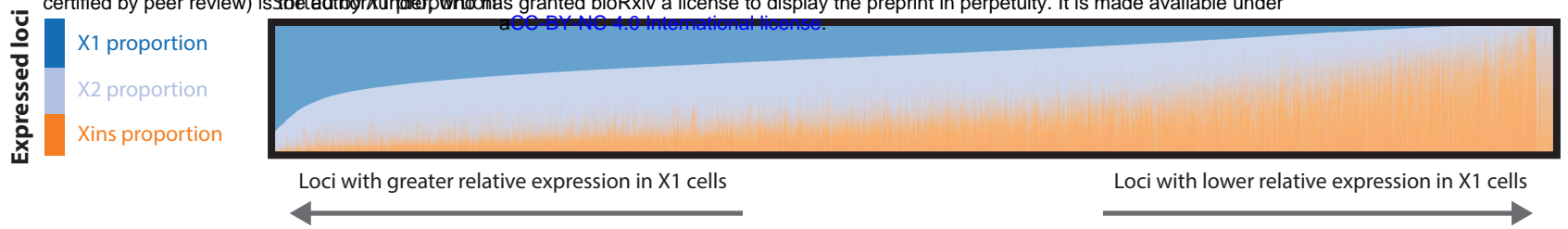
**A****B**

### Overlap among defined neoblast transcripts

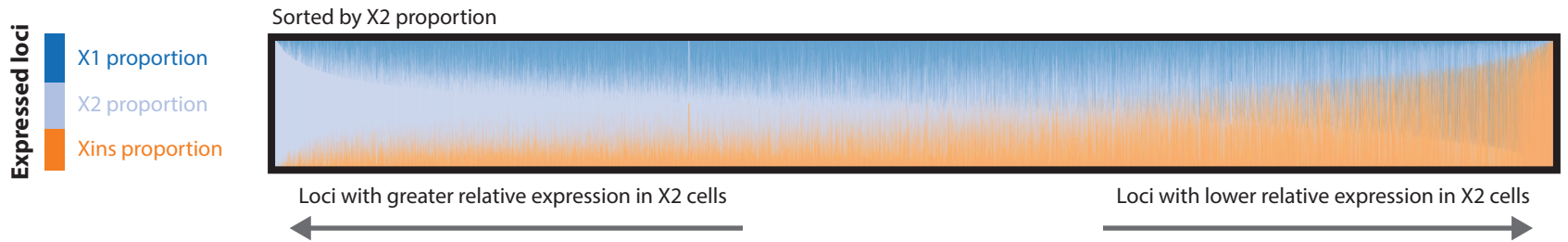


A

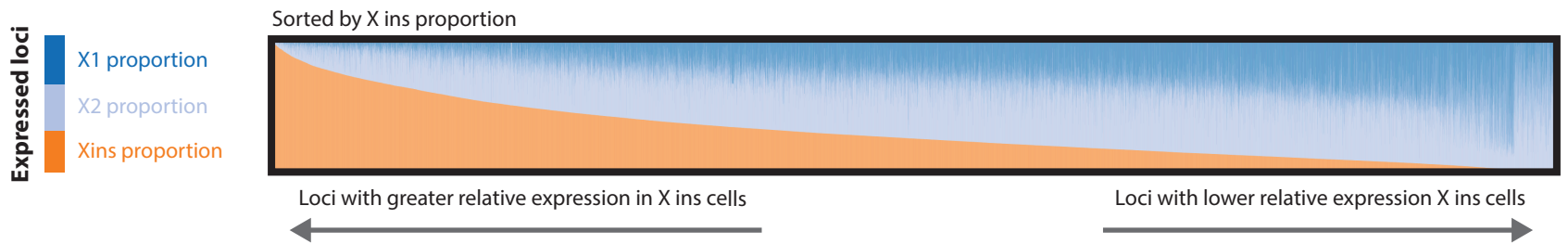
bioRxiv preprint doi: <https://doi.org/10.1101/122135>; this version posted March 29, 2017. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC 4.0 International license.



B



C



D

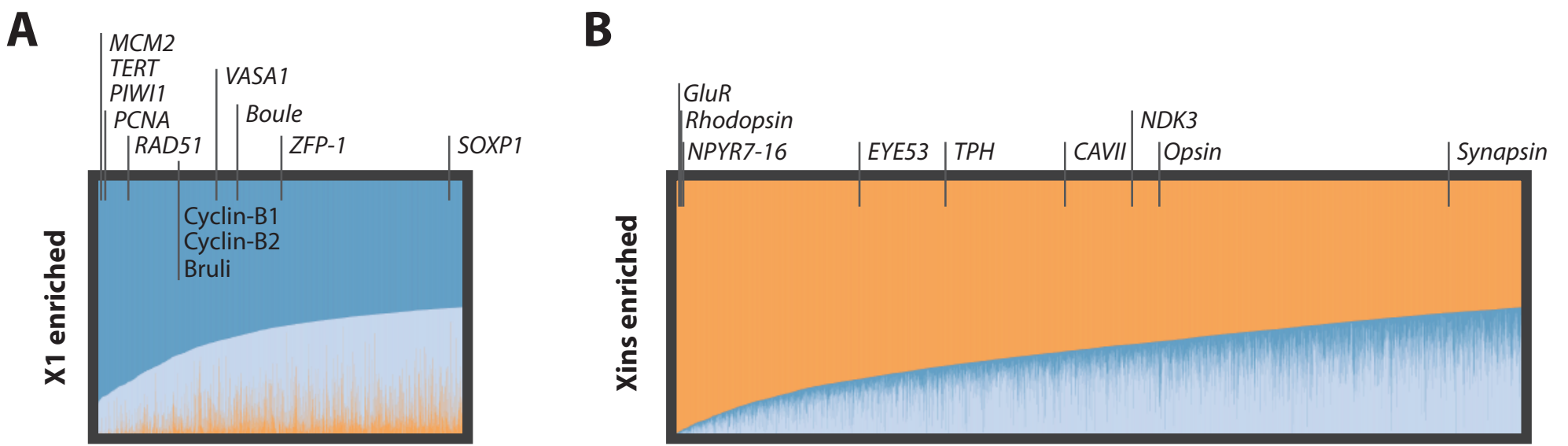
### 38,771 loci total

Categories defined by FACS RNA-seq data

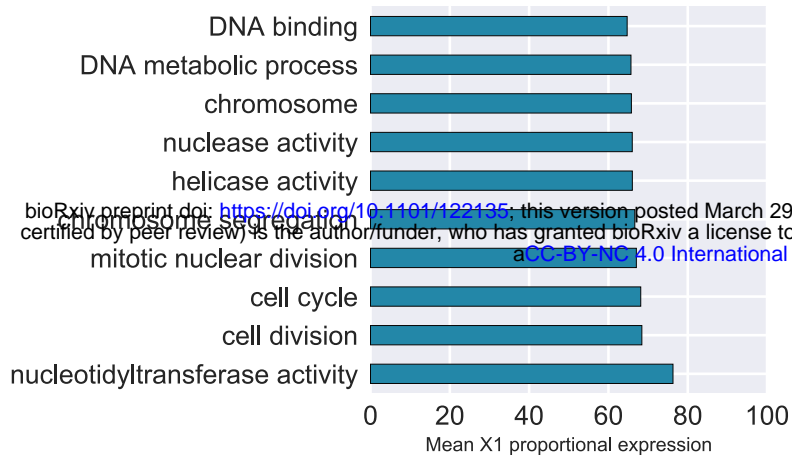


Category	Criteria	# Loci	# Coding* loci (% of category)
X1 enriched	X1 proportional expression $\geq 50\%$	2,253	1,544 (68%)
X2 enriched	X2 proportional expression $\geq 50\%$	8,444	4,781 (57%)
Xins enriched	Xins proportional expression $\geq 50\%$	5,119	3,887 (76%)
X1/X2 enriched	X1 + X2 proportional expression $\geq 75\%$ and not in X1 enriched nor X2 enriched	4,538	3,107 (68%)
X2/Xins enriched	X2 + Xins proportional expression $\geq 75\%$ and not in X2 enriched nor Xins enriched	3,652	2,688 (74%)
X1/xins enriched	X2 + Xins proportional expression $\geq 75\%$ and not in X2 enriched nor Xins enriched	303	0 (0%)
unenriched	Remaining loci with roughly equal proportions among X1, X2, and Xins	2,897	2,003 (69%)
unclassified	Loci with # of reads less than 10 in all FACS RNA-seq libraries	11,565	3,762 (33%)

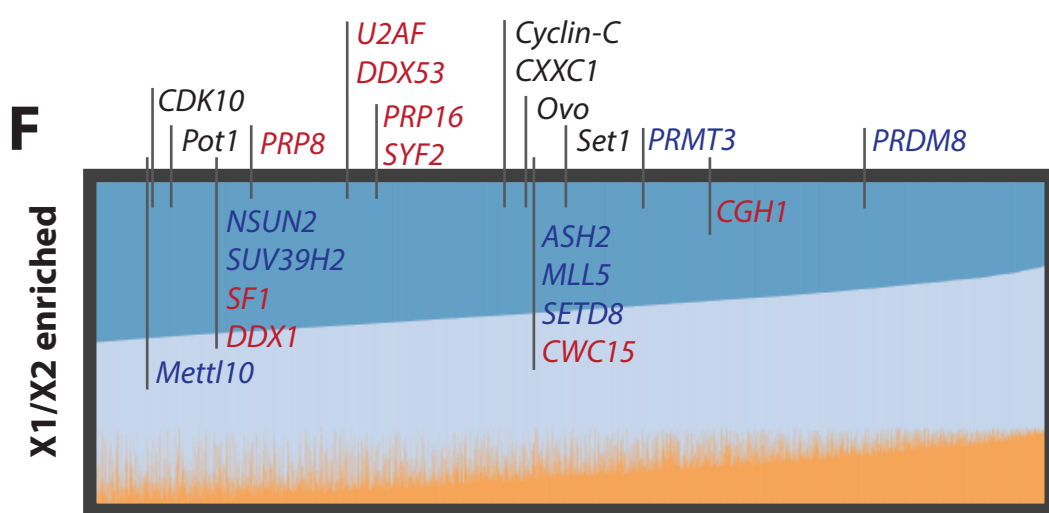
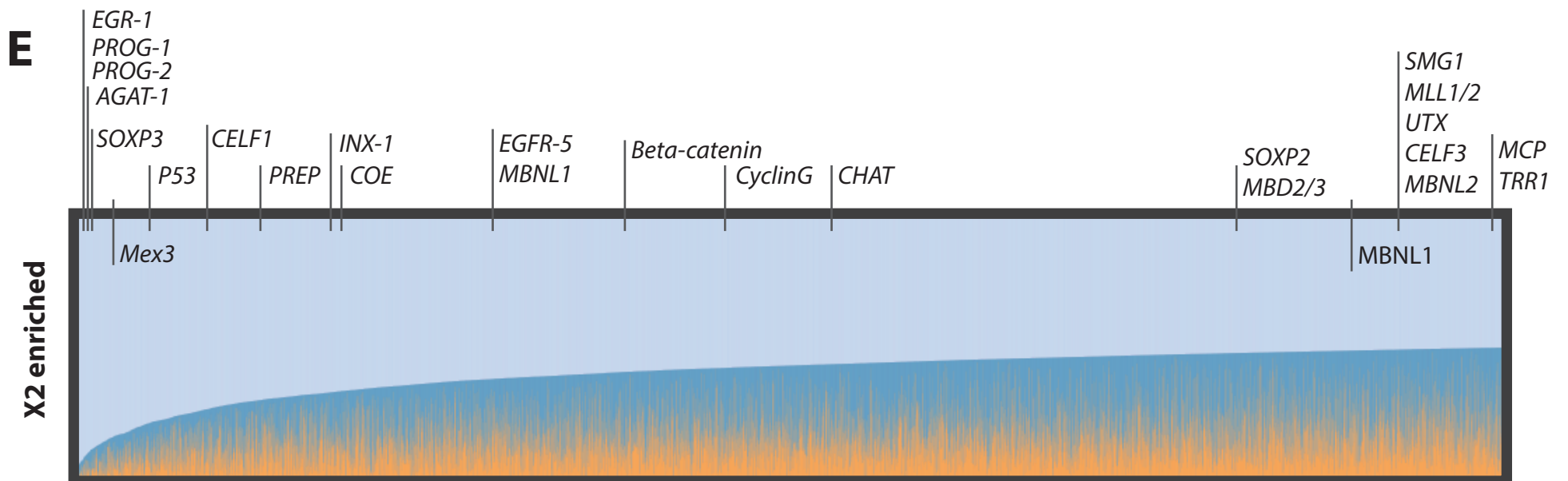
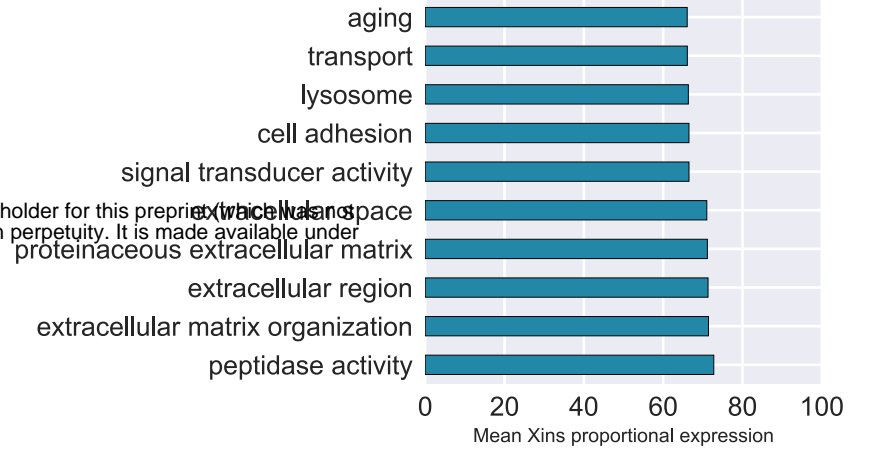
\* Coding is defined by having TransDecoder evidence which includes homology (Uniprot, PFAM), hexamer frequency, and ORF length



**C** X1 enriched genes GO enrichment

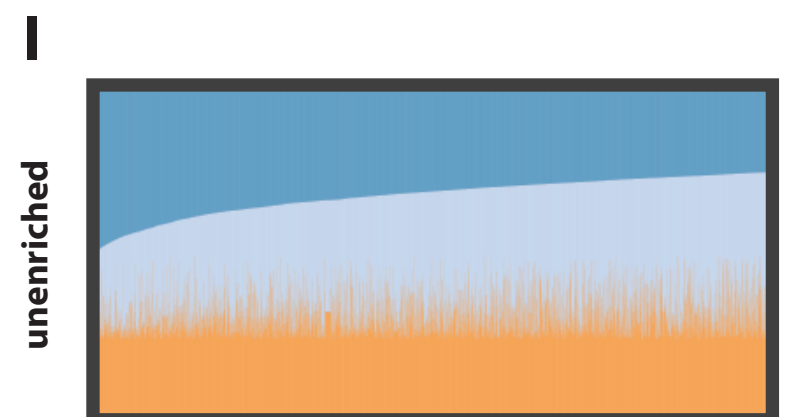
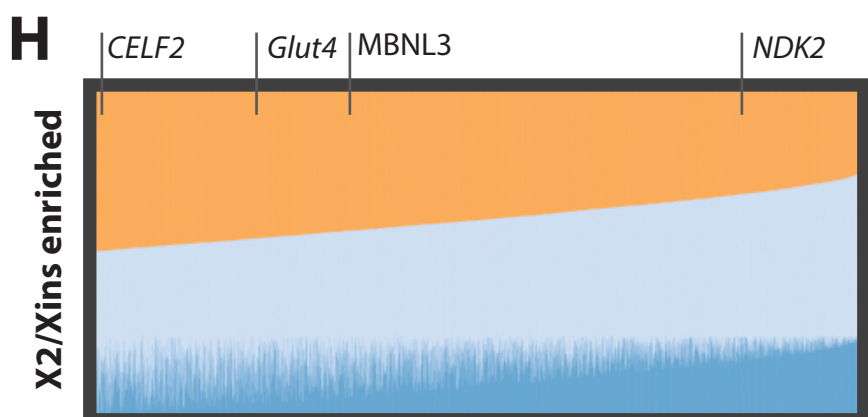
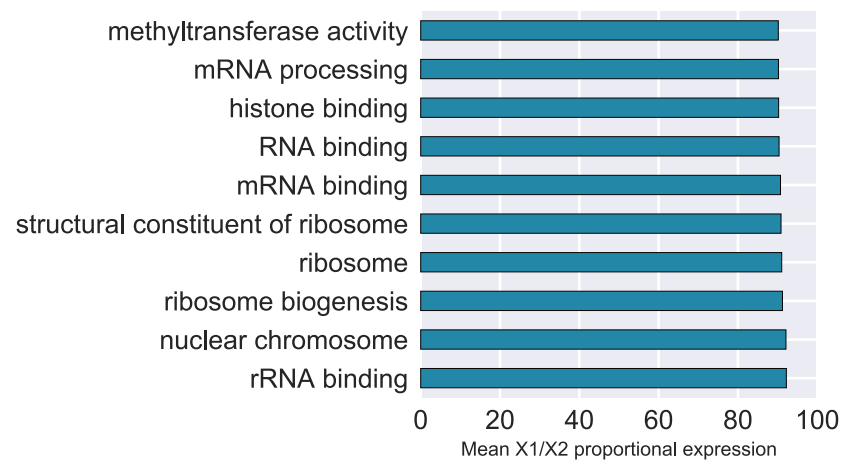


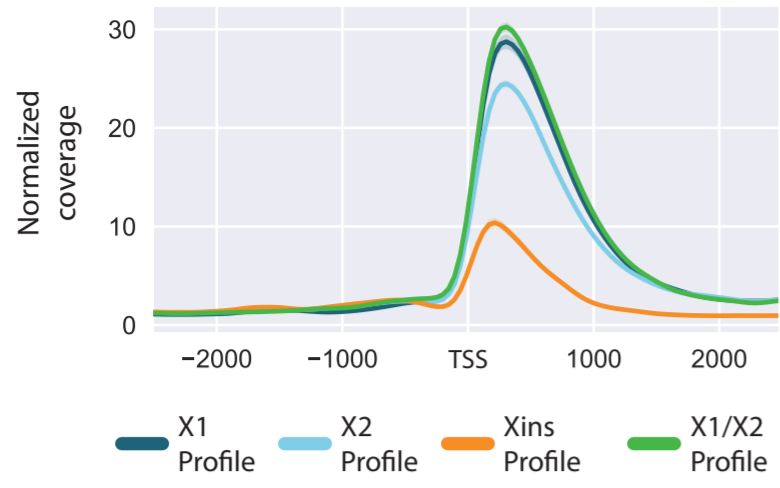
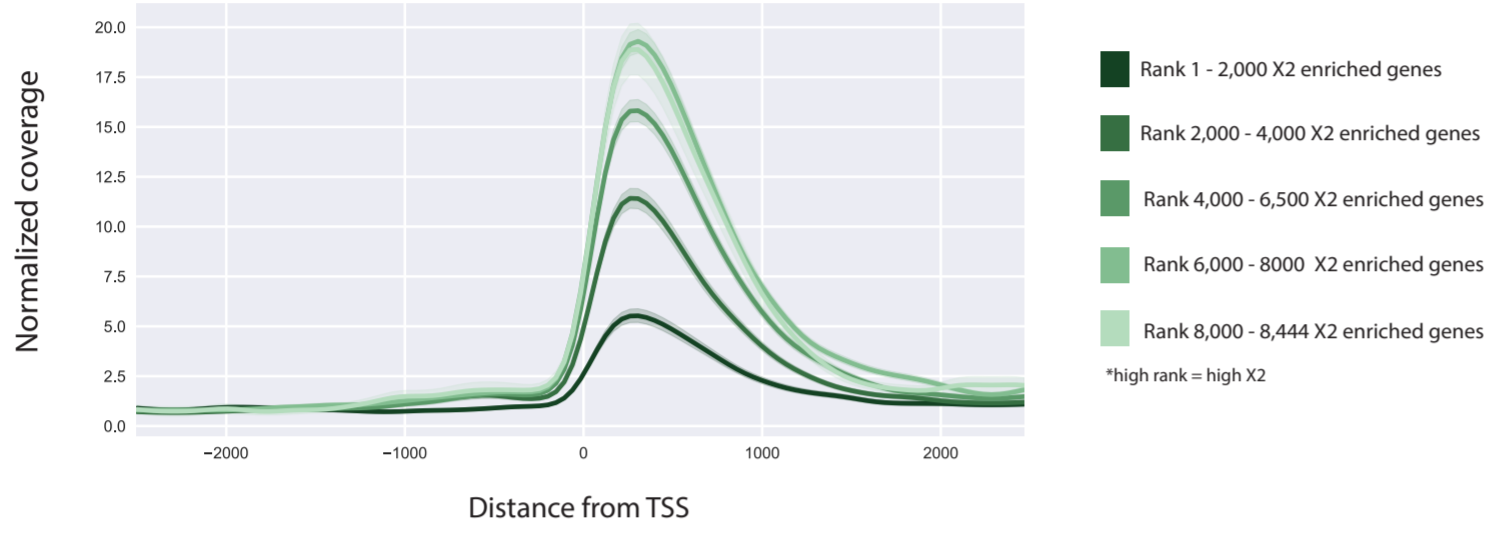
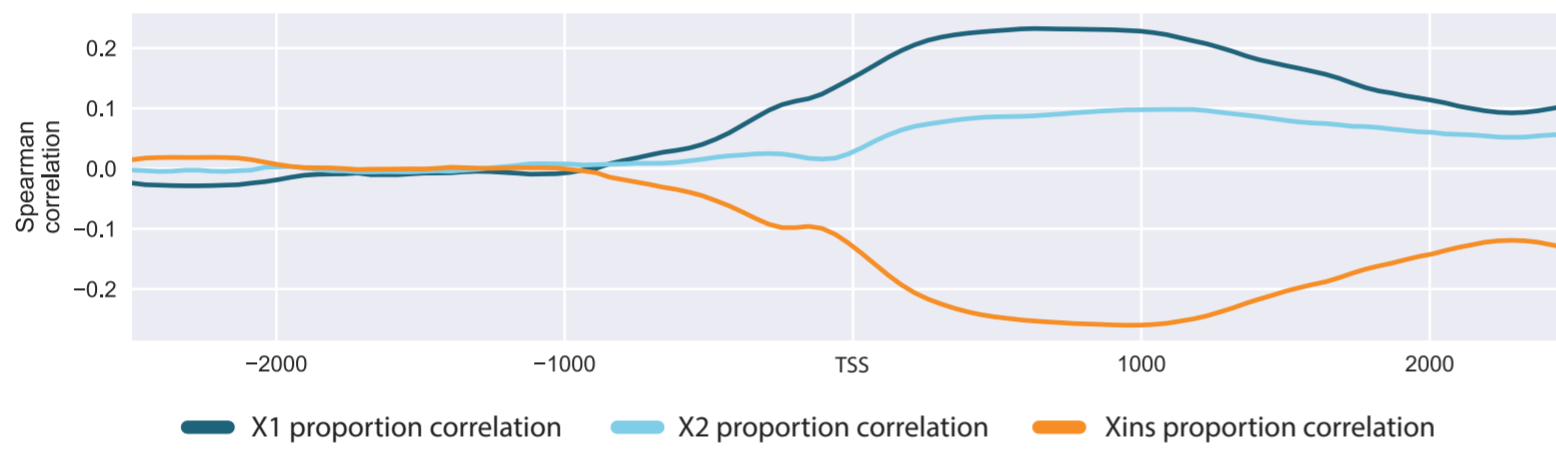
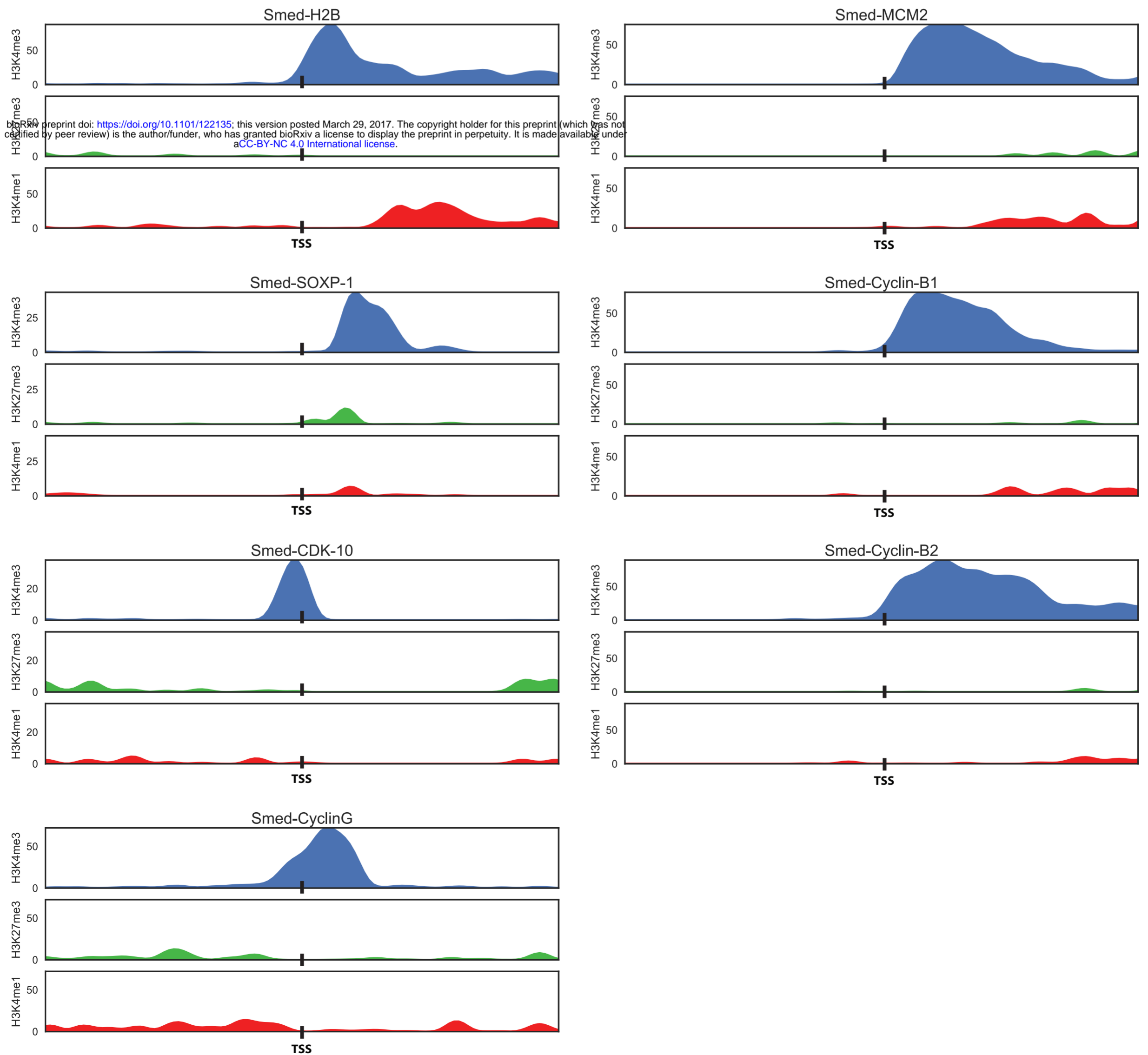
**D** Xins enriched genes GO enrichment



■ Genes associated with methyltransferase activity  
 ■ Genes associated with mRNA processing

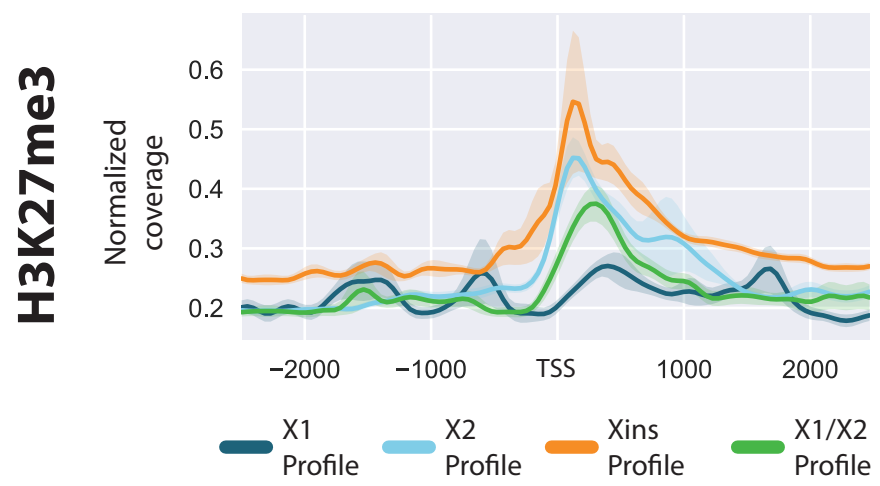
**G** X1/X2 enriched genes GO enrichment



**A****ChIP-seq Profile 2.5KB around TSS****H3K4Me3****B****Increase in H3k4me3 signal in X2 enriched genes as X2 proportion decreases****C****Correlation with proportional expression across 5kb around TSS****D****ChIP-seq profiles of known planarian neoblast genes**

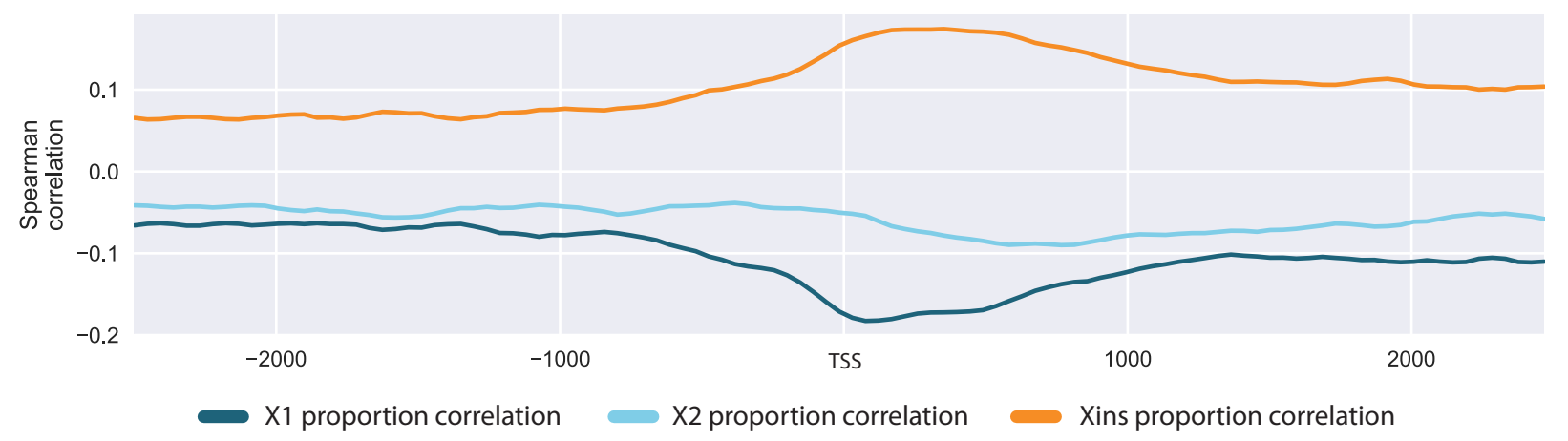
### A

#### ChIP-seq Profile 2.5KB around TSS



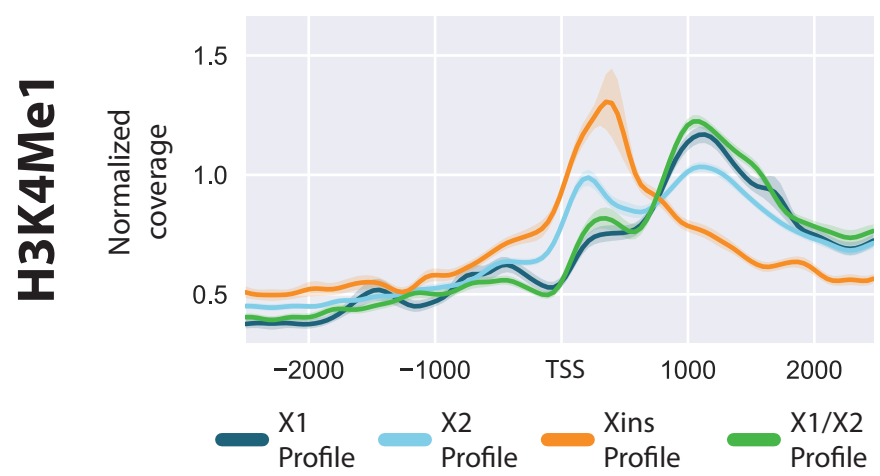
### B

#### Correlation with proportional expression across 5kb around TSS



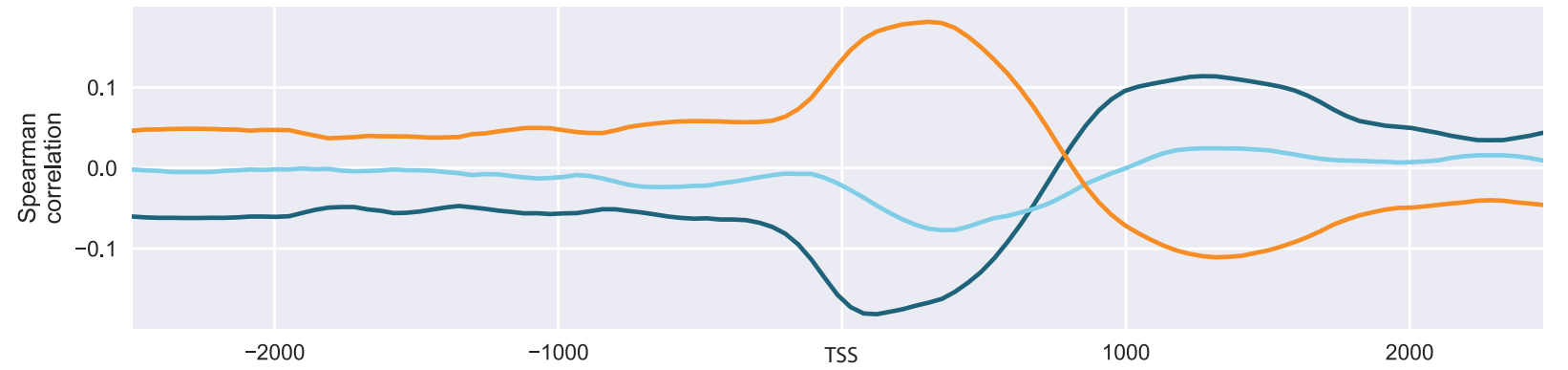
### C

#### ChIP-seq Profile 2.5KB around TSS



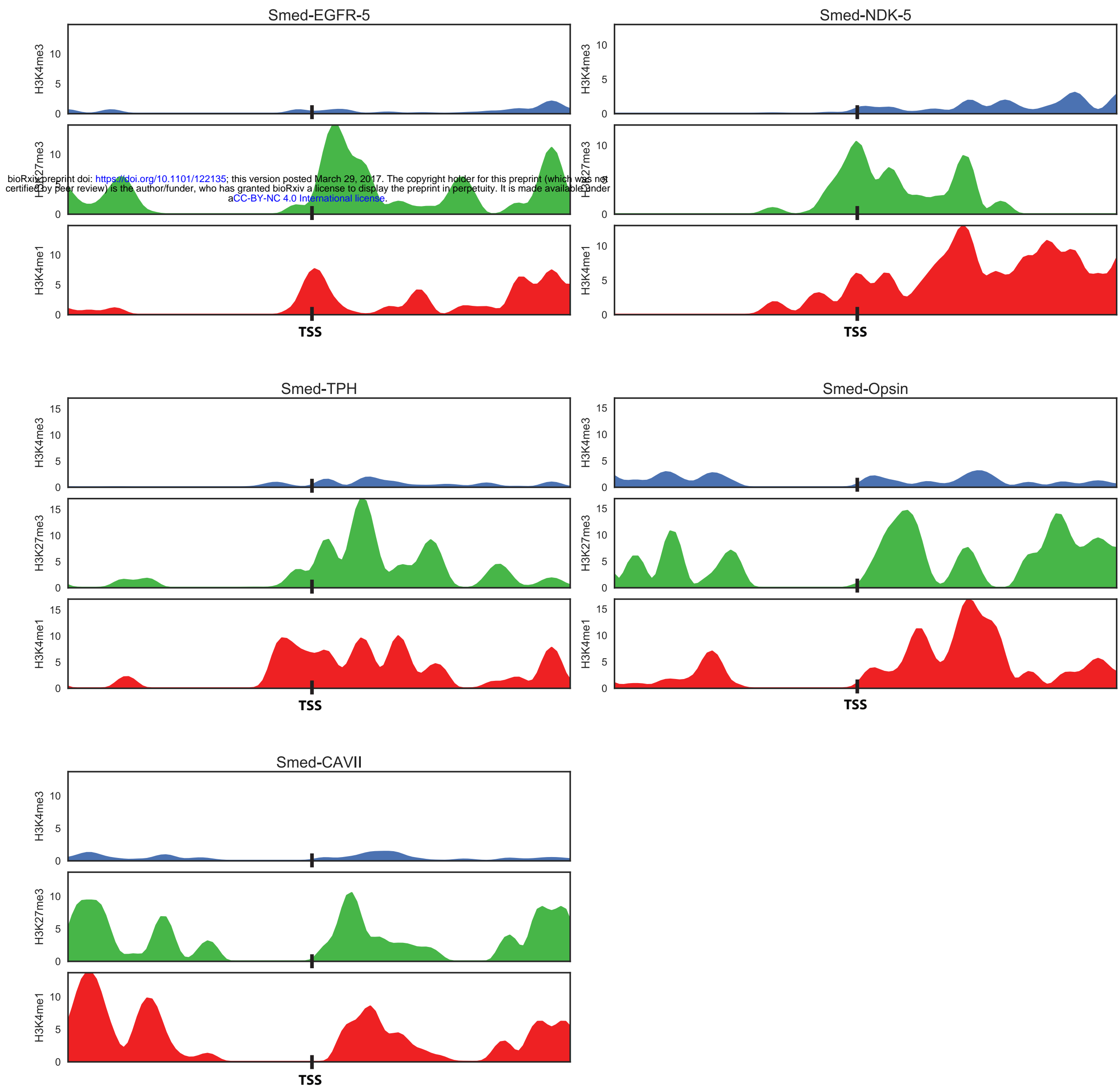
### D

#### Correlation with proportional expression across 5kb around TSS

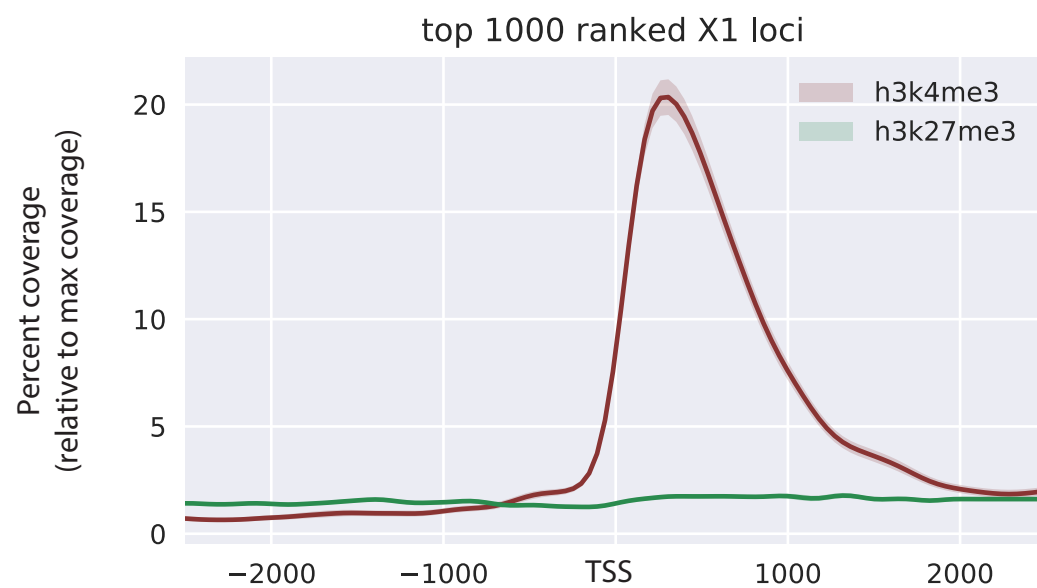


### E

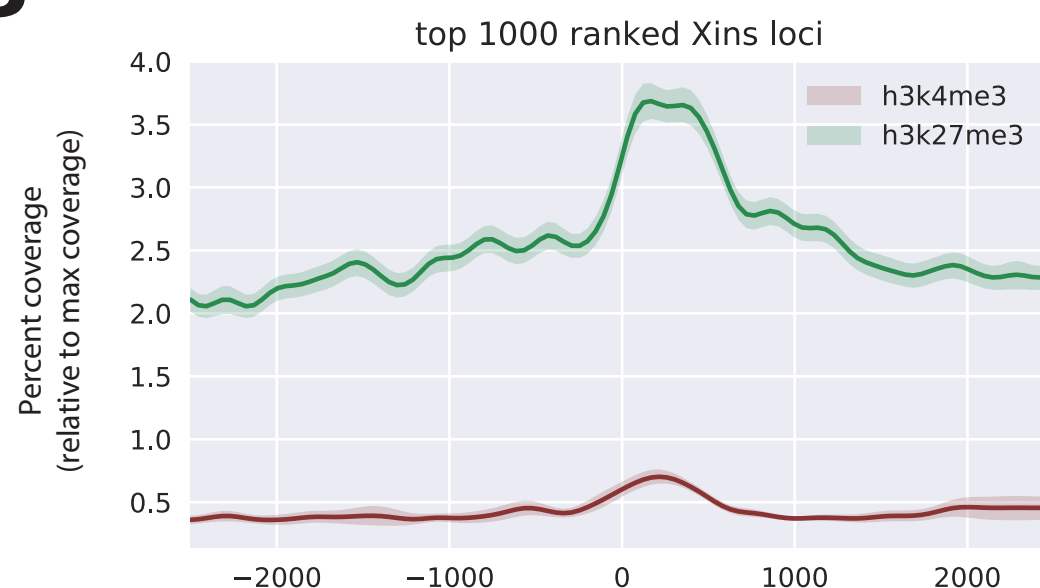
#### ChIP-seq profiles of known planarian differentiated genes



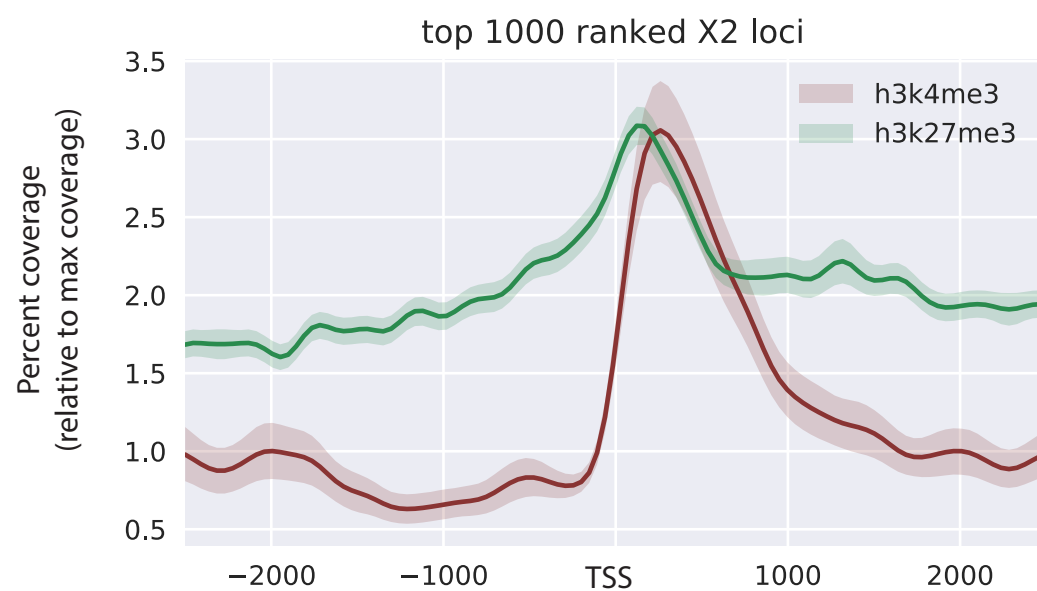
**A**



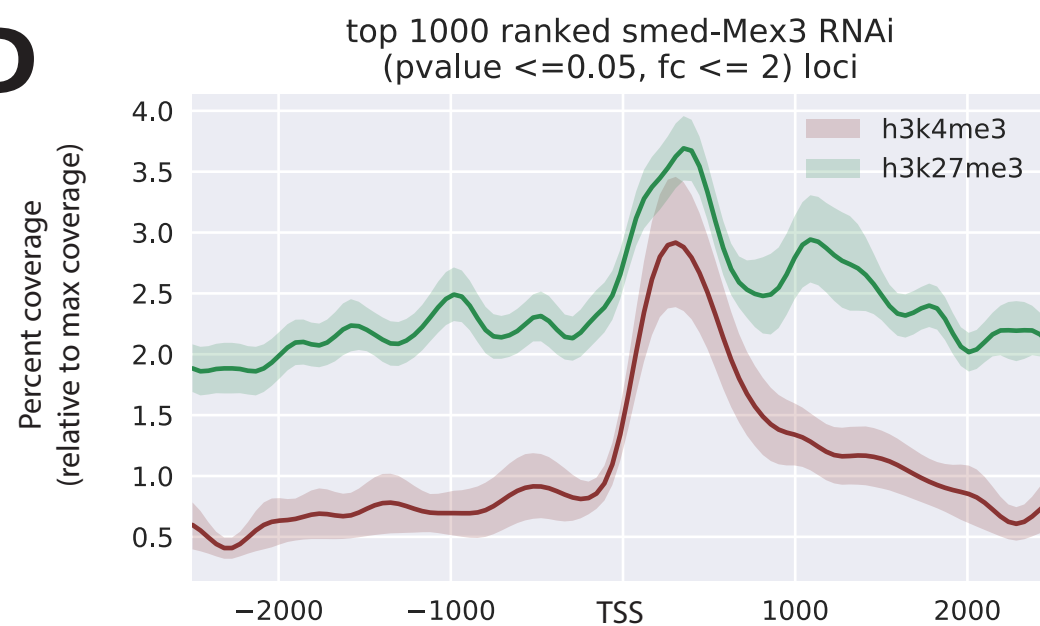
**B**



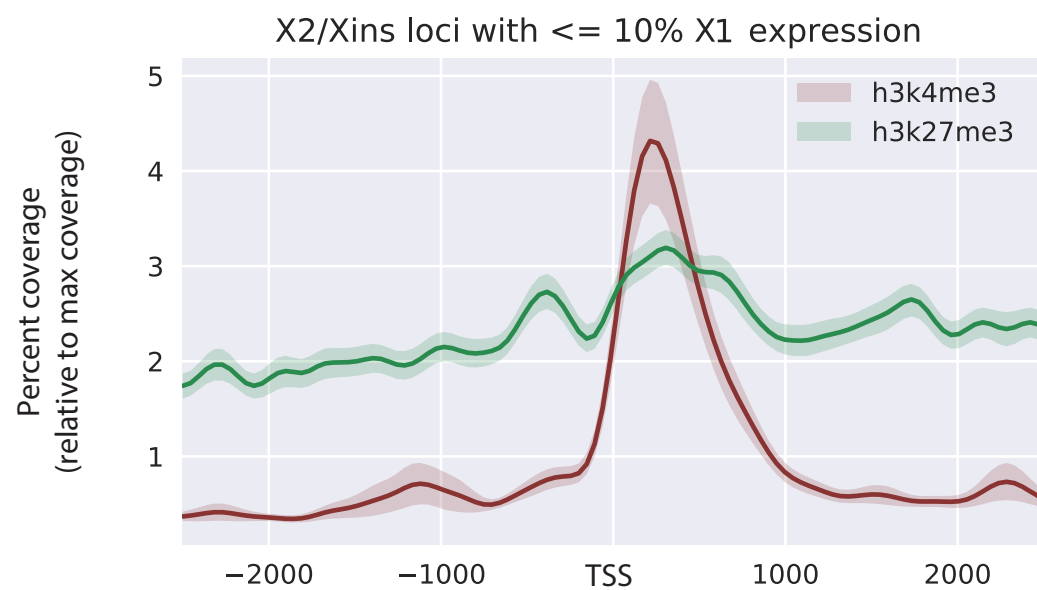
**C**



**D**

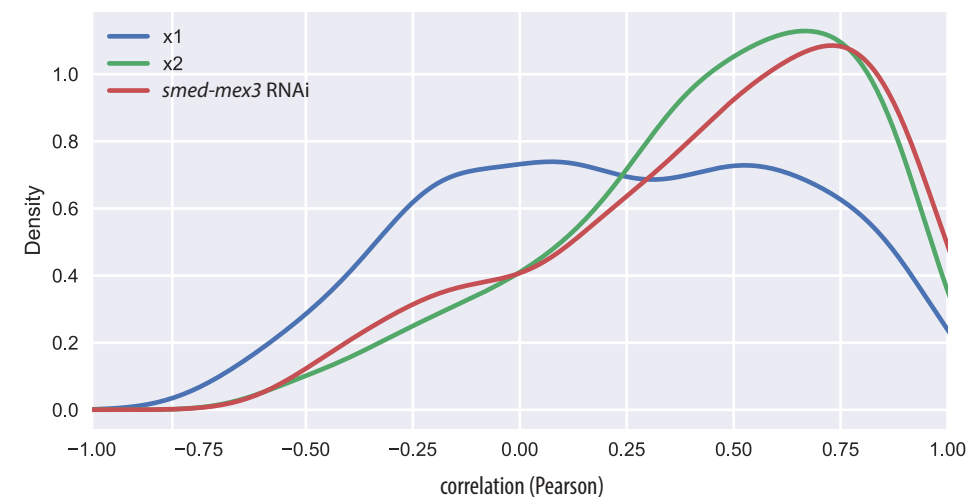


**E**

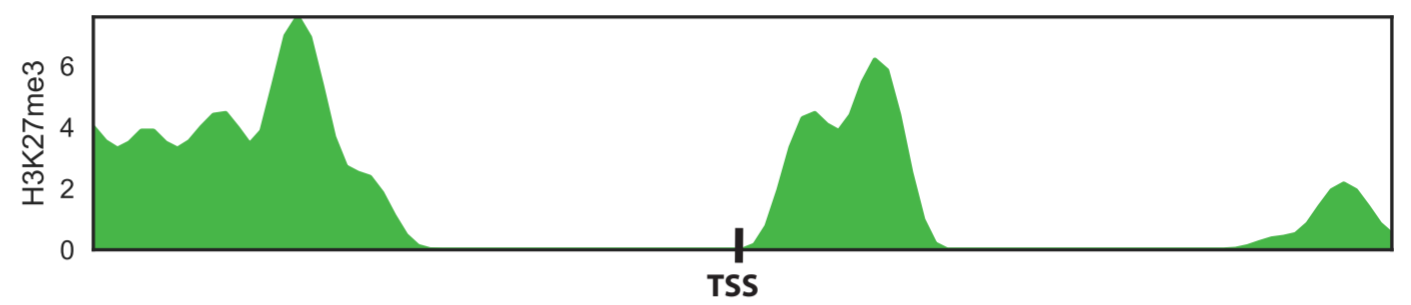
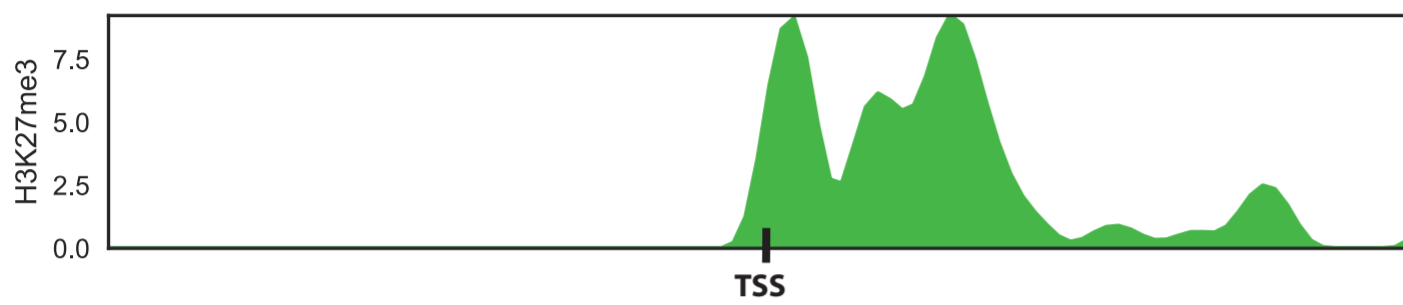
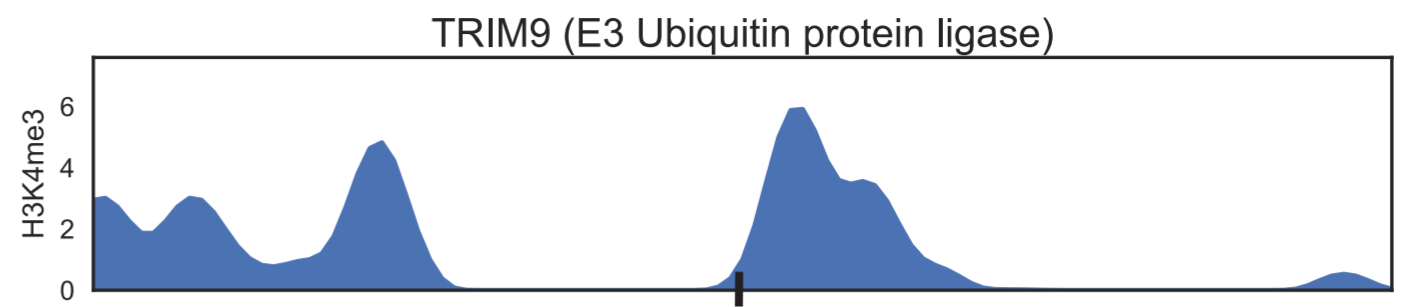
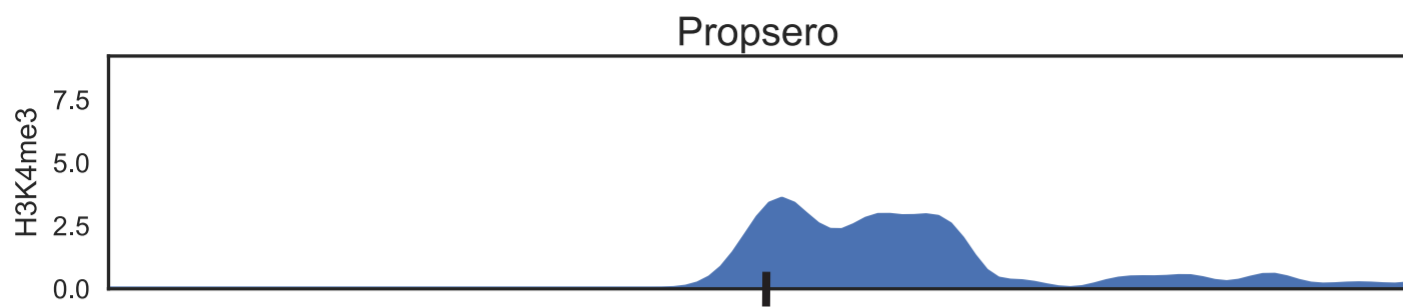


**F**

Distribution of correlations between H3K4me3 and H3K27me3 profiles of top 500 ranked loci in X1 enriched, X2 enriched, and smed-mex3



# ChIP-seq profiles of high ranked X2 loci



bioRxiv preprint doi: <https://doi.org/10.1101/122135>; this version posted March 29, 2017. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.

