

1 Chromatin and RNA Maps Reveal Regulatory Long Noncoding

2 RNAs in Mouse

3 Gireesh K. Bogu^{1,2,3,4,#}, Pedro Vizán^{2,4}, Lawrence W. Stanton^{5,6}, Miguel Beato^{2,4},
4 Luciano Di Croce^{2,4,7}, and Marc A. Marti-Renom^{1,2,4,7,#}

5

- 6 1. CNAG-CRG, Centre for Genomic Regulation (CRG), Barcelona Institute of
7 Science and Technology (BIST), Baldiri i Reixac 4, 08028 Barcelona, Spain
- 8 2. Gene Regulation, Stem Cells and Cancer Program, Centre for Genomic
9 Regulation (CRG), Dr. Aiguader 88, 08003 Barcelona, Spain
- 10 3. Bioinformatics and Genomics Programme, Centre for Genomic Regulation
11 (CRG) and UPF, Doctor Aiguader, 88, Barcelona 08003, Spain.
- 12 4. Universitat Pompeu Fabra (UPF), Barcelona, Spain
- 13 5. Department of Biological Sciences, National University of Singapore,
14 Singapore.
- 15 6. Stem Cell and Developmental Biology Group, Genome Institute of Singapore,
16 Singapore.
- 17 7. Institució Catalana de Recerca i Estudis Avançats (ICREA), Barcelona. Spain

18

19

20 Keywords: ChromHMM, ChIP-Seq, RNA-seq, lncRNA

21 #Address correspondence to Gireesh K. Bogu, gireesh.bogu@crq.eu, or Marc A.
22 Marti-Renom, martirenom@cnag.crg.eu

23

24 **ABSTRACT**

25 Discovering and classifying long noncoding RNAs across all mammalian tissues
26 and cell lines remains a major challenge. Previously, mouse lncRNAs were
27 identified using RNA-seq data from a limited number of tissues or cell lines.
28 Additionally, associating a few hundred lncRNA promoters with chromatin states
29 in a single mouse cell line has identified two classes of chromatin-associated
30 lncRNA. However, the discovery and classification of lncRNAs is still pending in
31 many other tissues in mouse. To address this, we built a comprehensive catalog
32 of lncRNAs by combining known lncRNAs with highly-confident novel lncRNAs
33 identified by mapping and *de novo* assembling billions of RNA-seq reads from
34 eight tissues and a primary cell line in mouse. Next, we integrated this catalog of
35 lncRNAs with multiple genome-wide chromatin-state maps and found two
36 different classes of chromatin state-associated lncRNAs, including promoter-
37 associated (plncRNAs) and enhancer-associated (elncRNAs) ones across
38 various tissues. Experimental knockdown of an elncRNA resulted in the down-
39 regulation of the neighboring protein-coding gene Kdm8, a histone-demethylase.
40 Our findings provide 2,803 novel lncRNAs and a comprehensive catalog of
41 chromatin-associated lncRNAs across different tissues in mouse.

42

43 INTRODUCTION

44 Previous large-scale transcriptome sequencing studies have confirmed that
45 ~80% of the human genome is transcribed, yet only a minor fraction of it (~3%)
46 codes for protein (1, 2). It is now known that a major fraction of the transcriptome
47 consists of RNAs from intergenic noncoding regions of the genome, which have
48 been termed as intergenic lncRNAs. Comprehensive lncRNA catalogs were
49 recently established for various cell lines and tissues in human, mouse, *C.*
50 *elegans*, *Drosophila*, and zebrafish (3-8). In addition, we now know the functions
51 of a limited number of the discovered lncRNAs, such as Xist in X chromosome
52 inactivation (9), HOTAIR in cancer metastasis (10), lnc-DC in dendritic cell
53 differentiation (11), Braveheart in heart development (12), Megamind and Cyrano
54 in embryonic development (13), Fendrr in cardiac mesoderm differentiation (14),
55 Malat1 in alternative splicing (15), and a few others including one from our
56 previous work showing that RMST lncRNA regulates neurogenesis by physically
57 interacting with Sox2 transcription factor (16).

58
59 Even though thousands of lncRNAs have been cataloged, it is still unclear how to
60 characterize regulatory lncRNAs. Very recently, regulatory lncRNAs were shown
61 to associate preferentially with promoter and enhancer chromatin states in a
62 single mouse cell line (17). While this observation is highly interesting, it is not
63 clear whether there were more lncRNAs associate with these two chromatin
64 states since the lncRNA associations were not tested in multiple tissues. In
65 addition, the lncRNA or chromatin state datasets used in the previous study (17)

66 were selected only in one single cell line, which technically limits testing
67 thousands of lncRNAs. Finally, it is also unknown whether these lncRNAs
68 associate with similar chromatin states across different tissues or not.

69

70 To build a comprehensive chromatin-associated mouse lncRNA dataset, we first
71 used billions of mapped RNA-seq reads to identify highly confident novel
72 lncRNAs and then combined this with thousands of known lncRNAs. Second, we
73 used more than a billion mapped ChIP-seq reads of various histone marks to
74 identify chromatin state maps. Finally, we integrated all these mouse lncRNAs
75 with chromatin state maps, resulting in a comprehensive catalog consisting of
76 thousands of chromatin state-associated lncRNAs. The analysis across multiple
77 tissues also revealed a novel set of lncRNAs that are significantly enriched with
78 promoter and enhancer chromatin states. Interestingly, the majority of the
79 lncRNAs chromatin states switch from one state to another state across all the
80 tissues or cell lines we tested. To our knowledge, this is the most comprehensive
81 dataset of chromatin state-associated lncRNAs in mouse, and we expect this will
82 be a valuable resource to help researchers select candidate lncRNAs for further
83 experimental studies.

84

85 RESULTS

86 Transcriptome mapping, assembly and quantification.

87 About 3 billion raw sequence reads of RNA-seq experiments were downloaded
88 from the ENCODE project (18) and analyzed using a computational pipeline
89 consisting of TopHat (v2.0.9) (19), Cufflinks (v2.1.1) (20), and Scripture (v4) (21)
90 (**Fig. 1A**). We constructed a map of RNA expression in mouse by first collecting
91 RNA sequencing reads using long (76–108 nucleotides), paired-end,
92 polyadenylated, strand-specific high-throughput RNA sequencing data from 8-
93 week-old adult brain, heart, kidney, small intestine, liver, spleen, testes, thymus
94 and a paired-end embryonic stem (ES) cell line (**Table S1**). Next, the collected
95 reads were mapped to the reference mouse genome using TopHat, which
96 uniquely mapped 85% (2,631,897,546) of the sequence reads with 2 mismatches
97 allowed. Of the mapped sequences, ~73% aligned with known transcript loci, and
98 the remaining 27% aligned to either intergenic loci or to coding genes in an
99 antisense direction, which suggested that novel transcripts might exist. To test
100 this, we assembled the mapped mouse transcriptome data in a *de novo*
101 approach using Scripture and Cufflinks to reconstruct transcripts and quantified
102 the expression by masking regions, including those containing snoRNAs, tRNAs,
103 miRNAs and pseudogenes. Transcripts that were significantly covered ($P < 0.01$)
104 were selected to avoid noisy transcripts (**Methods**). In total, Scripture identified
105 593,102 multi-exonic transcripts and Cufflinks, 539,775 transcripts, with an
106 overlap of 500,530 transcripts between the two methods. Of those overlapping
107 transcripts, ~86% (429,818) overlapped with known coding transcripts (annotated

108 in either RefSeq, UCSC, or Ensembl) and 10.2% (51,134) overlapped with
109 known noncoding transcripts (annotated either as snoRNA, tRNA, miRNA, or
110 pseudogenes). This shows the quality of transcripts and their ability to recover
111 known noncoding transcripts. The remaining 3.9% of transcripts (20,018) had no
112 overlap with any known coding or noncoding transcripts.

113

114 **Genome-wide identification and annotation of lncRNAs in mouse.**

115 We applied a computational pipeline to identify putative intergenic lncRNAs along
116 with other types of lncRNAs (e.g., antisense, intronic) (4, 5, 22). We identified
117 16,185 multi-exonic lncRNAs longer than 200 bp and with an expression ≥ 1
118 FPKM (fragments per kilobase of exonic length per million) in at least one given
119 tissue. Importantly, these lncRNAs did not contain transcripts with coding
120 potential as measured by the two independent methods including conservation-
121 independent CPAT (23) and conservation-dependent PhyloCSF (24) (**Methods**).
122 About 85% of this dataset overlapped with previously identified lncRNAs (17, 21,
123 25-29) (**Fig. S1**), supporting the accuracy of our prediction pipeline with a total of
124 34% of all known lncRNAs recovered (**Fig. 1B**). The remaining 2,803 identified
125 lncRNAs were considered as novel lncRNAs in mouse. Further, based on the
126 genomic location of lncRNAs relative to nearest protein-coding gene promoters,
127 we annotated 2,174 antisense (e.g., overlapping the protein-coding gene in an
128 antisense direction), 382 intergenic (e.g., located within 10 kb to the nearest
129 protein-coding gene), and 247 strictly intergenic lncRNAs (e.g., located more

130 than 10 kb away from the nearest protein-coding gene) (**Fig. S2**, and **Fig. 1C** for
131 an example of a novel lncRNA identified in testes).

132

133 **Properties of the 2,803 lncRNAs.**

134 It has been previously shown that lncRNAs comprise few exons, are shorter in
135 length, and are expressed at low levels in a highly tissue- or cell-specific nature
136 (3-5). The 2,803 lncRNAs reported here are consistent with these previous
137 studies. On average, our lncRNA transcripts have fewer exons (3 exons), are
138 shorter (6,336 nucleotides), and are expressed at lower levels (1.56 FPKM) than
139 the average for the 27,259 RefSeq protein-coding transcripts, which (on average)
140 have 10 exons, a length of 50,453 nucleotides, and expression levels of 4.68
141 FPKM (**Fig. S3**). To gain more insight, we combined our novel lncRNAs with all
142 the known lncRNAs and reanalyzed the genomic features by considering the
143 ones with an expression greater than 0.1 FPKM in at least one out of 8 tissues
144 and in a cell line, and the ones that are far from protein-coding genes (e.g., 10 kb
145 away from either a transcription start site [TSS] or a transcriptional end site [TES]
146 of a protein-coding gene). This resulted in 3,759 lncRNAs. On average, these
147 transcripts have an exon size of 482 nucleotides, a transcript size of 9,710
148 nucleotides, an expression level of 1.87 FPKM, and a conservation score of 0.1
149 phastCons. These results further confirmed the genomic features of lncRNA,
150 such as lower expression and conservation levels as compared to protein-coding
151 genes.

152

153 In mammals, lncRNAs are expressed in a tissue-specific manner (3-5). To
154 assess for any tissue specificity of our dataset of lncRNA, we compared each
155 lncRNA expression in a given tissue to its expression in the remaining 8 tissues
156 (**Fig. 2A; Table S2**). We observed that 62% of our novel intergenic lncRNAs are
157 tissue-specific, which is comparable to known intergenic lncRNAs (68% tissue-
158 specific). Moreover, protein-coding genes resulted in 36.4% tissue specificity
159 across the eight tissues and the ES cell line (**Fig. S4**). Overall, the results clearly
160 show that lncRNAs are highly tissue specific in nature. Next, we selected the
161 tissue-specific lncRNAs from our list as previously defined (e.g., with an entropy
162 > 0.4) (4). To experimentally validate a pair of these selected tissue-specific
163 lncRNAs, we measured the expression levels by RT-PCR of the heart (H-lnc1
164 and H-lnc2), liver (L-lnc1 and L-lnc2), and kidney (K-lnc1 and K-lnc2) lncRNAs
165 with respect to the housekeeping gene GAPDH (**Fig. 2B**), which confirmed their
166 tissue-specificity.

167
168 To assess whether our novel lncRNAs have active TSS and regulatory marks,
169 we overlapped CAGE tags and DNase I tags from the FANTOM and ENCODE
170 projects with the promoters of our lncRNA (30, 31). We observed an enrichment
171 of CAGE tags around our lncRNA promoters, as compared to random lncRNA
172 promoters (**Fig. S5A**). We also observed an enrichment of tissue-specific DNase
173 I tags in lncRNA promoters from the brain, kidney, liver, spleen, and thymus
174 tissues as well as for the ES cell line (**Fig. S5B**). Finally, we performed *de novo*
175 motif analysis using lncRNA promoters to explore whether any transcription

176 factors could be regulating these lncRNAs. Indeed, we found several significant
177 transcription factor binding motifs enriched near lncRNA promoters (**Fig. S5C**).
178 These results show that the 2,803 lncRNA promoters are enriched with various
179 regulatory marks in the mouse genome and could potentially have regulatory
180 roles.

181

182 **Genome-wide identification of chromatin state maps in mouse.**

183 Chromatin marks mapping across different cell lines in mammals have been
184 previously used to detect and annotate novel regulatory regions in the genome,
185 including for putative lncRNAs (5, 17, 32). We hypothesized that integrating
186 chromatin state maps with the promoters of the transcripts identified here using
187 RNA-seq expression could guide us in annotating the potential transcripts and in
188 predicting their mode of regulation. A map of chromatin marks was constructed
189 from ~1.4 billion mapped reads obtained from 72 pooled ENCODE genome-wide
190 ChIP-seq datasets in eight tissues (brain, heart, liver, small intestine, kidney,
191 spleen, testes, and thymus) and the one primary ES cell line. The ChIP-seq
192 datasets used included regulatory histone modifications, such as H3 lysine 4
193 mono-methylation (H3K4me1), H3 lysine 4 tri-methylation (H3K4me3), H3 lysine
194 36 tri-methylation (H3K36me3), H3 lysine 27 tri-methylation (H3K27me3), and H3
195 lysine 27 mono-acetylation (H3K27ac), as well as CCCTC-binding factor (CTCF)
196 marks and RNA polymerase II marks.

197

198 We applied the ChromHMM program (32) to create a chromatin state model at
199 200 bp resolution, which resulted in six major chromatin state maps (**Fig. 3A**),
200 including promoter (active and poised), enhancer (strong and poised/weak),
201 transcribed (transcription transition, elongation, and weak transcription),
202 insulator, repressed, and heterochromatin states (**Table S3**). In total, we mapped
203 261,175 promoter states (covering ~1% of the mouse genome), 863,677
204 enhancer states (~3%), 1,133,166 transcribed states (~12%), 150,752 repressed
205 states (~1%), 322,521 insulator states (~1%) and 995,562 heterochromatin
206 states (~82%). To validate the accuracy of the predicted chromatin states or
207 maps, we mapped (at ± 10 kb) our 206,045 unique non-overlapping active
208 promoter maps to known promoters of 23,431 RefSeq protein-coding genes and
209 3,190 RefSeq noncoding genes from TSSs. Our analysis recalled 82% (19,280)
210 of the protein-coding promoters and 75% (2,401) of the noncoding ones. We
211 repeated the above mapping using the poised promoter map and mapped an
212 additional 709 protein-coding and 92 noncoding gene promoters. All together, we
213 successfully mapped 85% of the known protein-coding and 78% noncoding gene
214 promoters. These results indicate that using combinatorial promoter chromatin
215 states to retrieve promoters results in ~6% higher recall than when using only
216 H3K4me3 as an active promoter chromatin mark (33).

217

218 **Classification of lncRNAs using chromatin state maps.**

219 Previously chromatin state maps at promoters were used to define two distinct
220 classes of lncRNAs (17). For example, enhancer-associated lncRNA (eIncRNA)

221 promoters or transcription start sites (TSSs) are depleted of H3K4me3 and
222 enriched with H3K4me1, and promoter-associated lncRNAs (plncRNAs) are
223 enriched with H3K4me3 and depleted of H3K4me1. Using a similar promoter-
224 overlapping approach for our chromatin state maps, we defined these two
225 classes of chromatin-associated lncRNAs across 8 tissues and an ES cell line.
226 For this classification, we first listed ~30,000 unique protein-coding promoter loci
227 and ~19,000 intergenic lncRNA promoter loci (200 bp long), which were then
228 passed through an expression filter (requiring >1 FPKM in a given tissue) and an
229 intergenic filter (requiring them to be 5 kb away from both TSS and TES of
230 protein-coding genes). We found a few thousand lncRNAs that passed the above
231 expression and intergenic filters (namely, 1,385 lncRNAs in whole brain, 1,236 in
232 ES cells, 903 in heart, 870 in kidney, 787 in liver, 435 in small intestine, 878 in
233 spleen, 2,083 in testes, and 932 in thymus). Overall, less than 10% (852) of
234 these intergenic lncRNAs significantly overlapped with an active promoter or a
235 strong enhancer chromatin state ($P < 0.001$, Fisher-exact test) (**Fig. 3B**).

236
237 We next focused our analysis on these significant chromatin state-associated
238 lncRNAs. In total, we identified 852 unique intergenic lncRNA transcripts
239 associated with either an active promoter or a strong enhancer chromatin state
240 (**Table S4, Fig. 3C and D**). This result apparently contradicts a previous study
241 (17), in which 52% of lncRNAs were found to be associated with an enhancer
242 chromatin state, and 48%, with a promoter chromatin state. These differences
243 could arise from several parameters used in the previous study that are distinct

244 to ours: specifically, the previous study considered single exonic transcripts, used
245 CAGE tags to define 5' ends, and used DNase-seq peaks to identify active
246 promoters. However, to check the consistency, we also used CAGE peaks from
247 FANTOM5 and DNase-seq peaks from ENCODE, along with RNA-seq
248 expression to identify active promoter lncRNAs in liver, spleen, and thymus. This
249 re-analysis resulted in more than 40% of the lncRNAs associated with enhancer
250 chromatin state in thymus (~50% with promoter chromatin state) and around 20%
251 in liver and spleen. (Table S5 and Figure 3D, Methods). Finally, we did not
252 notice any enrichment in the number of elncRNAs over plncRNAs in most of the
253 tissue we analyzed except brain and thymus. A total of 852 unique intergenic
254 lncRNAs were thus annotated as chromatin-associated, including 514 plncRNAs
255 and 433 elncRNAs.

256 Our approach successfully identified known enhancer-associated coding RNAs,
257 such as *Fos*, *Rgs2*, *Nr4a2*, and *Elf5* (34), and elncRNAs such as *lincRNA-Cox2*,
258 *lincRNA-Spasm*, and *lincRNA-Haunt* (35) (Fig. S6). Moreover, we also found
259 known promoter-associated coding RNAs in our analysis, such as *Sox2*, *Oct4*,
260 and *Nanog*, and plncRNAs, such as *linc1405* and *linc1428* (5) (Fig. S7).
261 Additionally, by pooling all promoter chromatin state maps into one major
262 promoter chromatin-state map, and enhancers into an enhancer chromatin-state
263 map, we were able to recall 71% of published enhancer-associated lncRNAs
264 (36). Our approach successfully recalled 64% of plncRNAs (74 out of 115) and
265 56% of elncRNAs (69 out of 124) from another study (17). We also
266 experimentally tested histone modifications around the lncRNA promoters, both

267 in mouse ES cells and heart cells (**Fig. S8**) using *Klf4* as negative control and
268 *Zic1* as positive control. All together, our study provides a confident list of
269 chromatin-associated lncRNAs across wide range of tissues in mouse.

270

271 **Properties of the chromatin-associated lncRNAs.**

272 To investigate whether the two types of chromatin-associated lncRNAs have
273 different properties, we calculated their sequence length and expression levels
274 (**Fig. 4A and B**). plncRNAs with a median length of ~6 kb were not significantly
275 different from elncRNAs. However, our finding of a ~6 kb for both elncRNAs and
276 plncRNAs differs from a previous study, which reported them to be ~1 kb long
277 (17). plncRNAs are highly expressed compared to elncRNAs, as previously
278 observed (17). We asked whether these chromatin-associated lncRNAs were
279 enriched in any biological processes by using nearest gene approach and whole-
280 genome background with a GREAT software (37). Indeed, they showed
281 enrichment of various biological processes (**Fig. S9**). Interestingly, we also
282 observed the changes in the status of chromatin-associated lncRNAs based on
283 their respective tissue or cell line. In total, ~17% chromatin-associated lncRNAs
284 (144 out of 852) tend to switch from one chromatin state to another in multiple
285 tissues (**Table S6**). plncRNAs are more likely to switch to plncRNAs and also the
286 percentage of this type of transition is higher than the plncRNAs-to-elncRNAs or
287 the elncRNAs-to-plncRNAs transition (**Fig. 4C, D and Table S6**).

288

289 We hypothesized that if a lncRNA is expressed in a specific tissue and also
290 associated with tissue-specific epigenetic modifications in the same tissue but
291 not in others, it could be associated with regulatory functions. To test this, we
292 selected for lncRNAs with the following characteristics: (1) associated with a
293 specific chromatin state only in ES cells, (2) expressed only in ES cells, (3)
294 associated with DNase I peaks only in ES cell, (4) associated with pluripotent
295 transcription factors in ES cells, and (5) close to a protein-coding gene
296 associated with pluripotency in ES cells. In total, 12 lncRNAs passed the above
297 filters.

298

299 For validation, we focused on a ES cell-specific, predicted regulatory enhancer-
300 associated lncRNA (chr7:132560406-132561472 (-)) located approximately 20 kb
301 away from the protein-coding gene *Kdm8*, which encodes a histone lysine
302 demethylase and regulates embryonic cell proliferation (**Fig. 5A** and **5D**) (37).
303 We named this as lncRNA-*Kdm8*, based on its proximity to the *Kdm8* protein-
304 coding gene. Using the RACE technique, we experimentally characterized the
305 lncRNA-*Kdm8* genomic structure; this revealed at least 3 variants (RACE-a, b
306 and c) in the 5' end of lncRNA-*Kdm8*, and also defined the exon-intron
307 boundaries (**Fig. 5B** and **5C**). We then knocked-down lncRNA-*Kdm8* with two
308 different siRNAs and checked the expression of the *Kdm8* transcript and the
309 positive control gene *Taf3*. As predicted, upon lncRNA knockdown, the
310 expression of the *Kdm8* gene significantly decreased as compared to *Taf3*, which
311 further supported the *cis* mode of enhancer-associated lncRNA gene regulation

312 (Fig. 5E) (38, 39). Together, our results show that chromatin-associated lncRNAs
313 annotated by its chromatin marks could have regulatory roles.
314

315 **DISCUSSION**

316

317 Our study identified novel lncRNAs in mouse by using deep RNA sequencing
318 data from eight tissues and an ES cell line. Public ENCODE large-scale RNA-seq
319 data allowed us to *de novo* reconstruct confident novel lncRNA transcripts. The
320 transcriptome data used in this study to discover lncRNAs go beyond previous
321 lncRNA studies in terms of depth (18). The tissue-specific nature of these
322 lncRNAs is in agreement with previous findings (3-5). The 2,803 lncRNAs
323 included 2,174 antisense and 629 intergenic transcripts. Antisense lncRNAs
324 have been shown to be key regulators and interestingly, many of the antisense
325 lncRNA transcripts we observed were from ES cells. We used intersection of
326 transcripts assembled by using two different *de novo* assemblers and also a
327 stringent expression threshold to filter out the spurious transcripts. Further, we
328 validated the expression of the lncRNA transcripts identified in this study by RT-
329 PCR, thus confirming the quality of the transcripts identified in this study as well
330 as their expression.

331

332 By using ChromHMM, we further characterized combinatorial chromatin state
333 maps in mouse, using more than 70 ChIP-seq datasets across the same tissues
334 used for lncRNA discovery. In previous studies, promoter, enhancer, and
335 insulator maps were identified using a specific set of ChIP-seq datasets, like
336 H3K4me3 (promoter), H3K4me1 with P300 (enhancer), and CTCF (insulator)
337 (33). We built upon that work by further including additional histone marks

338 allowing us to produce more detailed chromatin state maps. For example, the
339 Fendrr lncRNA, which was previously annotated as enhancer-associated, has
340 enhancer histone (p300/H3K4me1) marks (35) at the promoter but is also
341 enriched in H3K27me3 in brain. We conclude that its chromatin status is likely to
342 be poised or to switch to other states rather than to be enhancer-associated,
343 which emphasizes the importance of taking chromatin states into account when
344 classifying chromatin-associated lncRNAs.

345

346 By integrating chromatin state maps and promoters of lncRNAs across eight
347 tissues and an ES cell line, we were able to classify lncRNAs into two classes:
348 promoter-associated lncRNAs and enhancer-associated lncRNAs. Our study
349 provides a comprehensive catalog of chromatin-associated lncRNAs across
350 several mouse tissues. We also observed that plncRNAs were highly expressed,
351 shorter in length compared to other chromatin-associated lncRNAs, and retained
352 their embryonic promoter chromatin status in adult tissues. Experimental
353 knockdown of an enhancer-associated lncRNA partially validated the regulatory
354 behavior of chromatin state-associated lncRNAs in mouse.

355

356 Many of the bi-directional lncRNAs and enhancer-associated RNAs have been
357 shown to be non-polyadenylated (34, 40). However, recent findings (2, 17), along
358 with our study, suggest the existence of poly-adenylated bi-directional transcripts
359 and chromatin-associated RNAs. Still, because of the polyA-based RNA

360 sequencing, we could be missing a large fraction of non-polyadenylated
361 lncRNAs.
362
363 In the future, even more comprehensive catalogs of chromatin-associated
364 lncRNAs should be possible to obtain by associating of chromatin states and
365 lncRNA promoters across all tissues and cell lines in mammals. In addition, using
366 techniques like CRISPR against regulatory lncRNAs would reveal more valuable
367 information. All together, our study provides a novel set of classified lncRNAs,
368 which presents a valuable resource for future genomic experimental studies in
369 mouse.
370

371 **MATERIALS AND METHODS**

372 **Computational procedures**

373 *Data sources.*

374 All data used in the analysis were obtained from public databases. The links from
375 where the data was obtained are listed in **Table S7**.

376 *RNA-seq mapping and transcriptome assembly.*

377 TopHat-2.0.9 (19) was used to map RNA-seq reads against mouse reference
378 genome (mm9), using default parameters unless specified (**Table S8**). Cufflinks
379 (20) was used to assembled mapped reads to transcripts *de novo*, and
380 Cuffmerge was used against high-confidence *de novo* transcripts to generate a
381 single transcript annotation file, using default parameters unless specified (**Table**
382 **S6**). Scripture-v4 (21) was also used to assemble transcripts, using uniquely
383 mapped reads with default parameters unless specified (**Table S8**). Finally,
384 Qualimap-v.08 (41) was used with default parameters to count the number of
385 strand-specific reads overlapping with lncRNAs.

386 *Identification and genomic annotation of lncRNAs.*

387 We filtered out transcripts from 8 tissues and a primary ES cell line pooled by
388 Cuffmerge by using an in-house computational pipeline. Our pipeline relies in
389 previously published software and protocols for identifying lncRNAs from
390 transcriptomics data. The pipeline selects transcripts as lncRNAs by their size
391 (≥ 200 nucleotides), number of exons (≥ 2 exons), expression levels (> 1 FPKM in
392 at least one tissue or cell line that we used), overlap with coding regions (no
393 overlap with a known gene set from RefSeq, Ensembl, or UCSC on a similar

strand), overlap with noncoding regions (no overlap with known snoRNAs, tRNAs, miRNAs, lncRNAs, or pseudogenes), and noncoding potential (<0.44 CPAT and <100 PhyloCSF score). PhyloCSF (24) was used to calculate the coding potential of transcripts. First, we stitched mouse lncRNA exonic sequences into 18 mammals, using mm9-multiz30way alignments from UCSC. Second, we ran PhyloCSF against stitched sequences using default parameter unless specified (**Table S8**). We then removed the transcripts with open reading frames with a PhyloCSF score greater than 100, as previously suggested (36). The final lncRNA PhyloCSF score is the average decibans score of all its exons based on their strand direction and all possible frames. The transcripts that passed PhyloCSF and CPAT coding potential filters were further selected as potential lncRNAs.

lncRNAs that did not overlap with any known protein-coding gene (within a 10 kb window from both TSS and TES) were classified as intergenic lncRNAs or lncRNAs. lncRNAs that overlapped a transcript but on opposite strands were classified as antisense lncRNAs. lncRNAs that were close to a coding gene (within 10 kb from both TSS and TES) were annotated as either convergent (same strand as the nearest coding) or divergent (opposite strand as the nearest coding) lncRNAs.

Tissue specificity calculations.

To calculate tissue specificity of lncRNAs, we normalized the raw FPKM expression values as suggested in previous studies (4, 5). First, we added pseudo-count 1 to every raw FPKM value and second applied log2 normalization

417 to each value, to obtain a non-negative expression vector. Finally, we normalized
418 the expression vector by dividing it by the total expression counts. The resulting
419 matrix of lncRNA normalized expression levels in each of the replica experiments
420 per tissue or cell line was clustered by k-means.

421 *TFBS, CAGE tags, and DNase I sites enrichment analyses.*

422 To identify transcription factor binding sites, we first performed a *de novo* motif
423 analysis on the 2,803 lncRNA 1 kb promoters, using the HOMER software with
424 default parameters unless specified (**Table S8**). Second, the significant ($P < 1e$ -
425 5) *de novo* motifs from HOMER were used as input to the TOMTOM program to
426 search against the JASPAR CORE and UNIPROBE databases (42). Next, we
427 combined all identified motifs from both searches into a final list of transcription
428 factor motifs. We then checked the expression of genes in the master list and
429 required that the candidate transcription factor be expressed in the tissue.
430 Finally, we used the PWMEnrich program (43) to perform motif enrichment
431 analysis.

432 CAGE peak based annotations for mouse samples were downloaded from the
433 FANTOM5 database (30), and DNase I sites from ENCODE (31). We overlapped
434 these with the 2,803 lncRNA promoters and their corresponding random regions
435 using *sitepro* from the CEAS program (44) with default parameters. We used the
436 shuffledBed program (<https://code.google.com/p/bedtools/>) (45) with default
437 parameters to randomize the coding RNA and lncRNA promoters in the mm9
438 genome.

439 *Discovery of chromatin state maps.*

440 We first collected mapped ChIP-seq reads of H3K4me1, H3K4me3, H3K36me3,
441 H3K27me3, H3K27ac, CTCF, and RNA polymerase II from ENCODE. This data
442 was originally produced from mouse (C57BL/6-strain, E14, or 8-week-old) brain,
443 heart, kidney, liver, small intestine, spleen, testes, or thymus, or from an
444 embryonic stem (ES) cell line. Second, we used a Poisson-based multivariate
445 hidden Markov model²⁹ (ChromHMM, <http://compbio.mit.edu/ChromHMM/>) to
446 identify regions enriched in specific combinations of histone modifications as
447 previously described but without extending the read lengths. We ran the
448 ChromHMM software to produce classified maps containing from 2 to 50 states.
449 The 15-state model was rich enough and, at the same time, allowed us to
450 interpret the chromatin frequency observed across various tissues and cell lines.
451 Next, we classified the 15-state model into the final six major chromatin state
452 maps of active promoter, poised promoter, strong enhancer, poised or weak
453 enhancer, insulator, repressed, transcribed, or heterochromatin states. In total,
454 3,612,616 regions in the mouse genome were enriched in at least one of the six
455 major chromatin state maps.

456 *Collection of RNA promoters.*

457 We overlapped all 19,873 lncRNAs with protein-coding genes and removed the
458 ones that overlapped by at least one nucleotide on either strand. This resulted in
459 14,147 intergenic lncRNAs. We avoided protein-coding vicinities by removing the
460 lncRNAs that fall within 1 kb from either the TSS or the TES of any known
461 protein-coding gene. This resulted in 12,129 strictly intergenic lncRNAs. Further,
462 we selected lncRNAs with an expressed of more than 1 FPKM in a given tissue.

463 All together, the filters resulted in 1,385 lncRNAs in whole brain, 1,236 in ES
464 cells, 903 in heart, 870 in kidney, 787 in liver, 435 in small intestine, 878 in
465 spleen, 2,083 in testes and 932 in thymus. We created 200 bp promoters of
466 these expressed lncRNAs by extending the TSS 100 bp upstream and
467 downstream. We created random promoters by shuffling across intergenic space
468 and then overlapped these promoters with chromatin states in each tissue
469 separately. Next, we used ~30,000 RefSeq protein-coding gene promoters and
470 overlapped them with chromatin states in a similar fashion as above (>1 FPKM in
471 a given tissue).

472 *Overlapping chromatin state maps with RNA promoters.*

473 We used intersectBed from BEDtools package (45) to overlap RNA promoters
474 with chromatin state maps in each tissue or cell line. We considered the
475 chromatin association to be significant if the *P* value was less than 0.001
476 (Fischer-exact test) in all the tissues we tested. We found both active promoter
477 and strong enhancer chromatin states significantly associated with lncRNA
478 promoters (**Table S4** and **Fig. 3B**). We used CAGE peaks from FANTOM5 and
479 DNase-seq peaks from ENCODE, along with RNA-seq expression, to identify
480 active promoters lncRNA in liver, spleen, and thymus. We could not find both
481 CAGE and DNase-seq data for other tissues. We used the same 200 bp
482 promoter size for CAGE peaks (more than 1 tag) and overlapping DNase-seq
483 peaks (**Table S5**).

484 *Transition of chromatin-associated lncRNAs.*

485 We selected 200-bp-long promoters of expressed lncRNAs (>1 FPKM) in whole
486 brain and made sure that they did not overlap any protein-coding genes within a
487 5 kb distance (both from TSS and TES). We then overlapped the lncRNA
488 promoters with active promoter and strong enhancer chromatin states in whole
489 brain. The analysis resulted in 163 elncRNAs and 33 plncRNAs in whole brain.
490 We repeated the above steps in other tissues, resulting in hundreds of
491 chromatin-associated lncRNAs. This produced 41 ES-elncRNAs, 131 ES-
492 plncRNAs, 21 heart-elncRNAs, 61 heart-plncRNAs, 47 kidney-elncRNAs, 61
493 kidney-plncRNAs, 35 liver-elncRNAs, 77 liver-plncRNAs, 25 small intestine-
494 elncRNAs, 20 small intestine-plncRNAs, 20 spleen-elncRNAs, 65 spleen-
495 plncRNAs, 88 testes-elncRNAs, 258 testes-plncRNAs, 82 thymus-elncRNAs and
496 50 thymus-plncRNAs. Finally, we calculated the percentage of transition of
497 chromatin-associated lncRNA from one tissue to another (**Table S6**).

498 *Gene ontology analysis.*

499 We ran GREAT annotation tool on chromatin-associated lncRNA genomic
500 locations by taking the two nearest genes, using a default of a 1,000 kb distance
501 window. A whole-genome background was selected as a control.

502 **Experimental Procedures**

503 *Cell culture.*

504 Wild-type (E14Tg2A) ESCs were cultured feeder-free in plates coated with 0.1%
505 of gelatin in Glasgow minimum essential medium (Sigma) supplemented with β -
506 mercaptoethanol, sodium pyruvate, essential amino acids, GlutaMAX, 20% fetal
507 bovine serum (Hyclone), and leukemia inhibitory factor (LIF). Heart, liver, and

508 kidneys were isolated from 8-week-old C57BL/6J mice and snap-frozen before
509 RNA extraction for chromatin immunoprecipitation assays (only heart).
510 *Chromatin immunoprecipitation assay.*
511 ESCs were cross-linked in 1% formaldehyde (FA) for 10 min at room
512 temperature (RT). For ChIPs from heart, crosslinking was performed on 1- to 3-
513 mm³ fragments in a conical tube for 10 min rotating at RT in 1.5% FA.
514 Crosslinking was quenched with 0.125 M glycine for 5 min. Pelleted cells and
515 heart fragments were lysed and homogenized. Chromatin extraction and
516 immunoprecipitation was performed as previously described (Morey et al, 2012).
517 300 µg were used for immunoprecipitation. Antibodies used were: Suz12, Abcam
518 ab12073; histone H3, Abcam ab1791; histone H3K4me1, Abcam ab8895;
519 histone H3K27me3, Active-Motif 39155; and histone H3K27ac, Millipore 07-360.
520 The primers used in the qPCR assays are listed in **Table S2**.
521 *Expression and siRNA knockdown analyses.*
522 RNA from organs was extracted with Trizol (Life Technologies). cDNA was
523 generated from 1 mg of RNA with the First Strand cDNA Synthesis Kit
524 (Fermentas). The primers used in the RT-qPCR assays are listed in **Table S2**.
525 RT-PCR was performed in duplicates using GAPDH as a housekeeping gene for
526 normalization. For ES-specific lncRNA knock-downs, 50,000 cells/well in 6-well
527 plates were seeded and then transfected the next day with Lipofectamine
528 RNAiMAX Reagent and 75 pmol of siRNA duplexes (Invitrogen). Cells were
529 pelleted 24 h post-transfection, and RNA was extracted for RT-qPCR with the
530 RNA extraction kit (QIAGEN). cDNA was generated as explained above. The

531 primers used in the RT-qPCR assays and the siRNA duplexes used are listed in
532 **Table S9**. RT-PCR was performed in triplicates using GAPDH as a
533 housekeeping gene for normalization.

534 *Characterization of mouse lncRNA-Kdm8 using RACE.*

535 Total RNA extracted from mouse ES cells (E14) was used to generate RACE-
536 ready 3'- and 5'-cDNA using the SMARTer RACE cDNA Amplification Kit
537 (Clontech) following the manufacturer's protocol. cDNA ends were amplified with
538 universal primer mix and gene-specific primers (GSP), followed by a 'nested'
539 PCR with the nested universal primer and the nested gene-specific primers
540 (NGSP) (**Table S9**). RACE products were run on a 2% agarose gel, cloned in
541 pRACE (pUC19-based vector), and sequenced using M13 primers. Recovered
542 fragments were aligned to obtain the different full-length transcripts produced by
543 the lncRNA-Kdm8 (**Table S9**).

544 *Data access.*

545 All lncRNAs and chromatin state maps identified in this work for mouse (mm9)
546 are listed in the additional files lncRNAs.xlsx and ChromatinMaps.zip (**Table S7**).

547 **Competing Interests**

548 The authors declare that they have no competing interests.

549

550 **Authors' Contributions**

551 G.K.B. conceived the study, collected the data, analyzed the data, interpreted the
552 data, and wrote the manuscript. P.V. conducted qPCR and CHIP-PCR

553 experiments. L.W.S., M.B., L.D.C. and M.A.M-R. contributed ideas and wrote the
554 manuscript.

555

556 **ACKNOWLEDGMENTS**

557 We sincerely thank ENCODE consortium for publicly providing rich data. We are
558 thankful for the many productive discussions, especially with Rory Johnson
559 (lncRNAs), Jason Ernst and Guillaume Fillion (chromatin state maps), Irwin
560 Jungreis (PhyloCSF), Jochen Hecht (RACE), Sabah Kadri (Scripture), and
561 Veronica Raker (manuscript edition). We also thank the three anonymous
562 reviewers for their critical insights. The project was supported by a grant from la
563 Caixa to G.K.B. and by the Spanish MINECO to M.A.M-R. (BFU2010-19310 and
564 BFU2013-47736-P). We also acknowledge support of the Spanish Ministry of
565 Economy and Competitiveness, 'Centro de Excelencia Severo Ochoa 2013-
566 2017', SEV-2012-0208.

567

568 **FIGURES**

569 **Figure 1. Overview of the lncRNA discovery and chromatin state map**
570 **computational pipeline.** (A) Overview of the lncRNA discovery and chromatin
571 state map based classification pipeline that was employed using both RNA-seq
572 and ChIP-seq data from 8 tissues and one primary cell line (ES) in mouse. RNA-
573 seq reads from all the tissues and the cell line were mapped using TopHat2
574 against mouse reference genome (mm9), and transcriptomes were assembled in
575 *de novo* using Cufflink2 and ScriptureV4 assemblers. Common transcripts that
576 were assembled by both Cufflinks2 and ScriptureV4 were scanned for lncRNA
577 features like size, length, exon number, expression and coding score. A library of
578 intergenic lncRNAs was constructed by pooling lncRNAs identified in this study
579 and previous studies. In total, 10,728 unique lncRNAs were overlapped with
580 chromatin state maps discovered by using ChromHMM by pooling various ChIP-
581 Seq datasets and classified chromatin-associated lncRNAs in mouse. (B)
582 Overlap between lncRNAs identified in this study (light grey, left) and previously
583 published lncRNAs (dark grey, right). 2,803 non-annotated lncRNAs were
584 identified, and 34% (13,382) of the known lncRNAs were recovered in this study.
585 (C) RNA-seq coverage tracks showing the expression of a novel lncRNA
586 identified in this study (black). Transcription in testes is shown. "+" and "-"
587 indicate sense and antisense directions, respectively, and experimental replicas
588 are numbered as "1" and "2".

590 **Figure 2. Tissue- and cell-specific expression of lncRNAs.** (A) Heatmap
591 representing normalized FPKM expression values of the 2,803 lncRNAs (rows)
592 across eight tissues and a primary cell line (columns). Rows and columns were
593 ordered based on k-means clustering. Legend color intensity represents the
594 fractional density across the row of log10-normalized FPKM expression values
595 as estimated by ScriptureV4. Each tissue has 2 columns, representing their
596 replicates, and the ES cell line has 5 columns. (B) Experimentally validated
597 examples of lncRNAs with tissue-specific expression across heart, liver, and
598 kidney. Shown are RT-PCR duplicates normalized (against housekeeping gene
599 GAPDH) expression of heart-specific lncRNAs (H-lnc1 and H-lnc2), liver-specific
600 lncRNAs (L-lnc1 and L-lnc2), and kidney-specific lncRNAs (K-lnc1 and K-lnc2)
601 (**Table S9**).

602
603 **Figure 3. Discovery of chromatin state maps and their association with**
604 **lincRNAs.** (A) Emission parameters learned *de novo* with ChromHMM on the
605 basis of combinations recurring in chromatin. Each point in the table denotes the
606 frequency with which a given mark is found at genomic positions corresponding
607 to a specific chromatin state. The observation frequency of various chromatin
608 marks, including H3K36me3, H3K4me1, H3K27ac, Pol II, H3K4me3, CTCF, and
609 H3K27me3, as well as respective input showing 6 major chromatin states,
610 including active promoter (red), poised promoter (purple), enhancer (yellow),
611 Polycomb (grey), insulator (blue), and heterochromatin (white). (B) Percentage of
612 protein-coding TSS (top) and intergenic lncRNAs (bottom) significantly enriched

613 with both active promoter and strong enhancer (***, $P < 0.001$, Fisher-exact test).
614 “D” and “R” labels correspond to the observed data and randomized TSSs,
615 respectively. (C) Percentage of lncRNAs and protein-coding genes that are
616 associated with promoter and enhancer chromatin states. (D) The number of
617 plncRNAs and elncRNAs across 8 tissues and an ES cell line. (E) Percentage of
618 lncRNAs (overlapped with both CAGE peaks and DNase I hypersensitive sites)
619 associated with promoter and enhancer chromatin states.

620
621 **Figure 4. Transcript length, expression, and transition of chromatin-**
622 **associated lncRNAs in mouse.** (A) Transcript length of elncRNAs (median =
623 6565 nt) and plncRNAs (median = 6450 nt) across eight tissues and a cell line,
624 showing no difference in length (Mann-Whitney test; NS, not significant; $P =$
625 0.9848). (B) Log-normalized expression (FPKM) of elncRNAs (median = 0.08
626 FPKM) and plncRNAs (median = 0.33 FPKM) across eight tissues and an ES cell
627 line, showing a significant difference between them (Mann-Whitney test, *** $P =$
628 1.221e-10). (C) Circos plot showing the transition of plncRNA to elncRNA, or
629 elncRNA to plncRNA, across eight tissues and an ES cell line. Top bars indicate
630 the total number of chromatin-associated lncRNAs that undergo a transition per
631 tissue or cell line, which included whole brain (20 plncRNAs and 72 elncRNAs),
632 ES cells (62 plncRNAs and 8 elncRNAs), heart (44 plncRNAs and 4 elncRNAs),
633 small intestine (17 plncRNAs and 18 elncRNAs), kidney (50 plncRNAs and 24
634 elncRNAs), liver (46 plncRNAs and 10 elncRNAs), spleen (55 plncRNAs and 12
635 elncRNAs), testes (29 plncRNAs and 12 elncRNAs) and thymus (47 plncRNAs
636 and 40 elncRNAs). Links inside the bars indicate the number of lncRNAs that
637 switch their chromatin states from one tissue to another (red, plncRNAs; gold,
638 elncRNAs). The lncRNA transition table used to generate the circos plot is shown
639 in Table S6. (D) Percentage of chromatin-associated transitions across all the
640 mouse tissues, showing the high percentage of plncRNA-to-plncRNA transitions
641 as compared to elncRNA-to-elncRNA transitions.

642
643 **Figure 5. An enhancer-associated lncRNA, lncRNA-*Kdm8*, regulates the**
644 **expression of a neighboring protein-coding gene *Kdm8*.** (A) The lncRNA-
645 *Kdm8* locus promoter overlaps with a enhancer chromatin state and occurs
646 within 20 kb of the TSS of a protein-coding gene, *Kdm8* (e.g., it is an enhancer-
647 associated lncRNA). Gene tracks represent DNase I hypersensitive sites (HS)
648 and ChIP-seq data for H3K4me1, H3K27ac, and H3K4me3 from ENCODE. The
649 genomic scale is indicated on the top, and the scale of both DNase I HS and
650 ChIP-seq data on the top right. (B and C) The 5'- and 3'-ends and the exon-
651 intron boundaries of the enhancer-associated lncRNA, lncRNA-*Kdm8*, were
652 determined by RACE (see Supplemental Materials and Methods). Black arrows
653 depict TSSs and direction of transcription for respective genes. *Kdm8* mRNA and
654 the lncRNA-*Kdm8* are shown in green and red, respectively. Genomic DNA
655 sequence corresponding to the 5'- and 3'-ends of the cloned lncRNA are shown
656 in black at the bottom of the lncRNA-*Kdm8* gene track, defining accurate 5'-end
657 and exon-intron boundaries for exon 1, exon 3, exon 4, and exon 5 of lncRNA-
658 *Kdm8*. (D) Expression levels of lncRNA-*Kdm8* in mES cells and other tissues, as

659 measured by directional RNA-seq and expressed as fragments per kilobase of
660 exonic length per million (FPKM). (E) RT-PCR expression (triplicates, normalized
661 against housekeeping gene RPO) after siRNA-based knockdown of lncRNA-
662 *Kdm8* (chr7:132560406-132561472, -) resulted in a significant decrease of
663 neighboring gene *Kdm8* (t-test $*P \leq 0.05$, $**P \leq 0.01$), which was not observed for
664 the negative control of the distant coding gene *Taf3* (chr2:9836179-9970236,
665 +). Primers used for siRNA oligonucleotides of the lncRNA-*Kdm8* are given
666 in **Table S9**.
667

668 REFERENCES

- 669 1. Dunham I, Kundaje A, Aldred SF, Collins PJ, Davis CA, Doyle F,
670 Epstein CB, Frietze S, Harrow J, Kaul R, Khatun J, Lajoie BR, Landt
671 SG, Lee B-K, Pauli F, Rosenbloom KR, Sabo P, Safi A, Sanyal A,
672 Shores N, Simon JM, Song L, Trinklein ND, Altshuler RC, Birney E,
673 Brown JB, Cheng C, Djebali S, Dong X, Dunham I, Ernst J, Furey TS,
674 Gerstein M, Giardine B, Greven M, Hardison RC, Harris RS, Herrero J,
675 Hoffman MM, Iyer S, Kellis M, Khatun J, Kheradpour P, Kundaje A,
676 Lassmann T, Li Q, Lin X, Marinov GK, Merkel A, Mortazavi A, Parker
677 SCJ, Reddy TE, Rozowsky J, Schlesinger F, Thurman RE, Wang J,
678 Ward LD, Whitfield TW, Wilder SP, Wu W, Xi HS, Yip KY, Zhuang J,
679 Bernstein BE, Birney E, Dunham I, Green ED, Gunter C, Snyder M,
680 Pazin MJ, Lowdon RF, Dillon LAL, Adams LB, Kelly CJ, Zhang J,
681 Wexler JR, Green ED, Good PJ, Feingold EA, Bernstein BE, Birney E,
682 Crawford GE, Dekker J, Elnitski L, Farnham PJ, Gerstein M, Giddings
683 MC, Gingeras TR, Green ED, Guigó R, Hardison RC, Hubbard TJ,
684 Kellis M, Kent WJ, Lieb JD, Margulies EH, Myers RM, Snyder M,
685 Stamatoyannopoulos JA, Tenenbaum SA, Weng Z, White KP, Wold B,
686 Khatun J, Yu Y, Wrobel J, Risk BA, Gunawardena HP, Kuiper HC,
687 Maier CW, Xie L, Chen X, Giddings MC, Bernstein BE, Epstein CB,
688 Shores N, Ernst J, Kheradpour P, Mikkelsen TS, Gillespie S, Goren
689 A, Ram O, Zhang X, Wang L, Issner R, Coyne MJ, Durham T, Ku M,
690 Truong T, Ward LD, Altshuler RC, Eaton ML, Kellis M, Djebali S, Davis
691 CA, Merkel A, Dobin A, Lassmann T, Mortazavi A, Tanzer A, Lagarde
692 J, Lin W, Schlesinger F, Xue C, Marinov GK, Khatun J, Williams BA,
693 Zaleski C, Rozowsky J, Röder M, Kokocinski F, Abdelhamid RF, Alioto
694 T, Antoshechkin I, Baer MT, Batut P, Bell I, Bell K, Chakraborty S,
695 Chen X, Chrast J, Curado J, Derrien T, Drenkow J, Dumais E, Dumais
696 J, Duttagupta R, Fastuca M, Fejes-Toth K, Ferreira P, Foissac S,
697 Fullwood MJ, Gao H, Gonzalez D, Gordon A, Gunawardena HP,
698 Howald C, Jha S, Johnson R, Kapranov P, King B, Kingswood C, Li G,
699 Luo OJ, Park E, Preall JB, Presaud K, Ribeca P, Risk BA, Robyr D,
700 Ruan X, Sammeth M, Sandhu KS, Schaeffer L, See L-H, Shahab A,
701 Skancke J, Suzuki AM, Takahashi H, Tilgner H, Trout D, Walters N,
702 Wang H, Wrobel J, Yu Y, Hayashizaki Y, Harrow J, Gerstein M,
703 Hubbard TJ, Reymond A, Antonarakis SE, Hannon GJ, Giddings MC,
704 Ruan Y, Wold B, Carninci P, Guigó R, Gingeras TR, Rosenbloom KR,
705 Sloan CA, Learned K, Malladi VS, Wong MC, Barber GP, Cline MS,
706 Dreszer TR, Heitner SG, Karolchik D, Kent WJ, Kirkup VM, Meyer LR,
707 Long JC, Maddren M, Raney BJ, Furey TS, Song L, Grassegger LL,
708 Giresi PG, Lee B-K, Battenhouse A, Sheffield NC, Simon JM, Showers
709 KA, Safi A, London D, Bhinge AA, Shestak C, Schaner MR, Ki Kim S,
710 Zhang ZZ, Mieczkowski PA, Mieczkowska JO, Liu Z, McDaniell RM, Ni
711 Y, Rashid NU, Kim MJ, Adar S, Zhang Z, Wang T, Winter D, Keefe D,
712 Birney E, Iyer VR, Lieb JD, Crawford GE, Li G, Sandhu KS, Zheng M,

713 Wang P, Luo OJ, Shahab A, Fullwood MJ, Ruan X, Ruan Y, Myers RM,
714 Pauli F, Williams BA, Gertz J, Marinov GK, Reddy TE, Vielmetter J,
715 Partridge E, Trout D, Varley KE, Gasper C, Bansal A, Pepke S, Jain P,
716 Amrhein H, Bowling KM, Anaya M, Cross MK, King B, Muratet MA,
717 Antoshechkin I, Newberry KM, McCue K, Nesmith AS, Fisher-Aylor KI,
718 Pusey B, DeSalvo G, Parker SL, Balasubramanian S, Davis NS,
719 Meadows SK, Eggleston T, Gunter C, Newberry JS, Levy SE, Absher
720 DM, Mortazavi A, Wong WH, Wold B, Blow MJ, Visel A, Pennachio LA,
721 Elnitski L, Margulies EH, Parker SCJ, Petrykowska HM, Abyzov A,
722 Aken B, Barrell D, Barson G, Berry A, Bignell A, Boychenko V,
723 Bussotti G, Chrast J, Davidson C, Derrien T, Despacio-Reyes G,
724 Diekhans M, Ezkurdia I, Frankish A, Gilbert J, Gonzalez JM, Griffiths
725 E, Harte R, Hendrix DA, Howald C, Hunt T, Jungreis I, Kay M, Khurana
726 E, Kokocinski F, Leng J, Lin MF, Loveland J, Lu Z, Manthavadi D,
727 Mariotti M, Mudge J, Mukherjee G, Notredame C, Pei B, Rodriguez JM,
728 Saunders G, Sboner A, Searle S, Sisu C, Snow C, Steward C, Tanzer
729 A, Tapanari E, Tress ML, van Baren MJ, Walters N, Washietl S,
730 Wilming L, Zadissa A, Zhang Z, Brent M, Haussler D, Kellis M,
731 Valencia A, Gerstein M, Reymond A, Guigó R, Harrow J, Hubbard TJ,
732 Landt SG, Fietze S, Abyzov A, Addleman N, Alexander RP, Auerbach
733 RK, Balasubramanian S, Bettinger K, Bhardwaj N, Boyle AP, Cao AR,
734 Cayting P, Charos A, Cheng Y, Cheng C, Eastman C, Euskirchen G,
735 Fleming JD, Grubert F, Habegger L, Hariharan M, Harmanci A, Iyengar
736 S, Jin VX, Karczewski KJ, Kasowski M, Lacroute P, Lam H, Lamarre-
737 Vincent N, Leng J, Lian J, Lindahl-Allen M, Min R, Miotto B, Monahan
738 H, Moqtaderi Z, Mu XJ, O'Geen H, Ouyang Z, Patacsil D, Pei B, Raha
739 D, Ramirez L, Reed B, Rozowsky J, Sboner A, Shi M, Sisu C, Slifer T,
740 Witt H, Wu L, Xu X, Yan K-K, Yang X, Yip KY, Zhang Z, Struhl K,
741 Weissman SM, Gerstein M, Farnham PJ, Snyder M, Tenenbaum SA,
742 Penalva LO, Doyle F, Karmakar S, Landt SG, Bhanvadia RR,
743 Choudhury A, Domanus M, Ma L, Moran J, Patacsil D, Slifer T,
744 Vectorsen A, Yang X, Snyder M, White KP, Auer T, Centanin L,
745 Eichenlaub M, Gruhl F, Heermann S, Hoeckendorf B, Inoue D, Kellner
746 T, Kirchmaier S, Mueller C, Reinhardt R, Schertel L, Schneider S, Sinn
747 R, Wittbrodt B, Wittbrodt J, Weng Z, Whitfield TW, Wang J, Collins PJ,
748 Aldred SF, Trinklein ND, Partridge EC, Myers RM, Dekker J, Jain G,
749 Lajoie BR, Sanyal A, Balasundaram G, Bates DL, Byron R, Canfield
750 TK, Diegel MJ, Dunn D, Ebersol AK, Frum T, Garg K, Gist E, Hansen
751 RS, Boatman L, Haugen E, Humbert R, Jain G, Johnson AK, Johnson
752 EM, Kutayavin TV, Lajoie BR, Lee K, Lotakis D, Maurano MT, Neph SJ,
753 Neri FV, Nguyen ED, Qu H, Reynolds AP, Roach V, Rynes E, Sabo P,
754 Sanchez ME, Sandstrom RS, Sanyal A, Shafer AO, Stergachis AB,
755 Thomas S, Thurman RE, Vernot B, Vierstra J, Vong S, Wang H,
756 Weaver MA, Yan Y, Zhang M, Akey JM, Bender M, Dorschner MO,
757 Groudine M, MacCoss MJ, Navas P, Stamatoyannopoulos G, Kaul R,
758 Dekker J, Stamatoyannopoulos JA, Dunham I, Beal K, Brazma A,

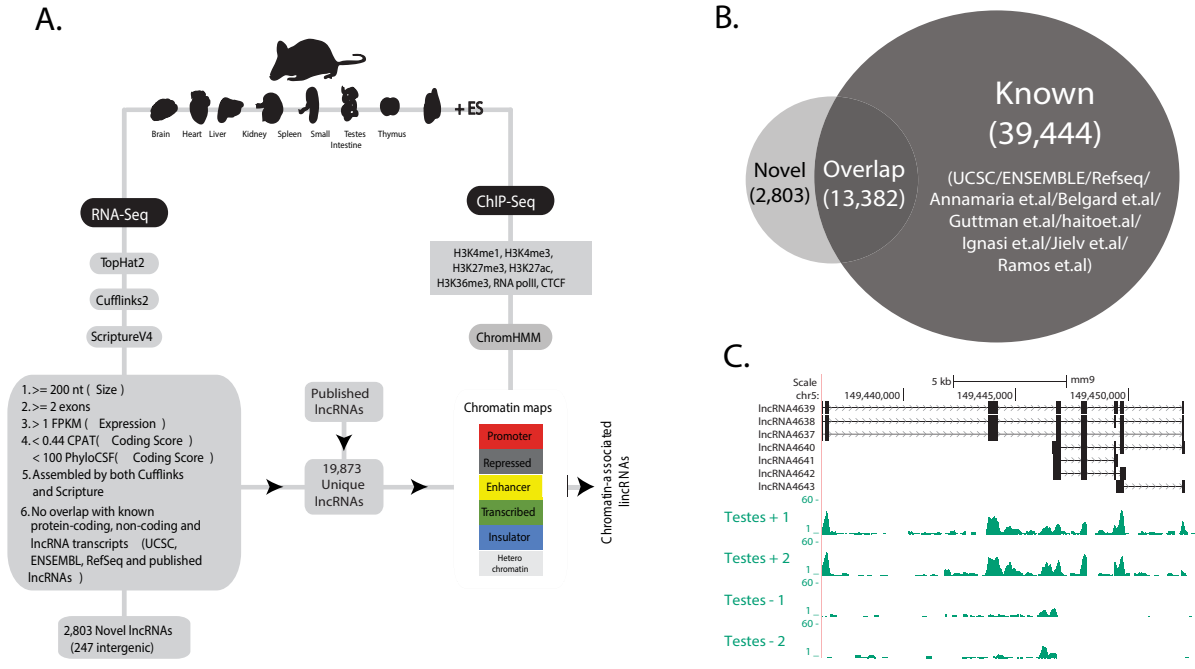
- 759 Flicek P, Herrero J, Johnson N, Keefe D, Lukk M, Luscombe NM,
760 Sobral D, Vaquerizas JM, Wilder SP, Batzoglou S, Sidow A, Hussami
761 N, Kyriazopoulou-Panagiotopoulou S, Libbrecht MW, Schaub MA,
762 Kundaje A, Hardison RC, Miller W, Giardine B, Harris RS, Wu W,
763 Bickel PJ, Banfai B, Boley NP, Brown JB, Huang H, Li Q, Li JJ, Noble
764 WS, Bilmes JA, Buske OJ, Hoffman MM, Sahu AD, Kharchenko PV,
765 Park PJ, Baker D, Taylor J, Weng Z, Iyer S, Dong X, Greven M, Lin X,
766 Wang J, Xi HS, Zhuang J, Gerstein M, Alexander RP,
767 Balasubramanian S, Cheng C, Harmanci A, Lochovsky L, Min R, Mu
768 XJ, Rozowsky J, Yan K-K, Yip KY, Birney E. 2012. An integrated
769 encyclopedia of DNA elements in the human genome. *Nature* **489**:57–74.
- 770 2. Djebali S, Davis CA, Merkel A, Dobin A, Lassmann T, Mortazavi A,
771 Tanzer A, Lagarde J, Lin W, Schlesinger F, Xue C, Marinov GK,
772 Khatun J, Williams BA, Zaleski C, Rozowsky J, Röder M, Kokocinski
773 F, Abdelhamid RF, Alioto T, Antoshechkin I, Baer MT, Bar NS, Batut P,
774 Bell K, Bell I, Chakraborty S, Chen X, Chrast J, Curado J, Derrien T,
775 Drenkow J, Dumais E, Dumais J, Duttagupta R, Falconnet E, Fastuca
776 M, Fejes-Toth K, Ferreira P, Foissac S, Fullwood MJ, Gao H, Gonzalez
777 D, Gordon A, Gunawardena H, Howald C, Jha S, Johnson R, Kapranov
778 P, King B, Kingswood C, Luo OJ, Park E, Persaud K, Preall JB, Ribeca
779 P, Risk B, Robyr D, Sammeth M, Schaffer L, See L-H, Shahab A,
780 Skancke J, Suzuki AM, Takahashi H, Tilgner H, Trout D, Walters N,
781 Wang H, Wrobel J, Yu Y, Ruan X, Hayashizaki Y, Harrow J, Gerstein
782 M, Hubbard T, Reymond A, Antonarakis SE, Hannon G, Giddings MC,
783 Ruan Y, Wold B, Carninci P, Guigó R, Gingeras TR. 2013. Landscape of
784 transcription in human cells. *Nature* **488**:101–108.
- 785 3. Derrien T, Johnson R, Bussotti G, Tanzer A, Djebali S, Tilgner H,
786 Guernec G, Martin D, Merkel A, Knowles DG, Lagarde J, Veeravalli L,
787 Ruan X, Ruan Y, Lassmann T, Carninci P, Brown JB, Lipovich L,
788 Gonzalez JM, Thomas M, Davis CA, Shiekhata R, Gingeras TR,
789 Hubbard TJ, Notredame C, Harrow J, Guigo R. 2012. The GENCODE
790 v7 catalog of human long noncoding RNAs: Analysis of their gene
791 structure, evolution, and expression. *Genome Research* **22**:1775–1789.
- 792 4. Cabili MN, Trapnell C, Goff L, Koziol M, Tazon-Vega B, Regev A, Rinn
793 JL. 2011. Integrative annotation of human large intergenic noncoding
794 RNAs reveals global properties and specific subclasses. *Genes &
795 Development* **25**:1915–1927.
- 796 5. Guttman M, Amit I, Garber M, French C, Lin MF, Feldser D, Huarte M,
797 Zuk O, Carey BW, Cassady JP, Cabili MN, Jaenisch R, Mikkelsen TS,
798 Jacks T, Hacohen N, Bernstein BE, Kellis M, Regev A, Rinn JL, Lander
799 ES. 2009. Chromatin signature reveals over a thousand highly conserved
800 large non-coding RNAs in mammals. *Nature* **458**:223–227.
- 801 6. Nam J-W, Bartel DP. 2012. Long noncoding RNAs in *C. elegans*. *Genome
802 Research* **22**:2529–2540.
- 803 7. Young RS, Marques AC, Tibbit C, Haerty W, Bassett AR, Liu JL,
804 Ponting CP. 2012. Identification and Properties of 1,119 Candidate

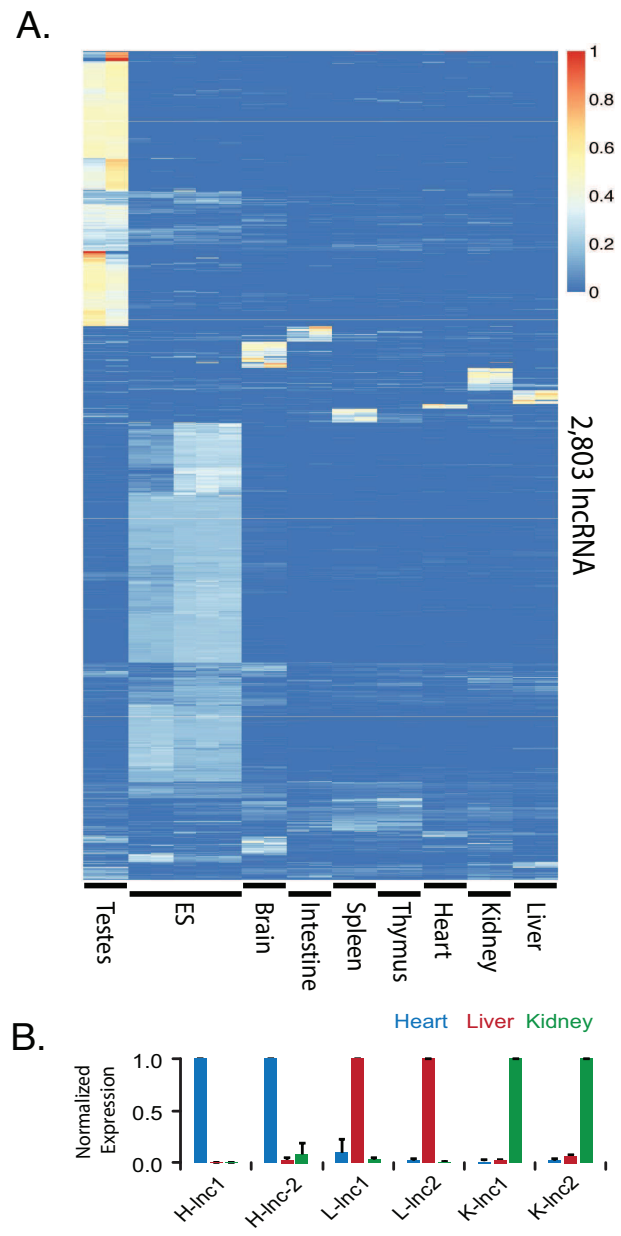
- 805 LincRNA Loci in the *Drosophila melanogaster* Genome. *Genome Biology*
806 and *Evolution* **4**:427–442.
- 807 8. **Pauli A, Valen E, Lin MF, Garber M, Vastenhouw NL, Levin JZ, Fan L,**
808 **Sandelin A, Rinn JL, Regev A, Schier AF.** 2012. Systematic identification
809 of long noncoding RNAs expressed during zebrafish embryogenesis.
810 *Genome Research* **22**:577–591.
- 811 9. **Panning B, Dausman J, Jaenisch R.** 1997. X Chromosome Inactivation
812 Is Mediated by Xist RNA Stabilization. *Cell* **90**:907–916.
- 813 10. **Gupta RA, Shah N, Wang KC, Kim J, Horlings HM, Wong DJ, Tsai M-C,**
814 **Hung T, Argani P, Rinn JL, Wang Y, Brzoska P, Kong B, Li R, West**
815 **RB, van de Vijver MJ, Sukumar S, Chang HY.** 2010. Long non-coding
816 RNA HOTAIR reprograms chromatin state to promote cancer metastasis.
817 *Nature* **464**:1071–1076.
- 818 11. **Wang P, Xue Y, Han Y, Lin L, Wu C, Xu S, Jiang Z, Xu J, Liu Q, Cao X.**
819 2014. The STAT3-Binding Long Noncoding RNA Inc-DC Controls Human
820 Dendritic Cell Differentiation. *Science* **344**:310–313.
- 821 12. **Klattenhoff CA, Scheuermann JC, Surface LE, Bradley RK, Fields PA,**
822 **Steinhauser ML, Ding H, Butty VL, Torrey L, Haas S, Abo R,**
823 **Tabebordbar M, Lee RT, Burge CB, Boyer LA.** 2013. Braveheart, a Long
824 Noncoding RNA Required for Cardiovascular Lineage Commitment. *Cell*
825 **152**:570–583.
- 826 13. **Ulitsky I, Shkumatava A, Jan CH, Sive H, Bartel DP.** 2011. Conserved
827 Function of lincRNAs in Vertebrate Embryonic Development despite Rapid
828 Sequence Evolution. *Cell* **147**:1537–1550.
- 829 14. **Grote P, Wittler L, Hendrix D, Koch F, Währisch S, Beisaw A, Macura**
830 **K, Bläss G, Kellis M, Werber M, Herrmann BG.** 2013. The Tissue-
831 Specific lncRNA Fendrr Is an Essential Regulator of Heart and Body Wall
832 Development in the Mouse. *Developmental Cell* **24**:206–214.
- 833 15. **Tripathi V, Ellis JD, Shen Z, Song DY, Pan Q, Watt AT, Freier SM,**
834 **Bennett CF, Sharma A, Bubulya PA, Blencowe BJ, Prasanth SG,**
835 **Prasanth KV.** 2010. The Nuclear-Retained Noncoding RNA MALAT1
836 Regulates Alternative Splicing by Modulating SR Splicing Factor
837 Phosphorylation. *Molecular Cell* **39**:925–938.
- 838 16. **Ng S-Y, Bogu GK, Soh BS, Stanton LW.** 2013. The Long Noncoding
839 RNA RMST Interacts with SOX2 to Regulate Neurogenesis. *Molecular Cell*
840 **51**:349–359.
- 841 17. **Marques AC, Hughes J, Graham B, Kowalczyk MS, Higgs DR, Ponting**
842 **CP.** 2013. Chromatin signatures at transcriptional start sites separate two
843 equally populated yet distinct classes of intergenic long noncoding RNAs.
844 *Genome Biol* **14**:R131.
- 845 18. **Pervouchine DD, Djebali S, Breschi A, Davis CA, Barja PP, Dobin A,**
846 **Tanzer A, Lagarde J, Zaleski C, See L-H, Fastuca M, Drenkow J, Wang**
847 **H, Bussotti G, Pei B, Balasubramanian S, Monlong J, Harmanci A,**
848 **Gerstein M, Beer MA, Notredame C, Oacutete RG, Gingeras TR.** 1AD.
849 Enhanced transcriptome maps from multiple mouse tissues reveal
850 evolutionary constraint in gene expression. *Nat Comms* **6**:1–11.

- 851 19. **Trapnell C, Pachter L, Salzberg SL.** 2009. TopHat: discovering splice
852 junctions with RNA-Seq. *Bioinformatics* **25**:1105–1111.
- 853 20. **Trapnell C, Roberts A, Goff L, Pertea G, Kim D, Kelley DR, Pimentel H,**
854 **Salzberg SL, Rinn JL, Pachter L.** 2012. Differential gene and transcript
855 expression analysis of RNA-seq experiments with TopHat and Cufflinks.
856 *Nat Protoc* **7**:562–578.
- 857 21. **Guttman M, Garber M, Levin JZ, Donaghey J, Robinson J, Adiconis X,**
858 **Fan L, Koziol MJ, Gnirke A, Nusbaum C, Rinn JL, Lander ES, Regev A.**
859 2010. Ab initio reconstruction of cell type-specific transcriptomes in mouse
860 reveals the conserved multi-exonic structure of lincRNAs. *Nat Biotechnol*
861 **28**:503–510.
- 862 22. **Jia H, Osak M, Bogu GK, Stanton LW, Johnson R, Lipovich L.** 2010.
863 Genome-wide computational identification and manual annotation of
864 human long noncoding RNA genes. *RNA* **16**:1478–1487.
- 865 23. **Wang L, Park HJ, Dasari S, Wang S, Kocher J-P, Li W.** 2013. CPAT:
866 Coding-Potential Assessment Tool using an alignment-free logistic
867 regression model. *Nucleic Acids Research*.
- 868 24. **Lin MF, Jungreis I, Kellis M.** 2011. PhyloCSF: a comparative genomics
869 method to distinguish protein coding and non-coding regions.
870 *Bioinformatics* **27**:i275–i282.
- 871 25. **Luo H, Sun S, Li P, Bu D, Cao H, Zhao Y.** 2013. Comprehensive
872 Characterization of 10,571 Mouse Large Intergenic Noncoding RNAs from
873 Whole Transcriptome Sequencing. *PLoS ONE* **8**:e70835.
- 874 26. **Morán I, Akerman I, van de Bunt M, Xie R, Benazra M, Nammo T,**
875 **Arnes L, Nakić N, García-Hurtado J, Rodríguez-Seguí S, Pasquali L,**
876 **Sauty-Colace C, Beucher A, Scharfmann R, van Arensbergen J,**
877 **Johnson PR, Berry A, Lee C, Harkins T, Gmyr V, Pattou F, Kerr-Conte**
878 **J, Piemonti L, Berney T, Hanley N, Gloyn AL, Sussel L, Langman L,**
879 **Brayman KL, Sander M, McCarthy MI, Ravassard P, Ferrer J.** 2012.
880 Human β Cell Transcriptome Analysis Uncovers lncRNAs That Are Tissue-
881 Specific, Dynamically Regulated, and Abnormally Expressed in Type 2
882 Diabetes. *Cell Metabolism* **16**:435–448.
- 883 27. **Lv J, Cui W, Liu H, He H, Xiu Y, Guo J, Liu H, Liu Q, Zeng T, Chen Y,**
884 **Zhang Y, Wu Q.** 2013. Identification and Characterization of Long Non-
885 Coding RNAs Related to Mouse Embryonic Brain Development from
886 Available Transcriptomic Data. *PLoS ONE* **8**:e71152.
- 887 28. **Ramos AD, Diaz A, Nellore A, Delgado RN, Park K-Y, Gonzales-Roybal**
888 **G, Oldham MC, Song JS, Lim DA.** 2013. Integration of Genome-wide
889 Approaches Identifies lncRNAs of Adult Neural Stem Cells and Their
890 Progeny In Vivo. *Cell Stem Cell* **12**:616–628.
- 891 29. **Belgard TG, Marques AC, Oliver PL, Abaan HO, Sirey TM, Hoerder-**
892 **Suabedissen A, García-Moreno F, Molnár Z, Margulies EH, Ponting**
893 **CP.** 2011. NeuroResource. *Neuron* **71**:605–616.
- 894 30. **DGT TFCATRPAC.** 2015. A promoter-level mammalian expression atlas.
895 *Nature* **507**:462–470.
- 896 31. **Mouse ENCODE Consortium, Stamatoyannopoulos JA, Snyder M,**

- 897 Hardison R, Ren B, Gingeras T, Gilbert DM, Groudine M, Bender M,
898 Kaul R, Canfield T, Giste E, Johnson A, Zhang M, Balasundaram G,
899 Byron R, Roach V, Sabo PJ, Sandstrom R, Stehling AS, Thurman RE,
900 Weissman SM, Cayting P, Hariharan M, Lian J, Cheng Y, Landt SG, Ma
901 Z, Wold BJ, Dekker J, Crawford GE, Keller CA, Wu W, Morrissey C,
902 Kumar SA, Mishra T, Jain D, Byrsk-Bishop M, Blankenberg D, Lajoie
903 BR, Jain G, Sanyal A, Chen K-B, Denas O, Taylor J, Blobel GA, Weiss
904 MJ, Pimkin M, Deng W, Marinov GK, Williams BA, Fisher-Aylor KI,
905 DeSalvo G, Kiralusha A, Trout D, Amrhein H, Mortazavi A, Edsall L,
906 McCleary D, Kuan S, Shen Y, Yue F, Ye Z, Davis CA, Zaleski C, Jha S,
907 Xue C, Dobin A, Lin W, Fastuca M, Wang H, Guigó R, Djebali S,
908 Lagarde J, Ryba T, Sasaki T, Malladi VS, Cline MS, Kirkup VM,
909 Learned K, Rosenbloom KR, Kent WJ, Feingold EA, Good PJ, Pazin
910 M, Lowdon RF, Adams LB. 2012. An encyclopedia of mouse DNA
911 elements (Mouse ENCODE). *Genome Biol* **13**:418.
- 912 32. Ernst J, Kellis M. 2012. ChromHMM: automating chromatin-state
913 discovery and characterization. *Nat Meth* **9**:215–216.
- 914 33. Shen Y, Yue F, McCleary DF, Ye Z, Edsall L, Kuan S, Wagner U, Dixon
915 J, Lee L, Lobanenko VV, Ren B. 2012. A map of the cis-regulatory
916 sequences in the mouse genome. *Nature* **488**:116–120.
- 917 34. Kim T-K, Hemberg M, Gray JM, Costa AM, Bear DM, Wu J, Harmin DA,
918 Laptewicz M, Barbara-Haley K, Kuersten S, Markenscoff-
919 Papadimitriou E, Kuhl D, Bito H, Worley PF, Kreiman G, Greenberg
920 ME. 2010. Widespread transcription at neuronal activity-regulated
921 enhancers. *Nature* **465**:182–187.
- 922 35. Sauvageau M, Goff LA, Lodato S, Bonev B, Groff AF, Gerhardinger C,
923 Sanchez-Gomez DB, Hacisuleyman E, Li E, Spence M, Liapis SC,
924 Mallard W, Morse M, Swerdel MR, D'Ecclesiss MF, Moore JC, Lai V,
925 Gong G, Yancopoulos GD, Frendewey D, Kellis M, Hart RP,
926 Valenzuela DM, Arlotta P, Rinn JL. 2013. Multiple knockout mouse
927 models reveal lincRNAs are required for life and brain development. *eLife*
928 **2**:e01749–e01749.
- 929 36. Alvarez-Dominguez JR, Hu W, Yuan B, Shi J, Park SS, Gromatzky AA,
930 Oudenaarden AV, Lodish HF. 2014. Global discovery of erythroid long
931 noncoding RNAs reveals novel regulators of red cell maturation. *Blood*
932 **123**:570–581.
- 933 37. McLean CY, Bristor D, Hiller M, Clarke SL, Schaar BT, Lowe CB,
934 Wenger AM, Bejerano G. 2010. GREAT improves functional interpretation
935 of cis-regulatory regions. *Nat Biotechnol* **28**:495–501.
- 936 38. Lai F, Orom UA, Cesarini M, Beringer M, Taatjes DJ, Blobel GA,
937 Shiekhattar R. 2013. Activating RNAs associate with Mediator to enhance
938 chromatin architecture and transcription. *Nature* **494**:497–501.
- 939 39. Qrom UA, Derrien T, Beringer M, Gumireddy K, Gardini A, Bussotti G.
940 2010. Long Noncoding RNAs with Enhancer-like Function. *Cell* **143**:46–58.
- 941 40. Wu X, Sharp PA. 2013. Perspective. *Cell* **155**:990–996.
- 942 41. Garcia-Alcalde F, Okonechnikov K, Carbonell J, Cruz LM, Gotz S,

- 943 **Tarazona S, Dopazo J, Meyer TF, Conesa A.** 2012. Qualimap: evaluating
944 next-generation sequencing alignment data. *Bioinformatics* **28**:2678–2679.
945 42. **Gupta S, Stamatoyannopoulos JA, Bailey TL, Noble W.** 2007.
946 Quantifying similarity between motifs. *Genome Biol* **8**:R24.
947 43. **Stojnic R.** 2014. PWMEnrich: PWM enrichment analysis. R package
948 version 3.6.1 1–46.
949 44. **Shin H, Liu T, Manrai AK, Liu XS.** 2009. CEAS: cis-regulatory element
950 annotation system. *Bioinformatics* **25**:2605–2606.
951 45. **Quinlan AR, Hall IM.** 2010. BEDTools: a flexible suite of utilities for
952 comparing genomic features. *Bioinformatics* **26**:841–842.
953



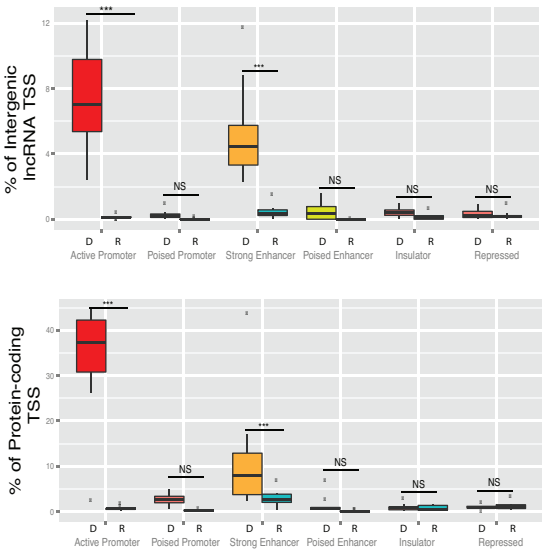


A.

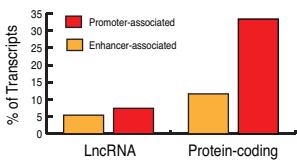
	H3K36me3	H3K4me1	H3K27ac	Pol2	Input	H3K4me3	CTCF	H3K27me3	Coverage (Mean)	Length (Mean)	
1	81	1	1	4	1	0	3	3	4.2	2.5	Transcription Elongation
2	16	1	0	1	1	0	1	2	6.7	2.7	Weakly Transcribed
3	84	60	33	12	1	1	6	7	0.8	0.9	Transcriptional Transition
4	90	57	62	34	2	88	12	13	0.6	0.7	Weak/poised Enhancer
5	5	23	7	8	1	74	4	6	0.4	0.5	Active Promoter
6	10	89	70	62	5	96	36	41	0.2	0.6	Strong Enhancer
7	4	11	85	61	6	97	23	20	0.4	0.9	Active Promoter
8	6	62	81	15	2	3	9	8	0.7	0.7	Strong Enhancer
9	3	37	4	5	1	0	2	6	2.1	0.8	Weak/poised Enhancer
10	6	54	7	17	3	53	17	89	0.3	0.9	Poised Promoter
11	2	2	0	1	1	0	2	49	1.3	1.8	Repressed
12	1	0	0	0	0	0	0	6	18.5	17.7	Heterochromatin
13	0	0	0	0	0	0	0	1	44.1	120.7	Heterochromatin
14	0	1	0	2	1	0	1	2	19.7	5.7	Heterochromatin
15	4	12	3	21	2	1	41	12	0.7	0.4	Insulator

Chromatin mark observation frequency (%)

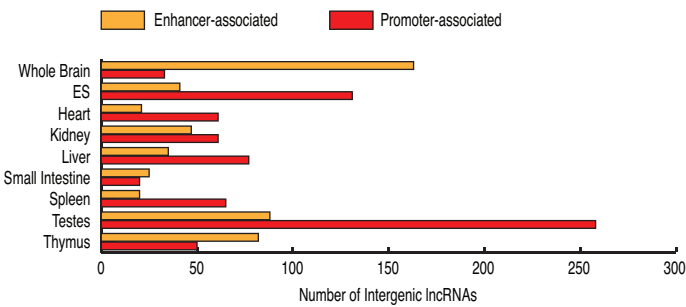
B.



C.



D.



E.

