

Conserved noncoding transcription and core promoter regulatory code in early *Drosophila* development

Philippe J. Batut^{1,2} & Thomas R. Gingeras^{1,*}

¹ Watson School of Biological Sciences, Cold Spring Harbor Laboratory, NY, USA

² Present address: Lewis-Sigler Institute, Princeton University, Princeton, NJ, USA

* Corresponding author: gingeras@cshl.edu

Multicellular development is largely determined by transcriptional regulatory programs that orchestrate the expression of thousands of protein-coding and noncoding genes. To decipher the genomic regulatory code that specifies these programs, and to investigate globally the developmental relevance of noncoding transcription, we profiled genome-wide promoter activity throughout embryonic development in 5 *Drosophila* species. We show that core promoters, generally not thought to play a significant regulatory role, in fact impart broad restrictions on the developmental timing of gene expression on a genome-wide scale. We propose a hierarchical model of transcriptional regulation during development in which core promoters define broad windows of opportunity for expression, by defining a limited range of transcription factors from which they are able to receive regulatory inputs. This two-tiered mechanism globally orchestrates developmental gene expression, including noncoding transcription on a scale that defies our current understanding of ontogenesis. Indeed, noncoding transcripts are far more prevalent than ever reported before, with ~4,000 long noncoding RNAs expressed during embryogenesis. Over 1,500 are functionally conserved throughout the *melanogaster* subgroup, and hundreds are under strong purifying selection. Overall, this work introduces a hierarchical model for the developmental regulation of transcription, and reveals the central role of noncoding transcription in animal development.

1 INTRODUCTION

2
3 Development in metazoans is orchestrated by complex gene regulatory programs encoded in the
4 sequence of the genome¹⁻⁴. The expression of thousands of protein-coding and noncoding genes, in
5 precise spatiotemporal patterns, progressively refines the organization of embryonic structures and
6 specifies the differentiation of specialized cell types. In *Drosophila*, many of the master genes governing
7 early development^{1,2} encode regulators of transcription⁵⁻⁸, and transcriptional regulation largely accounts
8 for embryo patterning^{5,9}.

9
10 The regulatory code that specifies these programs, however, remains poorly understood.
11 Sequences that bind transcriptional activators and repressors, known as enhancers¹⁰⁻¹⁶, are generally
12 thought to be the primary determinants of gene expression specificity. Sequences that serve as docking
13 sites for the basal transcription machinery, the core promoters¹⁷, are usually assumed to be structural
14 elements that contribute little or no regulatory information.

15
16 Indeed, core promoters contain sequence motifs (e.g., TATA boxes) that serve as a platform for
17 RNA Pol II initiation at transcription start sites (TSSs), but are not by themselves sufficient to induce
18 transcription^{14,17}. Sequence-specific activators and repressors, collectively designated as transcription
19 factors (TFs), bind to enhancers and foster the assembly of the basal transcription machinery at
20 associated core promoters^{14,17}. Beyond these general principles, the syntax of the code, and in particular
21 the functional relationship between these two classes of elements, remains obscure. Interacting core
22 promoter-enhancer pairs may be directly adjacent^{12,16,18}, or may be located hundreds of kilobases apart in
23 metazoan genomes^{14,19}, and the rules enforcing specific interactions are unknown^{14,17}. There is evidence
24 that core promoters can influence the expression specificity of some genes^{18,20}, but so far this has not
25 been systematically studied in the context of development. Understanding the basis of transcriptional
26 control requires parsing out these complex interactions.

27
28 In addition to delineating the rules of gene regulation, it is necessary to expand the concept of
29 gene^{21,22} to include all the noncoding loci that may control developmental processes. Indeed, it has
30 become increasingly clear that noncoding transcription is very prevalent in Eukaryotes²³⁻²⁸, and both
31 genetic and biochemical studies have unambiguously established long noncoding RNAs (lncRNAs) as
32 functional components of the cellular machinery²⁹⁻³². Exhaustively annotating noncoding transcripts, and
33 identifying those with biologically relevant functions, is crucial to our understanding of development.

34
35 Deciphering the meaning of regulatory sequences, or assessing the biological relevance of
36 lncRNA genes, are daunting tasks that require innovative strategies. The use of genome-wide functional
37 assays in a phylogenetic framework is a powerful and general approach to such questions³³⁻³⁶. Indeed,
38 direct measurements of complex genome functions in multiple species provide a sort of genomic Rosetta
39 Stone from which the underlying code can begin to be parsed out.

40
41 Here, we used high-throughput TSS mapping in tightly resolved time series to establish genome-
42 wide promoter activity profiles throughout embryonic development in 5 *Drosophila* species spanning 25-
43 50 million years (MY) of evolution. Combining TSS identification at single-nucleotide resolution³⁷ with
44 quantitative measurements of developmental expression patterns, we uncovered unique features of
45 expression timing and core promoter structure to generate novel insights into transcriptional regulation.

46
47 We report that distinct types of core promoters are selectively active in three broad phases of
48 embryonic development: specific combinations of core motifs mediate transcription during Early,
49 Intermediate and Late embryogenesis. Each individual class of core promoters is functionally associated

50 with distinct sets of transcription factors. We propose a two-tier model of transcriptional control in which
51 core promoters and enhancers mediate, respectively, coarse-grained and fine-grained developmental
52 regulation.

53
54 We also show that noncoding transcription is far more widespread than anticipated in *Drosophila*,
55 with 3,973 promoters driving the expression of lncRNAs during embryogenesis. The analysis of these
56 core promoters, most of which are currently unannotated, shows that they largely share the structural and
57 functional properties of their counterparts at protein-coding genes. Through the analysis of their fine
58 structure and sequence conservation, we demonstrate that evolutionarily conserved lncRNAs promoters
59 are under strong purifying selection at the levels of primary sequence and expression specificity. We
60 functionally characterize the *schmurri-like RNA (slr)* locus, which expresses a lncRNA in a highly
61 conserved spatiotemporal pattern suggestive of a role in early dorsoventral patterning.

62
63 In summary, these results uncover a major active role for core promoters in regulating
64 transcription, by defining windows of opportunity for activation by enhancer sequences. They also reveal
65 a vastly underappreciated aspect of developmental transcriptomes, by showing that noncoding
66 transcription is extremely prevalent, tightly regulated and, crucially, deeply conserved.

67
68
69
70

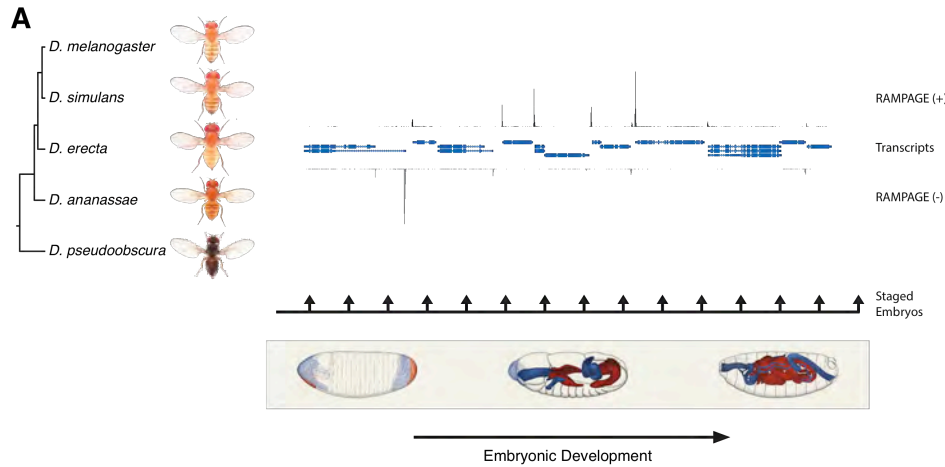
71 RESULTS

72 73 **Global multispecies profiling of developmental promoter activity**

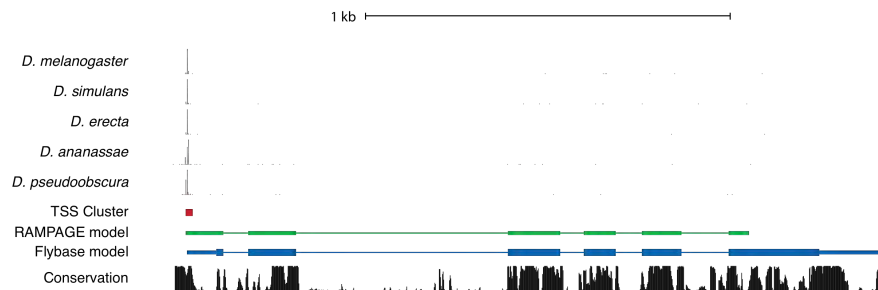
74
75 To explore the genome-wide dynamics of transcriptional regulation and their evolution, we
76 generated developmental transcriptome profiles at both very high temporal resolution (1 hour) and high
77 sequence coverage (137-180 million read pairs per species) for 5 *Drosophila* species spanning 25-50
78 million years of evolution: *D. melanogaster*, *D. simulans*, *D. erecta*, *D. ananassae* and *D. pseudoobscura*
79 (Fig. 1A; total of 120 samples). We focused on embryonic development, a crucial period during which
80 the body plan is established and all larval organs are generated. This data allows direct, genome-wide
81 comparisons of promoter activity in a phylogenetic framework (Fig. 1B).

82
83 In order to map transcription start sites (TSSs) with single-base resolution and accurately measure
84 the activity of individual promoters, we used a high-fidelity method called RAMPAGE³⁷ based on high-
85 throughput sequencing of 5'-complete complementary DNA (cDNA) molecules. It offers unprecedented
86 specificity and detection sensitivity for TSSs, and its multiplexing capabilities allow for the seamless
87 acquisition of high-resolution developmental time series³⁷. Since eukaryotic promoters often allow
88 productive transcription initiation from multiple positions³⁷⁻⁴⁰, we use a dedicated peak-calling algorithm
89 to group neighboring TSSs into TSS clusters (TSCs) corresponding to individual promoters^{37,41}.

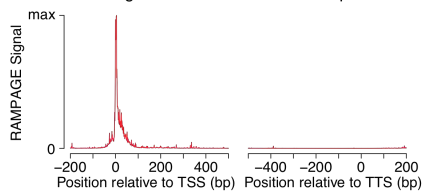
90
91 For each species, we identified 2.2×10^4 to 2.7×10^4 high-confidence TSCs. The narrow distribution
92 of raw RAMPAGE signal (Fig. 1C & S1) and of TSCs (Fig. S2) over annotated loci confirms our very
93 high specificity for true TSSs. The quantification of promoter expression levels is highly reproducible
94 across biological replicates (Fig. 1D-E). Importantly, paired-end sequencing of cDNAs allows for
95 evidence-based assignment of novel TSCs to existing gene annotations, and provides valuable
96 information about overall transcript structure (Fig. 1B). We can thus attribute 82% of *D. melanogaster*
97 TSCs to annotated genes, the remaining 18% potentially driving the expression of unannotated non-
98 protein-coding transcripts.



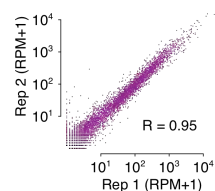
B Multispecies RAMPAGE: *NLaz* locus



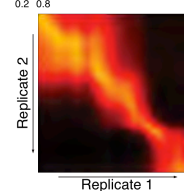
C RAMPAGE signal over annotated transcripts



D Quantification reproducibility



E Pearson R^2



99

100 **Figure 1: Comparative profiling of embryonic promoter activity**

101 (a) Genome-wide TSS usage maps were generated by RAMPAGE assays for developmental series in 5 species (central panel,
102 locus chr2L:1151500-1185900). Blue symbols: FlyBase transcript annotations; black density plot: density of RAMPAGE read
103 5' ends on the positive strand; grey density plot: read 5' ends on the negative strand (inverted y-axis). Fly photographs (N.

104 Gompel) and embryo drawings (V. Hartenstein, CSHL Press 1993) reproduced with permission from FlyBase. (b)

105 Distribution of RAMPAGE signal over the *NLaz* gene in 5 *Drosophila* species. From top: RAMPAGE signal intensity tracks,
106 *D. melanogaster* TSC, transcript model inferred from RAMPAGE data, transcript model from FlyBase, sequence conservation

107 (phastCons). For visualization of non-*melanogaster* data, sequencing reads were mapped to the appropriate genomes and
108 projected onto orthologous *D. melanogaster* positions based on whole-genome alignments. (c) Metaprofile of RAMPAGE

109 read 5' ends over FlyBase-annotated mature transcripts (introns excluded). (d) Reproducibility of RAMPAGE signal

110 quantification for individual TSCs (n=9,299) for two biological replicates of the first *D. melanogaster* time point (0-1h). TSCs

111 with no signal in either replicate were excluded. (e) Correlation matrix for the *D. melanogaster* time series biological

112 replicates. We plotted the correlation of TSC expression (n=24,832 TSCs) after smoothing and alignment of the time series.

113

114 The comparison of biological replicates for the *D. melanogaster* time series confirms our ability
115 to quantitatively measure promoter expression dynamics throughout development. Indeed, this analysis
116 (Fig. S3) shows excellent reproducibility (Pearson $R^2 = 0.95$) for TSCs with maximum expression ≥ 25
117 reads per million (RPM, Fig. 2A & S3), and very good reproducibility (Pearson $R^2 = 0.92$) with
118 maximum expression ≥ 10 RPM (Fig. S3).

119
120 Post-synchronization of developmental series was achieved by global alignments of all time
121 series to the *D. melanogaster* reference to maximize the overall similarity between genome-wide
122 expression profiles^{42,43}. The resulting alignments for well-known developmental genes, used here as
123 diagnostic markers, confirm the very high quality of the global alignments (Fig. 2B & S4). For all genes
124 with one-to-one orthologs, the expression profiles are overall tightly conserved across species (Fig. 2C &
125 S5A), but with substantial gene-to-gene variability: *hunchback*, for instance, displays considerable
126 conservation in the expression of both of its promoters (Pearson R^2 of 0.88 and 0.97; Fig. S5B), whereas
127 the *RpL19* promoter shows rapid divergence ($R^2=0.08$; Fig. S5C).

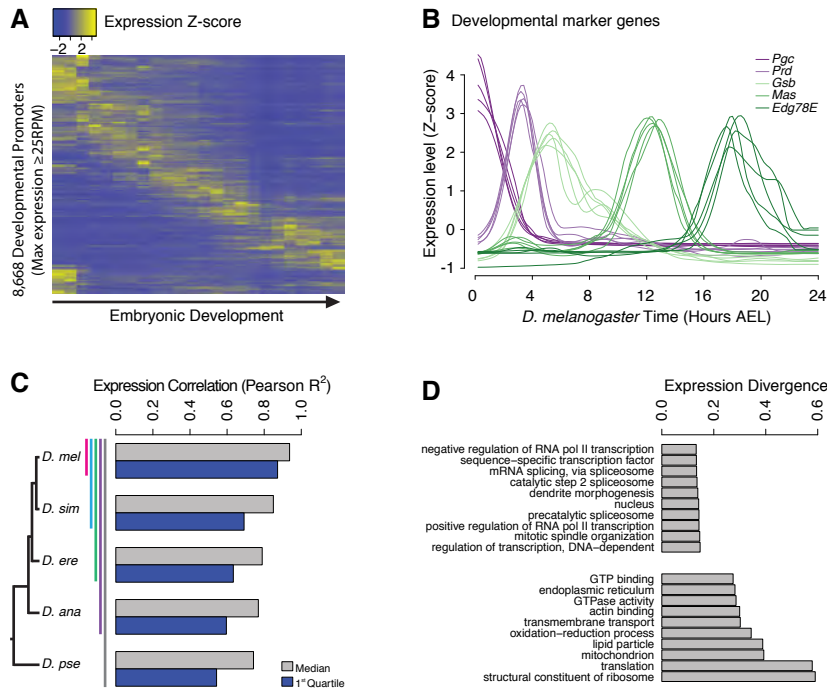
128
129 We found a strong relationship between selective constraints on expression specificity and gene
130 function: indeed, the degree of expression divergence differs substantially between Gene Ontology (GO)
131 annotation categories (Fig. 2D). Functions related to the regulation of transcription and splicing display
132 the strongest conservation of expression, in accordance with the known molecular function of many
133 master regulators of early development. Categories related to the core translational machinery and
134 cytoskeletal structures display much more plastic expression specificities.

135
136 The high similarity of biological replicates, the accuracy of inter-species alignments for well-
137 known developmental genes, and the biological features of evolutionary divergence patterns, together
138 confirm our ability to accurately quantify promoter expression and its variation across species. Our
139 observations also highlight the central importance of systems-level selective constraints, such as those
140 acting on gene function and developmental stages, in shaping the evolution of gene expression.

141
142
143
144
145 **Core promoter structure defines broad developmental phases of gene expression**

146
147 We leveraged our multispecies expression data to study promoter structure-function relationships,
148 and thus gain insights into the regulatory code that determines developmental gene expression. We
149 focused on a set of 3,462 promoters functionally active in all species that we classified either as
150 housekeeping (< 5 -fold variation throughout development) or as developmentally regulated ($\geq 60\%$
151 of total expression within an 8-hour window). The developmentally regulated group was further clustered
152 based on temporal correlation, yielding a total of 9 distinct expression clusters (Fig. 3A, see Methods).
153 These thresholds were chosen to maximize the total number of promoters included and the similarity of
154 profiles within each cluster, while still yielding clusters large enough for statistical analysis.

155
156
157



158
159
160
161
162
163
164
165
166
167
168
169

Figure 2: Evolutionary divergence of the developmental timing of promoter activity

(a) Hierarchically clustered expression profiles for individual *D. melanogaster* promoters throughout embryogenesis (24 time points; replicate 1). Only promoters with a maximum expression level ≥ 25 RPM ($n=8,668$), for which quantification reproducibility is very high, are included. (b) Expression profiles for 5 developmental marker genes, after global alignment of all time series to *D. melanogaster* (see Methods). The 5 curves for each gene correspond to the 5 species. (c) Conservation of the temporal expression profiles of individual promoters. For each subclade, we computed the average correlation coefficient between all pairs of species for each individual gene. The graph shows the median and first quartile over all genes with orthologs throughout the subclade. (d) The evolutionary divergence of expression specificity varies widely between Gene Ontology (GO) categories. For each gene with orthologs in all 5 species, maximum expression ≥ 25 RPM and expression changes ≥ 5 -fold, we computed a measure of overall divergence across the clade (see Methods). The bar plot shows the average divergence by GO category, for the 20 categories with the lowest (top) and greatest (bottom) divergence.

170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185

These 9 expression clusters fall into 3 main groups, defined by the core motif composition of the promoters (Fig. 3B). Indeed, we unexpectedly observed robust enrichments for specific sets of motifs in the promoters of all individual expression clusters. Three major classes emerge, within which motif enrichments are relatively homogeneous (Fig. 3B).

Remarkably, these three classes define distinct temporal phases of embryonic development (Fig. 3A-B). The Early expression class, enriched for DRE and Ohler-1/5/6/7 motifs, consists of the promoters for maternally deposited transcripts, including the housekeeping cluster. The Intermediate class, enriched for Initiator (INR), MTE and DPE motifs, mediates transcription throughout a broad phase of mid-embryogenesis, from the onset of zygotic expression to the end of organogenesis. The Late class, enriched for TATA boxes, drives transcription around the transition to the first larval stage. Notably, expression clusters with vastly different specificities (e.g., D1 and D2) share the same promoter structure trends.

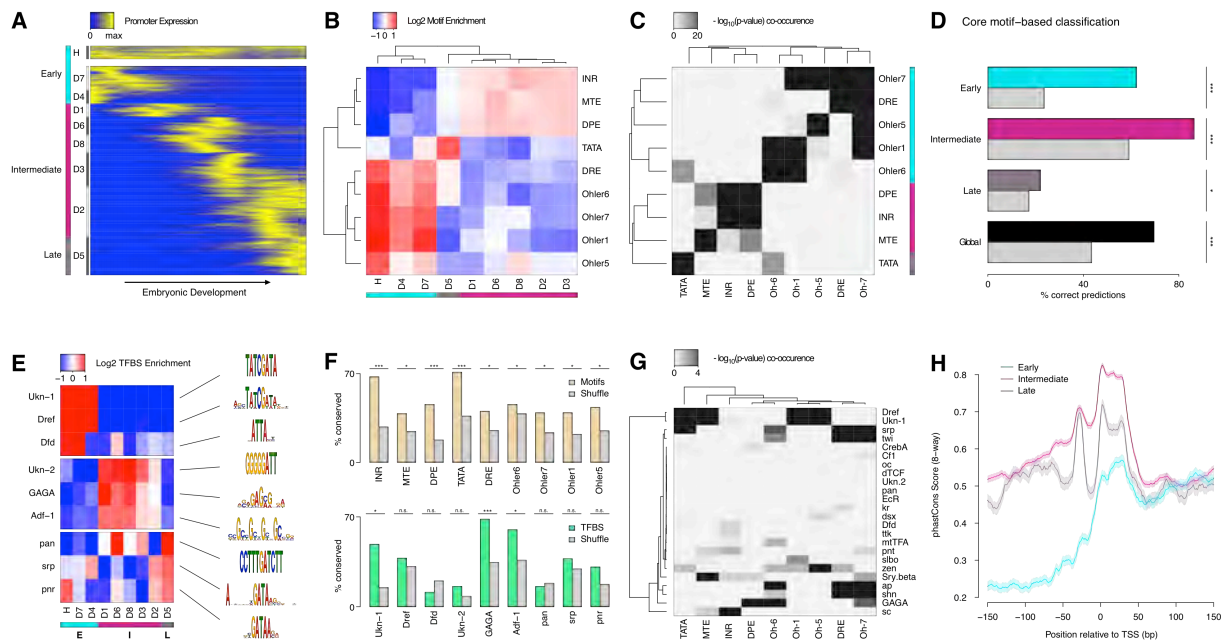
186 Importantly, clustering core motifs by their tendency to co-occur within the same promoters,
187 regardless of expression timing, recapitulates the same three main motif groups (Fig. 3C). It does,
188 however, also uncover a finer data structure, which suggests that there are a variety of promoter subtypes.
189 In addition, we were able to train a classifier with substantial predictive power to distinguish the three
190 promoter classes based solely on core motif scores (Fig. 3D). Taken together, our observations suggest
191 that core promoter elements play a significant role in restricting windows of opportunity for expression
192 during distinct periods of development.

193
194 Sequences proximal to the three core promoter classes are also enriched for different sets of
195 TFBSs (Fig. 3E). The Early class preferentially harbors binding sites for Dref and Dfd, while the
196 Intermediate class favors Trl (GAGA motif) and Adf-1. The Late class is enriched in pan, srp and pnr
197 sites. The presence of most core promoter motifs and TFBSs is conserved between species far beyond
198 random expectations (Fig. 3F), which validates the overall quality of our motif predictions. Interestingly,
199 TFBSs are often specific for only a subset of expression clusters within a class (e.g., Dfd or GAGA).
200 This suggests a model in which core promoter structure defines broad developmental periods of
201 expression potential, and precise expression timing is then refined by sequence-specific transcription
202 factors. Combinatorial encoding through stereotypical sets of TFBSs is likely to sharpen expression
203 patterns even further (Fig. S6A-B).

204
205 The analysis of favored pairings between individual TFBSs and core promoter motifs suggests a
206 possible mechanism to mediate this 2-step specification of expression patterns. Indeed, some TFBSs
207 appear to be strongly associated with specific sets of core promoter motifs (Fig. 3G). Dref sites, for
208 instance, are preferentially found along DRE and Ohler-1/6/7 core motifs. Likewise, twi and srp TFBSs,
209 which often tend to be found together (Fig. S6A), have a robust association with INR and MTE. These
210 strong affinities suggest that core motifs may tune the ability of a promoter to respond to specific sets of
211 transcription factors. They may do so by recruiting different sets of general transcription factors (GTFs)
212 that functionally interact with distinct groups of activators. Such a mechanism may channel various
213 regulatory inputs to limited subsets of promoters and thus limit crosstalk between parallel pathways.

214
215 We found that the three promoter classes display markedly different profiles of sequence
216 conservation (Fig. 3H & S6C). Importantly, this analysis only includes those promoters for which we
217 have detected transcriptional activity in all 5 species, and we can therefore categorically rule out
218 differences in rates of promoter gain/loss as an explanation. These observations suggest that the 3
219 promoter classes indeed have intrinsic structural differences that make them subject to distinct regimes of
220 natural selection and sequence evolution.

221
222
223
224



225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243

Figure 3: Core promoter structure defines 3 broad phases of embryonic development.

(a) Clustering of expression profiles for 3,462 promoters (maximum expression ≥ 10 RPM) defined as housekeeping (H) or developmentally regulated (D). Developmentally regulated promoters are divided into 8 groups based on hierarchical clustering of expression profiles (see Methods). Core promoter structure defines 3 broad developmental phases (color sidebar). (b) Relative enrichment of core promoter motifs in each expression cluster. Three major clusters can be defined (bottom color bar), which correspond to 3 phases of embryonic development (see (a); Early: 817 TSCs; Intermediate: 2,047; Late: 598). (c) Clustering of core promoter motifs based on their co-occurrence in the same promoters. This approach yields the same 3 major sets of motifs previously defined based on expression profiles (see (b)). (d) Predictions of expression timing from core promoter structure. Quadratic discriminant analysis on log-transformed motif scores was used to predict the developmental phase in which promoters are expressed. The performance of the classifier in leave-one-out cross-validation (color bars) is compared to random expectation (grey bars; FDR-corrected chi-square tests applied to individual contingency tables as appropriate). (e) Distinct sets of TFBSs are enriched near the core promoters active in the 3 developmental phases. The top 3 motifs for each core promoter class are shown. (f) The conservation of core promoter motif composition between *D. melanogaster* and *D. pseudoobscura* confirms the biological relevance of the motifs. (Grey bars: conservation of shuffled motifs; FDR-corrected chi-square tests). (g) Many TFBSs are strongly associated with specific sets of core promoter motifs. Results are shown for the 3 motifs most enriched near the promoters of each expression cluster. (h) The 3 major core promoter classes display distinct profiles of sequence conservation. Lines: median phastCons scores across promoters of the class. Envelope: standard error in the estimate of the median, computed by bootstrapping.

244
245
246
247
248
249
250
251
252
253
254
255

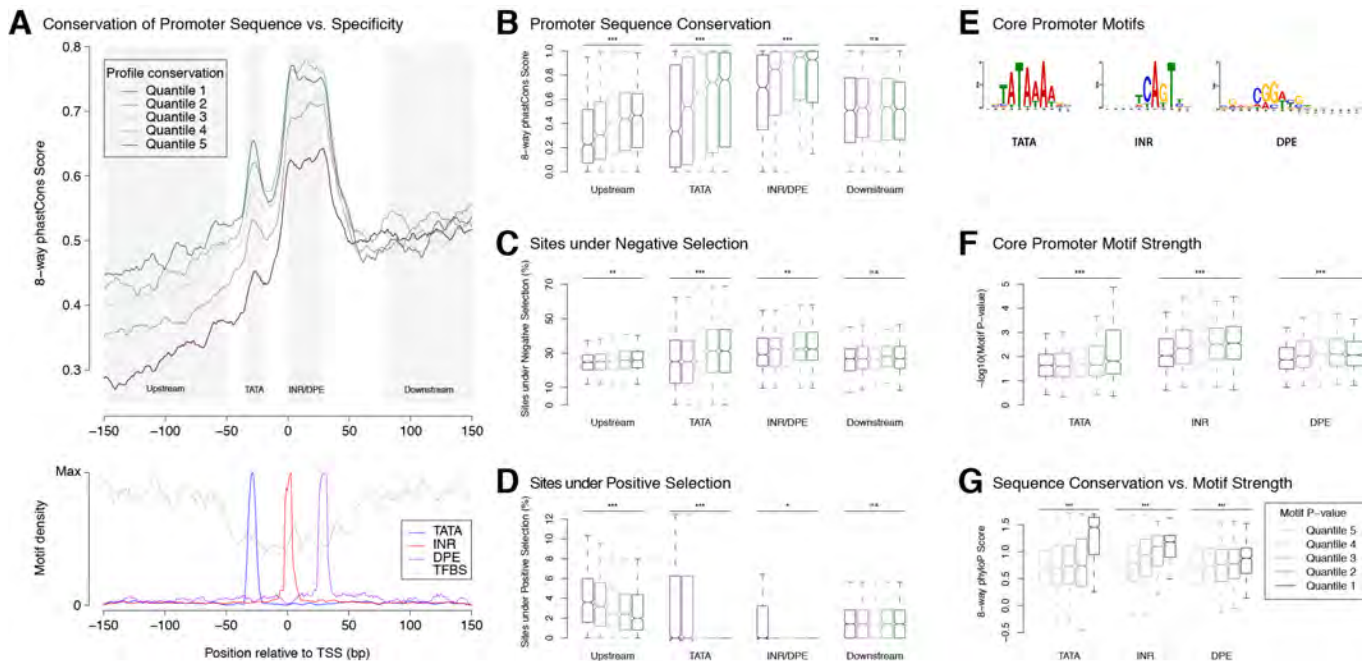
Selection on expression specificity shapes promoter sequence evolution

To further explore the selective pressures acting on regulatory sequences, we investigated quantitative relationships between the evolution of promoter structure (primary sequence) and function (expression specificity). We report a subtle but highly significant correlation between the conservation of promoter sequence and that of expression profiles (Fig. 4A-B). This shows that the effects of selection on expression specificity are reflected in the evolution of promoter sequences. The main areas of differential sequence conservation overlap regions preferentially occupied by precisely positioned core promoter elements (TATA, INR, DPE) and transcription factor binding sites (Fig. 4A). Importantly, this correlation does not hold for sequences >50 bp downstream of the TSS, which are likely subject to additional selective pressures on 5'-UTR and protein-coding sequences. Promoters with highly divergent

256 expression profiles are depleted of sites under purifying selection, and enriched for sites under positive
 257 selection (Fig. 4C-D).

258 *Ab initio* motif discovery within the regions of differential sequence conservation returned the
 259 canonical consensus sequences for TATA, INR and DPE (Fig. 4E). We found that promoters with highly
 260 conserved expression profiles tend to have core promoter elements whose sequence is closer to the motif
 261 consensus (Fig. 4F), and those in turn tend to be more conserved at the sequence level (Fig. 4G).
 262 Upstream flanking sequences also tend to show higher conservation of individual TFBSs (Fig. S7A).
 263 Importantly, it is possible to detect such a correlation for the binding sites of some individual
 264 transcription factors (Fig. S7B). This is a rather striking fact, considering that promoter-proximal
 265 enhancers generally bind more than one factor, and that many promoters are additionally regulated by
 266 distal enhancers not taken into account here. Finally, the conservation of promoter TFBS composition, as
 267 expected, also correlates with interspecific divergence (Fig. S7C).

268
 269
 270
 271



272
 273
 274 **Figure 4: Selection on expression specificity shapes promoter sequence evolution.**

275 (a) Upper panel: average sequence conservation for *D. melanogaster* promoters, by expression profile conservation quantile.
 276 All *D. melanogaster* promoters with maximum expression ≥ 25 RPM and functionally conserved in all 5 species were included
 277 ($n=4,973$). Lower panel: density of core promoter elements and TFBSs over all promoters. TATA box, INR, DPE: respective
 278 consensus sequences STATAAA, TCAGTY or TCATTCG, KCGGTTSK or CGGACGT⁴⁴; TFBS motifs from Jaspar. (b)
 279 Sequence conservation by profile conservation quantile, over promoter subregions depicted as shaded areas in A. Upstream
 280 region runs from -300 to -50bp. (***) p-value $< 10^{-8}$ for Pearson correlation test between profile conservation and sequence
 281 conservation; ** $p < 10^{-4}$; * $p < 10^{-2}$; *n.s.* not significant) (c) Proportion of bases under purifying selection (phyloP score > 0.1).
 282 (d) Proportion of bases under positive selection (phyloP score < -0.1). (e) Core motif position weight matrices derived *de novo*
 283 from promoter sequences (see Methods). (f) Core motif strength correlates with expression profile conservation. (g) Core
 284 motif sequence conservation correlates with proximity to motif consensus.

285
 286

287 These observations establish a clear relationship between core promoter sequence and expression
288 specificity. Selective constraints on developmental expression timing act particularly strongly on core
289 promoter elements, most notably Initiator, DPE and TATA elements. This is consistent with the idea that
290 core promoter motifs are indeed key determinants of this specificity. Together, our results further support
291 the hypothesis that core motifs and general transcription factors play a crucial role in determining
292 promoter expression specificity.

293
294
295

296 **Promoter birth and death are widespread and dynamic**

297

298 In addition to changes in the specificity of individual promoters, gene expression programs
299 evolve through turnover of regulatory elements^{36,45,46}. And indeed, we observed widespread birth and
300 death of promoters throughout the clade: only 49% of *D. melanogaster* TSCs are functionally conserved
301 in all 5 species (Fig 5A). To rule out genome assembly artifacts, we restricted our analysis to those with
302 syntenic alignments to other genomes, and measured a functional conservation rate of 75% (Fig 5A). As
303 some peaks lack alignments owing to genuine insertions or deletions, we expect the true conservation
304 rate to be within the 49-75% range. Analyzing TSC conservation between species pairs, or from a *D.*
305 *simulans*-centric perspective, yields similar conclusions (Fig. S8).

306

307 The analysis of replicates for the *D. melanogaster* time series shows the false positive rate for
308 gain/loss event detection to be under 0.1% (Fig. 5A). Although TSCs with lower expression levels tend
309 to be less conserved, general trends are shared between TSCs of all expression levels (Fig. S9). Even
310 though a gain/loss of expression during embryogenesis may reflect a shift in specificity rather than a
311 complete gain/loss of function, the vast majority (91.4%) of *D. pseudoobscura* TSCs that were inferred
312 to be lost in *D. melanogaster* based on embryo data are never expressed at any other stage of the life
313 cycle (analysis of published data³⁷). Therefore, we conclude that our strategy accurately and robustly
314 detects promoter gain and loss events. Consistent with this, we can reconstruct the known species
315 phylogeny by treating the presence or absence of individual TSCs as binary discrete characters in a
316 standard parsimony framework (Fig. 5B). Functional conservation is reflected in sequence conservation:
317 promoters active in all species display far higher sequence conservation than species-specific ones, which
318 appear no more constrained than surrounding regions (Fig 5C & S10).

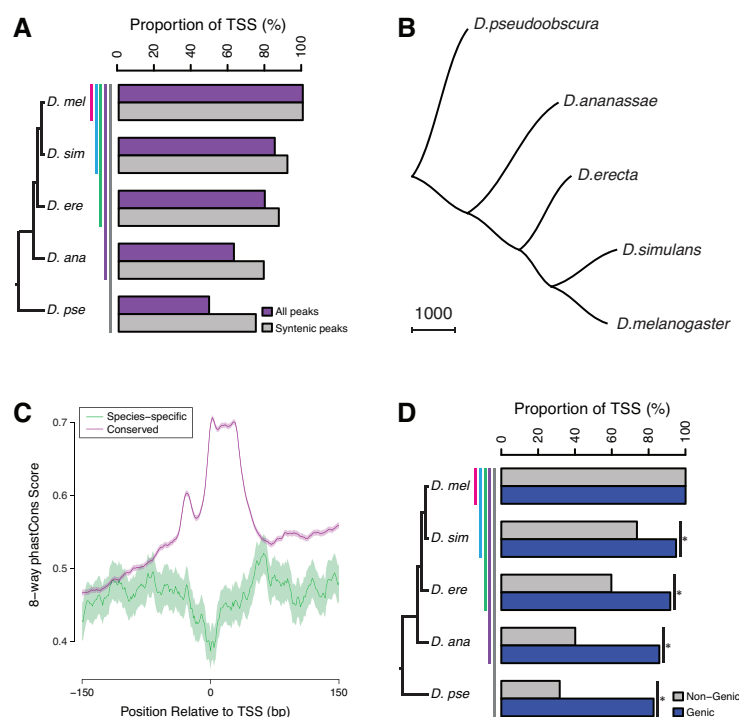
319

320 While purifying selection clearly plays a major role, a sizeable proportion of TSCs (25-50%) are
321 not shared between all species, underscoring the inherently fluid nature of the regulatory landscape.
322 Comparison to published data on the evolution of Twist binding sites³⁶ revealed that overall, promoters
323 evolve no more slowly, and possibly even faster, than TFBSs do (Fig. S11). Although a rigorous
324 comparison between such disparate data types is delicate, this shows that promoters and TFBSs do not
325 turn over on vastly different timescales.

326

327 Our data also reveals the existence of several thousand novel promoters that cannot be assigned to
328 any annotated genes. Many of those may drive the expression of long noncoding transcripts, and it has
329 long been a matter of debate to what extent this transcriptional activity plays meaningful biological roles.
330 Hence we analyzed their rates of gain and loss to explore the selective pressures they may be subject to.
331 We found a stark contrast in the degree of functional conservation of the two classes, with novel non-
332 genic promoters evolving at a substantially higher rate than genic ones (Fig. 5D). There is, however, a
333 very substantial proportion that is deeply conserved, suggesting the possibility of widespread
334 functionality. The discrepancy between the classes may be due to a larger proportion of noncoding
335 transcripts being devoid of biological roles and evolving neutrally. Alternatively, it may instead reflect a

336 more pronounced tendency for noncoding transcription to take on lineage-specific roles and thereby be a
 337 driver of adaptation, as has been suggested before. To gain a better understanding of these questions, we
 338 sought to investigate in more depth the evolution of noncoding transcription.
 339
 340



341
 342

Figure 5: Widespread evolutionary gain and loss of promoters.

344 (a) Proportion of *D. melanogaster* TSCs reproducibly detected in biological duplicates (first pair of bars) and functionally
 345 conserved in all species of subclades of increasing sizes. Subclades include all descendants of a common ancestor, and are
 346 designated by the species that is most distantly related to *D. melanogaster*. (b) The species phylogeny can be accurately
 347 reconstructed from patterns of TSC gain and loss. The presence/absence of each TSC was treated as a discrete character and
 348 the unrooted tree reconstructed using the Phylip software package. (c) Average profiles of sequence conservation over the
 349 TSCs functionally conserved across all 5 species and those specific to *D. melanogaster*. The shaded areas represent the
 350 standard deviation, estimated from 1,000 bootstraps. (e) TSCs driving the expression of annotated genes display a far higher
 351 degree of functional conservation than “orphan” TSCs ($p < 0.01$ for all pairwise comparisons; chi-square test with Bonferroni
 352 correction).

353

354

355

356

357

Deep conservation of long noncoding RNA promoters

358

359

360

361

362

363

364

365

366

367

The prevalence and biological relevance of noncoding transcription have long been major areas of
 contention. Attempts at resolving these issues using genomic sequence conservation have been largely
 inconclusive, probably due to minimal selective constraints on the primary sequence of these
 transcripts^{27,28,47}. Studying promoter activity experimentally in a phylogenetic framework provides a
 unique opportunity to rigorously address the question of lncRNA conservation and functionality.
 Furthermore, our ability to pinpoint TSSs with single-nucleotide accuracy gives us unprecedented
 leverage to elicit elusive sequence conservation patterns. Furthermore, beyond sequence conservation,
 we are also in a position to assess selective constraints on the expression specificity of these promoters.

368 We found 3,682 embryonic TSCs in *D. melanogaster* that could not be functionally linked to any
369 annotated protein-coding or small RNA gene, and could therefore represent putative lncRNA promoters.
370 We also identified TSCs for 291 annotated lncRNAs, bringing the total up to 3,973. Their developmental
371 expression kinetics appear to be diverse and exquisitely stage-specific (Fig. 6A).

372
373 The analysis of published genome-wide DNaseI hypersensitivity data⁴⁸ confirmed that these
374 putative lncRNA TSCs are likely to correspond to genuine promoters (Fig. 6B & S12). In addition, to
375 verify that the transcripts expressed from these TSCs are indeed independent and devoid of any
376 significant protein-coding potential, we built transcript models from a recently published RNA-seq
377 developmental time course²⁵. We successfully generated transcript models for 16,105 TSCs, including
378 1,475 lncRNA TSCs. Most of them appear to correspond to full-length transcripts, and the vast majority
379 of putative lncRNAs do not overlap annotated protein-coding sequences (Fig. S13). Analyses of protein-
380 coding potential confirm that the overwhelming majority of transcripts are unlikely to encode proteins, or
381 even peptides as short as 10 amino acids (Fig. 6C & S13). Transcripts from 18 loci are likely to encode
382 short open reading frames (sORFs, <100 residues). We conclude that the vast majority of candidate
383 transcripts are likely to be genuine lncRNAs.

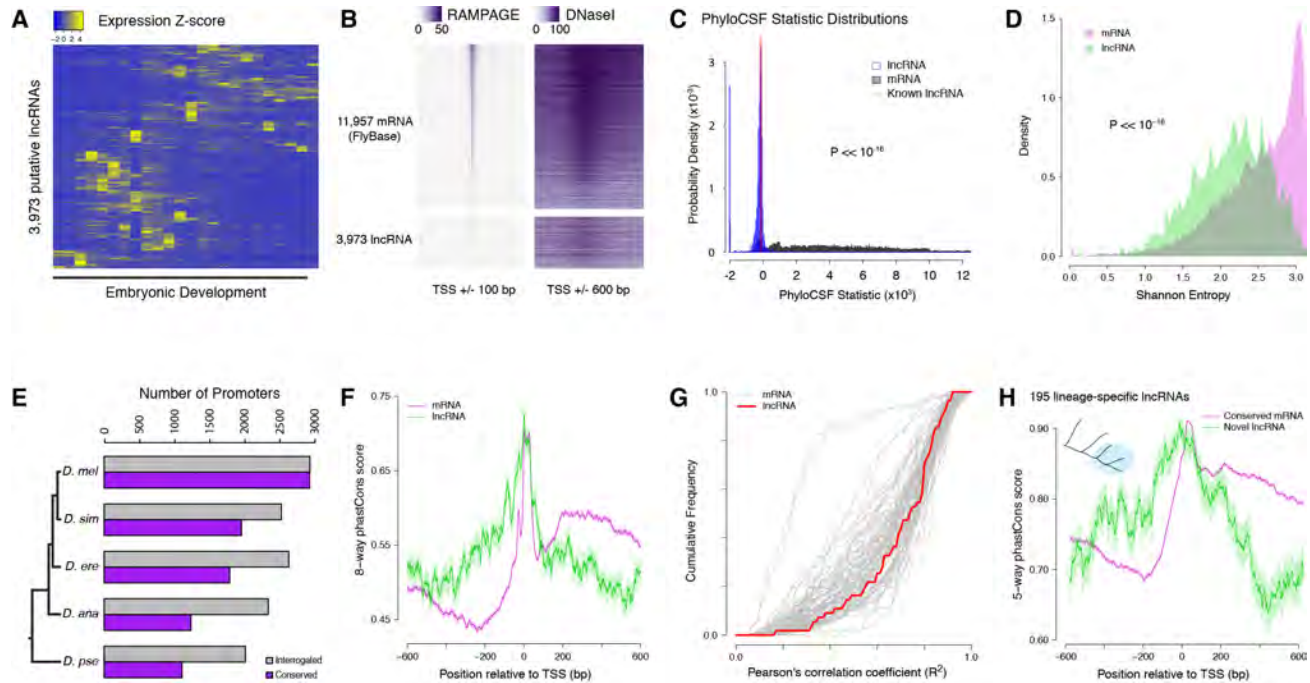
384
385 The expression profiles of lncRNA promoters have unique properties that set them apart. As a
386 class, they tend to be substantially more stage-specific than their protein-coding gene counterparts, as
387 measured by the Shannon entropy of their activity profiles (Fig. 6D). Comparing their developmental
388 expression timing to that of protein-coding gene promoters suggests a broad diversity of potential
389 developmental functions (Fig. S14 & Table S1).

390
391 The detection of lncRNA TSCs is highly reproducible across *D. melanogaster* biological
392 replicates (Fig. 6E). Of all *D. melanogaster* lncRNA TSCs, 2,016 can be aligned to the *D. pseudoobscura*
393 genome assembly and 1,111 are functionally conserved (Fig. 6E), suggesting that they have been
394 maintained since the last common ancestor of these two species. In order to investigate whether lncRNA
395 promoters are under purifying selection, we focused on an extremely stringently selected set of 631 TSCs
396 that are active in all 5 species. This set includes well-known essential noncoding transcription units, such
397 as *bithoraxoid*^{49,50} (Fig. S15) and *roX1*⁵¹. Overall, the level of sequence conservation at these functionally
398 preserved lncRNA promoters is similar to that observed at protein-coding gene promoters (Fig. 6F).
399 Their developmental expression specificity is also more constrained than that of many protein-coding
400 gene promoters (Fig. 6G). Both observations taken together argue strongly for sustained selective
401 pressure on these 631 lncRNA promoters for 25-50 million years. Furthermore, many TSCs of interest
402 were excluded from this analysis simply because of the poor quality of genome assemblies, and this is
403 therefore an extremely conservative set. Therefore, to place a more reasonable lower bound on the true
404 number of conserved lncRNA promoters, we focused on those shared between the 3 species of the
405 *melanogaster* subgroup. These 1,529 promoters similarly display a high degree of sequence conservation
406 within the subgroup, and their expression specificity also appears constrained (Fig. S16).

407
408 Still, it remains that lncRNA promoters as a class evolve much faster than those of protein-coding
409 genes (Fig. 5D), and it has been a matter of debate whether this reflects lineage-specific functions or
410 merely neutral evolution. To address this question, we focused on 195 lncRNA TSCs that are specific to
411 the *melanogaster* subgroup, despite the orthologous sequences being present in the genomes of the two
412 outgroups (Fig. 6H). Surprisingly, they display the same degree of sequence conservation as the protein-
413 coding gene promoters that are shared throughout the subgroup (Fig. 6H). The assessment of
414 conservation at orthologous sequences in outgroup species confirms that the selective constraints are
415 indeed lineage-specific (Fig. S17). This argues that these evolutionarily recent lncRNAs have come
416 under purifying selection after acquiring lineage-specific functions.

417
418
419
420
421
422
423
424

Taken together, our observations show that a vast proportion of lncRNAs are indeed under purifying selection for biological functions relevant to embryonic development.



425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448

Figure 6: Strong purifying selection on long noncoding RNA promoters.

(a) Developmental expression profiles of putative lncRNA promoters (n=3,973). (b) Heatmaps of RAMPAGE signal (left; number of reads) and DNase-seq signal (right; arbitrary units) over individual TSCs. We are comparing TSCs that overlap FlyBase-annotated mRNA transcription start sites (top), which we use as positive controls, to the TSCs of putative lncRNAs (bottom). In each group, TSCs are sorted by total RAMPAGE signal intensity. (c) Distribution of phyloCSF scores for transcript models corresponding to putative lncRNAs (n=1,475) and mRNAs (n=16,105). Transcript models were built from publicly available RNA-seq data using Cufflinks. The phyloCSF metric quantifies the protein-coding potential of transcripts, based on the presence and conservation of ORFs. (d) Shannon entropy of the temporal expression profiles for lncRNA (n=2,397) and mRNA (n=18,067) promoters with maximum expression ≥ 2 RPM. Overall, the profiles of lncRNA promoters have lower entropy, reflecting more acutely stage-specific expression. (e) Number of *D. melanogaster* lncRNA promoters functionally preserved in other species. The grey bars represent the number of promoters for which the multiple sequence alignments passed our filtering criteria, and therefore could be interrogated. (f) Promoter sequence conservation. Considering all promoters that are functionally preserved in all 5 species, the sequences of lncRNA promoters (n=631) are under comparable selective pressure to those of protein-coding genes. (g) The developmental expression profiles of functionally conserved lncRNA promoters are far more constrained than those of many categories of protein-coding genes. We considered all promoters with maximum expression ≥ 25 RPM and expression changes ≥ 5 -fold (n=55 lncRNA promoters). (h) We identified 195 *D. melanogaster* lncRNA promoters that are functionally preserved within the *melanogaster* subgroup, but not in the 2 outgroup species, and are therefore likely to have been recently acquired specifically in this lineage (inset, top left). Lineage-specific lncRNA promoters display a level of sequence conservation within the subgroup similar to that of conserved protein-coding gene promoters.

449 ***schnurri-like RNA: A deeply conserved, developmentally regulated lncRNA gene***
450

451 To validate our findings in one specific case, we focused on the *FBgn0264479* locus, which
452 displays one the most tightly conserved expression patterns among all the lncRNA genes in our dataset.
453 This is an intriguing embryonic transcript that, although it has been annotated based on expressed
454 sequence tag (EST) data, has to our knowledge never been characterized. The 0.5kb *FBgn0264479* RNA
455 is extremely unlikely to encode functional peptides, as assessed by phyloCSF analysis (score of -217.3)
456 and manual curation (Fig. S18). It is highly expressed in all 5 species surveyed, in a strikingly conserved
457 temporal pattern restricted to a ~3-hour period encompassing the onset of gastrulation (Fig. 7A-B &
458 S19). Northern-blot analysis confirmed the size and expression dynamics of the transcript (Fig. 7C &
459 S20).

460
461 The body of the transcription unit displays hallmarks of robust purifying selection within the
462 *melanogaster* subgroup (Fig. 7A). In addition, publicly available chromatin immunoprecipitation data^{52,53}
463 reveals the binding of several transcription factors to the promoter region (Fig. S19), and their putative
464 binding sites identified by sequence motif search also show evidence of purifying selection (Fig. S21).
465 The expression dynamics of the transcription factors are consistent with their regulating the
466 *FBgn0264479* promoter (Fig. S21).

467
468 Fluorescent *in situ* hybridization (FISH) in early embryos revealed expression along the ventral
469 and dorsal midlines at the late blastoderm stage, with the exclusion of lateral regions, the primordial
470 germ cells and the prospective head (Fig. 7D). Based on this distinctive expression pattern, highly
471 reminiscent of the well-characterized *schnurri* gene⁵⁴, we renamed this lncRNA gene *schnurri-like RNA*
472 (*slr*). This early expression domain subsequently evolves into a complex segmented pattern by the end of
473 germband extension (Fig. 7D). In the late blastoderm, the RNA is found almost exclusively in the
474 cytoplasm at the apical pole of the cells (Fig. 7D). It appears at that stage to be generally concentrated in
475 a single major focus per cell (Fig. 7E). Secondary structure predictions show that the *slr* transcript is
476 likely to be highly structured (Fig. 7F), suggesting a high potential for RNA-protein interactions.

477
478 Although further characterization will be required to decipher the precise biological roles of this
479 lncRNA, our observations confirm its noncoding nature, definitively establish its developmental
480 expression specificity, and provide clear evidence that it has been under strong purifying selection over at
481 least 25-50 million years.

482
483
484
485
486

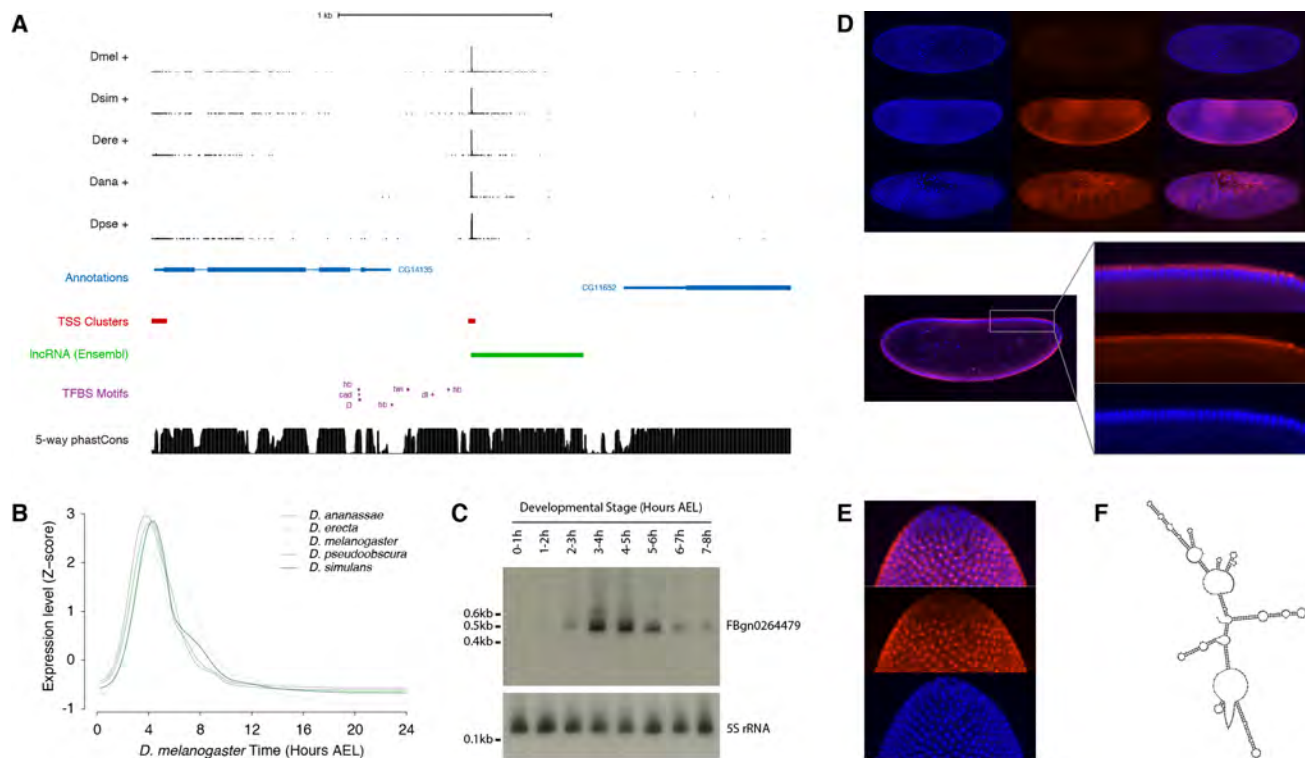


Figure 7: Functional characterization of *schnurri-like* RNA.

(a) The *schnurri-like* RNA locus (*FBgn0264479*; UCSC Genome Browser). Data tracks, from top to bottom: RAMPAGE signal in 5 species (5 upper tracks), FlyBase gene annotations (blue), *D. melanogaster* TSCs (red), transcript annotation from Ensembl (green), TFBS motif predictions around the *slr* promoter (purple), sequence conservation within the *melanogaster* subgroup (black). (b) *slr* transcript expression profiles in 5 species. (c) Northern-blot against the *slr* RNA. The 5S ribosomal RNA was used as a loading control (lower panel). (d) RNA-FISH for the *slr* transcript. Upper panel, top to bottom: maximum-intensity projections of confocal series for embryos at stages 4, 5 and 7 (Blue: DAPI, red: FISH). Bottom panel: Single confocal section from the embryo in the middle of the top panel. Controls with sense probes showed very little background. (e) *slr* RNA-FISH; lateral view of the posterior pole of a stage 5 embryo. (f) Mfold-predicted secondary structure of the RNA.

DISCUSSION

This work provides, to our knowledge, the first genome-wide overview of promoter evolution in *Drosophila*, and we leveraged this comparative framework to study the sequence determinants of developmental expression specificity. Through nucleotide-resolution mapping of TSSs and quantitative measurements of expression kinetics, our study yields new insights into the transcriptional regulatory code. We find that distinct classes of core promoters drive transcription in three broad phases of embryonic development. Each class is defined by a characteristic set of core motifs, and is associated with regulation by specific groups of transcription factors. Of note, we successfully detected *in vivo* some functional associations between Dref and housekeeping promoters, and between Trl and some developmentally regulated promoters, that were recently described *in vitro*¹⁸ (Fig. 3). Our analysis generalizes the concept of specific interactions between core promoters and enhancer sequences, and demonstrates for the first time its global relevance in a developmental context.

517 We propose a hierarchical model for transcriptional regulation in which core promoter syntax
518 defines broad temporal windows of opportunity for activation, and precise expression timing is
519 subsequently refined by the binding of sequence-specific activators and repressors at enhancers. Core
520 promoters may restrict regulatory inputs by recruiting different sets of general transcription factors
521 (GTFs) that functionally interact with distinct groups of transcription factors²⁰. Some GTFs have been
522 shown to shape the expression specificity of individual promoters, and it is known that different
523 activators and repressors have distinct requirements for cofactors and GTFs^{18,20}. Such a mechanism may
524 channel regulatory inputs to limited subsets of promoters, and thus limit crosstalk between promoters and
525 enhancers across the genome. Notably, what applies to time may also apply to space, and it is possible
526 that similar core promoter/enhancer interplay hierarchically specifies gene expression in broad
527 developmental lineages and individual cell types.

528
529 Evolutionary analysis of developmental expression specificity further supports this model. First,
530 the three major classes of core promoters defined here show drastically different patterns of sequence
531 evolution, suggesting substantial differences in their underlying structure and functional interactions with
532 the transcriptional machinery. At a finer scale, the conservation of expression specificity across species
533 correlates with the degree of sequence conservation at canonical core promoter elements, such as TATA
534 boxes and Initiator motifs. This is highly suggestive of an instructive role for these sequence elements,
535 once thought generic, in defining developmental gene expression patterns.

536
537 Approximately 4,000 promoters were found to drive the expression of lncRNAs during
538 embryogenesis, a strikingly high number for this very brief developmental period. Our findings likely
539 apply to other developmental stages as well, as we previously reported the existence of 7,421 putative
540 lncRNA promoters in an analysis of the whole *D. melanogaster* life cycle³⁷. In addition, we detected the
541 expression of only 205 of 1,119 recently identified lncRNAs²⁷. This suggests that we are only beginning
542 to scratch the surface of lncRNA biology in *Drosophila*. Importantly, we show here that vast numbers of
543 these promoters are under strong selective pressure, at the levels of both promoter sequence and
544 expression specificity. A *melanogaster* subgroup core set of at least 1,529 is under substantial selective
545 constraint, and most of those are therefore highly likely to have biologically relevant activities.

546
547 In agreement with previous reports^{28,47,55}, we find that lncRNA genes evolve faster than their
548 protein-coding counterparts. It has been an unresolved debate so far whether this reflects neutral
549 evolution or lineage-specific functions. Our observation that lineage-specific lncRNAs are also under
550 substantial selective pressure reveals that noncoding transcription may be a major driver of phenotypic
551 diversification and organismal adaptation. We recently showed that transposable elements play an
552 important role in the evolutionary gain of promoters, and in particular of lncRNA promoters³⁷. We
553 propose that transposon proliferation is a major mechanism favoring the neofunctionalization of
554 intergenic regions as sources of biologically active noncoding transcripts.

555
556 To open a window into the biology of developmentally regulated lncRNAs, we focused on
557 *schmurri-like RNA*, a deeply conserved yet never-before characterized gene. We experimentally validate
558 the existence and expression kinetics of the slr transcript, and demonstrate that both its promoter
559 sequence and its developmental expression specificity are deeply conserved across drosophilids.
560 Interestingly, this lncRNA is expressed in a spatial profile highly similar to that of the *schmurri* gene,
561 which is part of the Dpp (TGF- β /Smad) signaling pathway and plays an essential role in the
562 establishment of dorsoventral embryo polarity^{54,56}. The punctate cytoplasmic localization pattern of the
563 slr lncRNA is reminiscent of the targeting of multiple Dpp pathway components to endosomes in larval
564 wing discs⁵⁷⁻⁵⁹. Endosomes localize apically in the late embryonic blastoderm⁶⁰, and mutants for *Sara*,
565 the Smad endosome-targeting factor^{57,59}, die early in embryogenesis⁵⁷, suggesting that endosome-based

566 signaling is also essential at that stage. Taken together, our observations are suggestive of a possible role
567 for the slr lncRNA in TGF- β signaling, a crucial pathway in all animals that plays a major role in human
568 disease, including cancer.

569

570 In recent years, it has become clear that noncoding transcription serves a myriad of molecular
571 functions in Eukaryotes, and plays a part in virtually every known biological process²⁹⁻³². LncRNAs have
572 been shown to regulate transcription and chromatin structure, as well as mRNA stability and protein
573 localization. Sometimes it is the transcription of the locus itself that plays a mechanistic role, rather than
574 the resulting transcript – as in the case of the upstream *bx*d promoter⁵⁰, which has one of the most highly
575 conserved expression profiles that we have observed. Our work unambiguously demonstrates the
576 biological relevance of noncoding transcription to developmental processes, and establishes *D.*
577 *melanogaster* as an excellent model for exploring the diverse functions of lncRNA genes. We expect that
578 systematic efforts on a larger scale will illuminate the biology of a long-ignored class of genes that has
579 proven its worth.

580

581

582

583

584 **AUTHOR CONTRIBUTIONS**

585

586 P.J.B and T.R.G. conceived the project and designed experiments. P.J.B carried out experiments and data
587 analysis. P.J.B and T.R.G. wrote the manuscript.

588

589

590 **ACKNOWLEDGEMENTS**

591

592 The authors would like to thank Alexander Dobin, Felix Schlesinger, Chris Zaleski, Carrie Davis and all
593 other members of the Gingeras group at CSHL for their assistance and advice, as well as Richard
594 McCombie and the CSHL sequencing facility for their services. We also thank Alexander Gann, Gregory
595 Hannon, Zachary Lippman, Joshua Dubnau, Adrian Krainer, Brenton Graveley and Mike Levine for
596 helpful discussions and advice. We are grateful to Thomas Kaufman and the FlyBase team for their
597 permission to reproduce images. Work supported in part by the National Human Genome Research
598 Institute, modENCODE Project, contract U01HG004271.

599

600

601 **DATA AVAILABILITY**

602

603 The primary data for this study is available through the GEO database, under accession numbers
604 GSE36212 and GSE89335. GSE36212 is public, and reviewers can access GSE89335 here:
605 <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?token=onklcqedfohvkt&acc=GSE89335>

606

607

608 **ONLINE METHODS**

609

610 **Fly stocks & Embryo collections:**

611 All Drosophila strains were obtained from the Drosophila Species Stock Center at UC San Diego, CA
612 (<https://stockcenter.ucsd.edu/info/welcome.php>). For each species considered we worked with the
613 reference genome strain. Stock numbers: *D. melanogaster* #14021-0231.36, *D. simulans* #14021-
614 0251.195, *D. erecta* #14021-0224.01, *D. ananassae* #14024-0371.13, *D. pseudoobscura* #14011-
615 0121.94. Stocks were maintained on standard cornmeal medium. Embryo collections were performed in
616 population cages (Flystuff, #59-116). 2- to 7-day-old flies were left to acclimatize to the cage for at least
617 48h and regularly fed with grape juice-agar plates (Flystuff, #47-102) generously loaded with yeast paste.
618 After two 2-hour pre-lays, embryos were collected in 1-hour windows and aged appropriately (24 time
619 points, 0-24h). Embryos were washed with deionized water, dechorionated for 90 sec with 50% bleach,
620 rinsed abundantly with water, and snap-frozen in liquid nitrogen.

621

622 **RNA Extraction & RAMPAGE Library preparation:**

623 Total RNA was extracted from embryos using a Beadbeater (Biospec, Cat. #607) with 1.0 mm zirconia
624 beads (Biospec, #11079110zx) and the RNAdvance Tissue kit (Agencourt #A32649) according to the
625 manufacturer's instructions, including DNaseI treatment. We systematically checked on a Bioanalyzer
626 RNA Nano chip (Agilent) that the RNA was of very high quality. Libraries were prepared as described
627 before^{37,41}. 5'-monophosphate transcripts were depleted by TEX digest (Epicentre #TER51020). For
628 every time series, each sample was labeled with a different sequence barcode during reverse-
629 transcription, and all samples for the series were then pooled and processed together as a single library.
630 Quality control and library quantification were carried out on a Bioanalyzer DNA High Sensitivity chip.
631 Each library was sequenced on one lane of an Illumina HiSeq 2000.

632

633 **Genome references & annotations:**

634 All reference sequences and annotations were obtained from Flybase (<http://flybase.org>). *D.*
635 *melanogaster* release 5.49, *D. simulans* r1.4, *D. erecta* r1.3, *D. ananassae* r1.3, *D. pseudoobscura* r2.9.

636

637 **Primary data processing:**

638 Reads were mapped to the appropriate reference genomes using the STAR aligner⁶¹. Peaks were called
639 on the pooled data from whole time series, using a custom peak-caller described previously^{37,41}. We used
640 parameters optimized to yield good TSS specificity with respect to annotations and comparable numbers
641 of peaks for all species. All peaks overlapping FlyBase-annotated rDNA repeats were filtered out.

642

643 **TSC conservation:**

644 Functional conservation was assessed for all peaks with ≥ 15 RAMPAGE tags that did not map to
645 heterochromatic regions or chr4 in *D. melanogaster*, or orthologous regions in other species. We
646 translated the genomic coordinates of each peak in each species to coordinates in the multiple sequence
647 alignment of all genomes (15-way MultiZ alignment from UCSC,
648 <http://hgdownload.soe.ucsc.edu/goldenPath/dm3/multiz15way>). To be considered for analysis, each peak
649 was required to have a unique syntenic alignment in all other species considered, defined as follows: both
650 ends of an 800-bp window centered on the middle of the peak had to map to the same strand of the same
651 chromosome or scaffold, 50% of bases had to be aligned (*i.e.*, not in assembly gaps), and 25% of bases
652 had to align to orthologous bases (not alignment gaps). Raw 5' signal for each genome was also
653 translated into multiple alignment coordinates. For each peak from each species, functional conservation
654 was assessed by counting the number of RAMPAGE tags in each species. A peak was considered absent
655 in a target species if it had at least a 100-fold lower signal than in the reference species. Peaks with <100
656 tags in the reference species were considered absent if they had no detectable signal in a target species.

657

658 **Phylogeny reconstruction:**

659 The peaks from all species were merged and collapsed in multiple alignment space to generate a non-
660 redundant set of all peaks in the clade. The conservation of these peaks was assessed as described above.
661 The phylogenetic tree was inferred by treating the presence/absence of each peak as a 2-state discrete
662 character, sequentially using the MIX and PARS program of the PHYLIP suite according to the
663 recommendations of the software documentation (<http://evolution.genetics.washington.edu/phylip.html>).

664

665 **Sequence conservation:**

666 Per-base conservation scores were computed by running the phastCons and phyloP programs of the
667 PHAST suite v1.1 on the MultiZ alignment according to the recommendations of the software
668 documentation (<http://compgen.bscb.cornell.edu/phast>). Depending on the subclade of interest, some
669 species were excluded from the alignment for certain analyses. Pre-computed phastCons scores for the
670 full 15-way alignment were downloaded from UCSC
671 (<http://hgdownload.soe.ucsc.edu/goldenPath/dm3/phastCons15way>).

672

673 **Core promoter motifs:**

674 For analyses of motif composition, we only considered *D. melanogaster* TSCs that were functionally
675 conserved across all 5 species. We used pairwise chained alignments downloaded from UCSC
676 (<http://hgdownload.soe.ucsc.edu/downloads.html#fruitfly>) to align the most heavily used position of each
677 TSC (*i.e.*, the main TSS) to all other genomes. Peaks for which the maximum position could not be
678 aligned to all genomes were excluded from the analysis. A custom script was used to search for matches
679 to previously characterized core promoter motifs⁴⁴ within a 301-bp window centered on the main TSS.
680 Consensus sequences for sets of peaks with matches to individual motifs were computed using MEME
681 v4.9.0 (<http://meme.nbcr.net/meme>).

682 **Time series alignment:**

683 Z-score transformed gene expression time series from all species were registered to one another using the
684 GTEM suite⁴² according to the recommendations of the software documentation
685 (<http://flydev.berkeley.edu/cgi-bin/GTEM/index.htm>). One-to-one orthology calls from Flybase (2012
686 release 2) were used to match gene expression profiles between species. We pre-processed pairs of
687 datasets (*D. melanogaster* and another species) to compensate for differences in annotation quality and
688 peak calling between species. We identified orthologs of TSCs that had detectable expression (≥ 10 tags)
689 but initially failed to be called in one species. In addition, when a functionally conserved TSC had been
690 attributed to an annotated gene in one species but not the other, we corrected this discrepancy by
691 attributing it to the gene in both species. For the *D. ananassae* dataset, the 8th time point failed to yield
692 acceptable data, and was excluded from the analysis. All time series were upsampled 5-fold and
693 smoothed with a 2-hour window size using RZ-Smooth v4.1. Optimal global alignment paths between *D.*
694 *melanogaster* and the other datasets were computed with T-Warp v3.2 with Pearson distance matrices (3-
695 hour window). M-Align v2.8 was used to align each series to the *D. melanogaster* reference and smooth
696 the final aligned series (1-hour window). The expression profiles of individual TSCs were registered to
697 one another with M-Align, using the optimal alignment path computed for gene expression profiles. Prior
698 to alignment, we used the UCSC liftOver tool
699 (http://hgdownload.cse.ucsc.edu/admin/exe/linux.x86_64/liftOver) to identify *D. melanogaster* TSCs that
700 aligned well ($\geq 50\%$ of bases aligned) to all other genomes. The temporal expression profiles of those
701 orthologous genomic positions only were aligned.

702

703 **Expression profile conservation:**

704 We defined the conservation of individual expression profiles (TSCs or genes) across a clade as the
705 average Pearson R^2 for all pairwise comparisons of species within the clade.

706

707 **Clustering of TSC expression profiles:**

708 We classified all *D. melanogaster* TSCs with maximum expression ≥ 10 RPM (n=11,900) as either
709 Housekeeping (< 5 -fold variation throughout the time series, n=587) or Developmentally Regulated
710 ($\geq 60\%$ of total expression within an 8-hour window, n=6,015). We further selected TSCs functionally
711 conserved in all 5 species (n=240 and n=3,824, respectively). Developmentally regulated TSC profiles
712 were hierarchically clustered (R *hclust*, distance metric $1 - cor(t(expr), method="pearson")$) and initially
713 grouped into 12 clusters (R *cutree*, k=12). After filtering out excessively small clusters (< 200 TSCs,
714 n=4), further analysis was conducted on the remaining 8 regulated clusters (n=3,222 TSCs). See Fig. 3A
715 for clustering results.

716

717 **lncRNA transcript reconstruction & phyloCSF ORF analysis:**

718 We ran Cufflinks (v2.2.1) independently on each dataset of a published RNA-seq developmental series.
719 Cuffmerge was used to generate a consensus annotation set. Transcript models were attributed to a
720 RAMPAGE TSC if their 5' end lay within 150 of that TSC. Models without a matching TSC were
721 excluded from further analyses. phyloCSF was run on these annotation sets and the 15-way multiZ
722 whole-genome alignments.

723

724 **Analysis of sequence motifs:**

725 We used the MEME Suite v4.9.0 (primarily FIMO and MAST) to search promoter regions for a
726 previously published compendium of motifs⁶² corresponding to core promoter motifs and TFBSs.
727 MEME was used to generate consensus motifs from specific subregions of RAMPAGE-defined
728 promoters (Fig. 4).

729

730

731 **Other software:**

732 Custom analysis scripts were written in Python 2.7 (<http://www.python.org>). R was used for graphics
733 generation (<http://www.r-project.org>).

734

735 **FBgn0264479 promoter sequence analysis:**

736 We identified potential regulators of the FBgn0264479 promoter based on the BDTNP transcription
737 factor ChIP-chip data available from the UCSC genome browser website^{52,53}. For 10 factors that bind the
738 promoter in embryos (bcd, cad, D, dl, ftz, gt, h, hb, Kr, twi), we searched the 600bp upstream of the main
739 TSS for Jaspar TFBS motifs (<http://jaspar.genereg.net>), using FIMO. We identified 7 motifs at a p-value
740 cutoff of 2×10^{-4} .

741

742 **Northern-blot:**

743 We ran 8 μ g of total RNA per sample on an 8% acrylamide 8M urea gel with a Invitrogen Novex minigel
744 system. Transfer to a nylon membrane was carried out in 0.5X TBE in a Novex XCell II module,
745 followed by UV-crosslinking (1,200J). For detection, we used a combination of 6 oligonucleotide probes
746 targeting FBgn0264479 (5'-gaacatcgcttgcaagtgcag, 5'-cgatggatgtgtcggtcgg, 5'-ctctcgttctttgattctc, 5'-
747 caggatgtgtggtgtccac, 5'-agattggatccttatggttg, 5'-atatgctgacactgcatggt). 30 pmol of oligo mix were
748 radioactively labeled with γ^{32} -ATP and PNK. Following phenol-chloroform extraction, the labeled
749 probes were hybridized for 2 hours at 42°C in 40mL of ULTRAHyb buffer (ThermoFischer). After serial
750 washes with decreasing concentrations of SSC buffer (final stringency 0.5X), the membrane was exposed
751 on Kodak BioMax autoradiography film. After stripping and control re-exposure, a similar protocol was
752 used to detect the 5S rRNA on the same membrane, using a single probe (5'-caacacgcggtgtccaagccg).

753

754 **Fluorescent *in situ* hybridization:**

755 Templates for probe synthesis were generated by amplification of FBgn0264479 cDNAs with primers 5'-
756 CGATGTTCTCCGACCGACAA and 5'-TGCACTACTTAGACTAAATTGGCT. In separate reactions,
757 a T3 promoter sequence was added at either one end or the other, for the generation of sense and
758 antisense probes. Amplicons were cloned into a TOPO-T/A vector (Life Technologies #K4575-01) and
759 checked by Sanger sequencing. RNA-FISH was performed on 0-5 hours AEL *y; cn b sp* embryos as
760 described before⁶³, and imaged on a Perkin-Elmer UltraVIEW VoX confocal microscope. Biotin-
761 conjugated mouse monoclonal anti-DIG (Jackson ImmunoResearch Laboratories Inc., Cat. No. 200-062-
762 156) previously validated for this application⁶³.

763

764 **Sample size estimates:**

765 In this work, relevant comparisons are between groups of promoters or groups of genes. All comparisons
766 were designed to include all TSCs or genes of interest throughout the genome (*e.g.*, all lncRNA TSCs *vs.*
767 all genic TSCs), while applying expression level thresholds calibrated on the analysis of *D. melanogaster*
768 replicates to ensure measurement reproducibility.

769

770 **Code availability:**

771 RAMPAGE analysis pipeline and custom analysis scripts (Python & R) available upon request.

772

773 **REFERENCES**

774

775

- 776 1. Lewis, E.B. A gene complex controlling segmentation in *Drosophila*. *Nature* **276**, 565-70 (1978).
- 777 2. Nusslein-Volhard, C. & Wieschaus, E. Mutations affecting segment number and polarity in
- 778 *Drosophila*. *Nature* **287**, 795-801 (1980).
- 779 3. Levine, M. & Davidson, E.H. Gene regulatory networks for development. *Proc Natl Acad Sci U S*
- 780 *A* **102**, 4936-42 (2005).
- 781 4. Peter, I.S. & Davidson, E.H. Evolution of gene regulatory networks controlling body plan
- 782 development. *Cell* **144**, 970-85 (2011).
- 783 5. St Johnston, D. & Nusslein-Volhard, C. The origin of pattern and polarity in the *Drosophila*
- 784 embryo. *Cell* **68**, 201-19 (1992).
- 785 6. Levine, M. & Hoey, T. Homeobox proteins as sequence-specific transcription factors. *Cell* **55**,
- 786 537-40 (1988).
- 787 7. Hoey, T. & Levine, M. Divergent homeo box proteins recognize similar DNA sequences in
- 788 *Drosophila*. *Nature* **332**, 858-61 (1988).
- 789 8. Desplan, C., Theis, J. & O'Farrell, P.H. The sequence specificity of homeodomain-DNA
- 790 interaction. *Cell* **54**, 1081-90 (1988).
- 791 9. Segal, E., Raveh-Sadka, T., Schroeder, M., Unnerstall, U. & Gaul, U. Predicting expression
- 792 patterns from regulatory sequence in *Drosophila* segmentation. *Nature* **451**, 535-40 (2008).
- 793 10. Banerji, J., Rusconi, S. & Schaffner, W. Expression of a beta-globin gene is enhanced by remote
- 794 SV40 DNA sequences. *Cell* **27**, 299-308 (1981).
- 795 11. Banerji, J., Olson, L. & Schaffner, W. A lymphocyte-specific cellular enhancer is located
- 796 downstream of the joining region in immunoglobulin heavy chain genes. *Cell* **33**, 729-40 (1983).
- 797 12. Zinn, K., DiMaio, D. & Maniatis, T. Identification of two distinct regulatory regions adjacent to
- 798 the human beta-interferon gene. *Cell* **34**, 865-79 (1983).
- 799 13. Small, S., Kraut, R., Hoey, T., Warrior, R. & Levine, M. Transcriptional regulation of a pair-rule
- 800 stripe in *Drosophila*. *Genes Dev* **5**, 827-39 (1991).
- 801 14. Spitz, F. & Furlong, E.E. Transcription factors: from enhancer binding to developmental control.
- 802 *Nat Rev Genet* **13**, 613-26 (2012).
- 803 15. Schaffner, W. Enhancers, enhancers - from their discovery to today's universe of transcription
- 804 enhancers. *Biol Chem* **396**, 311-27 (2015).
- 805 16. Benoist, C. & Chambon, P. Deletions covering the putative promoter region of early mRNAs of
- 806 simian virus 40 do not abolish T-antigen expression. *Proc Natl Acad Sci U S A* **77**, 3865-9 (1980).
- 807 17. Lenhard, B., Sandelin, A. & Carninci, P. Metazoan promoters: emerging characteristics and
- 808 insights into transcriptional regulation. *Nat Rev Genet* **13**, 233-45 (2012).
- 809 18. Zabidi, M.A. *et al.* Enhancer-core-promoter specificity separates developmental and
- 810 housekeeping gene regulation. *Nature* **518**, 556-9 (2015).
- 811 19. Lettice, L.A. *et al.* A long-range *Shh* enhancer regulates expression in the developing limb and fin
- 812 and is associated with preaxial polydactyly. *Hum Mol Genet* **12**, 1725-35 (2003).
- 813 20. Goodrich, J.A. & Tjian, R. Unexpected roles for core promoter recognition factors in cell-type-
- 814 specific transcription and gene regulation. *Nature Reviews Genetics* **11**, 549-558 (2010).
- 815 21. Gerstein, M.B. *et al.* What is a gene, post-ENCODE? History and updated definition. *Genome*
- 816 *Res* **17**, 669-81 (2007).
- 817 22. Gingeras, T.R. Origin of phenotypes: genes and transcripts. *Genome Res* **17**, 682-90 (2007).
- 818 23. Kapranov, P. *et al.* RNA maps reveal new RNA classes and a possible function for pervasive
- 819 transcription. *Science* **316**, 1484-1488 (2007).
- 820 24. Djebali, S. *et al.* Landscape of transcription in human cells. *Nature* **489**, 101-8 (2012).

- 821 25. Graveley, B.R. *et al.* The developmental transcriptome of *Drosophila melanogaster*. *Nature*
822 (2010).
- 823 26. Gerstein, M.B. *et al.* Integrative analysis of the *Caenorhabditis elegans* genome by the
824 modENCODE project. *Science* **330**, 1775-87 (2010).
- 825 27. Young, R.S. *et al.* Identification and properties of 1,119 candidate lincRNA loci in the *Drosophila*
826 *melanogaster* genome. *Genome Biol Evol* **4**, 427-42 (2012).
- 827 28. Derrien, T. *et al.* The GENCODE v7 catalog of human long noncoding RNAs: analysis of their
828 gene structure, evolution, and expression. *Genome Res* **22**, 1775-89 (2012).
- 829 29. Augui, S., Nora, E.P. & Heard, E. Regulation of X-chromosome inactivation by the X-
830 inactivation centre. *Nat Rev Genet* **12**, 429-42 (2011).
- 831 30. Ulitsky, I. & Bartel, D.P. lincRNAs: genomics, evolution, and mechanisms. *Cell* **154**, 26-46
832 (2013).
- 833 31. Guttman, M. & Rinn, J.L. Modular regulatory principles of large non-coding RNAs. *Nature* **482**,
834 339-46 (2012).
- 835 32. Ponting, C.P., Oliver, P.L. & Reik, W. Evolution and functions of long noncoding RNAs. *Cell*
836 **136**, 629-41 (2009).
- 837 33. Tsankov, A., Yanagisawa, Y., Rhind, N., Regev, A. & Rando, O.J. Evolutionary divergence of
838 intrinsic and trans-regulated nucleosome positioning sequences reveals plastic rules for chromatin
839 organization. *Genome Res* **21**, 1851-62 (2011).
- 840 34. Schmidt, D. *et al.* Five-vertebrate ChIP-seq reveals the evolutionary dynamics of transcription
841 factor binding. *Science* **328**, 1036-40 (2010).
- 842 35. Stefflova, K. *et al.* Cooperativity and rapid evolution of cobound transcription factors in closely
843 related mammals. *Cell* **154**, 530-40 (2013).
- 844 36. He, Q. *et al.* High conservation of transcription factor binding and evidence for combinatorial
845 regulation across six *Drosophila* species. *Nat Genet* **43**, 414-20 (2011).
- 846 37. Batut, P., Dobin, A., Plessy, C., Carninci, P. & Gingeras, T.R. High-fidelity promoter profiling
847 reveals widespread alternative promoter usage and transposon-driven developmental gene
848 expression. *Genome Res* (2012).
- 849 38. Carninci, P. *et al.* Genome-wide analysis of mammalian promoter architecture and evolution.
850 *Nature Genetics* **38**, 626-635 (2006).
- 851 39. Valen, E. *et al.* Genome-wide detection and analysis of hippocampus core promoters using
852 DeepCAGE. *Genome Research* **19**, 255-265 (2009).
- 853 40. Hoskins, R.A. *et al.* Genome-wide analysis of promoter architecture in *Drosophila melanogaster*.
854 *Genome Res* (2011).
- 855 41. Batut, P. & Gingeras, T.R. RAMPAGE: Promoter Activity Profiling by Paired-End Sequencing
856 of 5'-Complete cDNAs. *Curr Protoc Mol Biol* **104**, 25B 11 1-25B 11 16 (2013).
- 857 42. Goltsev, Y. & Papatsenko, D. Time warping of evolutionary distant temporal gene expression
858 data based on noise suppression. *Bmc Bioinformatics* **10**, - (2009).
- 859 43. Kalinka, A.T. *et al.* Gene expression divergence recapitulates the developmental hourglass model.
860 *Nature* **468**, 811-U102 (2010).
- 861 44. FitzGerald, P.C., Sturgill, D., Shyakhtenko, A., Oliver, B. & Vinson, C. Comparative genomics
862 of *Drosophila* and human core promoters. *Genome Biol* **7**, R53 (2006).
- 863 45. Odom, D.T. *et al.* Tissue-specific transcriptional regulation has diverged significantly between
864 human and mouse. *Nature Genetics* **39**, 730-732 (2007).
- 865 46. Villar, D. *et al.* Enhancer evolution across 20 mammalian species. *Cell* **160**, 554-66 (2015).
- 866 47. Haerty, W. & Ponting, C.P. Mutations within lincRNAs are effectively selected against in fruitfly
867 but not in human. *Genome Biol* **14**, R49 (2013).
- 868 48. Thomas, S. *et al.* Dynamic reprogramming of chromatin accessibility during *Drosophila* embryo
869 development. *Genome Biol* **12**, R43 (2011).

- 870 49. Lipshitz, H.D., Peattie, D.A. & Hogness, D.S. Novel transcripts from the Ultrabithorax domain of
871 the bithorax complex. *Genes Dev* **1**, 307-22 (1987).
- 872 50. Petruk, S. *et al.* Transcription of bxd noncoding RNAs promoted by trithorax represses Ubx in cis
873 by transcriptional interference. *Cell* **127**, 1209-21 (2006).
- 874 51. Franke, A. & Baker, B.S. The rox1 and rox2 RNAs are essential components of the
875 compensasome, which mediates dosage compensation in *Drosophila*. *Mol Cell* **4**, 117-22 (1999).
- 876 52. Li, X.Y. *et al.* Transcription factors bind thousands of active and inactive regions in the
877 *Drosophila* blastoderm. *PLoS Biol* **6**, e27 (2008).
- 878 53. MacArthur, S. *et al.* Developmental roles of 21 *Drosophila* transcription factors are determined
879 by quantitative differences in binding to an overlapping set of thousands of genomic regions.
880 *Genome Biol* **10**, R80 (2009).
- 881 54. Arora, K. *et al.* The *Drosophila* schnurri gene acts in the Dpp/TGF beta signaling pathway and
882 encodes a transcription factor homologous to the human MBP family. *Cell* **81**, 781-90 (1995).
- 883 55. Kutter, C. *et al.* Rapid turnover of long noncoding RNAs and the evolution of gene expression.
884 *PLoS Genet* **8**, e1002841 (2012).
- 885 56. Dai, H. *et al.* The zinc finger protein schnurri acts as a Smad partner in mediating the
886 transcriptional response to decapentaplegic. *Dev Biol* **227**, 373-87 (2000).
- 887 57. Bokel, C., Schwabedissen, A., Entchev, E., Renaud, O. & Gonzalez-Gaitan, M. Sara endosomes
888 and the maintenance of Dpp signaling levels across mitosis. *Science* **314**, 1135-9 (2006).
- 889 58. Coumailleau, F., Furthauer, M., Knoblich, J.A. & Gonzalez-Gaitan, M. Directional Delta and
890 Notch trafficking in Sara endosomes during asymmetric cell division. *Nature* **458**, 1051-5 (2009).
- 891 59. Gonzalez-Gaitan, M. Signal dispersal and transduction through the endocytic pathway. *Nat Rev*
892 *Mol Cell Biol* **4**, 213-24 (2003).
- 893 60. Fabrowski, P. *et al.* Tubular endocytosis drives remodelling of the apical surface during epithelial
894 morphogenesis in *Drosophila*. *Nat Commun* **4**, 2244 (2013).
- 895 61. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21 (2012).
- 896 62. Stark, A. *et al.* Discovery of functional elements in 12 *Drosophila* genomes using evolutionary
897 signatures. *Nature* **450**, 219-32 (2007).
- 898 63. Legendre, F. *et al.* Whole mount RNA fluorescent in situ hybridization of *Drosophila* embryos. *J*
899 *Vis Exp*, e50057 (2013).
- 900

SUPPLEMENTARY FIGURES

CORE PROMOTER STRUCTURE AND NONCODING TRANSCRIPTION SHAPE EARLY DEVELOPMENT IN DROSOPHILA

PHILIPPE J. BATUT & THOMAS R. GINGERAS

List of Figures

1	Distribution of raw RAMPAGE signal over transcript annotations	2
2	Distribution of RAMPAGE peaks over transcript annotations	3
3	Reproducibility of expression time series	4
4	Time series alignment by time-warping of gene expression profiles	5
5	Evolutionary conservation of TSC expression profiles	6
6	Core promoter types: Co-occurrence of TFBSs & Sequence conservation profiles	7
7	Conservation of transcription factor binding sites	8
8	Alternative analyses of TSC conservation	9
9	TSC conservation by expression quantiles	10
10	Phylogeny of sequenced species	11
11	Evolutionary rates of gain and loss for TSCs and Twist TFBSs	12
12	DNase I hypersensitivity at RAMPAGE TSCs	13
13	Independence and protein-coding potential of putative lncRNAs	14
14	Clustering of <i>D. melanogaster</i> developmental expression profiles	15
15	<i>Bithoraxoid</i> locus	16
16	Conservation of <i>melanogaster</i> subgroup lncRNA TSCs	17
17	Sequence conservation over <i>melanogaster</i> subgroup-specific lncRNA TSC	18
18	<i>FBgn0264479</i> protein-coding potential	19
19	<i>FBgn0264479</i> locus and expression	20
20	<i>FBgn0264479</i> Northern-blot	21
21	<i>FBgn0264479</i> transcriptional regulation	22

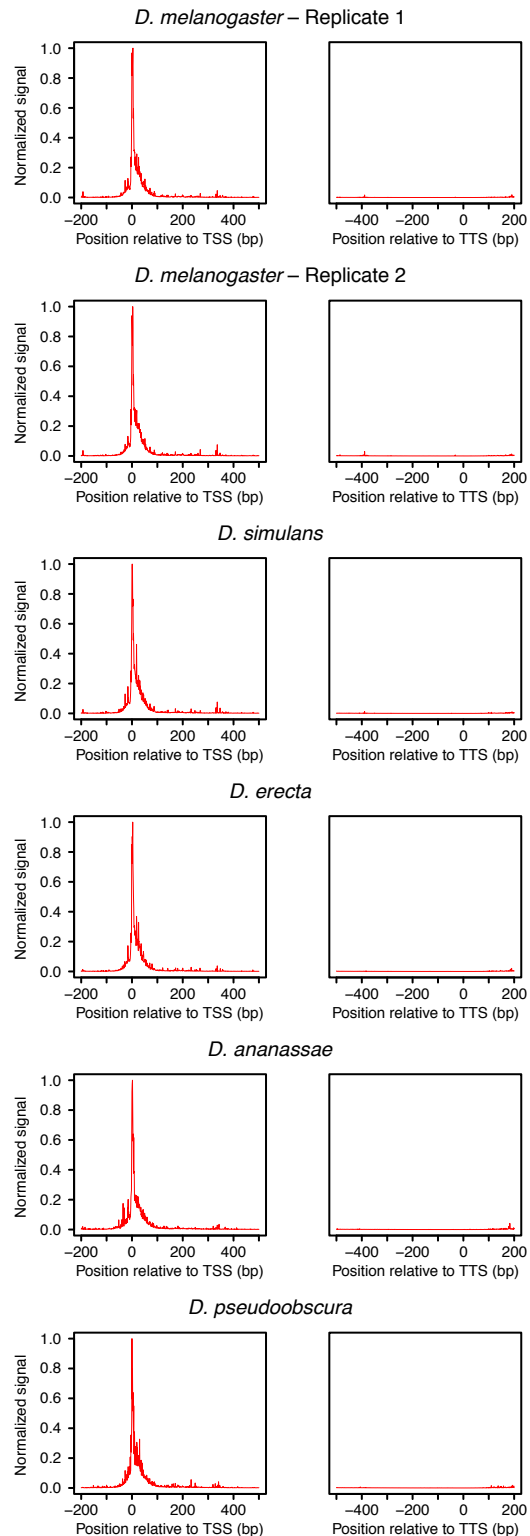


Figure 1. Distribution of raw RAMPAGE signal over transcript annotations. For each species, RAMPAGE reads were mapped to the appropriate genome. The raw 5' signal was then converted to orthologous *D. melanogaster* coordinates using chained pairwise alignments from UCSC. Metaprofiles were constructed by summing signal intensity over FlyBase r5.49 mRNA annotations.

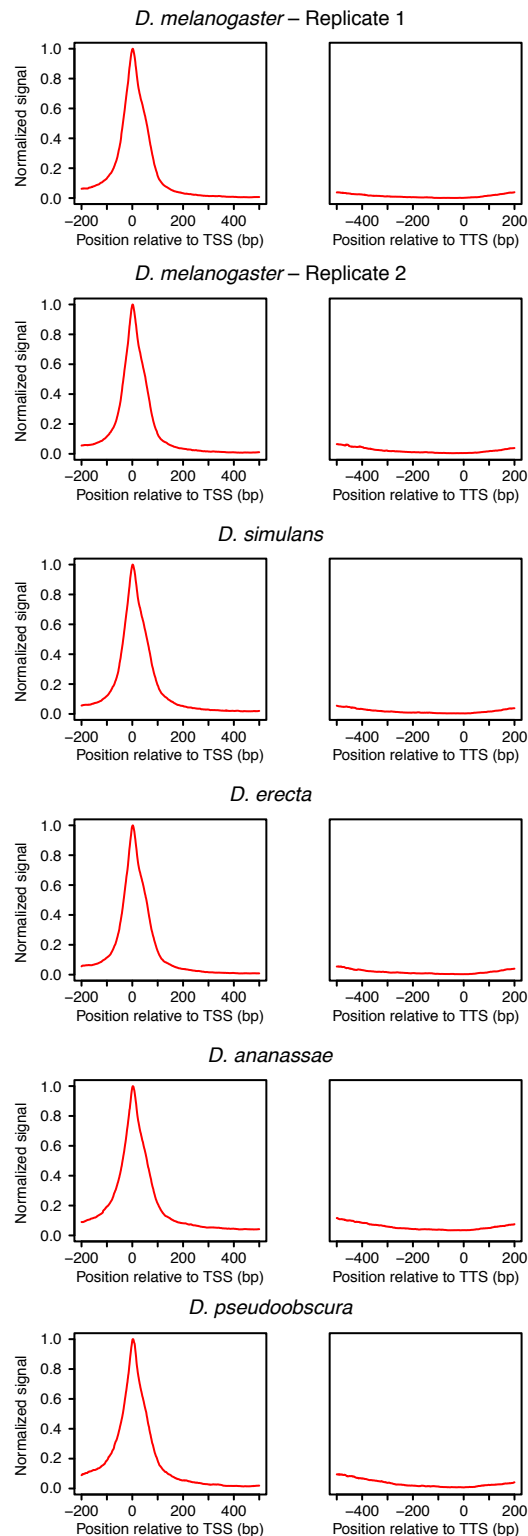


Figure 2. Distribution of RAMPAGE peaks over transcript annotations

For each species, RAMPAGE reads were mapped to the appropriate genome and peaks called as described in Methods. The peak coordinates were then converted to orthologous *D. melanogaster* coordinates using chained pairwise alignments and the liftOver tool from the UCSC Genome Browser. Metaprofiles were constructed by summing signal intensity over FlyBase r5.49 mRNA annotations.

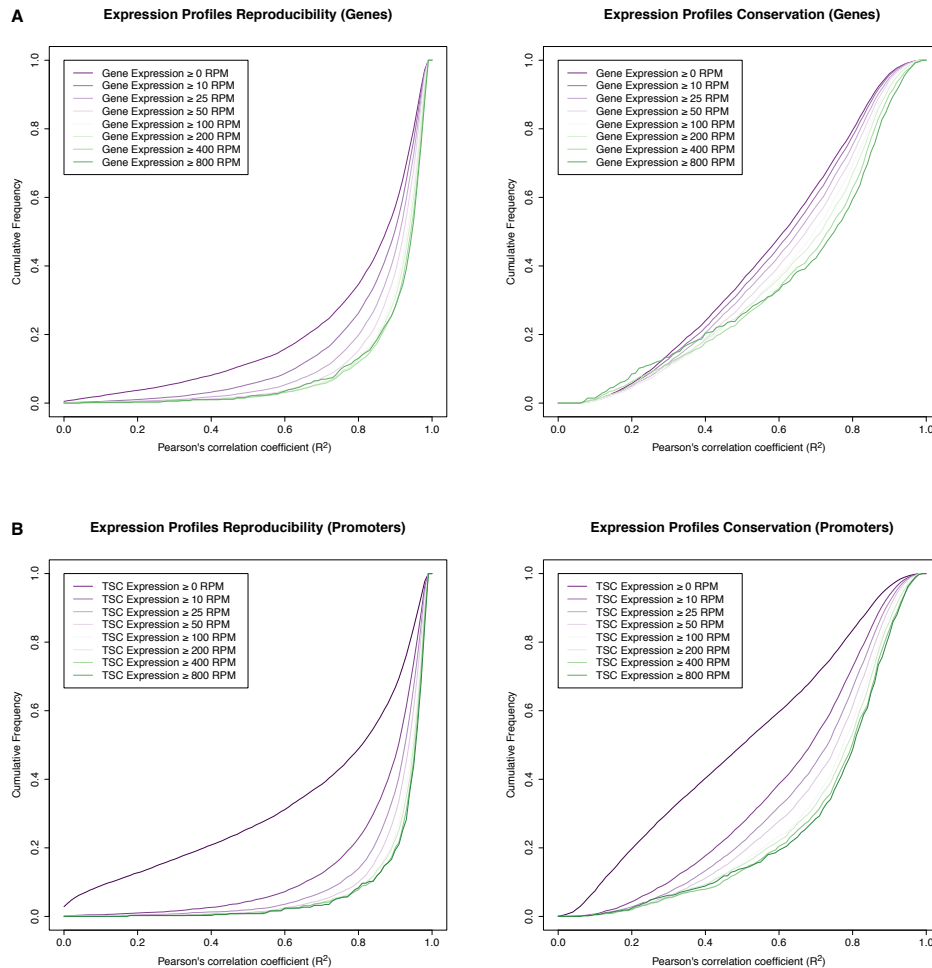


Figure 3. Reproducibility of expression time series

Reproducibility of expression profile measurements for genes (A) and individual TSCs (B). The graphs represent cumulative distributions of expression profile correlations across *D. melanogaster* biological replicates (left), and across all species (right). We only considered genes whose maximum expression level throughout the *D. melanogaster* time series (replicate 1) exceeded a given threshold (see legends; RPM: reads per million). Note that the variation across species vastly exceeds the variation across replicates, at all expression thresholds considered.

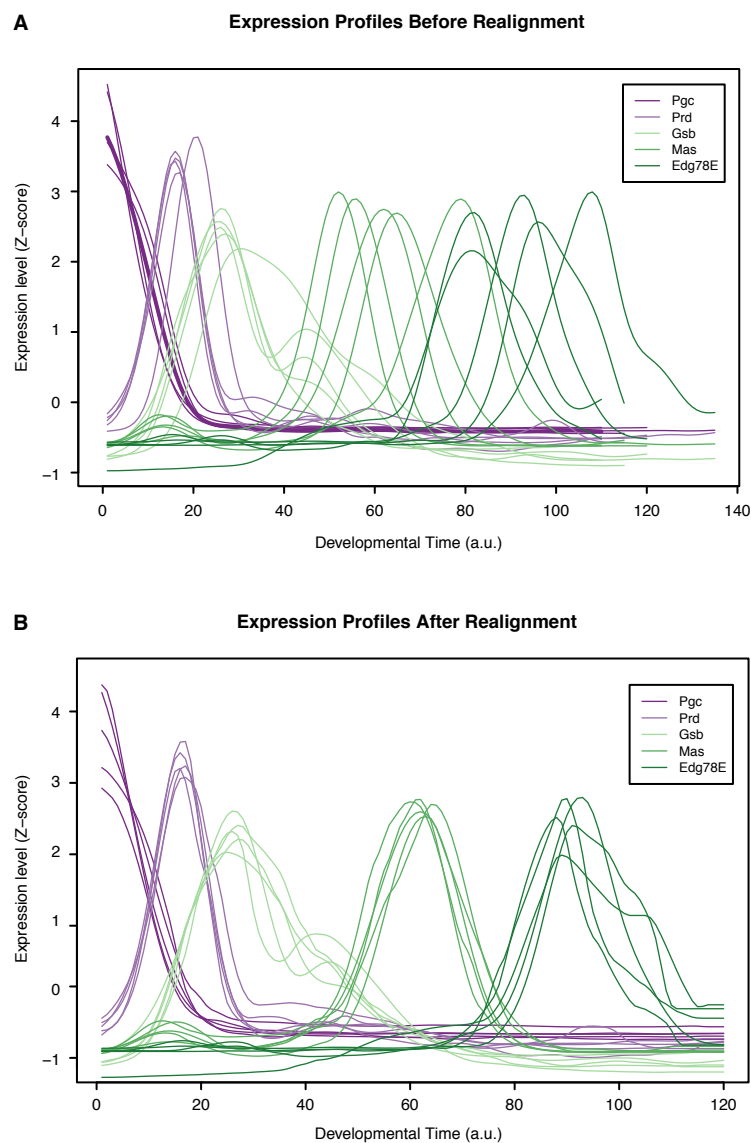


Figure 4. Time series alignment by time-warping of gene expression profiles
Global gene expression profiles from all species were aligned to the *D. melanogaster* time series as described in Methods. This figure shows the expression profiles for well-characterized developmental regulators before (A) and after (B) alignment. The time scale corresponds to the absolute *D. melanogaster* developmental time (24 hours divided into 120 units by upsampling).

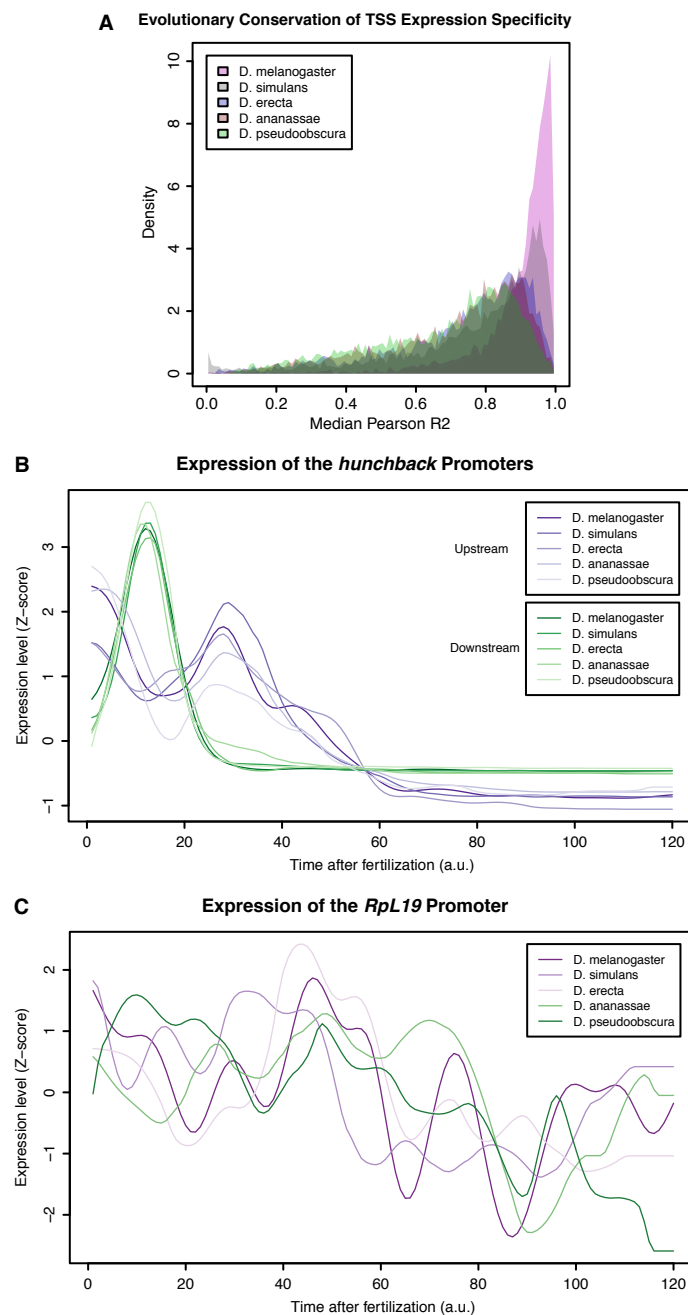


Figure 5. Evolutionary conservation of TSC expression profiles

(A) Distribution of average correlation coefficients for all orthologous TSCs between pairs of species. (B) Aligned expression profiles for the 2 promoters of the *hunchback* gene. (C) Aligned expression profiles for the *RpL19* gene.

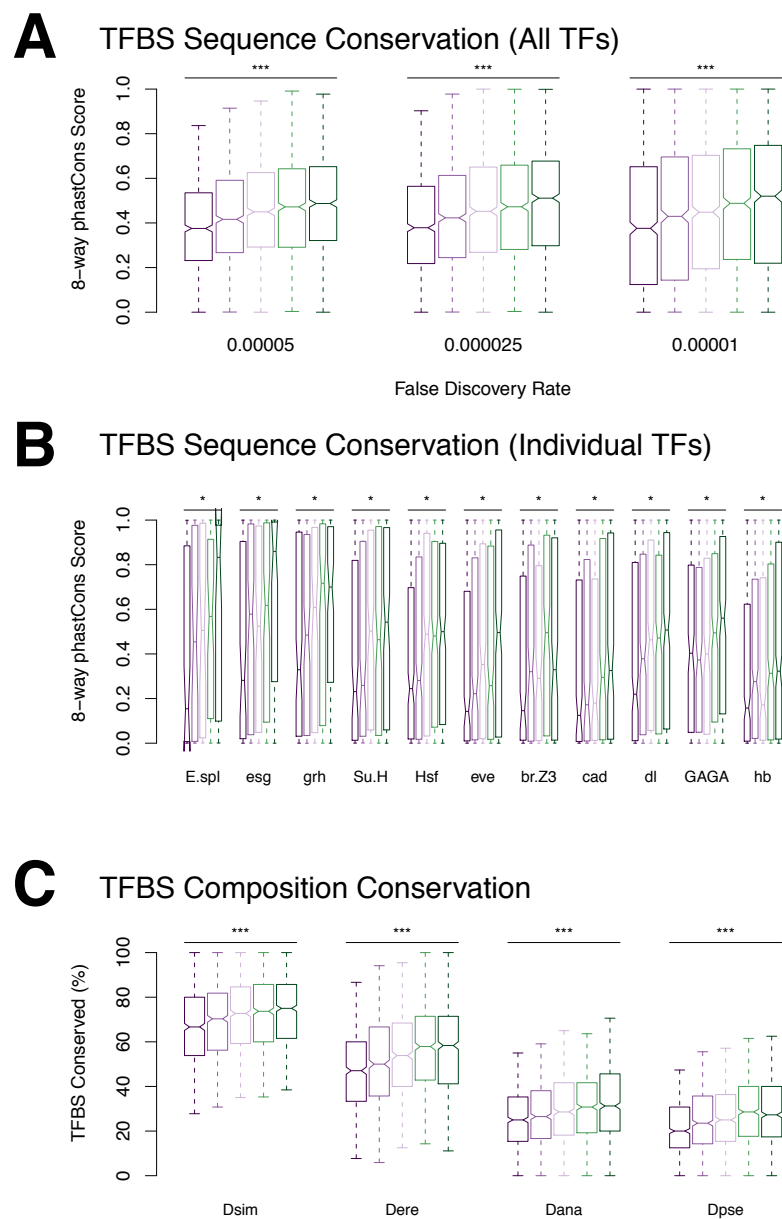


Figure 7. Conservation of transcription factor binding sites

(A) TFBS sequence conservation versus profile conservation. (B) Conservation of individual motif types. We found a significant correlation between motif sequence conservation and promoter profile conservation for 20 motif types (Bonferroni-corrected p-value < 0.01). (C) Conservation of promoter TFBS composition, regardless of TFBS position.

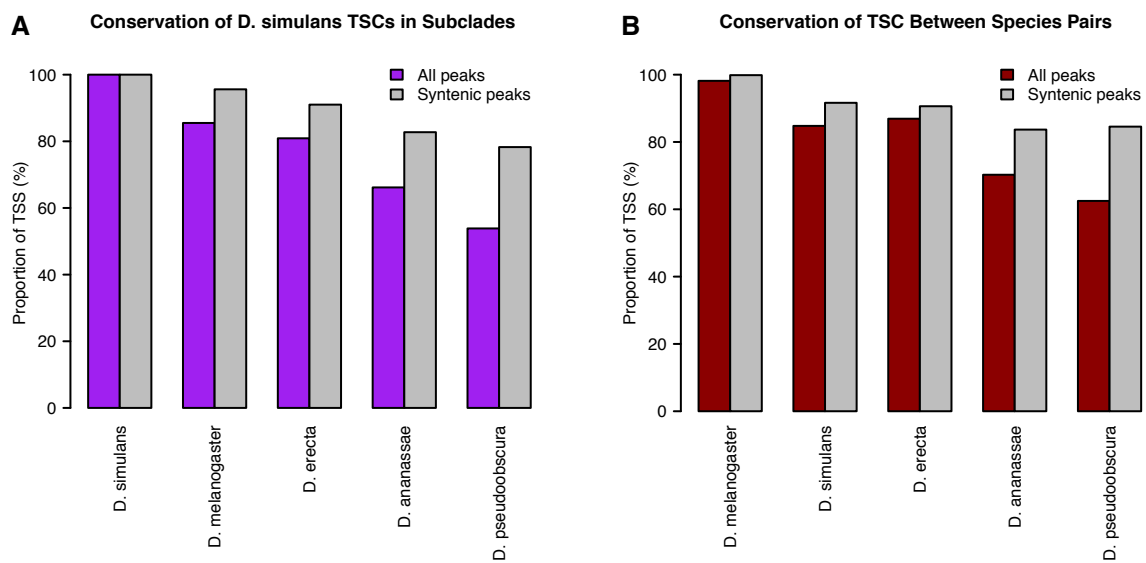


Figure 8. Alternative analyses of TSC conservation

TSC conservation was quantified as described in Methods. (A) *D. simulans*-centric analysis. (B) Quantification of TSC conservation between species pairs, as opposed to subclades.

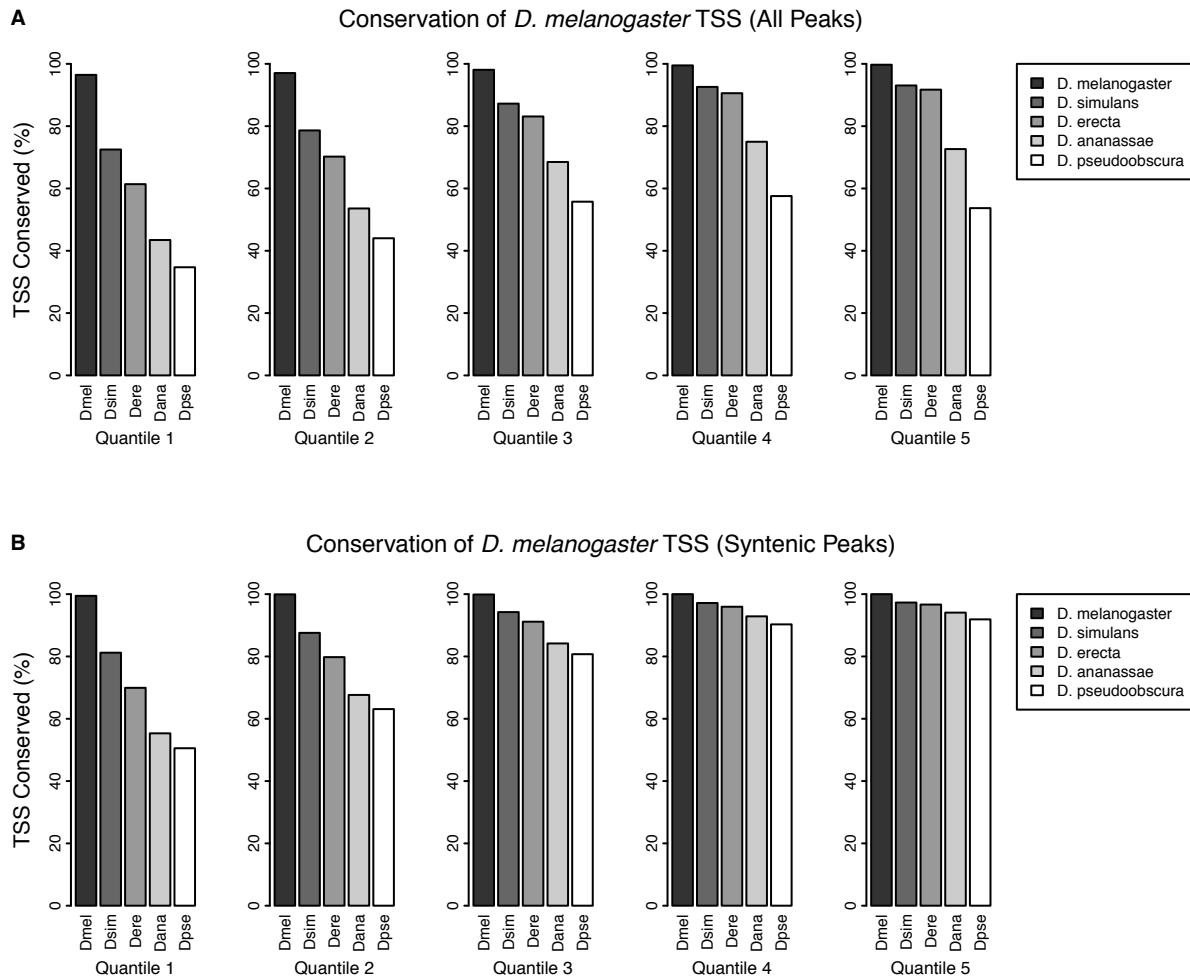


Figure 9. TSC conservation by expression quantiles

D. melanogaster TSCs were categorized into 5 expression quantiles based on total raw signal for the full time series. Functional conservation was assessed as described in Methods. (A) Conservation of all TSCs. (B) Conservation of TSCs with syntenic alignments in all 5 species.

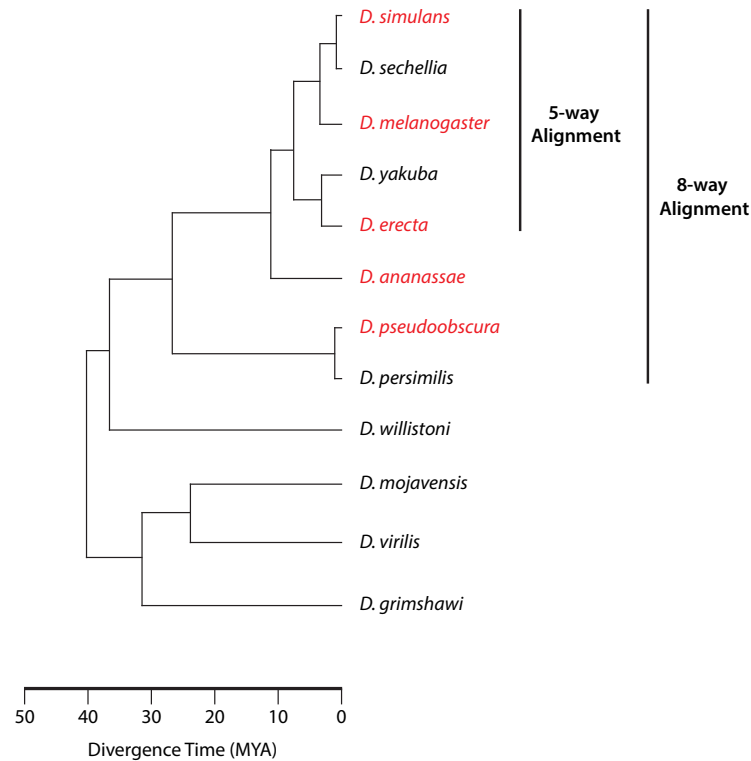


Figure 10. Phylogeny of sequenced species

The species for which we gathered data appear in red. When assessing sequence conservation for features conserved across all 5 species studied, we included the genome sequences of the 5 species studied and the 3 additional sequenced species from the same monophyletic group (8-way alignment). When assessing conservation throughout the *melanogaster* subgroup, we used the genomes of our 3 species and the other 2 from the subgroup (5-way alignment).

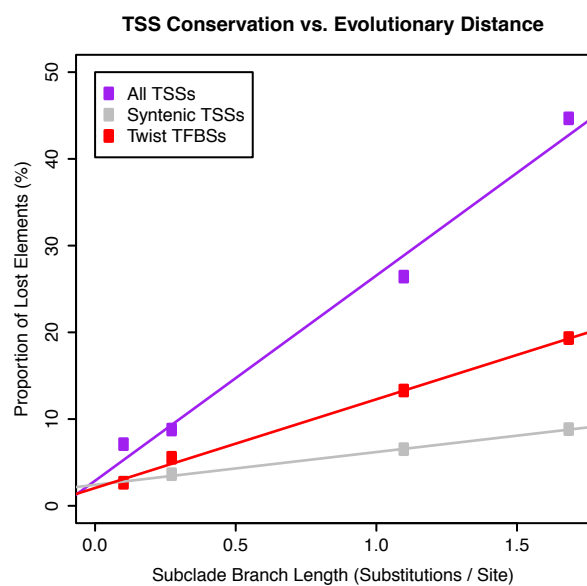


Figure 11. Evolutionary rates of gain and loss for TSCs and Twist TFBSs
Twist TFBS data from <Twist paper>.

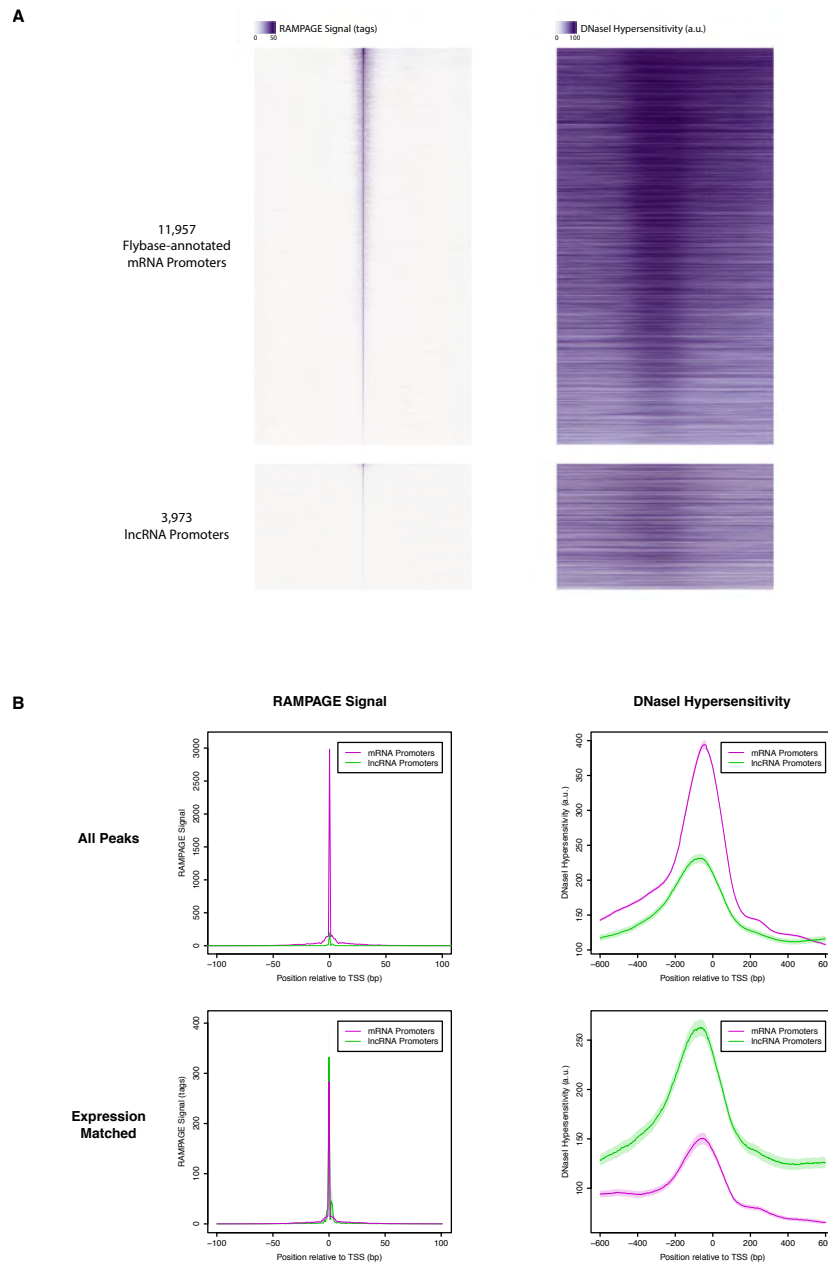


Figure 12. DNase I hypersensitivity at RAMPAGE TSCs

(A) Heatmaps of RAMPAGE signal (left) and DNase-seq signal (right) over individual TSCs. We are comparing TSCs that overlap FlyBase-annotated mRNA transcription start sites (top), which we use as positive controls, to the TSCs of putative lncRNAs (bottom). In each group, TSCs are sorted by total RAMPAGE signal intensity. (B) Class-wise average profiles of RAMPAGE signal (left) and DNase-seq signal (right). The lines represent median profiles, the shaded areas cover +/- 1 standard deviation as estimated by bootstrapping. When considering all peaks (top), lncRNA promoters show weaker DNase sensitivity than FlyBase-annotated controls, but the latter also have considerably stronger RAMPAGE signal. When matching RAMPAGE signal distributions (bottom), this trend is reversed. Shaded areas represent +/-1 standard deviation, as estimated by downsampling and bootstrapping.

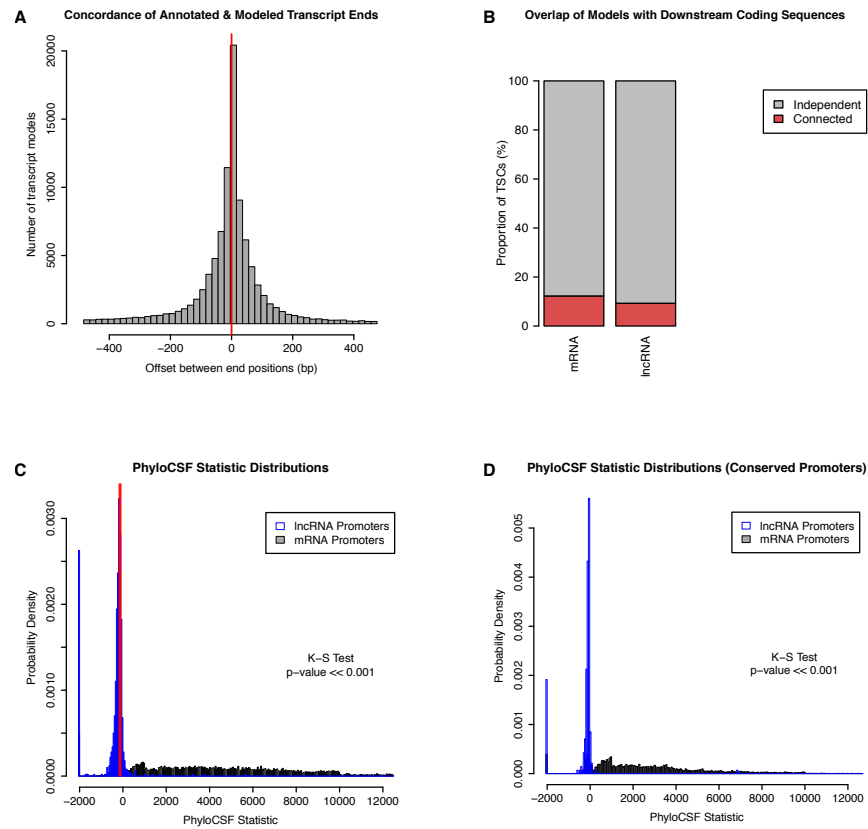
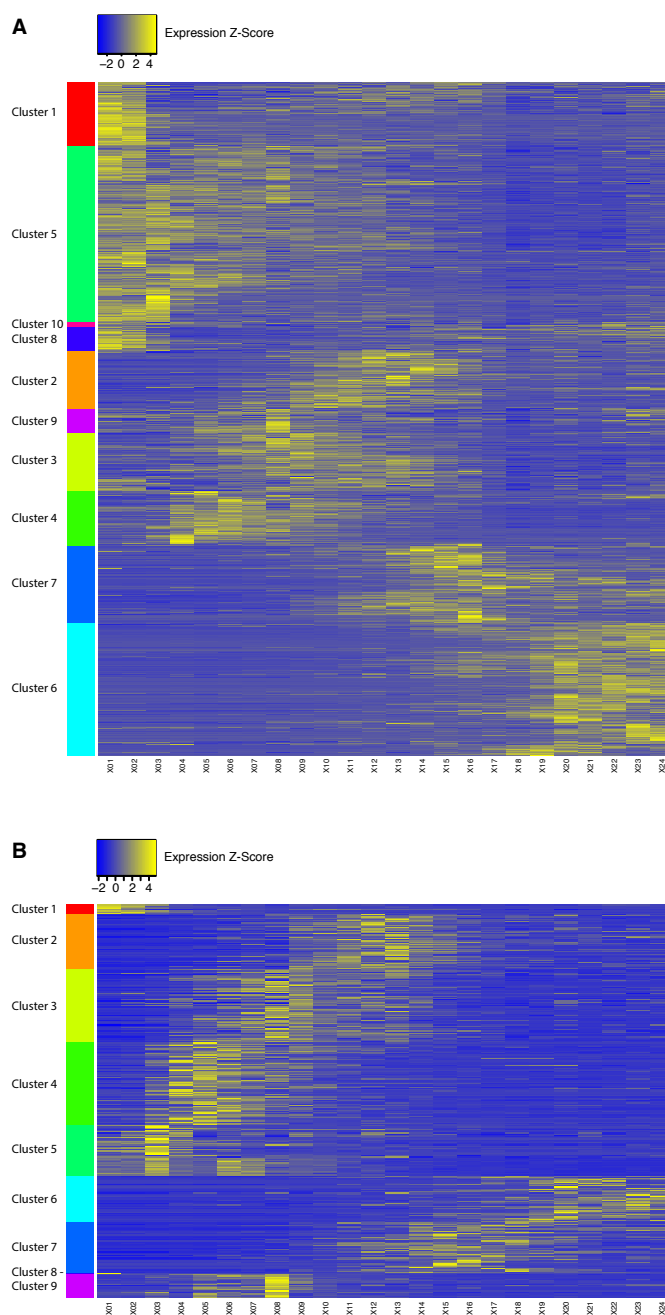


Figure 13. Independence and protein-coding potential of putative lncRNAs

We used published *D. melanogaster* embryo RNA-seq data and Cufflinks to generate transcript models, and only considered those starting within 150bp of a TSC. (A) Concordance of annotated and Cufflinks-modeled 3' ends. For each Cufflinks model starting at a protein-coding gene TSC, we are representing the distance to the closest annotated transcript 3' end. (B) Fraction of TSCs for which at least one transcript model overlaps a downstream annotated CDS. Note that while a minority of lncRNA TSC transcript models overlap a downstream protein-coding genes, we observed a similar propensity of Cufflinks models to fuse together consecutive protein-coding genes. (C) PhyloCSF score distributions of protein-coding and putative lncRNA transcript models, as a measure of protein-coding potential. For each TSC, we only considered the transcript model with the highest score. Transcripts with no ORF were assigned a default score of -2,000. Red lines: scores for known lncRNAs (roX1, roX2, bithoraxoid). (D) PhyloCSF analysis restricted to TSCs that are functionally conserved across all 5 species.



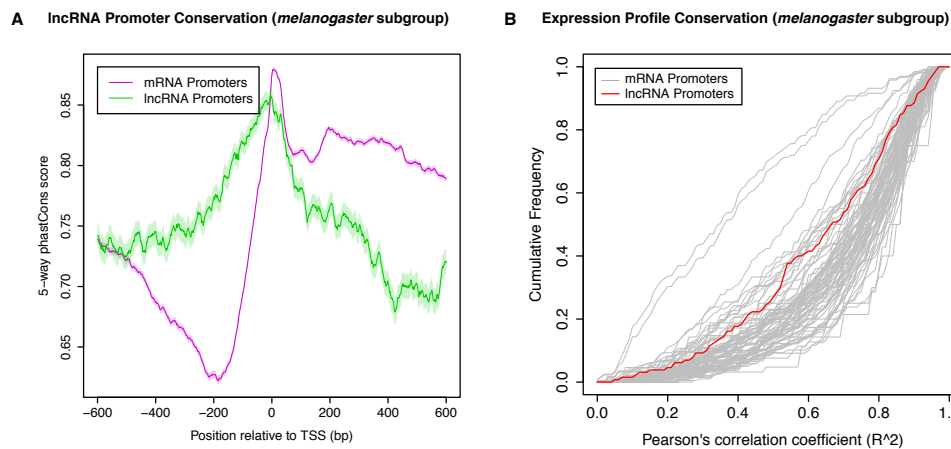
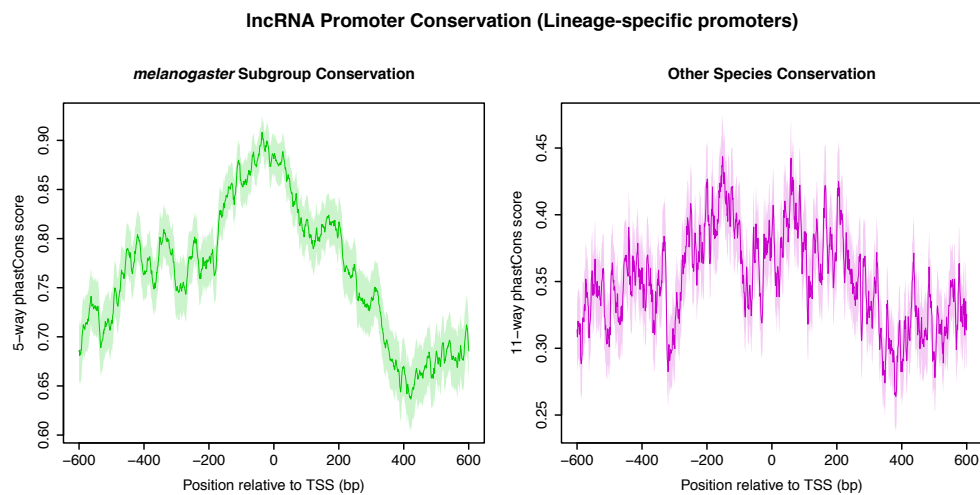


Figure 16. Conservation of *melanogaster* subgroup lncRNA TSCs

(A) Promoter sequence conservation. *melanogaster* subgroup-specific phastCons scores over lncRNA or protein-coding gene promoters that are shared across the 3 species of the subgroup. These phastCons scores were computed by including exclusively the genomes of the 5 sequenced *melanogaster* subgroup species in the input multiple sequence alignment. (B) Expression specificity conservation. All functionally shared promoters with maximum ≥ 25 RPM and ≥ 5 -fold expression changes in *D. melanogaster* were included in the analysis (130 lncRNA promoters).



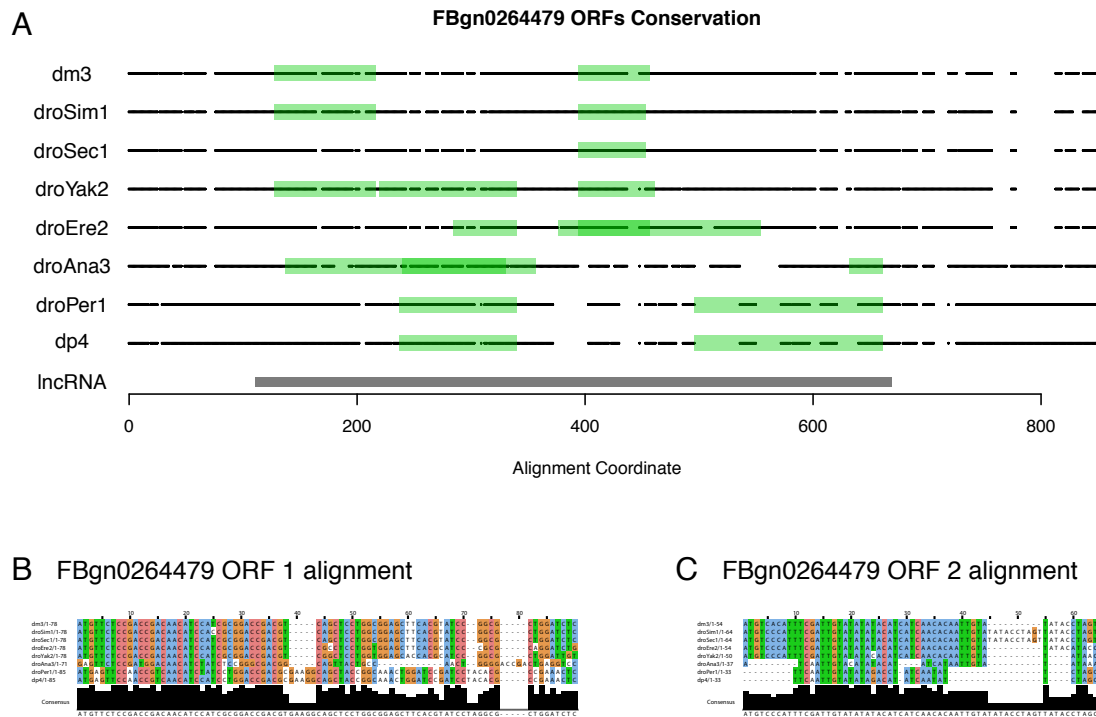
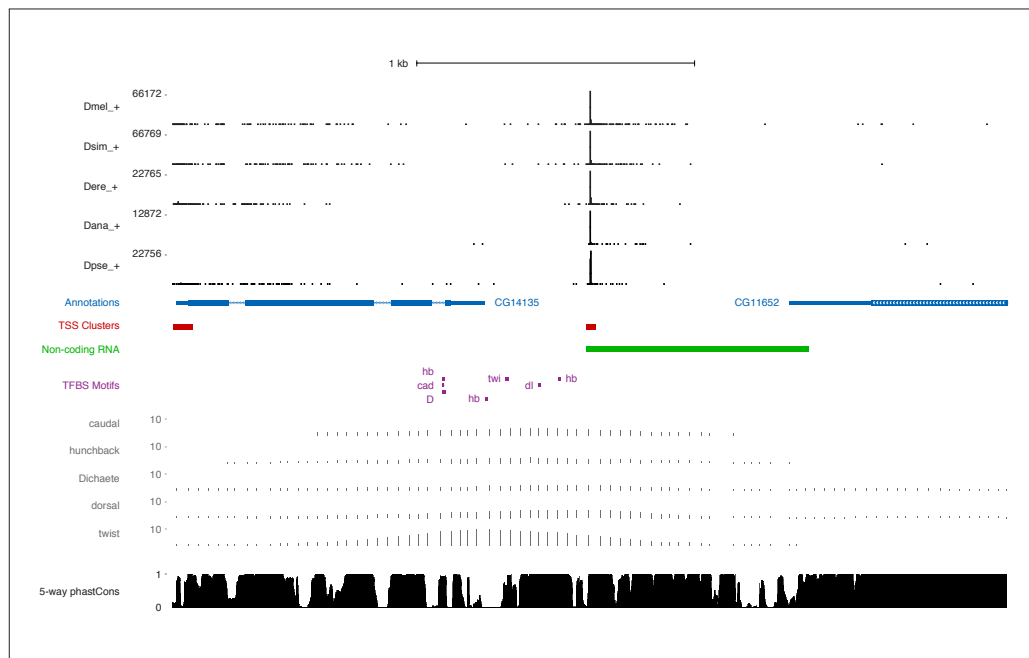


Figure 18. *FBgn0264479* protein-coding potential

(A) Projection of the putative ORFs in all orthologs onto the multiZ multiple sequence alignment (UCSC). Black points: bases aligned to *D. melanogaster* (match or mismatch). White spaces: alignment gaps. lncRNA: *D. melanogaster* transcript. Note the absence of any ORF conserved in all species that express the *FBgn0264479* transcript. (B) Multiple sequence alignment for *D. melanogaster* ORF #1. (C) Multiple sequence alignment for *D. melanogaster* ORF #2.



FBgn0264479 developmental expression profiles

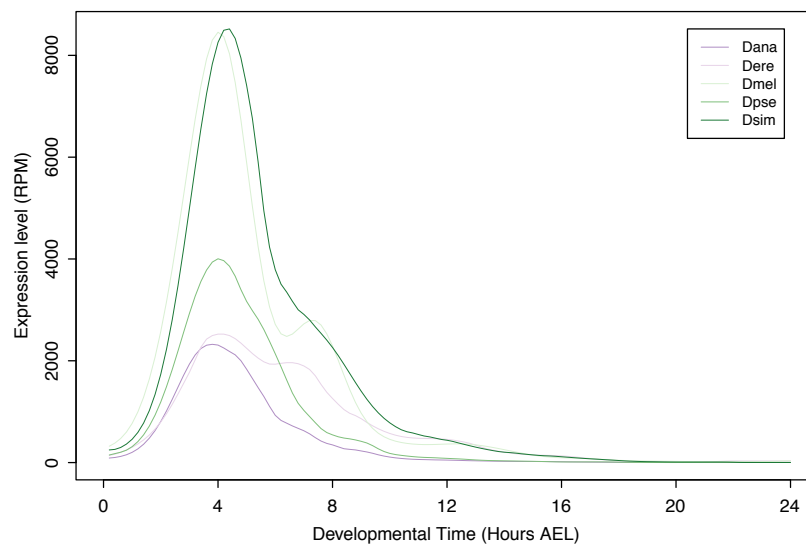


Figure 19. *FBgn0264479* locus and expression

Top: *FBgn0264479* locus organization. Includes all data tracks from Fig. 6A. In addition, the 5 tracks just above the phastCons scores represent the chromatin immunoprecipitation $\hat{\text{A}}$ microarray data (ChIP-chip) for 5 transcription factors: caudal, hunchback, Dichaete, dorsal, twist (modENCODE data visualized on the UCSC Genome Browser). Bottom: Expression profiles in reads per million (RPM).

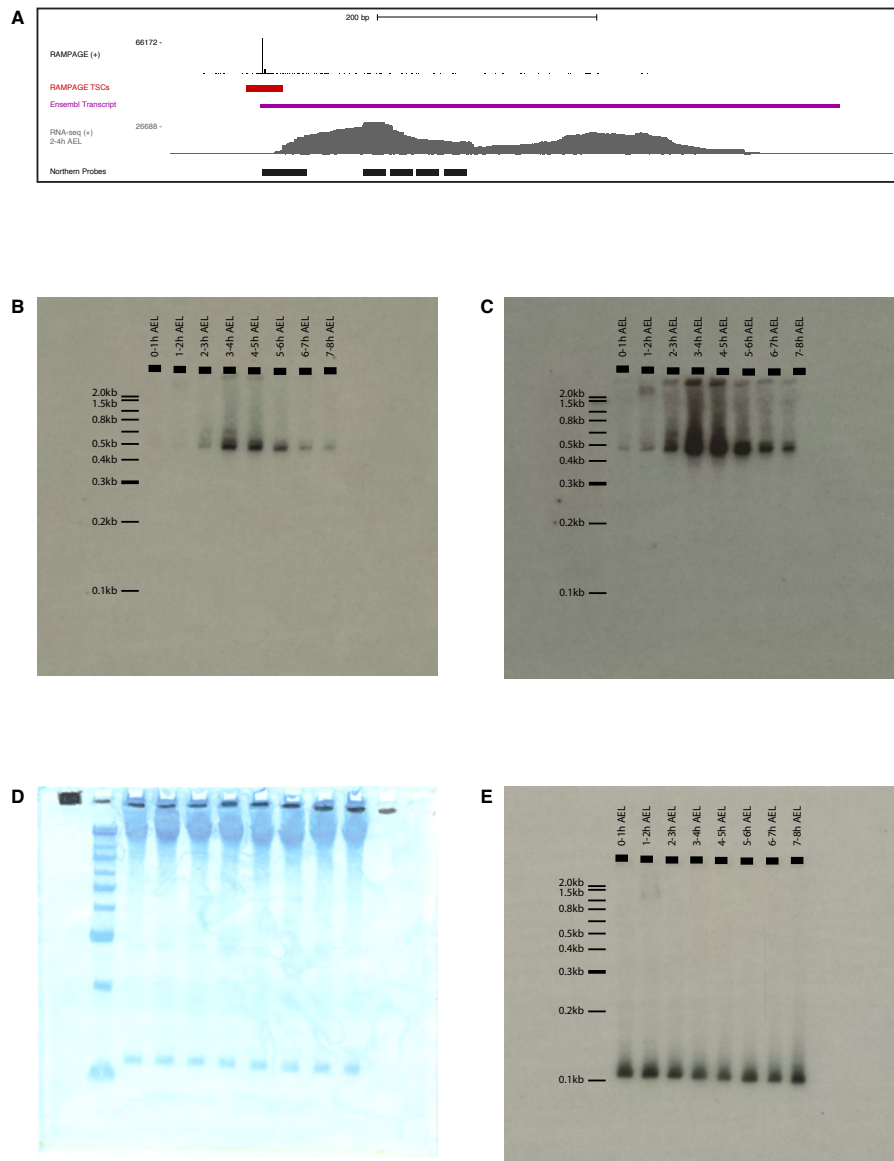


Figure 20. FBgn0264479 Northern-blot

(A) Description of the probes. Tracks from top to bottom: RAMPAGE signal, TSC, Ensembl transcript annotation, RNA-seq signal (modENCODE, 2-4 hours time point), FBgn0264479 probes (antisense to target). We used a mix of all 6 radiolabeled oligonucleotide probes. (B) Exposure of the whole membrane, hybridized with anti-FBgn0264479 probes. (C) Same blot, longer exposure. (D) Methylene blue staining of the membrane prior to hybridization. First lane: Invitrogen 0.1-2.0kb RNA ladder. (E) Hybridization with anti-5S rRNA probe.

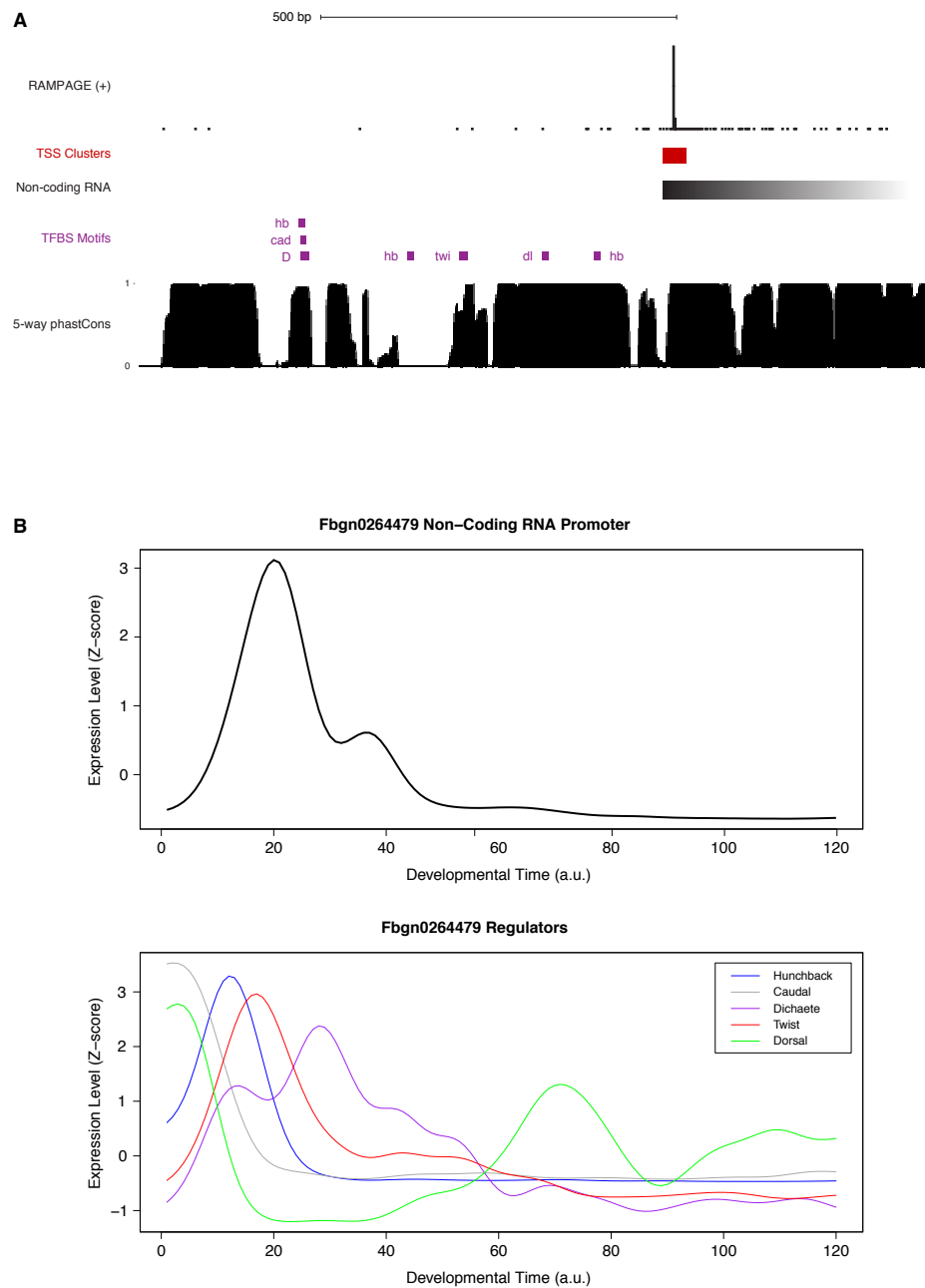


Figure 21. *FBgn0264479* transcriptional regulation

(A) Sequence conservation over the putative binding sites for the factors with ChIP-chip signal over the promoter (see Suppl. Fig. 18). We used FIMO to search for Jaspas-defined motifs within 600bp upstream of the main TSS. Note that all putative TFBSs but one are under strong purifying selection within the *melanogaster* subgroup. (B) Expression profiles of the genes encoding the putative regulators of *FBgn0264479*.