

MIT Open Access Articles

*Local regulation of gene expression by
lncRNA promoters, transcription and splicing*

The MIT Faculty has made this article openly available. **Please share** how this access benefits you. Your story matters.

Citation: Engreitz, Jesse M. et al. "Local Regulation of Gene Expression by lncRNA Promoters, Transcription and Splicing." *Nature* 539, 7629 (October 2016): 452–455 © 2016 Macmillan Publishers Limited, part of Springer Nature

As Published: <https://doi.org/10.1038/nature20149>

Publisher: Nature Publishing Group

Persistent URL: <http://hdl.handle.net/1721.1/117775>

Version: Author's final manuscript: final author's manuscript post peer review, without publisher's formatting or copy editing

Terms of Use: Article is made available in accordance with the publisher's policy and may be subject to US copyright law. Please refer to the publisher's site for terms of use.



1 **Local regulation of gene expression by**
2 **lncRNA promoters, transcription, and splicing**
3
4

5 **Authors:** Jesse M. Engreitz^{1,2}, Jenna E. Haines^{1,6}, Elizabeth M. Perez¹, Glen Munson¹, Jenny
6 Chen^{1,2}, Michael Kane¹, Patrick E. McDonel^{1,7}, Mitchell Guttman³, Eric S. Lander^{1,4,5*}

7 **Affiliations:**

8 ¹Broad Institute of Harvard and MIT, Cambridge MA 02142

9 ²Division of Health Sciences and Technology, MIT, Cambridge MA 02139

10 ³Division of Biology and Biological Engineering, California Institute of Technology, Pasadena
11 CA 91125

12 ⁴Department of Biology, MIT, Cambridge MA 02139

13 ⁵Department of Systems Biology, Harvard Medical School, Boston MA 02114

14 ⁶Present address: Department of Molecular & Cell Biology, University of California Berkeley,
15 Berkeley CA 94720

16 ⁷Present address: University of Massachusetts Medical School, Worcester MA 01655

17 *Correspondence to: eric@broadinstitute.org (ESL)

18
19
20
21 **Mammalian genomes are pervasively transcribed^{1,2} to produce thousands of long**
22 **noncoding RNAs (lncRNAs)^{3,4}. A few of these lncRNAs have been shown to recruit**
23 **regulatory complexes through RNA-protein interactions to influence the expression of**
24 **nearby genes⁵⁻⁷, and it has been suggested that many other lncRNAs similarly act as local**
25 **regulators^{8,9}. Such local functions could explain the observation that lncRNA expression is**
26 **often correlated with the expression of nearby genes^{2,10,11}. However, such correlations have**
27 **been challenging to dissect¹² and could alternatively result from processes that are not**
28 **mediated by the lncRNA transcripts themselves. For example, some gene promoters have**
29 **been proposed to have dual functions as enhancers¹³⁻¹⁶, and the process of transcription *per***
30 ***se* has been proposed to contribute to gene regulation by recruiting activating factors or**
31 **remodeling nucleosomes^{10,17,18}. Here we used genetic manipulations to dissect 12 genomic**
32 **loci that produce lncRNAs and found that 5 of these loci influence the expression of a**
33 **neighboring gene in *cis*. Surprisingly, none of these effects required the specific lncRNA**

34 transcripts themselves and instead involved general processes associated with their
35 production, including enhancer-like activity of gene promoters, the process of
36 transcription, and the splicing of the transcript. Importantly, such effects were not limited
37 to lncRNA loci: we found that 4 of 6 protein-coding loci similarly influenced the expression
38 of a neighbor. These results demonstrate that ‘crosstalk’ among neighboring genes is a
39 prevalent phenomenon that can involve multiple mechanisms and *cis* regulatory signals,
40 including a novel role for RNA splice sites. These mechanisms may explain the function
41 and evolution of some genomic loci that produce lncRNAs and broadly contribute to the
42 regulation of both coding and noncoding genes.

43 We analyzed 12 lncRNA loci whose RNA transcripts in mouse embryonic stem cells (mESCs)
44 show preferential localization to the nucleus and span a range of abundance levels (Methods,
45 **Extended Data Fig. 1**). For each locus, we looked for direct regulatory effects on local gene
46 expression by using a genetic approach based on classical *cis-trans* tests (**Fig. 1a, Note S1**).
47 Specifically, we generated clonal cell lines carrying heterozygous knockouts of the promoter
48 (~600-1,000 bp deletions) (**Fig. 1b**) and compared the expression of nearby genes within 1
49 megabase on the *cis* and *trans* alleles (*i.e.*, on the modified and unmodified homologous
50 chromosomes in the same cells) (**Note S2**). Changes in neighboring gene expression that involve
51 only the *cis* allele likely result from direct, local functions of the lncRNA locus, while changes
52 that involve both the *cis* and *trans* alleles likely result as indirect, downstream consequences of
53 the lncRNA acting elsewhere (**Note S1**). We performed genetic modifications in 129/Castaneus
54 F1 hybrid mESCs that contain a polymorphic site every ~140 basepairs (bp), enabling us to
55 distinguish the two alleles using RNA sequencing (**Extended Data Fig. 2, Note S3**).

56 At 5 of these 12 lncRNA loci, promoter knockouts significantly affected the expression of a
57 nearby gene in an allele-specific manner (false discovery rate <10%), including both activating
58 and repressive effects (**Fig. 1c,d, Note S4, Extended Data Fig. 3**). For each locus, the affected
59 gene was located immediately adjacent to, and within 5-71 kb of, the knocked-out promoter (**Fig.**
60 **1c, Extended Data Fig. 4**). This indicates that a substantial fraction of lncRNA loci influence
61 the expression of a neighboring gene.

62 To test whether such effects were specific to lncRNA loci, we deleted the promoters of 6 protein-
63 coding genes (**Extended Data Fig. 1**). Surprisingly, knockouts at 4 of these loci also affected the

64 expression of a neighbor in *cis* (**Fig. 1c,d, Extended Data Fig. 5**). Thus, both noncoding and
65 coding loci can directly influence local gene expression. These regulatory connections likely
66 contribute to the observed correlations in the expression of neighboring genes, which have been
67 reported both for lncRNAs and for mRNAs^{10,11,19,20}.

68 Because in these experiments we deleted gene promoters, the mechanisms underlying such *cis*
69 effects could in principle involve (i) DNA regulatory elements in gene promoters¹³⁻¹⁶; (ii) the
70 process of transcription^{10,17,18}; or (iii) the RNA transcripts themselves⁵⁻⁹ (**Extended Data Fig.**
71 **6a**). To begin to distinguish among these possible mechanisms, we inserted early
72 polyadenylation signals (pAS), 0.5-3 kb downstream of each transcription start site (TSS), that
73 eliminated the production of most of the RNA while leaving the promoter sequence intact (**Fig.**
74 **2, Extended Data Fig. 6b,c**, see Methods). We examined 4 lncRNA loci and 2 mRNA loci
75 where promoter deletion affected the expression of a neighboring gene (see **Note S5**).

76 As one example, we describe the linc1536 locus, hereafter called Bendr (Bend4-regulating
77 Effects Not Dependent on the RNA, **Fig. 2a**). Whereas deleting the Bendr promoter reduced the
78 expression of the adjacent Bend4 gene by 57%, inserting a pAS into the first intron of Bendr
79 (~570 bp downstream of the TSS in this ~13-kb locus) had no effect on Bend4 expression
80 despite eliminating the spliced Bendr RNA (**Fig. 2b,c**). Furthermore, global run-on sequencing
81 (GRO-seq) did not detect any transcriptionally engaged polymerase upstream of the pAS
82 insertion (**Fig. 2c, Extended Data Fig. 7a**) — perhaps because the pAS prevents RNA splicing,
83 which may dramatically reduce transcriptional activity in the modified locus^{21,22}. Therefore, *cis*
84 activation of Bend4 requires neither the mature Bendr RNA transcript nor significant Bendr
85 transcription. Instead, this effect is likely mediated by DNA regulatory elements in the ~750 bp
86 knocked-out promoter-proximal region.

87 In total, at 5 of the 6 loci examined with pAS insertions (including 3 lncRNAs and 2 mRNAs),
88 DNA regulatory elements in the promoter-proximal sequences appeared to be responsible for
89 activating a neighboring gene (**Extended Data Fig. 7b**). Although the promoters in these loci
90 would not be classified as “enhancers” based on H3K4me3/H3K4me1 ratios²³, they are bound by
91 mESC transcription factors (**Extended Data Fig. 7c**) and are located in close proximity to their
92 neighboring target genes (**Fig. 1c, Extended Data Fig. 7d,e**), suggesting that these promoters
93 may affect local gene expression through mechanisms similar or identical to enhancers^{13,24,25}.

94 We also identified one locus, linc1319 (renamed Blustr: Bivalent Locus (Sfmbt2) is Up-
95 regulated by the Splicing and Transcription of an RNA), where both promoter deletions and pAS
96 insertions substantially reduced the expression of a neighboring gene, Sfmbt2, located 5 kb
97 upstream (**Fig. 3a**). To dissect the regulatory mechanism, we tested whether the activation of
98 Sfmbt2 is mediated by (i) a sequence-specific function of the Blustr transcript or (ii) the process
99 of transcription (by which we mean one or more sequence-independent functions associated with
100 transcription, such as changes in chromatin state or recruitment of co-factors). To test the first
101 possibility, we knocked out each of the 3 downstream exons and 3 introns. None of these
102 deletions impaired Sfmbt2 activation (**Fig. 3b, Note S6**), suggesting that the activation of Sfmbt2
103 does not require unique sequences or structures in the Blustr transcript itself. To test the second
104 possibility, we engineered pAS insertions at five different locations in the first exon or intron
105 (+40 bp to +15 kb downstream of the TSS) and found that increasing the length of the Blustr
106 transcribed region led to increased activation of Sfmbt2 (**Fig. 3b, Extended Data Fig. 8a,b**). We
107 note that changing the length of the transcribed region affected the total amount of engaged
108 polymerase in the Blustr locus (**Fig. 3c**). Thus, Sfmbt2 activation responds to changes in the
109 length/amount of transcriptional activity in the Blustr locus but does not appear to require
110 specific sequence elements in the mature Blustr transcript (**Note S7**).

111 Because promoter-proximal splice sites and the process of splicing can enhance transcription —
112 in some cases by as much as 100-fold^{21,22} — we tested whether the splicing of Blustr is involved
113 in Sfmbt2 activation. Upon deleting the 5' splice site of the first intron of Blustr (**Extended Data**
114 **Fig. 8c**), we observed a 94% reduction in Blustr transcription (as assayed by GRO-seq), a 92%
115 reduction in the levels of the mature Blustr transcript, and an 85% reduction in Sfmbt2
116 expression (**Fig. 3b,c, Extended Data Fig. 8a,b**), demonstrating that the first 5' splice site of
117 Blustr has a critical role in activating Blustr and Sfmbt2 transcription. In contrast, downstream
118 splice sites were dispensable: upon deleting downstream Blustr exons, splicing skipped over the
119 removed exon to the next available 3' splice site (**Extended Data Fig. 8d**) and Sfmbt2
120 expression was unaffected (**Fig. 3b**).

121 Together, these data demonstrate that the 5' splice site and the process of transcription in the
122 Blustr locus are important for its ability to regulate Sfmbt2. This indicates that the Blustr RNA is
123 in fact required for Sfmbt2 activation (splicing involves direct interactions between the

124 spliceosome and the nascent transcript), although this mechanism does not appear to depend on
125 the precise sequence of the RNA beyond the presence of initial splice signals. One possibility is
126 that the 5' splice site promotes transcriptional activity in the *Blustr* locus, which in turn recruits
127 components of the transcriptional machinery that act on the nearby *Sfmbt2* promoter (**Fig. 3d,**
128 **Note S7**). Consistent with this model, altering transcription or splicing in the *Blustr* locus led to
129 changes in chromatin state at the *Sfmbt2* promoter (including reductions in H3K4me3 and
130 spreading of H3K27me3) and reduced occupancy of engaged RNA polymerase in the paused
131 position just downstream of the *Sfmbt2* TSS (**Extended Data Fig. 8b,e,f**). Thus, changes in
132 *Blustr* transcription and splicing may affect *Sfmbt2* expression in part by altering chromatin state
133 and RNA polymerase occupancy at the *Sfmbt2* promoter (**Fig. 3d, Note S7**).

134 In summary, genetic dissection of 12 lncRNA loci and 6 mRNA loci found that 9 loci (50%)
135 regulate the expression of a neighboring gene (**Extended Data Fig. 9**). In most of these loci,
136 including *Bendr*, local effects are mediated by enhancer-like functions of DNA elements in
137 promoters. In one locus, *Blustr*, the processes of transcription and splicing also contribute to *cis*
138 regulatory functions, perhaps by increasing the local concentration of transcription-associated
139 factors. We did not identify any lncRNA loci in which local effects are mediated by sequence-
140 specific functions of the lncRNA transcript. Because there exist thousands of other loci that fit
141 our selection criteria, we expect that similar mechanisms broadly contribute to gene regulation in
142 many loci (**Note S8**).

143 The frequent ‘crosstalk’ between neighboring genes observed in our study indicates that gene
144 loci can encode multiple independent categories of functions. Category I involves functions of
145 the RNA product: mRNAs template protein synthesis, and some noncoding transcripts (*e.g.*,
146 *XIST*) act as functional lncRNAs. Category II involves the effects of transcription-related
147 processes — including mechanisms mediated by promoters, transcription, and splicing — on the
148 regulation of other nearby genes.

149 The fact that many lncRNA loci have category II functions does not necessarily mean that they
150 do not also have category I functions, and we note that our experiments do not rule out the
151 possibility that the lncRNAs dissected in this study have RNA-mediated functions other than on
152 local gene regulation. However, the prevalence of category II functions suggests a model for the
153 evolutionary origins of some lncRNAs. In loci where a promoter acts as an enhancer, RNA

154 transcripts may arise as non-functional byproducts¹⁶. In loci where co-transcriptional processes
155 have *cis* regulatory functions, the nascent transcripts might contribute through mechanisms like
156 splicing that require little RNA-sequence specificity. These possibilities are particularly
157 intriguing in light of the patterns of evolutionary conservation of lncRNA loci²⁶⁻²⁸. For example,
158 although most lncRNA transcripts expressed in mESCs are not conserved (no RNA detected in
159 syntenic loci in other mammals, see Methods), the promoters in some of these loci correspond to
160 conserved DNA sequences that have an enhancer chromatin signature in human ESCs (**Fig. 4,**
161 **Extended Data Fig. 10, Note S9**). These sequences may have conserved functional roles as *cis*
162 regulatory elements, rather than as lncRNA promoters. Thus, mechanisms associated with *cis*
163 functions by promoters, transcription, and/or RNA processing may contribute to the functions
164 and evolution of an important subset of noncoding loci in mammalian genomes (**Extended Data**
165 **Fig. 10c**).

166 Beyond the implications for lncRNAs, these *cis* regulatory connections between neighboring
167 genes occur in both protein-coding and noncoding loci and thus appear to represent a
168 fundamental property of mammalian gene regulatory networks. The properties of these *cis*
169 regulatory connections — including mechanisms for specificity and the potential for cooperative
170 dynamics of gene activation — represent key areas for future investigation.

171 **Fig. 1. Many lncRNA and mRNA loci influence the expression of neighboring genes.** (a)
 172 Knocking out a promoter (black) could affect a neighboring gene (blue) directly (local) or
 173 indirectly (downstream). (b) Knockout of the linc1536 promoter. Left: genotypes. Right: allele-
 174 specific RNA expression for 129 and Castaneus (Cast) alleles normalized to 81 control clones
 175 (+/+). Error bars: 95% confidence interval (CI) for the mean (**Table S1**). (c) Gene neighborhoods
 176 oriented so each knocked-out gene (black) is transcribed in the positive direction. Blue
 177 neighboring genes show allele-specific changes in expression. ^See Note S3. (d) Average RNA
 178 expression on promoter knockout compared to wild-type alleles in 2+ clones (**Table S1**). *: FDR
 179 < 10%. ***: FDR < 0.1%.

180

181 **Fig. 2. Enhancer-like function of the Bendr promoter.** (a) Transcriptionally engaged RNA
 182 polymerase (GRO-Seq) and H3K4me3 occupancy (ChIP-Seq). (b) p(A)+ RNA expression upon
 183 deleting the Bendr promoter or inserting a pAS on modified versus unmodified alleles. Error
 184 bars: 95% CI for the mean of 2+ clones (see Methods, Table S1). (c) Allele-specific GRO-seq
 185 signal for clones carrying the indicated modifications. Both clones are modified on the 129
 186 allele, and only reads specifically mapping that allele are shown. Y-axis: normalized read count.
 187 Bar plot quantifies signal at Bend4, including 7 additional wild-type controls not shown on left.

188

189 **Fig. 3. Transcription and splicing of Blustr activates Sfmbt2 expression.** (a) p(A)+ RNA-seq,
 190 GRO-seq, and H3K4me3 ChIP-Seq in the Blustr locus. Sfmbt2 has two alternative TSSs. (b)
 191 p(A)+ RNA expression on knocked-out alleles compared to controls (arrows). Error bars: 95%
 192 CI for the mean for 2+ clones (pAS at +15 kb has 1 clone only, **Table S1**). Sfmbt2 pAS
 193 comparisons: two-sided *t*-test $P < 0.05$ (*) or < 0.01 (**). (c) Allele-specific GRO-seq signal for
 194 clones carrying indicated modifications. Only reads mapping to the modified allele are shown
 195 (Cast for pAS +2 kb; 129 for others). (d) Model for how transcription in the Blustr locus
 196 activates Sfmbt2.

197

198 **Fig. 4. Evolutionary conservation of mESC lncRNAs and their promoters.** (a) Classification
 199 of a subset of lncRNAs expressed in mESCs (see **Note S9**, Methods). (b) 11 have promoters
 200 whose syntenic sequence corresponds to putative DNA regulatory elements (REs) marked by
 201 DNase I hypersensitivity (HS) in human ESCs. (c) Example: linc1494. (d) Enhancers and
 202 lncRNA promoters are significantly enriched for corresponding to human REs (pie chart, ***: P
 203 $< 10^{-10}$, Chi-squared test versus GC-matched random regions) and show elevated sequence
 204 conservation compared to GC-matched regions (bar plot, **: $P < 0.01$, ***: $P < 0.001$, Mann-
 205 Whitney test versus ii+iii).

206

207 **References:**

- 208 1. Okazaki, Y. *et al.* Analysis of the mouse transcriptome based on functional annotation of 60,770 full-length
 209 cDNAs. *Nature* **420**, 563–573 (2002).
- 210 2. Kapranov, P. *et al.* RNA maps reveal new RNA classes and a possible function for pervasive transcription.
 211 *Science* **316**, 1484–1488 (2007).
- 212 3. Guttman, M. *et al.* Chromatin signature reveals over a thousand highly conserved large non-coding RNAs in
 213 mammals. *Nature* **458**, 223–227 (2009).
- 214 4. Carninci, P. *et al.* The transcriptional landscape of the mammalian genome. *Science* **309**, 1559–1563 (2005).
- 215 5. Lee, J. T. Lessons from X-chromosome inactivation: long ncRNA as guides and tethers to the epigenome.
 216 *Genes Dev* **23**, 1831–1842 (2009).
- 217 6. Nagano, T. *et al.* The Air noncoding RNA epigenetically silences transcription by targeting G9a to
 218 chromatin. *Science* **322**, 1717–1720 (2008).
- 219 7. Wang, K. C. *et al.* A long noncoding RNA maintains active chromatin to coordinate homeotic gene
 220 expression. *Nature* **472**, 120–124 (2011).
- 221 8. Ørom, U. A. *et al.* Long noncoding RNAs with enhancer-like function in human cells. *Cell* **143**, 46–58
 222 (2010).
- 223 9. Guil, S. & Esteller, M. Cis-acting noncoding RNAs: friends and foes. *Nat Struct Mol Biol* **19**, 1068–1075
 224 (2012).
- 225 10. Ebisuya, M., Yamamoto, T., Nakajima, M. & Nishida, E. Ripples from neighbouring transcription. *Nat. Cell*
 226 *Biol.* **10**, 1106–1113 (2008).
- 227 11. Cabili, M. N. *et al.* Integrative annotation of human large intergenic noncoding RNAs reveals global
 228 properties and specific subclasses. *Genes Dev* **25**, 1915–1927 (2011).
- 229 12. Bassett, A. R. *et al.* Considerations when investigating lncRNA function in vivo. *Elife* **3**, e03058 (2014).
- 230 13. Li, G. *et al.* Extensive promoter-centered chromatin interactions provide a topological basis for transcription
 231 regulation. *Cell* **148**, 84–98 (2012).
- 232 14. Rajagopal, N. *et al.* High-throughput mapping of regulatory DNA. *Nat Biotechnol* **34**, 167–174 (2016).
- 233 15. Yin, Y. *et al.* Opposing Roles for the lncRNA Haunt and Its Genomic Locus in Regulating HOXA Gene
 234 Activation during Embryonic Stem Cell Differentiation. *Cell Stem Cell* **16**, 504–516 (2015).
- 235 16. Paralkar, V. R. *et al.* Unlinking an lncRNA from Its Associated cis Element. *Mol Cell* **62**, 104–110 (2016).
- 236 17. Martens, J. A., Laprade, L. & Winston, F. Intergenic transcription is required to repress the *Saccharomyces*
 237 *cerevisiae* SER3 gene. *Nature* **429**, 571–574 (2004).
- 238 18. Shearwin, K. E., Callen, B. P. & Egan, J. B. Transcriptional interference--a crash course. *Trends Genet* **21**,
 239 339–345 (2005).
- 240 19. Purmann, A. *et al.* Genomic organization of transcriptomes in mammals: Coregulation and cofunctionality.
 241 *Genomics* **89**, 580–587 (2007).
- 242 20. Kosak, S. T. *et al.* Coordinate gene regulation during hematopoiesis is related to genomic organization.
 243 *PLoS Biol* **5**, e309 (2007).
- 244 21. Brinster, R. L., Allen, J. M., Behringer, R. R., Gelinas, R. E. & Palmiter, R. D. Introns increase
 245 transcriptional efficiency in transgenic mice. *Proc Natl Acad Sci U S A* **85**, 836–840 (1988).
- 246 22. Fong, Y. W. & Zhou, Q. Stimulatory effect of splicing factors on transcriptional elongation. *Nature* **414**,
 247 929–933 (2001).
- 248 23. Calo, E. & Wysocka, J. Modification of enhancer chromatin: what, how, and why? *Mol Cell* **49**, 825–837
 249 (2013).
- 250 24. Andersson, R., Sandelin, A. & Danko, C. G. A unified architecture of transcriptional regulatory elements.
 251 *Trends Genet* **31**, 426–433 (2015).
- 252 25. Kim, T.-K. & Shiekhattar, R. Architectural and Functional Commonalities between Enhancers and
 253 Promoters. *Cell* **162**, 948–959 (2015).
- 254 26. Necsulea, A. *et al.* The evolution of lncRNA repertoires and expression patterns in tetrapods. *Nature* **505**,
 255 635–640 (2014).
- 256 27. Hezroni, H. *et al.* Principles of long noncoding RNA evolution derived from direct comparison of
 257 transcriptomes in 17 species. *Cell Rep* **11**, 1110–1122 (2015).
- 258 28. Chen, J. *et al.* Evolutionary analysis across mammals reveals distinct classes of long non-coding RNAs.
 259 *Genome Biol* **17**, 19 (2016).
- 260
- 261

262 **Acknowledgements:** We thank S. Grossman, J. Rinn, M. Yassour, P. Sharp, L. Boyer, M. Ray,
263 C. Fulco, M. Munschauer, T. Wang, and N. Friedman for discussions; A. Goren and Broad
264 Technology Labs for ChIP; J. Lis, D. Mahat, and A. Shishkin for technical advice and reagents;
265 and J. Flannick for computational tools. J.M.E. is supported by the Fannie and John Hertz
266 Foundation and the National Defense Science and Engineering Graduate Fellowship. M.G. is
267 supported the NIH Director's Early Independence Award (DP5OD012190), the Edward
268 Mallinckrodt Foundation, the Sontag Foundation, and the Searle Scholars Program. Work in the
269 Lander Lab is supported by the Broad Institute. The Broad Institute, which E.S.L. directs, holds
270 patents and has filed patent applications on technologies related to other aspects of CRISPR.

271
272 **Author contributions:** J.M.E., M.G., and E.S.L. conceived and designed the study. J.M.E.,
273 J.E.H., G.M., M.K., and P.E.M. developed knockout protocols and performed genetic
274 manipulations. E.M.P. and J.M.E. performed all other experiments. J.M.E. developed
275 computational tools and analyzed data. J.M.E. and J.C. performed evolutionary analysis. J.M.E.
276 and E.S.L. wrote the manuscript with input from all authors. E.S.L. supervised the work and
277 obtained funding.

278

279 **Methods**

280

281 **Cell lines and cell culture.** F1 hybrid 129/Castaneus female mouse embryonic stem cells (gift
282 from Kathrin Plath) were cultured in serum-free N2B27-based medium (250 ml Neurobasal
283 media (Gibco), 250 ml DMEM/F12 (Gibco), 5 ml 100× N2 supplement (Gibco), 5 ml 50× B27
284 supplement (Gibco), 5 ml 200 mM L-Glutamine (Gibco), 3.6 µl 2-mercaptoethanol, 50 µg
285 human leukemia initiation factor (5×10^5 units, EMD Millipore), 7.4 µg Progesterone, 10 mg
286 Bovine Insulin (Sigma), 350 µl 7.5% BSA Fraction V (Gibco), supplemented with MEK
287 inhibitor PD0325901 (50 µl 10 mM, SelleckChem), and GSK3b inhibitor CHIR99021 (150 µl 10
288 mM, SelleckChem)). Prior to plating cells, tissue culture dishes were pretreated with PBS +
289 0.2% gelatin (Sigma) and 1.75 µg/ml laminin (Sigma) for 2-10 hours at 37°C. At each passage,
290 cells were trypsinized for 3-5 minutes in TVP Solution (0.025% trypsin, 1% Chicken Serum
291 (Sigma), and 1 mM EDTA in PBS pH 7.4) at room temperature. Cells tested negative for
292 mycoplasma contamination and were authenticated by comparing polymorphisms to 129S1 and
293 Castaneus genomes.

294 **Cellular fractionation.** To estimate the relative abundance of lncRNAs in different cellular
295 compartments and to characterize transcriptional activity in Blustr knockouts, we performed
296 cellular fractionation to isolate chromatin-associated, soluble nuclear, and cytoplasmic fractions
297 essentially as described²⁹. Briefly, we first lysed 5 million cells in 200 µl cold cell lysis buffer
298 (10 mM Tris-HCl pH 7.5, 0.05% IGEPAL CA-630, 150 mM NaCl), incubating on ice for 5
299 minutes. We layered the cell lysate over 2.5 volumes of chilled sucrose cushion (24% sucrose in
300 cell lysis buffer) and centrifuged at $15,000 \times g$ for 10 minutes. The supernatant from this spin
301 became the cytoplasmic fraction. After washing the pellet of nuclei with PBS (pH 7.5) + 1 mM
302 EDTA, we resuspended the pellet in 100 µl of cold glycerol buffer (20 mM Tris-HCl pH 7.5, 75
303 mM NaCl, 0.5 mM EDTA, 0.85 mM DTT, 0.125 mM PMSF, 50% glycerol) by gently flicking
304 the tube. We added 100 µl of cold nuclei lysis buffer (10 mM HEPES pH 7.5, 1 mM DTT, 7.5
305 MgCl₂, 0.2 mM EDTA, 0.3 M NaCl, 1 M urea, 1% IGEPAL CA-630), then vortexed for four
306 seconds. After 2 minutes on ice, we spun the nuclear lysate at $15,000 \times g$ for 2 minutes. This
307 supernatant was collected as the soluble nuclear (nucleoplasm) fraction. We rinsed the remaining
308 pellet (chromatin fraction) in PBS + 1 mM EDTA, then resuspended the chromatin in 300 µl

309 chromatin DNase buffer (20 mM Tris-HCl pH 7.5, 50 mM KCl, 4 mM MgCl₂, 0.5 mM CaCl₂, 2
310 mM TCEP, 0.5 mM PMSF, 0.4% sodium deoxycholate, 1% IGEPAL CA-630, 0.1% N-
311 lauroylsarcosine) plus 15 µl murine RNase inhibitor (NEB) and 30 µl TURBO DNase (Ambion).
312 The DNase digestion proceeded for 20 minutes at 37°C and was halted by adding 10 mM EDTA
313 and 5 mM EGTA. Protein was digested with proteinase K for 1 hour at 37°C. RNA was isolated
314 using Zymo RNA Concentrator-25 columns (two columns for the cytoplasmic fraction). With
315 this method, nuclear-associated endoplasmic reticulum is known to fractionate with the
316 nucleoplasm²⁹, and we observed that nucleolar RNAs fractionated with chromatin (data not
317 shown). From each cellular fraction, we sequenced total RNA and polyadenylated RNA (selected
318 using oligo d(T)₂₅ magnetic beads, NEB) using a strand-specific RNA-sequencing protocol for
319 Illumina instruments described previously³⁰.

320 **Selection criteria for knocked-out lncRNAs.** We selected lncRNA loci initially identified and
321 defined by a chromatin signature of H3K4me3 at promoters and H3K36me3 through gene
322 bodies³. We further required that lncRNAs selected for knockout analysis have TSSs, as defined
323 by capped analysis of gene expression (CAGE), located >5 kb from other genes (for epigenomic
324 annotation of each locus, see <http://pubs.broadinstitute.org/neighborng-genes/>). To prioritize
325 intergenic lncRNA loci that may regulate local gene expression, we focused on lncRNAs that
326 have subcellular localization biased toward the nucleus versus the cytoplasm (**Extended Data**
327 **Fig. 1**). We performed cellular fractionation experiments in V6.5 male mESCs as described
328 above and sequenced RNA from chromatin-associated, soluble nuclear, and cytoplasmic
329 fractions (GEO Accession GSE80262). We calculated a relative nuclear-to-cytoplasmic ratio
330 (chromatin RPKM + soluble nuclear RPKM divided by cytoplasmic RPKM) and focused on
331 lncRNAs with ratios above the median (1.5): these lncRNAs are preferentially localized to the
332 nucleus compared to other lncRNAs and mRNAs. We selected nuclear-biased lncRNAs that
333 span a range of abundance levels (**Extended Data Fig. 1**). We also included some lncRNAs that
334 are conserved across mammalian evolution (Snhg3, Snhg17, Meg3, and linc2025).

335 **Selection criteria for knocked out mRNAs.** We selected 6 mRNAs for promoter knockouts
336 based on the following criteria. We knocked out 2 mRNAs that are moderately expressed and are
337 not expected to be essential for mESC growth (Dicer1 and Crlf3). We knocked out 2 mRNAs
338 that are located adjacent to knocked-out lncRNAs (Sfmbt2 and Rcc1), in order to look for

339 reciprocal regulatory effects between the lncRNA and the affected mRNA. We knocked out 2
340 mRNAs that are located adjacent to a gene that is itself adjacent to a lncRNA (Gpr19 and
341 Slc30a9), in order to determine whether affected genes are specifically responsive to lncRNA
342 promoters or are generally responsive to other promoters in the locus. Similar to the lncRNAs
343 selected, the TSSs of these selected mRNAs are located >5 kb from other genes.

344 **CRISPR sgRNA design.** To design single-guide RNAs (sgRNAs), we built custom software to
345 calculate a specificity score (based on potential off-target sites using the algorithm described at
346 crispr.mit.edu³¹) and an efficacy score (based on a sequence model for sgRNA efficiency as
347 previously described³²) for each 20-nt targeting sequence. We removed guides with specificity
348 scores <20 or efficacy scores >0.7. To avoid T-rich sequences that result in premature
349 termination of Pol III-mediated sgRNA transcription, we removed guides with more than 1 “T”
350 in the 4 bases closest to the seed region, guides with more than 3 consecutive T’s, and guides
351 with more than 8 T’s total. We removed guides with homopolymer stretches of 5 or more bases
352 and guides with GC content <20% or >90%. We removed guides that overlapped a known
353 129/Castaneus SNP³³. Within a given region, we typically chose the three remaining guides with
354 the highest specificity scores. The sequences of all sgRNAs used in this study are listed in **Table**
355 **S2**.

356 **Promoter deletion guide placement.** To knock out a lncRNA or mRNA promoter, we chose 2-3
357 sgRNAs located in windows 300-500 bp upstream and downstream of the TSS, leading to
358 deletions of approximately 600-1000 bp surrounding the TSS. We adjusted the precise deletion
359 boundaries outward if we could not successfully design guides in these regions (*e.g.*, because
360 they were located in repetitive sequences). We note that we often found that the “wild-type”
361 alleles in heterozygous knockouts were affected by scars from repair of sgRNA double-stranded
362 breaks. Accordingly, we adjusted the bounds if necessary to cut outside of the exons of the
363 mRNA or lncRNA and thus avoid damaging the exonic sequences on the “wild-type” alleles in
364 heterozygous knockouts. We note that the presence of these scars (and their lack of allele-
365 specific effects on the expression of neighboring genes) indicate that the *cis* effects observed
366 upon deleting promoters are not merely a result of CRISPR-mediated cutting and subsequent
367 DNA repair.

368 **Genetic deletions with CRISPR/Cas9.** To delete specific sequences, we co-transfected 100 ng
369 of Cas9-expressing plasmids (“PX330-NoGuide”), 300 ng of a pool of sgRNA-expressing
370 plasmids (“pZB-Sg3”), and 100 ng of a plasmid expressing EGFP and a puromycin selectable
371 marker from a CAG promoter (pS-pp7-GFPiP). To create PX330-NoGuide, we modified PX330
372 (gift from Feng Zhang, Addgene plasmid #44230³⁴) to remove the sgRNA expression cassette.
373 To generate pZB-Sg3, we cloned a human U6 promoter and optimized sgRNA scaffold
374 sequence³⁵ into a minimal vector with an ampicillin-selectable marker and a ColE1 replication
375 origin. We transfected batches of 250,000 mouse embryonic stem cells using the Neon
376 Transfection System (Invitrogen), using 1 pulse of 40 milliseconds at 1200 V and plated two
377 batches of cells (500,000 total) into a 96-well plate in 200 μ l media. As an internal control for
378 each set of transfections, we performed a transfection using 4 guides with no predicted target
379 sites in the mouse genome.

380 We verified efficient transfection by examining GFP expression after 24 hours. To select for
381 transfected cells, we replaced the media 24 hours after transfection with 200 μ l 2i + 1 μ g/ml
382 puromycin. One day later, we split the cells into a 10-cm plate with 8 ml of 0.5 μ g/ml
383 puromycin. One day later, we replaced the media with 10 ml of 2i with no puromycin. We
384 allowed cells to grow for 7-8 days, replacing the media every 2-3 days. We hand-picked 88
385 individual colonies and 8 control colonies for each transfection in 5 μ l media, added 20 μ l of
386 TVP for ~10-20 minutes at 37°C to dissociate the colonies, and then split the colonies into two
387 identical plates. We grew the cells in these plates for 4-5 days. We harvested one of the plates for
388 DNA and RNA extraction by removing most of the media and adding 3.5 \times volume Buffer RLT
389 (Qiagen) and froze the other plate for later recovery in Freezing Media (2i media + 10% fetal
390 bovine serum + 10% DMSO).

391 **Genotyping by PCR and sequencing.** To genotype each promoter knockout, we extracted
392 genomic DNA and performed PCR using primers spanning the deleted sequence. We genotyped
393 each clone by running the PCR products on agarose gels and comparing PCR amplicon sizes to
394 predicted wild-type and deletion band sizes. We confirmed the sequences of wild-type and
395 deletion bands by Sanger sequencing or high-throughput sequencing through barcoded amplicon
396 sequencing on an Illumina MiSeq (see **Table S2**). Where possible, we used known polymorphic
397 sites from 129S1 and Castaneus genomes³³ to determine the haplotype-resolved genotype of each

398 clone. Based on the genotyping data, we nominated clones for RNA sequencing. We eliminated
399 clones showing evidence of (i) polyclonal or subclonal mutations or (ii) complex mutations such
400 as inversion or duplication of the genomic sequence between the sgRNAs. The sequences of all
401 genotyping primers are listed in **Table S2**.

402 **RNA sequencing libraries.** We generated RNA sequencing libraries as previously
403 described^{30,36}, with some modifications for high sample throughput. We isolated RNA from
404 harvested mESCs using RNeasy 96 columns. We enriched for poly(A)+ RNA using oligo d(T)₂₅
405 magnetic beads (NEB) and eluted in 18 μ L H₂O. We fragmented RNA to an average of ~150-nt
406 by adding 2 μ L Ambion Fragmentation Buffer and incubating at 70°C for exactly 2.5 minutes.
407 After transferring quickly to ice, we added 40 μ L of a master mix containing 12 μ L 5 \times FNK
408 Buffer (50 mM Tris-HCl pH 7.5, 5 mM MgCl₂, 0.6 mM CaCl₂, 50 mM KCl, 10 mM DTT,
409 0.01% Triton X-100), 1 μ L Murine RNase Inhibitor (NEB), 3 μ L FastAP Thermosensitive
410 Alkaline Phosphatase (Thermo Scientific), 3 μ L T4 Polynucleotide Kinase (NEB), and 1 μ L
411 TURBO DNase (Life Technologies). We incubated this reaction for 37°C for 30 minutes, then
412 cleaned the reaction with MyOne SILANE magnetic beads³⁷ and eluted in 6 μ L of H₂O.

413 We proceeded with the library preparation as previously described³⁰, with one additional
414 modification. To simplify the library preparation for many samples, we added unique sample
415 barcodes (8 nt) during the first adapter ligation³⁶. We used 12 pools each with 4 barcodes in
416 order to mitigate differences in the efficiency of ligation for different adapter sequences.
417 Following the first adapter ligation, we pooled 12 samples together, including up to 9 clones
418 corresponding to a single target gene as well as 3 control clones, during the first 70% ethanol
419 wash of the SILANE-bead purification. We performed an extra SILANE purification using the
420 same beads to remove excess adapter and then proceeded with reverse transcription.

421 **Hybrid selection of RNA sequencing libraries.** To measure allele-specific expression for
422 hundreds of genes in a cost-effective manner, we developed a hybrid selection strategy to enrich
423 for allele-informative reads at target genes (**Extended Data Fig. 2**). We designed oligo pools to
424 capture allele-informative sequences in the ~1600 RNAs located in the genome within 1 Mb of
425 one of the knockout targets. These target RNAs were divided into two independent pools:
426 #140820 and #141203. We used RefSeq RNA annotations for mRNAs and our custom
427 annotations for most lncRNAs. We identified SNPs that would distinguish the 129S1 and

428 Castaneus genomes³³. We designed 120-bp capture oligos in the vicinity of each 129/Castaneus
429 polymorphic site, tiling every 15 bp across either 600 bp (pool #140820) or 240 bp (pool
430 #141203) centered on the SNP. We included probes targeting both alleles to minimize
431 differences in capture efficiency between the two alleles. We filtered capture probe sequences as
432 previously described³⁷. We included up to 10 oligos per targeted RNA, duplicating probes where
433 necessary to include the sequences corresponding to each allele. Empirically, this probe design
434 strategy in combination with the protocol described below enabled assessing allele-specific
435 expression for 84% (611 of 731) of the targeted expressed genes in mESCs (RPKM \geq 2) at a
436 sequencing depth of <5 million reads per sample. Target genes and oligos sequences for these
437 pools are listed in **Table S3**.

438 We synthesized pools of 12,000 capture oligos using CustomArray technology. Oligos in each
439 pool were flanked by unique primers (Left primer sequence: CTCCTACGAGCAGTTTGCC;
440 Right primer sequence: AGTTTACGCATTACGGGCAC). After one round of PCR to add a T7
441 promoter (GGATTCTAATACGACTCACTATAGGG), we generated biotinylated RNA probes
442 as described previously³⁸, adding in 20% Biotin-16-UTP (Roche) and 20% Biotin-14-CTP (Life
443 Technologies) to the *in vitro* transcription reactions. We generated RNA probes targeting both
444 strands by incorporating the T7 promoter into either side of the PCR product and performing two
445 separate *in vitro* transcription reactions per oligo pool.

446 To capture the allele-informative regions, we pooled the final, barcoded RNA sequencing
447 libraries from all samples in the batch and performed a modified version of solution hybrid
448 selection³⁹. We first combined 500 ng dsDNA library pool with 1 nmol of Illumina P5 and P7
449 primer mix in 21 μ l total. We denatured this mix at 94°C for 10 minutes and transferred
450 immediately to ice. We added 7.5 μ l 20 \times SSPE, 0.5 μ l Murine RNase Inhibitor (NEB), and 1 μ l
451 of 500 ng/ μ l biotinylated RNA probe, for a total volume of 30 μ l. We set up at least two
452 reactions per 10 libraries, including at least one reaction with each strand of probes. We
453 incubated the hybridization reaction at 65°C for 24-48 hours. For each capture sample, we
454 washed 30 μ l Streptavidin C1 MyOne magnetic beads (Invitrogen) in 5 \times SSPE and aliquoted
455 them into PCR tubes. After removing the wash from the beads, we added the hybridization
456 reaction and mixed to resuspend the beads. We captured the biotinylated probes by shaking at
457 65°C for 20 minutes. We washed the beads twice in 150 μ l Low Stringency Wash Buffer (1 \times

458 SSPE, 0.1% SDS, 1% NP-40, 4 M urea) at 62°C for 3-4 minutes, and twice in 150 µl High
459 Stringency Wash Buffer (0.1× SSPE, 0.1% SDS, 1% NP-40, 4 M urea). To elute, we removed
460 the final wash and resuspended beads in 10 µl 100 mM NaOH and heated to 70°C for 10
461 minutes. To complete the elution, we added 1 µl 1 M acetic acid and 14 µl NLS Elution Buffer
462 (20 mM Tris-HCl pH 7.5, 10 mM EDTA, 2% N-lauroylsarcosine, 2.5 mM TCEP) and heated to
463 94°C for 4 minutes. While hot, we placed samples on magnet, removed eluate, and then placed
464 the eluate on ice for at least 30 seconds. We cleaned the eluates with 20 µl MyOne SILANE
465 magnetic beads as described³⁷, using 75 µl RLT and 61 µl 100% ethanol for the initial
466 precipitation. We eluted in 23 µl H₂O, and used this as input for a 50 µl NEBNext High Fidelity
467 PCR reaction using 500 pmol each P5 and P7 Illumina primers (98°C for 30 s; 13 cycles of 98°C
468 for 15 s, 68°C for 30 s, 72°C for 30s; 72°C for 2 minutes, 4°C hold). We cleaned the PCR
469 reaction twice with 1× volume Agencourt Ampure XP magnetic beads and eluted in 20 µl H₂O.

470 **Allele-specific gene expression measurements from RNA sequencing.** We sequenced RNA
471 libraries on an Illumina HiSeq 2500 (Read 1: 38 cycles; Read 2: 30 cycles; Index: 8 cycles). The
472 first read includes the 8-nt barcode added during the first adapter ligation (see above). Following
473 processing to separate samples based on the inline barcodes, we filtered out sequencing reads
474 that aligned to highly abundant RNA transcripts, including ribosomal RNAs, snRNAs, and
475 repetitive elements, as defined by RefSeq and RepeatMasker. A FASTA file containing these
476 sequences is available at the Gene Expression Omnibus (GSE55914).

477 We developed a computational pipeline to estimate allele-specific expression from RNA-
478 sequencing data. We created two separate reference files for the 129S1 and Castaneus
479 haplotypes, starting with the mm9 genome build and layering on SNPs based on whole-genome
480 sequencing of each of the two mouse strains³³. We aligned RNA-sequencing data separately to
481 each of the two haplotypes using Tophat (version 2.0.8). We combined the results of the two
482 alignments using PySuspenders⁴⁰, which identifies reads that map specifically to one or the other
483 allele and splits them into separate BAM files. We discarded duplicate reads and reads with
484 MAPQ < 30. After generating separate BAM files containing the reads mapping to each allele,
485 we counted reads that mapped to each RefSeq transcript (including both spliced and unspliced
486 isoforms) using Scripture⁴¹ and calculated “allelic expression ratios” for each gene (counts from
487 129 allele divided by total counts from both 129 and Castaneus alleles). The distribution of

488 allelic expression ratios for all active genes in mESCs was centered on 0.5, indicating that on
489 average each gene is expressed equally from the 129 and Castaneus alleles (**Extended Data Fig.**
490 **2b**). This indicates that there is not systematic bias in our mapping procedure toward one allele
491 or the other.

492 **RNA-seq data analysis.** We processed RNA-sequencing datasets in batches corresponding to
493 sets of libraries made on the same day with the same hybrid selection probe pool. We removed
494 samples with fewer than 100,000 non-repetitive, unique, allele-informative reads. For within-
495 batch quality control, we performed hierarchical clustering on all samples by their allelic
496 expression ratios and removed the 2-5% of outlier samples, which were largely comprised of
497 clones that showed monoallelic expression from the X chromosome.

498 **Assessment of gene knockout by expression analysis.** The PCR genotyping procedure
499 described above provided putative genotypes for the cell clones. We confirmed the genotype of
500 cells by analyzing the allele-specific expression of the knocked out gene in each clone. We
501 required that clones show >80% reduction of expression of the knocked out gene on the
502 appropriate allele in order to include the clone in downstream analysis. Incomplete reduction of
503 expression in some cases appeared to result from use of alternative TSSs that were not included
504 in the deleted sequence. In other cases, incomplete reduction of expression appeared to result
505 from subclonal genetic mosaicism within the cell line, which likely resulted from deletions that
506 occurred after several cell divisions, leading to genetic differences between individual cells in a
507 colony. For further analysis, we focused on gene loci where we obtained at least 2 heterozygous
508 knockout clones.

509 **Identifying significant changes in allele-specific expression.** In developing a statistical
510 approach to identify local, *cis* effects of these genetic manipulations, we sought to distinguish
511 local effects of the genetic deletion from downstream effects that result as a consequence of
512 either lncRNA/mRNA functions elsewhere in the cell, off-target effects, or biological/technical
513 variation between clonal cell lines (**Note S1**). Our power to detect these effects varies between
514 different measured genes (due to their level of expression and availability of SNPs) and between
515 different knockout targets (due to differences in the numbers of knockout clones analyzed).

516 To account for these two variables, we developed a statistical approach to empirically estimate
517 the false discovery rate of allele-specific changes in the expression neighboring genes using
518 hundreds of genes on other chromosomes as controls. For each gene in the neighborhood of one
519 of our promoter deletions, we calculated three statistics: (i) a T-test statistic comparing the
520 average change in expression for each of the knockout alleles (including both heterozygous and
521 homozygous knockout clones), normalized to the expression of the gene on the wild-type allele
522 of the heterozygous clones; (ii) a z -score statistic comparing the expression of the knockout allele
523 in heterozygous clones to the expression of the wild-type allele in the same clone; and (iii) a T-
524 test statistic comparing the heterozygotes to the wild-type control clones using the allelic
525 expression ratio after applying a variance-stabilizing transformation (arcsin of the square root of
526 the allelic expression ratio). For a given gene, only samples with at least 20 allele-informative
527 reads were considered, in order to enable accurate estimates of allele-specific expression. These
528 three tests differ in whether they incorporate information from homozygous clones and how they
529 normalize between knockout and wild-type alleles. We required that a gene perform significantly
530 in each of the three tests in order to regard the gene as significant, as described below. We note
531 that each underlying measure was approximately normally distributed, with some apparent
532 outliers across hundreds of control clones; we conservatively included these outliers in
533 calculating each test statistic. We examined differences in variation between knockout and
534 control alleles with Levene's test. For estimates of the variance of distributions presented in
535 figures, see **Table S1**.

536 Because the distributions are only approximately normal, we assessed the significance of each of
537 these gene-level statistics by permutation, sampling other cell lines from the same experimental
538 batch and randomly assigning them as heterozygous or homozygous knockout clones to match
539 the distribution of genotypes of the real samples. We calculated an empirical false discovery rate
540 for the sum of these permutation ranks, testing each of the neighboring genes and using all of the
541 genes on other chromosomes as the background model. Neighboring genes with $FDR < 10\%$, a
542 transformed allelic expression ratio >0.03 , and an effect size of $>10\%$ in heterozygotes were
543 considered significant.

544 **Transcriptional read-through for Meg3 and Snhg3.** Promoter knockouts of Meg3 and Snhg3
545 led to reductions in one or more downstream genes oriented in the same direction as the

546 knockout target gene. We attributed these changes to transcriptional read-through based on the
547 following evidence (**Note S4, Extended Data Fig. 3**). For both Meg3 and Snhg3, we observed
548 evidence for transcription continuing past the annotated 3' end of the knockout target, through
549 intergenic regions, and into the downstream gene (as assayed by RNA sequencing of chromatin-
550 associated RNA). For the Meg3 locus, we did not observe H3K4me3 or CAGE reads at the 5'
551 ends of Rian and Mirg (downstream of Meg3), indicating that they are not expressed from their
552 own promoters. In the Snhg3 locus, the downstream affected gene (Rcc1) is in fact expressed
553 from its own promoter, but we found evidence for reads splicing from just downstream of Snhg3
554 into the first splice acceptor of Rcc1, indicating that at least some fraction of Rcc1 transcripts
555 begin at the Snhg3 promoter.

556 **Insertion of polyadenylation signals.** To halt transcription, we initially attempted to use a short
557 49-bp synthetic polyadenylation signal (spA) sequence⁴² to minimize the amount of genomic
558 sequence added (**Extended Data Fig. 6b**). For a given gene, we designed a guide 0.5-3 kb
559 downstream of the transcription start site. We designed 200-nt ssDNA oligos including the spA
560 sequence flanked by 75- and 76-bp homologous arms, centered on the sgRNA cut site (~4 bp
561 upstream of the PAM sequence), and ordered these as ultramers from Integrated DNA
562 Technologies (**Table S2**). To knock in polyadenylation signals, we transfected 100 ng PX330-
563 NoGuide, 100 ng pZB, 100 ng pS-pp7-GFPiP, and 100-200 ng of donor ssDNA oligo and
564 followed the selection procedure described for the promoter knockouts. To genotype these
565 insertions, we used a combination of PCR and high-throughput amplicon sequencing as
566 described above. We identified clones that had heterozygous insertions of the full 49-bp spA
567 sequence on one allele; we typically observed that the other allele had a short insertion or
568 deletion, consistent with non-homologous end joining (NHEJ)-mediated repair. This short pAS
569 sequence (spA) succeeded in halting the transcription of three RNAs: Blustr (pAS at +40bp and
570 +0.5 kb in Fig. 3), Gpr19, and Bendr. However, for other genes, transcription was unaffected
571 despite pAS knock-in, consistent with the location-dependent efficiency previously observed for
572 this pAS sequence⁴².

573 Accordingly, we built a larger construct containing three polyadenylation signals (p3PA,
574 **Extended Data Fig. 6c**). The structure of this construct upon insertion into the genome through
575 homologous recombination is as follows: spA – EFS promoter – Puromycin resistance gene

576 IRES thymidine kinase – WPRE – SV40 pAS – PGK pAS (“p3PA-Puro-iTk”). We co-
577 transfected 300 ng of this construct with 100 ng of pZB and 100 ng of PX330-NoGuide, waited
578 three days, and then selected for cells with integrations with 1 µg/mL puromycin for one week.
579 We picked individual colonies and used PCR to genotype clones, using primers spanning the
580 insertion junctions. We sequenced these PCR products to determine the allele of insertion.
581 Following genotyping, we expanded clonal cell lines and transfected with PX330 and a pool four
582 sgRNAs to delete the selection cassette, leaving behind three tandem pASs. Following selection
583 with 2 µg/mL ganciclovir, we again picked individual colonies, used PCR to confirm loss of the
584 cassette, and sequenced RNA from multiple clones. PCR primer sequences for cloning homology
585 arms and genotyping p3PA insertions are listed in **Table S2**.

586 **Knockouts of Blustr exons and introns.** To delete each exon and intron of Blustr, we
587 transfected cells with pools of guides as described for the promoter deletions, using 2 guides on
588 each side. We assessed the genotype of clonal cell lines as described above for promoter
589 deletions. To confirm exon knockout from RNA sequencing data, we examined SNPs in each of
590 the exons. Upon knockout of exon 2, for example, we observed loss of RNA sequencing reads
591 mapping to exon 2, while reads mapping to other exons were still present. We also identified
592 reads spanning a new splice junction between exon 1 and exon 3, further confirming that exon 2
593 was removed from the mature transcript. For barplots in Fig. 3 measuring Blustr expression, the
594 values represent the normalized read counts of the remaining exons that were not deleted in that
595 experiment. To confirm intron knockout, we used PCR primers spanning the deletion junction
596 and sequenced the resulting PCR products. We note that the intron knockouts, by design, do not
597 affect the sequence of the spliced Blustr RNA.

598 **5' splice site knockout.** To knock out the 5' splice site of Blustr, we co-transfected mESCs as
599 described above, using a single sgRNA pZB plasmid and 200 ng of ssDNA oligonucleotide
600 donor for homologous recombination (**Extended Data Fig. 8c**). The oligo was ordered as an
601 ultramer from Integrated DNA Technologies (**Table S2**). We genotyped these insertions through
602 amplicon sequencing using an Illumina MiSeq (primers in **Table S2**).

603 **Transcriptional activity with GRO-Seq.** We used precision run-on sequencing (PRO-seq)⁴³, a
604 variant of global run-on sequencing⁴⁴, to map transcriptionally engaged RNA polymerase for a
605 subset of clones. Clones for PRO-seq (as well as ChIP-Seq and ATAC-Seq) were chosen from

606 among the recoverable knockout cell lines with a preference for clones with homozygous
607 knockouts or knockouts on the 129 allele only. We performed PRO-seq as previously
608 described⁴⁵, with modifications. We harvested 10 million mESCs by scraping, washing in cold
609 PBS, and spinning at $330 \times g$ for 3 minutes. The cell pellet was resuspended in 1 ml cold
610 Douncing Buffer (10 mM Tris-HCl pH 7.4, 300 mM Sucrose, 3 mM CaCl₂, 2 mM MgCl₂, 0.1%
611 (v/v) Triton X-100, and 0.5 mM DTT) per 1 million cells. The cells were incubated on ice in the
612 cold room for 5 minutes and dounced 25 times. The nuclei were pelleted at $500 \times g$ for 2
613 minutes, washed twice in 5 ml Douncing Buffer, and centrifuged at $500 \times g$ for 2 minutes. The
614 nuclei were then gently resuspended in 100 μ l of cold Storage Buffer (10 mM Tris-HCl, pH 8.0,
615 25% (v/v) glycerol, 5 mM MgAc₂, 0.1 mM EDTA, and 0.5 mM DTT), immediately flash frozen,
616 and stored at -80°C until use.

617 A 28 μ l 2 \times Nuclear Run-On (NRO) mix was prepared as follows: 1 M Tris-HCl, pH 8.0, 1M
618 MgCl₂, 2M KCl, and 0.1 M DTT. 5 μ l of 1 mM Biotin-11-CTP (Perkin Elmer), 1 μ l of 0.05 mM
619 CTP, 2.5 μ l of 2 mM ATP, 2.5 μ l of 2 mM GTP, 2.5 μ l of 2 mM UTP (Sigma Aldrich), 6.5 μ l of
620 nuclease free water, and 2 μ l of SUPERaseIn (Ambion) were added to the 2 \times NRO mix and
621 mixed well prior to the addition of 50 μ l of 2% NLS. The NRO reaction mix was mixed well and
622 preheated to 37°C. 100 μ l of NRO mix was added to 100 μ l of nuclei in Storage Buffer. The
623 reaction was mixed gently by pipetting and incubated at 37°C for 3 minutes, mixing halfway
624 through. To halt the reaction 500 μ l of Trizol LS (Thermo Fisher) was added, mixed well, and
625 incubated at room temperature for 5 minutes. RNA was isolated through a chloroform extraction
626 and ethanol precipitation, and resuspended in 20 μ l of H₂O. The RNA was heat denatured at
627 65°C for 40 seconds and fragmented on ice for 10 minutes with 5 μ l of 1N NaOH. To stop the
628 reaction, 5 μ l of 1 M Acetic Acid and 20 μ l of 1 M Tris-HCl, pH 7.4 were added. To remove
629 unincorporated biotinylated nucleotides, the sample was passed through a P-30 exchange column
630 (BioRad). 1 μ l of RNase inhibitor was added to the \sim 50 μ l of RNA and the first biotin
631 enrichment was then performed.

632 Each biotin enrichment was performed as follows. To prepare the Streptavidin M280 Beads
633 (Invitrogen) for biotin enrichment, 100 μ l of beads were taken per sample and washed once in
634 0.1 N NaOH with 50 mM NaCl and twice in 100 mM NaCl. Beads were resuspended in 160 μ l

635 of Binding Buffer (10 mM Tris-HCl, pH 7.4, 300 mM NaCl, and 0.1% (v/v) Triton X-100). To
636 each sample an equal volume of Streptavidin M280 beads was added, mixed, and incubated on a
637 rotator for 20 minutes at room temperature. The beads were magnetically separated and washed
638 twice in 500 μ l of ice cold High Salt Wash Buffer (50 mM Tris-HCl, pH 7.4, 2 M NaCl, and
639 0.5% (v/v) Triton X-100), twice in 500 μ l of Binding Buffer, and once in 500 μ l of Low Salt
640 Wash Buffer (50 mM Tris-HCl, pH 7.4 and 0.1% (v/v) Triton X-100). To harvest the RNA, 300
641 μ l of Trizol (Thermo Fisher) was added to the beads, vortexed for 20 seconds, and incubated at
642 room temperature for 3 minutes. 60 μ l of chloroform was added and mixture was incubated at
643 room temperature for 3 minutes. The samples were centrifuged at $14,000 \times g$ for 5 minutes at
644 4°C . The aqueous phase was collected and transferred to a new tube; the remaining organic
645 phase was removed from the beads. The Trizol extraction was then repeated as above and the
646 two aqueous phases were combined. RNA was purified with a chloroform extraction and ethanol
647 precipitation, and resuspended in nuclease free water. RNA sequencing libraries were then
648 prepared as described above, except that SILANE clean-ups were replaced with Streptavidin-
649 biotin capture enrichments until after reverse transcription (a total of 3 enrichments).

650 We sequenced PRO-seq libraries to a depth of ~ 10 million 30-bp paired-end reads. To analyze
651 the data, we mapped and processed the RNA sequencing data as described above, including
652 aligning individually to the 129 and *Castaneus* genomes. Figures showing “Allele-specific GRO-
653 seq” depict coverage for reads that uniquely map to the specific allele indicated in the figure. To
654 assess the relative read density in the promoter-proximal region and gene body of *Sfmbt2*, we
655 counted reads in the 2 kb region downstream of the first *Sfmbt2* TSS and in the remainder of the
656 gene body⁴⁶. We calculated the pause index as the ratio of these two quantities, normalized to
657 total read count. We noticed that different PRO-seq libraries had subtle biases in the relative
658 fraction of reads aligning to the TSS versus the gene body, leading to slightly offset distributions
659 of pause indices across all genes, and so we corrected for these biases in each library by
660 normalizing TSS and gene body RPKMs to the median of the $\sim 5,000$ genes with coverage across
661 all samples.

662 **Chromatin accessibility with ATAC-Seq.** Libraries were generated as previously described⁴⁷
663 using 50,000 mESCs. We generated duplicate ATAC-Seq libraries for each clonal cell line
664 examined and sequenced each to a depth of ~ 40 million 30-bp paired end reads. We aligned

665 paired-end DNA sequencing reads using bowtie2⁴⁸ to each of the 129 and Castaneus genomes
666 with the following parameters: “--met-stderr --maxins 1000”, removed duplicate reads using
667 Picard (<http://picard.sourceforge.net>), and filtered to uniquely aligning reads using samtools
668 (MAPQ < 30, <https://github.com/samtools/samtools>). For plotting normalized read coverage at
669 the Blustr and Sfmbt2 promoters, we combined data from the two biological replicates (two
670 independent measures of the same cell line) and connected paired-end reads to generate
671 fragments. Fragment coverage was normalized by the total number of uniquely mapping reads.

672 **Chromatin immunoprecipitation.** ChIP-seq for H3K4me3 and H3K27me3 was performed
673 using monoclonal antibodies as previously described⁴⁹. Sequencing data was analyzed as for
674 ATAC-Seq described above.

675 **Validation of allele-specific RNA expression with ddPCR.** To validate our RNA-seq based
676 measurements of allele specific expression, we used a quantitative allele-specific PCR assay to
677 verify measurements for Blustr and Sfmbt2. We isolated RNA from harvested mESCs using
678 RNeasy 96 columns and performed a DNase treatment followed by reverse transcription of 500
679 ng of RNA (total reaction volume 20 µl). We performed droplet digital PCR (ddPCR) using Bio-
680 Rad Custom ddPCR Assays that involve qPCR primers flanking a polymorphic site and two
681 allele-specific fluorescent probes. For Blustr: Left primer sequence:
682 GACAAATACTCCCTTCAACA; Right primer sequence: GAACAGTTTGTCTGCTGCC; Probe
683 sequence: TAAGTGAGGTGAACTCCAAG (129 allele, FAM) or
684 AGTGAGGCGAACTTCAAG (Castaneus, HEX). For Sfmbt2: Left primer sequence:
685 TGTAAGTTTGCCTGATACTC; Right primer sequence: TCTAATGTACCTCAGCCC; Probe
686 sequence: TTTCTATGAGCAGTTCAAC (129 allele, FAM) or TCCTATGAACCGTTCAGC
687 (Castaneus, HEX). ddPCR was done with 2.2 µl of cDNA, 11 µl of Supermix (BioRad), 1.1 µl of
688 each probe, and 7.7 µl of water per reaction followed by droplet generation. PCR was performed
689 as follows: 95°C for 10 minutes; and cycling at 94°C for 30 s and 55°C for 1 minute for a total of
690 40 cycles; and 98°C for 10 minutes. Readout was done using the QX200 Droplet Reader and
691 QuantaSoft Software (BioRad) to determine the total number of droplets containing each allele.
692 We calculated allelic expression ratios from these values and compared it to values generated
693 through RNA-sequencing and hybrid selection of the same RNA samples (**Extended Data Fig.**
694 **2d,e**).

695 **External ChIP-Seq, RNA-Seq, and DNase HS data.** We utilized the following data from
696 ENCODE⁵⁰: H3K4me3, H3K4me1, H3K27ac, and CTCF ChIP-Seq in mESCs (ES-Bruce4);
697 DNase hypersensitivity sequencing in mESCs (E14); H3K4me3, H3K4me1, and CTCF ChIP-
698 Seq and DNase HS data in H1-hESCs; and RNA-sequencing data in H1-hESCs (nuclear p(A)+,
699 nuclear total). To assess transcription factor binding to mRNA and lncRNA promoters
700 (**Extended Data Fig. 7c**), we examined mESC ChIP-seq peaks available from Kagey *et al.* at the
701 Gene Expression Omnibus (GSE22562)⁵¹.

702 **DNA purification for examining proximity contacts.** To examine the proximity contacts of the
703 linc1405 locus, we used the RAP-DNA protocol, which we initially developed in order to map
704 RNA localization to chromatin, to capture linc1405 DNA³⁷. Briefly, we crosslinked live cells to
705 fix endogenous chromatin complexes, then purified a target DNA region using a pool of
706 oligonucleotides targeting the linc1405 locus (**Table S3**). Here, we used probes that are the same
707 strand as the linc1405 RNA – in this way, we specifically capture the linc1405 DNA and do not
708 directly capture the linc1405 RNA itself. We mapped the 3-D proximity contacts of the linc1405
709 locus through high-throughput sequencing of co-purified DNA and calculated the normalized
710 enrichment to an input DNA library in 1-kb windows (**Extended Data Fig. 7e**). Annotations for
711 topologically associated domains (TADs) were downloaded from the Ren Lab
712 (<http://chromosome.sdsc.edu/mouse/hi-c/download.html>)⁵².

713 **LncRNA transcript annotations.** For evolutionary conservation analysis, we used lncRNA
714 annotations and isoforms previously defined based on RNA sequencing in mouse embryonic
715 stem cells, combining annotations generated with multiple methods (Scripture⁴¹ and slncky²⁸).
716 We filtered the combined list using slncky²⁸ to eliminate transcripts predicted to encode proteins
717 or micropeptides by UCSC, transcripts that partially align to protein-coding genes (*e.g.*,
718 pseudogenes or incomplete reconstructions), and species-specific coding gene duplications.
719 Subsequently we performed several manual curation steps. We examined each isoform using a
720 combination of long-read RNA-sequencing data, total chromatin-associated RNA sequencing
721 data, capped analysis of gene expression (CAGE) data, and poly(A+) 3'-end sequencing data
722 from mESCs^{28,30,41,53}. We eliminated transcripts that appeared to result from an extended 3'UTR
723 of an upstream protein-coding transcript. Because the precise 5' ends of transcripts are
724 imprecisely assigned by based on RNA-sequencing data alone, we re-assigned 5' ends (TSSs)

725 using a sliding-window approach to find the 10-bp window with the highest number of same-
726 strand CAGE reads within 300-bp of the initial calculated TSS. We additionally manually
727 curated the TSS of each lncRNA, some of which were incorrectly assigned by more than 300 bp,
728 based on CAGE and H3K4me3 ChIP-Seq data, and eliminated any where we could not identify
729 the TSS (*e.g.*, due to unmappable sequence or very low abundance).

730 **Analysis of lncRNA and promoter conservation.** To categorize lncRNAs by their conservation
731 properties and promoter locations, we examined a set of 307 lncRNAs expressed in mESCs as
732 described above. We assessed the conservation of each lncRNA through a two-step approach.
733 We first used slncky to look in syntenic locations for evidence of lncRNA transcripts in deep
734 p(A)+ RNA-seq of rat, chimp, and human induced pluripotency stem cells (iPSCs)²⁸. LncRNAs
735 called “conserved” by this first filter have substantial evidence based on RNA-seq that allows for
736 independent reconstruction of the transcript in one or more of these other organisms. We
737 categorized the remaining lncRNAs by the location of their TSS: 71 lncRNAs originate within
738 500-bp of an mRNA TSS on the opposite strand (“divergent”); 59 lncRNAs originate within the
739 long-terminal repeats (LTRs) of endogenous retroelements; and 79 lncRNAs have their
740 promoters in intergenic regions that do not overlap with LTRs and do not emerge from a
741 bidirectional mRNA promoter (henceforth, “intergenic”).

742 Because some conserved lncRNAs might be too lowly expressed to assemble a transcript *de*
743 *novo* in a given species, we examined more closely the 79 intergenic lncRNAs that were called
744 “mouse-specific” in the initial slncky analysis. We applied a second, more stringent threshold to
745 remove lncRNAs misclassified as mouse-specific due to low abundance. For each intergenic
746 lncRNA locus, we used liftOver⁵⁴ to map the 10 bp surrounding the mouse TSS (mm9) to the
747 human genome (hg19) (minMatch=0.1, UCSC chain). 37 of these transcripts did not lift over at
748 this step, and thus were considered mouse-specific. For the 42 that did lift over, we examined the
749 syntenic region for evidence of p(A)+ RNA-seq data from human iPSCs²⁸ or p(A)+ nuclear-
750 fraction RNA-seq from hESCs (−100 to +900 bp relative to the TSS), or for evidence of p(A)+
751 nuclear-fraction or whole-cell CAGE from hESCs (−250 to +250 bp relative to the TSS), and
752 removed from consideration any lncRNAs that showed evidence for RNA-seq or CAGE above a
753 certain threshold. We chose this threshold based on a set of random intergenic regions, which
754 were matched to the set of intergenic mouse-specific lncRNAs based on GC content. We

755 eliminated from consideration the 10 lncRNAs that showed RNA-seq or CAGE signal greater
756 the 90th percentile of random regions, corresponding to approximately 2 CAGE or RNA-seq
757 reads in the windows described above. These 10 lncRNAs were added to the “conserved” section
758 of the pie chart in **Fig. 4a**. Several of these 10 lncRNAs correspond to substantially shortened,
759 single-exon p(A)+ transcripts that show minimal overlap with the syntenic exons in mouse;
760 although a majority of the exonic sequence of these transcripts are not in fact conserved between
761 human and mouse, we excluded these from consideration as putative mouse-specific lncRNAs.

762 For the purposes of examining the conservation properties of these intergenic mouse-specific
763 lncRNAs, we defined a matched set of “enhancer” elements. We first generated a list of
764 regulatory elements in mESCs using the DNase hotspots called by ENCODE-UW in ES-E14
765 cells. As an estimate of the activity of each element, we calculated the density of H3K27ac reads
766 in the region. From the set of intergenic elements that did not overlap a promoter, lncRNA
767 promoter, or LTR, we selected a random subset matched to the intergenic lncRNA promoters for
768 H3K27ac density (binned by 10 reads / bp) and distance to the TSS of the closest active gene
769 (binned by 5 kb). We call these elements “enhancers” because they are marked by DNase
770 hypersensitivity and H3K27ac but do not overlap a known gene promoter.

771 We compared the sequence conservation and functional conservation of three classes of
772 elements: intergenic mouse-specific lncRNAs, matched intergenic enhancer elements, and GC-
773 matched random intergenic elements. First, we computed the rate at which each set maps to
774 human sequence. We centered each element and used liftOver (--minMatch=0.1) to identify the
775 syntenic region in the human genome. Elements that did not lift over at this step correspond to
776 the white segment of the pie charts in **Fig. 4** (iii – “did not map”). For elements that did lift over
777 to human, we next defined the subset that map to putative regulatory elements in human. We
778 examined a 500-bp window centered on the lifted over region and counted reads in hESC
779 DNase-seq data from ENCODE. We defined regions showing DNase HS scores higher than 95%
780 of the mappable random intergenic regions as putative DNA regulatory elements. We note that
781 these random intergenic regions include some enhancers – they are matched to lncRNA
782 promoters for GC content, and thus frequently correspond to regulatory elements (which are GC-
783 rich) that happen to be active in hESCs. For both intergenic mouse-specific lncRNAs and
784 enhancers, ~33% of elements corresponded to putative DNA regulatory elements in human (**Fig.**

785 **4d**), representing a ~6.6-fold enrichment versus the random intergenic controls. To compare
786 sequence conservation of these classes of elements, we calculated the average SiPhy score⁵⁵
787 across each 500-bp region surrounding the mouse TSS or the center of the enhancer element,
788 using the 29 mammals alignment from the mouse perspective⁵⁶. We used a two-sided Mann-
789 Whitney U-test to look for changes in the distributions of SiPhy scores to the set of mappable
790 random intergenic regions (**Fig. 4d** – random ii+iii).

791 **Impact of expression level on conservation analysis.** Although the set of intergenic mESC
792 lncRNAs examined above does not show any significant evidence for p(A)+ RNA in the syntenic
793 locus in human, some of these transcripts may not be detected in human and yet still be truly
794 conserved. These transcripts might be misclassified as “mouse-specific” lncRNAs for several
795 reasons, including: (i) low expression level in hESCs and iPSCs such that the lncRNA, by
796 chance, is not detected based on the depth of sequencing data available; or (ii) the lncRNA is not
797 expressed in hESCs or iPSCs, but is expressed in a different human cell type and thus may have
798 a conserved function.

799 To estimate the false positives resulting from these and other scenarios, we examined the
800 properties of a set of 853 conserved mRNAs matched to the intergenic “mouse-specific”
801 lncRNAs based on expression in mESCs. We counted the frequency at which these mRNAs
802 would be called “not conserved” by the same procedures described above: we applied the nuclear
803 p(A)+ CAGE and RNA-seq filters to eliminate transcripts that show detectable transcription in
804 the 1-kb region near the TSS. While 87% of the intergenic lncRNAs described above passed
805 these filters (and thus appeared to be mouse-specific), only 22% of the expression-matched
806 mRNAs passed; this indicates that the set of 69 mouse-specific intergenic lncRNAs are
807 approximately 3.9-fold enriched for human elements that are not transcribed in hESCs. Thus, the
808 mouse-specific lncRNAs defined above appear to consist largely of transcripts that are not
809 conserved.

810 We performed the following additional analyses to ensure the robustness of our conclusions
811 regarding the existence of lncRNAs that evolved from ancestral regulatory elements. First, we
812 examined the conservation of the first 5' splice sites of this set of lncRNAs. In 7 of these 11 loci,
813 the “GT” dinucleotide in the first 5' splice site is not conserved, suggesting that a similar spliced
814 transcript cannot be produced from this locus. Second, we re-performed the entire conservation

815 analysis focusing on the 50% of mESC intergenic lncRNAs with the highest expression levels –
 816 these lncRNAs are less likely to be missed in hESCs due to low abundance. We also adjusted our
 817 p(A)+ RNA and CAGE filters to require a complete absence of reads in the corresponding
 818 regions in hESCs and iPSCs. Using these filters, 79% of the intergenic lncRNAs are not
 819 detectably expressed in human cells, representing a ~12-fold enrichment over mRNAs matched
 820 for expression level. Therefore we are confident that most of these lncRNAs are correctly
 821 classified as mouse-specific. Of the 30 intergenic lncRNAs called mouse-specific by this more
 822 conservative analysis, 5 do indeed correspond to putative DNA regulatory elements, including
 823 linc1494 (**Fig. 4c**), representing a >8-fold enrichment versus GC-matched random sequences
 824 (Chi-squared $P < 10^{-10}$). Thus, our conclusions that some lncRNAs appear to evolve from
 825 ancestral regulatory elements are robust even with stringent thresholds.

826 **Software for data analysis and graphical plots.** We used the following software for data
 827 analysis and graphical plots: R Bioconductor (version 3.0)⁵⁷, Gviz (version 1.10.11), gplots
 828 (version 2.17.0), GenomicRanges (version 1.18.4)⁵⁸, rtracklayer (version 1.26.3)⁵⁹, BEDTools⁶⁰,
 829 Integrative Genomics Viewer (version 2.3.26)⁶¹, and vcftools (version 0.1.12)⁶².

830 **Data availability.** Sequencing data for this study is available at the Gene Expression Omnibus
 831 (GSE80262 and GSE85798), and additional visualizations of the data are available at
 832 <http://pubs.broadinstitute.org/neighbor-genes/>.

833 **Additional references:**

- 834 29. Bhatt, D. M. *et al.* Transcript dynamics of proinflammatory genes revealed by sequence analysis of
 835 subcellular RNA fractions. *Cell* **150**, 279–290 (2012).
 836 30. Engreitz, J. M. *et al.* RNA-RNA interactions enable specific targeting of noncoding RNAs to nascent Pre-
 837 mRNAs and chromatin sites. *Cell* **159**, 188–199 (2014).
 838 31. Hsu, P. D. *et al.* DNA targeting specificity of RNA-guided Cas9 nucleases. *Nat Biotechnol* **31**, 827–832
 839 (2013).
 840 32. Wang, T., Wei, J. J., Sabatini, D. M. & Lander, E. S. Genetic screens in human cells using the CRISPR-
 841 Cas9 system. *Science* **343**, 80–84 (2014).
 842 33. Keane, T. M. *et al.* Mouse genomic variation and its effect on phenotypes and gene regulation. *Nature* **477**,
 843 289–294 (2011).
 844 34. Cong, L. *et al.* Multiplex genome engineering using CRISPR/Cas systems. *Science* **339**, 819–823 (2013).
 845 35. Chen, B. *et al.* Dynamic imaging of genomic loci in living human cells by an optimized CRISPR/Cas
 846 system. *Cell* **155**, 1479–1491 (2013).
 847 36. Shishkin, A. A. *et al.* Simultaneous generation of many RNA-seq libraries in a single reaction. *Nat Methods*
 848 **12**, 323–325 (2015).
 849 37. Engreitz, J., Lander, E. S. & Guttman, M. RNA antisense purification (RAP) for mapping RNA interactions
 850 with chromatin. *Methods Mol Biol* **1262**, 183–197 (2015).
 851 38. Engreitz, J. M. *et al.* The Xist lncRNA exploits three-dimensional genome architecture to spread across the
 852 X chromosome. *Science* **341**, 1237973 (2013).

- 853 39. Gnirke, A. *et al.* Solution hybrid selection with ultra-long oligonucleotides for massively parallel targeted
854 sequencing. *Nat Biotechnol* **27**, 182–189 (2009).
- 855 40. Huang, S., Holt, J., Kao, C.-Y., McMillan, L. & Wang, W. A novel multi-alignment pipeline for high-
856 throughput sequencing data. *Database (Oxford)* **2014**, bau057–bau057 (2014).
- 857 41. Guttman, M. *et al.* Ab initio reconstruction of cell type-specific transcriptomes in mouse reveals the
858 conserved multi-exonic structure of lincRNAs. *Nat Biotechnol* **28**, 503–510 (2010).
- 859 42. Levitt, N., Briggs, D., Gil, A. & Proudfoot, N. J. Definition of an efficient synthetic poly(A) site. *Genes Dev*
860 **3**, 1019–1025 (1989).
- 861 43. Kwak, H., Fuda, N. J., Core, L. J. & Lis, J. T. Precise maps of RNA polymerase reveal how promoters direct
862 initiation and pausing. *Science* **339**, 950–953 (2013).
- 863 44. Core, L. J., Waterfall, J. J. & Lis, J. T. Nascent RNA sequencing reveals widespread pausing and divergent
864 initiation at human promoters. *Science* **322**, 1845–1848 (2008).
- 865 45. Mahat, D. B. *et al.* Base-pair-resolution genome-wide mapping of active RNA polymerases using precision
866 nuclear run-on (PRO-seq). *Nat Protoc* **11**, 1455–1476 (2016).
- 867 46. Adelman, K. & Lis, J. T. Promoter-proximal pausing of RNA polymerase II: emerging roles in metazoans.
868 *Nat Rev Genet* **13**, 720–731 (2012).
- 869 47. Buenrostro, J. D., Giresi, P. G., Zaba, L. C., Chang, H. Y. & Greenleaf, W. J. Transposition of native
870 chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins and
871 nucleosome position. *Nat Methods* **10**, 1213–1218 (2013).
- 872 48. Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**, 357–359
873 (2012).
- 874 49. Busby, M. *et al.* *Systematic comparison of monoclonal versus polyclonal antibodies for mapping histone*
875 *modifications by ChIP-seq.* (2016). doi:10.1101/054387
- 876 50. Mouse ENCODE Consortium *et al.* An encyclopedia of mouse DNA elements (Mouse ENCODE). *Genome*
877 *Biol* **13**, 418 (2012).
- 878 51. Kagey, M. H. *et al.* Mediator and cohesin connect gene expression and chromatin architecture. *Nature* **467**,
879 430–435 (2010).
- 880 52. Dixon, J. R. *et al.* Topological domains in mammalian genomes identified by analysis of chromatin
881 interactions. *Nature* **485**, 376–380 (2012).
- 882 53. Fort, A. *et al.* Deep transcriptome profiling of mammalian stem cells supports a regulatory role for
883 retrotransposons in pluripotency maintenance. *Nat Genet* **46**, 558–566 (2014).
- 884 54. Kent, W. J. *et al.* The human genome browser at UCSC. *Genome Res* **12**, 996–1006 (2002).
- 885 55. Garber, M. *et al.* Identifying novel constrained elements by exploiting biased substitution patterns.
886 *Bioinformatics* **25**, i54–62 (2009).
- 887 56. Lindblad-Toh, K. *et al.* A high-resolution map of human evolutionary constraint using 29 mammals. *Nature*
888 **478**, 476–482 (2011).
- 889 57. Gentleman, R. C. *et al.* Bioconductor: open software development for computational biology and
890 bioinformatics. *Genome Biol* **5**, R80 (2004).
- 891 58. Lawrence, M. *et al.* Software for computing and annotating genomic ranges. *PLoS Comput Biol* **9**, e1003118
892 (2013).
- 893 59. Lawrence, M., Gentleman, R. & Carey, V. rtracklayer: an R package for interfacing with genome browsers.
894 *Bioinformatics* **25**, 1841–1842 (2009).
- 895 60. Quinlan, A. R. & Hall, I. M. BEDTools: a flexible suite of utilities for comparing genomic features.
896 *Bioinformatics* **26**, 841–842 (2010).
- 897 61. Robinson, J. T. *et al.* Integrative genomics viewer. *Nat Biotechnol* **29**, 24–26 (2011).
- 898 62. Danecek, P. *et al.* The variant call format and VCFtools. *Bioinformatics* **27**, 2156–2158 (2011).
- 899
900
901

902

903

904 **Extended Data Fig. 1. Expression and subcellular localization of knocked-out lncRNAs**
 905 **and mRNAs.** (a) Expression of lncRNAs and mRNAs in F1 129/Castaneus female mESCs,
 906 reported in fragments per kilobase per million (FPKM) in whole-cell p(A)+ RNA-seq.
 907 Cumulative fraction is plotted for all mRNAs expressed in mESCs. Large dots represent
 908 transcripts whose promoters we deleted in this study. LncRNAs and mRNAs span a >20-fold
 909 range of abundance levels. (b) Relative subcellular localization of lncRNAs and mRNAs. We
 910 sequenced p(A)+ RNA from chromatin, soluble nuclear, and cytoplasmic fractions (see
 911 Methods) and plotted the relative abundance of mature transcripts in each fraction. We selected
 912 lncRNAs that showed localization biased toward the nuclear fractions relative to most mRNAs.
 913 For comparison, we plotted 1,000 randomly selected mRNAs (light gray).

914
 915 **Extended Data Fig. 2. Generation of knockout clones and measurement of allele-specific**
 916 **RNA expression.** (a) Overview of knockout and measurement protocol. (b) Distribution of
 917 allelic expression ratios (number of informative reads mapping to 129S1 allele divided by the
 918 number mapping to either the 129S1 or the Castaneus allele) across active genes in mESCs. (c)
 919 Scatterplot of allelic expression ratios for genes with RPKM ≥ 2 that have more than 100 allele-
 920 informative reads across all libraries. Allelic expression ratios are consistent in RNA sequencing
 921 data before and after hybrid selection (HS). (d) Allelic expression ratios as measured by two
 922 independent methods for Blustr and (e) Sfmbt2 expression in 15 clonal cell lines containing
 923 genetic modifications in the Blustr locus. (f) Example locus showing hybrid selection strategy
 924 and RNA-seq coverage for cell lines with the indicated genotype for deletion of the Bendr
 925 promoter. Y-axis scales represent normalized read counts and are the same for all hybrid
 926 selection tracks. The absolute level of expression for any given gene varies among clonal cell
 927 lines; throughout this work, we instead consider the *relative* level of expression between the two
 928 alleles in heterozygous knockout cells. For similar plots of each gene studied, see
 929 <http://pubs.broadinstitute.org/neighboring-genes/>.

930
 931 **Extended Data Fig. 3. Read-through transcription at Meg3 and Snhg3 loci.** (a) Snhg3
 932 promoter knockout reduces the levels of Rcc1 mRNA by 23%. However, sequencing of
 933 chromatin-associated RNA shows that transcription continues past the annotated 3' end of Snhg3
 934 into the downstream Rcc1 gene (see Methods). This read-through transcription creates a fusion
 935 transcript containing exons of both Snhg3 and Rcc1, as well as intergenic RNA. We note that
 936 this fusion transcript is also annotated in the syntenic human locus as an alternative isoform of
 937 RCC1. Bars: relative p(A)+ RNA expression on modified versus unmodified alleles. Error bars:
 938 95% CI for the mean of 2+ clones (Table S1). (b) Meg3 promoter knockout eliminates the
 939 expression not only of Meg3 but also of two additional lncRNAs encoded downstream in a
 940 tandem orientation (Rian and Mirg). Although these three lncRNAs are annotated as separate
 941 genes, they appear to be derived from a single transcript driven by the Meg3 promoter. This is
 942 consistent with the presence of continuous chromatin-associated RNA throughout the locus and a
 943 lack of CAGE reads at the 5' ends of Rian and Mirg3.

944 **Extended Data Fig. 4. Promoter knockouts for 5 intergenic lncRNAs affect the expression**
 945 **of a neighboring gene.** Significance (z -score) of allele-specific expression ratios at all genes
 946 within 1 Mb of each of 5 lncRNA loci. Each dot represents a different heterozygous promoter
 947 knockout clone for a given gene. Dots are shown only for genes that are sufficiently highly

948 expressed to assess allele-specific expression (see Methods). The y-axis is capped at -10 to +10
 949 standard deviations from the mean. Black: knocked-out lncRNA. Blue: Gene with significant
 950 allele-specific change in gene expression (FDR < 10%). Independent clones are not expected to
 951 yield the same significance value (z -score), in part because read depth differs between samples.
 952
 953

954 **Extended Data Fig. 5. Promoter knockouts for 4 mRNAs affect the expression of a**
 955 **neighboring gene.** Significance (z -score) of allele-specific expression ratios at all genes within 1
 956 Mb of each of 4 mRNA loci. Each dot represents a different heterozygous promoter knockout
 957 clone for a given gene. Dots are shown only for genes that are sufficiently highly expressed to
 958 assess allele-specific expression (see Methods). The y-axis is capped at -10 to +10 standard
 959 deviations from the mean. Black: knocked-out lncRNA. Blue: Gene with significant allele-
 960 specific change in gene expression (FDR < 10%). Independent clones are not expected to yield
 961 the same significance value (z -score), in part because read depth differs between samples.
 962
 963

964 **Extended Data Fig. 6. Dissecting mechanisms for how gene loci regulate a neighbor. (a)**
 965 Three categories of possible mechanisms by which a gene locus might regulate the expression of
 966 a neighbor. **(b)** We used two strategies to insert pAS downstream of gene promoters. In the first
 967 strategy, we inserted a 49-bp synthetic pAS (“spA”) using a single-stranded DNA oligo with 75-
 968 bp homology arms (see Methods). **(c)** In the second pAS insertion strategy, we cloned a donor
 969 plasmid containing a selection cassette and three different pAS sequences (see Methods).
 970 Homology arms of 300-800 bp were used to integrate the cassette. After isolating clones with
 971 successful insertions, we used a second round of transfections to remove the selection cassette,
 972 leaving behind three tandem pASs. EFS = elongation factor 1 promoter. Puro = puromycin
 973 resistance gene (*pac*). HSV-tk = herpes simplex virus thymidine kinase.
 974

975 **Extended Data Fig. 7. Promoters of lncRNAs and mRNAs have enhancer-like functions.**
 976 **(a)** Allele-specific GRO-seq signal for clones with the indicated modifications at the Bendr
 977 locus. Only reads specifically mapping to one of the two alleles are shown. Y-axis scale
 978 represents normalized read count and is the same for all tracks. **(b)** Allele-specific p(A)+ RNA
 979 expression for genetic modifications at the linc1405, Snhg17, Gpr19, and Slc30a9 loci. Bars:
 980 Average RNA expression on modified compared to unmodified (wild-type) alleles. Error bars:
 981 95% CI for the mean of 2+ clones (Table S1). Gray arrows indicates distance from the targeted
 982 locus promoter to the affected neighboring gene. We note that, based on their location, the
 983 Snhg17 and Gpr19 pAS insertions likely allow more substantial splicing and transcription; for
 984 these loci, it is clear that the majority of the transcript is dispensable but it is possible that
 985 transcription close to the promoter may be involved in the *cis* regulatory function. **(c)** Presence
 986 (gray) or absence (white) of various chromatin marks and transcription factors in mESCs in a
 987 1.5-kb window centered on the TSS of each targeted gene. **(d)** Distance from each knocked-out
 988 gene to its neighboring target gene (x -axis) versus the magnitude of the effect on the expression
 989 of the neighboring gene (% compared to wild-type, y -axis). Blue genes represent those discussed
 990 in main text; gray genes are discussed in Note S5. **(e)** Proximity-based contacts between the
 991 linc1405 and Eomes loci (the pair of loci separated by the greatest linear distance). The y -axis
 992 shows enrichment in a sequencing-based proximity assay in which we used antisense oligos to
 993 capture linc1405 DNA and any interacting, crosslinked proximal DNA (see Methods). TAD

994 annotations are derived from Hi-C experiments in mESCs (see Methods). Blue arrow: focal
995 contact between the linc1405 and Eomes loci.

996

997 **Extended Data Fig. 8. Characterization of genetic modifications in the Blustr locus. (a)**

998 Allele-specific GRO-seq signal for clones with the indicated modifications at the Blustr locus.

999 Only reads specifically mapping to one of the two alleles are shown. Y-axis scale represents

1000 normalized read count and is the same for all tracks, and is magnified 5 times at the indicated

1001 location to better visualize the reads in the Sfm2t2 locus. **(b)** Quantification of allele-specific

1002 GRO-seq signal in the Sfm2t2 locus on alleles modified as indicated. TSS: region including the

1003 two alternative TSSs of Sfm2t2 and 2 kb downstream. Gene body: region containing the

1004 remainder of the Sfm2t2 gene locus. Pause index: ratio of TSS to gene body. Dashed gray lines

1005 indicate the 95% CI for the mean of 8 wild-type clones. **(c)** Schematic of the 5' end of the Blustr

1006 locus and genotypes of two knockout clones. The 5' splice site is located 78 bp downstream of

1007 the Blustr transcription start site (in this panel, Blustr is transcribed from left to right). One of the

1008 alleles from the two clones contains insertion of the oligo mediated by homologous

1009 recombination; the remaining three alleles contain insertions or deletions resulting from non-

1010 homologous end joining repair of sgRNA-mediated double-strand breaks, some of which also

1011 disrupt the 5' splice site. Barplots show allele-specific RNA expression for knockout clones and

1012 control clones (+/+). **(d)** Schematic of the observed splice structures of Blustr RNA transcripts in

1013 p(A)+ RNA sequencing of the exon deletion clones. Each deletion removes a region including

1014 ~50-200 bp on either side of the exon, thereby removing both the exon and its splice sites. The

1015 Exon 4 deletion removes the endogenous pAS, leading to new isoforms of the lncRNA transcript

1016 that splice into two cryptic splice acceptors downstream. **(e)** GRO-Seq, H3K4me3 ChIP-Seq, and

1017 chromatin accessibility (ATAC-Seq FPKM) at the Blustr and Sfm2t2 promoters in cell lines with

1018 the indicated genotypes. Deletion of the first 5' splice site leads to a significant reduction in

1019 H3K4me3, RNA polymerase occupancy, and chromatin accessibility at the Blustr promoter, as

1020 well as H3K4me3 and RNA polymerase occupancy (but not accessibility) at the Sfm2t2

1021 promoter. **(f)** H3K27me3 ChIP-seq at the Blustr and Sfm2t2 loci in cell lines with the indicated

1022 genotypes. Deletion of the Blustr promoter or 5' splice site leads to spreading of the repression-

1023 associated H3K27me3 modification across a ~30 kb region.

1024

1025 **Extended Data Fig. 9. Mechanisms for crosstalk between neighboring lncRNAs and**

1026 **mRNAs.** Proposed mechanisms based on pAS insertion experiments and other genetic

1027 manipulations (see text). †For proposed mechanisms, see Note S5.

1028

1029 **Extended Data Fig. 10. Classification of lncRNAs based on conservation and promoter**

1030 **location. (a)** Classification of 307 lncRNAs expressed in mESCs. "Conserved" transcripts are

1031 those that show significant evidence of capped analysis of gene expression (CAGE) data and/or

1032 p(A)+ RNA in syntenic loci (see Methods). Divergent: initiating within 500 bp of an mRNA

1033 TSS, on the opposite strand. ERV: endogenous retroviral repetitive element (see Note S9).

1034 Boxplot shows sequence-level conservation of the promoters of subsets of lncRNAs expressed in

1035 mESCs. Random intergenic regions are matched to lncRNA promoters by GC content. Positive

1036 SiPhy score indicates evolutionary constraint on functional sequences. Orange category

1037 corresponds to mouse-specific lncRNAs that appear to have evolved from ancestral regulatory

1038 elements (REs) and correspond to sequences that show evidence for DNase I hypersensitivity in

1039 human embryonic stem cells. Significance is calculated compared to random intergenic regions

1040 using a Mann-Whitney U-test. ***: $P < 0.001$. Whiskers represent data within $1.5\times$ the
1041 interquartile range of the box. **(b)** Chromatin and RNA data for 11 mouse-specific lncRNAs that
1042 appear to have evolved from ancestral regulatory elements. In mouse, these elements show
1043 evidence for CAGE, H3K4me3, and DNase I hypersensitivity, consistent with their roles as
1044 promoters. The syntenic sequences in human do not show evidence for CAGE but nonetheless
1045 are DNase I hypersensitive and are frequently marked by H3K4me1 and/or CTCF. **(c)** Model for
1046 evolution of lncRNAs from pre-existing enhancers, which often initiate weak bidirectional
1047 transcription to produce eRNA. Spliced transcripts may neutrally appear through the appearance
1048 of splice signals and loss of polyadenylation signals. In some cases, transcription, splicing, or
1049 other RNA processing mechanisms may feed back and contribute to the *cis* regulatory function
1050 of the promoter, producing a lncRNA as a byproduct.







