

1 DNA Transposons favour *de novo* transcript emergence through  
2 enrichment of transcription factor binding motifs

3 Running title. *De novo* transcripts emergence through TEs and motifs

4 Marie Kristin Lebherz<sup>1</sup>, Bertrand Fouks<sup>1</sup>, Julian Schmidt<sup>1</sup>, Erich Bornberg-Bauer<sup>1,2</sup>, Anna  
5 Grandchamp<sup>1,\*</sup>

6 <sup>1</sup>Institute for Evolution and Biodiversity, University of Münster, Münster, Germany

7 <sup>2</sup>Department of Protein Evolution, Max Planck Institute for Biology, Tübingen, Germany

8 \*Corresponding author: [a.grandchamp@uni-muenster.de](mailto:a.grandchamp@uni-muenster.de)

9 **Abstract**

10 *De novo* genes emerge from non-coding regions of genomes via succession of mutations.  
11 Among others, such mutations activate transcription and create a new open reading frame (ORF).  
12 Although the mechanisms underlying ORFs emergence are well documented, relatively little is  
13 known about the mechanisms enabling new transcription events. Yet, in many species a continuum  
14 between absent and very prominent transcription has been reported for essentially all regions of  
15 the genome.

16 In this study we searched for *de novo* transcripts by using newly assembled genomes and  
17 transcriptomes of seven inbred lines of *Drosophila melanogaster*, originating from six European  
18 and one African population. This setup allowed us to detect line specific *de novo* transcripts,  
19 and compare them to their homologous non-transcribed regions in other lines, as well as genic  
20 and intergenic control sequences. We studied the association with transposable elements and the  
21 enrichment of transcription factor motifs upstream of *de novo* emerged transcripts and compared  
22 them with regulatory elements.

23 We found that *de novo* transcripts overlap with TEs more often than expected by chance.  
24 The emergence of new transcripts correlates with high CpG islands and regions of TEs activity.  
25 Moreover, upstream regions of *de novo* transcripts are highly enriched with regulatory motifs.  
26 Such motifs abound in new transcripts overlapping with TEs, particularly DNA TEs, and are  
27 more conserved upstream *de novo* transcripts than upstream their non-transcribed homologs.  
28 Overall, our study demonstrates that TEs insertion is important for transcript emergence, partly  
29 by introducing new regulatory motifs from DNA TE families.

## 30 Introduction

31 For long, new genes were thought to exclusively arise from pre-existing genes (Guerzoni and McLysaght,  
32 2011). However recent studies showed that a non-negligible proportion of new genes also emerge  
33 *de novo* from non-coding regions of the genome (Schlötterer, 2015; Bornberg-Bauer et al., 2015;  
34 Rödelsperger et al., 2019; Tautz and Domazet-Lošo, 2011; McLysaght and Hurst, 2016; Van Oss and  
35 Carvunis, 2019; Bornberg-Bauer et al., 2021). Several *de novo* genes have been shown to become  
36 essential, bearing important organismal functions ,e.g. male fertility (Gubala et al., 2017) and cold  
37 resistance (Baalsrud et al., 2018). For a *de novo* gene to arise, it requires both the gain of an open  
38 reading frame (ORF) and the acquisition of transcription (Durand et al., 2019; Schlötterer, 2015).  
39 While the gain of ORFs in the emergence of *de novo* genes has been well studied (Zhuang and Cheng,  
40 2021; Delihias, 2022; Rödelsperger et al., 2019; Wang et al., 2020b; Carvunis et al., 2012; Grandchamp  
41 et al., 2023b), how transcription is acquired remains poorly understood.

42 The transcription of a gene is initiated at the core promoter which is located upstream the gene's 5'  
43 untranslated region (UTR) (Haberle and Stark, 2018; Butler and Kadonaga, 2002). Core promoters  
44 contain specific binding motifs, such as the *TATA box* or the *Initiator sequence*, that are recognized  
45 by transcription factors (tFs) (Boeva, 2016). Binding motifs with low identity to the consensus  
46 sequence are referred as minimal motif (Wang et al., 2020a). Transcription factors then recruit the  
47 protein complexes required for transcription (Butler and Kadonaga, 2002). However, transcription  
48 of low amounts of transcripts can also be initiated by a core promoter alone (reviewed in Haberle  
49 and Stark (2018); Small and Arnosti (2020)). Promoters can also produce antisense transcripts by  
50 initiating transcription in both direction (Scruggs et al., 2015). Furthermore, proximal and distal  
51 enhancers regulate the levels of transcription. Proximal enhancers (also called proximal promoters)  
52 are located directly upstream of core promoters, while distal enhancers influence transcription over  
53 long distances (Kim and Shiekhhattar, 2015; Haberle and Stark, 2018). Both contain tF binding motifs  
54 and can increase the amount of transcription initiated by the promoter (Haberle and Stark, 2018),  
55 independently of their locations and directions (Haberle and Stark, 2018). Enhancers often carry out  
56 bi-directional transcription, producing short but unstable transcripts in both directions (Small and  
57 Arnosti, 2020; Meers et al., 2018). Enhancers and promoters can also occasionally be converted into  
58 each other (Majic and Payne, 2020), and promoters can be interconnected by successive mutations  
59 without completely losing their activity (Kurafeiski et al., 2019)

60 In a non-coding region, the gain of transcription can result from random point mutations in a  
61 minimal motif and lead to stable transcription (Palazzo and Lee, 2015; Kapusta and Feschotte, 2014),

62 as genomes generally contain many cryptic functional sites with minimal promoters (Kapusta and  
63 Feschotte, 2014). Genomic mutations can also be initiated via the insertion of transposable elements  
64 (TEs). TEs are mobile DNA sequences that can move and amplify in genomes. They can be divided  
65 into two classes, based on their transposition mode: RNA and DNA transposons, which are further  
66 divided into sub classes and families based on their sequence characteristics (McCullers and Steiniger,  
67 2017). Several studies reported major reshuffling of genomic architectures due to TEs, as well as  
68 their role in adaptive evolution (Bourque et al., 2018; Delprat et al., 2009; Thybert et al., 2018). For  
69 example, syncytin genes, enabling cell-cell fusion in mammalian placenta, are derived from TEs (Malik,  
70 2012). TEs have also aided the evolution of the placenta in mammals, by acting on enhancers activity  
71 (Chuong et al., 2013). Other epigenetic mechanisms can influence transcription levels, such as DNA  
72 methylation, which represses genes transcription in vertebrates via the modulation of tFs activity (Law  
73 and Jacobsen, 2010). In invertebrates, methylation patterns are also associated with the regulation of  
74 transcription (Dixon and Matz, 2021), but the correlation between transcription and methylation is  
75 less clear than in vertebrates (Dunwell and Pfeifer, 2014; Lyko et al., 2000). Transcription is a highly  
76 dynamic and plastic process with high rates of transcripts gain and loss in closely related species, as  
77 well as among populations and individuals (Zhao et al., 2014; Grandchamp et al., 2022, 2023a; Neme  
78 and Tautz, 2016; Iyengar and Bornberg-Bauer, 2023), suggesting fast transcripts turnover. However,  
79 the mechanisms promoting *de novo* transcripts, i.e. transcription initiation from non-coding regions,  
80 remains elusive.

81 In this study, we investigate the mechanisms underlying novel transcript emergence at short evolutionary  
82 time scales by studying *de novo* transcripts in seven lines of *Drosophila melanogaster*, originating  
83 from different geographical locations (Grandchamp et al., 2022). By using long-read sequencing and  
84 a common annotation methodology across all genomes, our genomes and transcriptomes present a  
85 unique opportunity to precisely categorize *de novo* transcripts in each *Drosophila* line and investigate  
86 the molecular basis underlying the gain of transcription. Indeed, our dataset is allowing us to compare  
87 directly the related DNA sequences that are transcribed in one or several lines but not in others. In  
88 particular, we studied the role of transposable element insertions and motif enrichment upstream of *de*  
89 *novo* transcripts that emerged in each *Drosophila* line. Overall, our analyses reveal that the emergence  
90 of transcription is aided by an enrichment of motifs upstream of a DNA sequence, motif enrichment  
91 which is itself favored by nearby insertion of DNA transposons.

## 92 Results

### 93 General characteristics of *de novo* transcripts

94 To characterize the molecular basis underlying gains of transcription, we used a conservative approach  
95 to define *de novo* transcripts, ensuring detection of strictly *de novo* transcript (see Methods). Such  
96 definition and filtering led to the discovery of between 403 (Sweden [SE]) to 628 (Ukraine [UA]) *de*  
97 *nov* transcripts across *D. melanogaster* lines (mean =  $504 \pm 28.04$  (SE), Figure 1 A, [Supplemental](#)  
98 [deposit](#)). *De novo* transcripts were unevenly distributed among and along chromosomes, with the  
99 highest numbers of *de novo* transcripts in 3L and 3R chromosome arms ([SI-S1](#)). Most of *de novo*  
100 transcripts were found in only one *D. melanogaster* line (2389 / 3528) and only a few (38) were  
101 shared among all lines, suggesting a high birth / death rate of *de novo* transcripts, Figure 1 B).  
102 This high birth/death rate of *de novo* transcripts is likely the result of gain / loss of transcription,  
103 as most *de novo* transcripts (14058 blast hit out of 18903 blast searches in a maximum of 6 lines)  
104 had a 'non-transcribed' homolog in at least one other *D. melanogaster* line ([Supplemental deposit](#)).  
105 Moreover, *de novo* transcripts show different patterns from annotated transcripts (both genes and  
106 non-coding RNAs), with *de novo* transcripts having lower expression level, GC content, exon number,  
107 and spliced length compared to annotated transcripts (GLMM: TPM:  $p < 0.001$ , GC content:  $p < 0.001$ ,  
108 exon number:  $p < 0.001$ , spliced length:  $p < 0.001$ , ([SI-S2](#))).

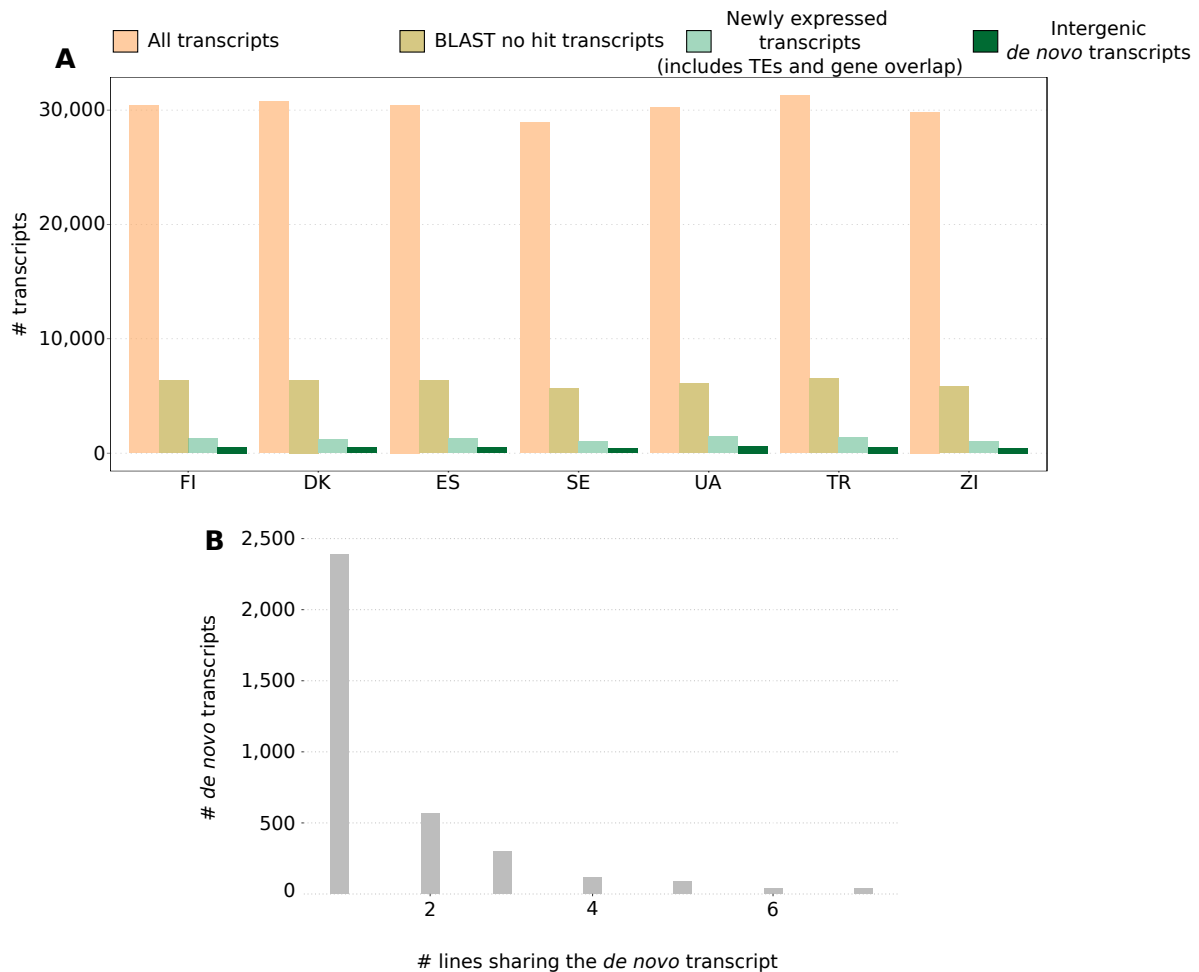


Figure 1: ***De novo* transcripts.**(A): Number of transcripts after filtering steps. The beige bar represents all transcripts detected with transcriptome assembly. The grey bar represents all transcripts without a BLAST hit. The green bar represents *de novo* transcripts after filtering for TPM and splicing. The dark green bar represents only the intergenic *de novo* transcripts after filtering out transcribed TE. (B): Number of *de novo* transcripts shared by lines

## 109 DNA transposon insertions favour the gain of transcription

110 For each genome, we performed a *de novo* annotation of TEs, using the TransposonUltimate pipeline  
 111 (Riehl et al., 2022) (method, [Supplemental deposit](#)). To understand how TEs can favour the gain of  
 112 transcription, we first assessed the relationship between TEs and *de novo* transcripts at the chromosome  
 113 scale (Figure 2). While *de novo* transcripts were evenly distributed along chromosomes, inactive  
 114 TEs and expressed TEs, were in higher density in the telomere regions of chromosomes (GLMM:  $p$   
 115  $<0.001$ ; Figure 2 A, [SI-S3](#)). Nevertheless, *de novo* transcript densities were positively correlated with  
 116 TE densities at a 100 kb scale (GLMM:  $p <0.001$ ). An important mechanism by which TE impact  
 117 gene expression is the import of epigenetic marks, such as DNA methylation (Zhou et al., 2020). We

118 therefore calculated the CpGoe, as an estimate for DNA methylation status, with high CpGoe values  
119 corresponding to low level of methylation. *De novo* transcripts displayed low level of methylation  
120 (CpGoe: mean & median = 0.902 +- sd 0.222) and their methylation status were negatively correlated  
121 with TE density (GLMM:  $p < 0.001$ ), highlighting the role of TEs in importing epigenetic marks  
122 (Figure 2 B, Supplemental deposit, SI-S4).

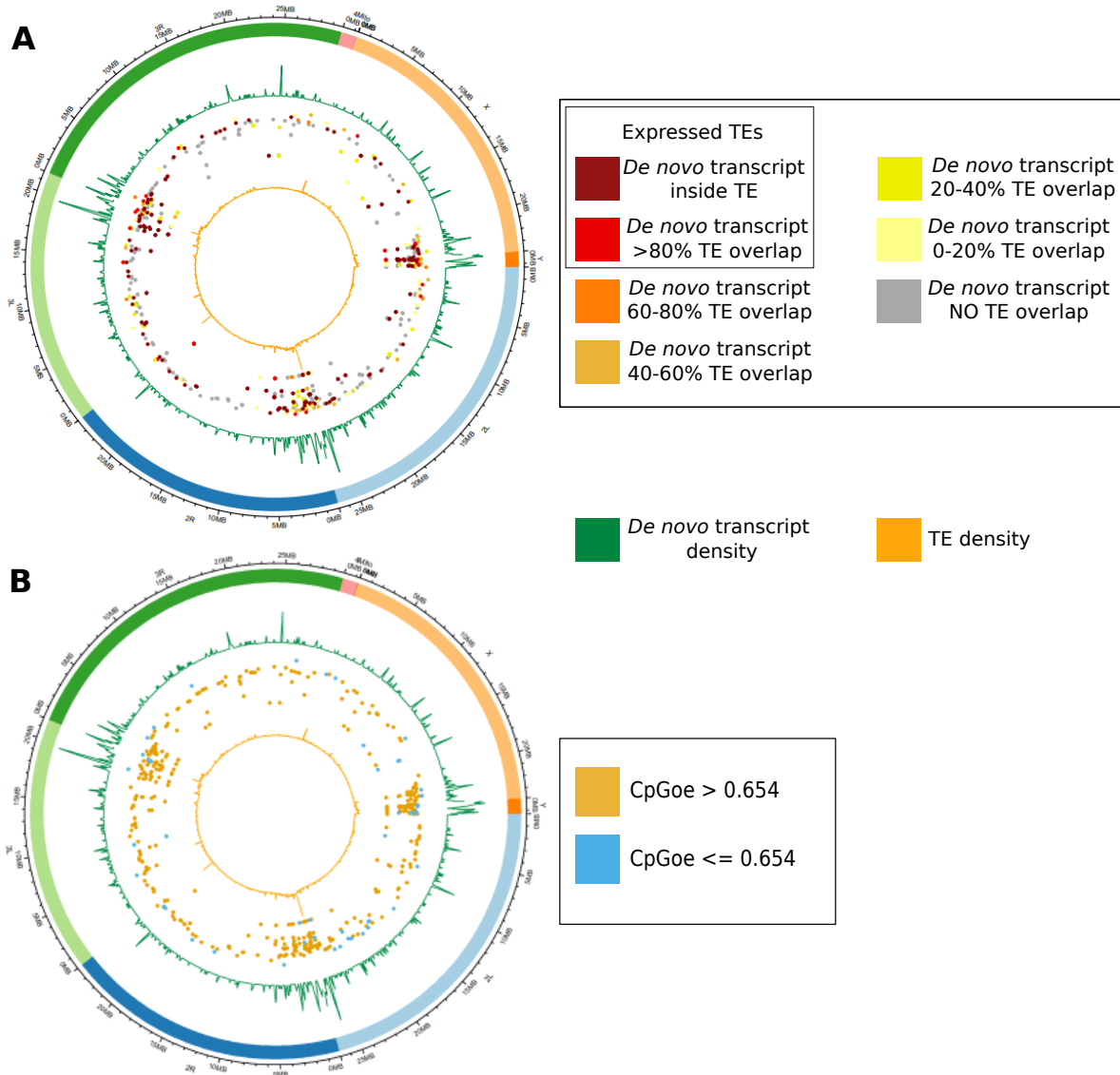


Figure 2: ***De novo* transcripts and TE density among the chromosomes.** The circular plots represent the *D.melanogaster* line collected in Denmark (DK). Plots with similar distributions can be found for all other lines in the *supplemental data*. The 8 chromosome arms are represented with specific colours. In the 2 circle plots, the green circles represent *de novo* transcripts, and the yellow lines represent TEs distribution. (A) The coloured dot represents expressed TEs and *de novo* transcripts distribution according to their relative overlap with TEs. (B) The coloured dots represent the CpGoe values of *de novo* transcripts according to their genomic distribution.

123 In addition to our chromosome scale analyses, we also investigated the impact of TE insertions on *de*  
124 *novo* transcripts by comparing the number of TE overlapping with these transcripts, as well as their  
125 down- and upstream regions, with random intergenic regions as a negative control. *De novo* transcripts  
126 displayed a higher amount of TE insertions compared to other sequences, however with a lower length  
127 of TE overlap (GLMM:  $p < 0.001$ ; Figure 3 A,B, [Supplemental deposit](#)). Furthermore to be able to  
128 pinpoint precisely the role of TE insertions on the gain of transcription, we directly compared *de novo*  
129 transcripts with their 'non-transcribed' homolog sequences present in other *D. melanogaster* lines. Our  
130 analyses revealed that TE insertions did not differ between *de novo* transcripts and 'non-transcribed'  
131 homologs, however *de novo* transcripts displayed shorter overlaps with TEs as well as a lower number  
132 of TE insertions compared to 'non-transcribed' homologs (GLMM,  $p < 0.001$ , [SI-S5](#)). Moreover, RNA  
133 TEs were less abundant in *de novo* transcripts compared to 'non-transcribed' homologs (GLMM,  $p$   
134  $< 0.001$ , Figure 3 C, [SI-S5](#)). On the contrary, DNA TEs were more abundant in *de novo* transcripts  
135 compared to 'non-transcribed' homologs (GLMM,  $p < 0.001$ ). Our results highlight a different impact  
136 between TE classes (DNA vs. RNA) on transcription gain.

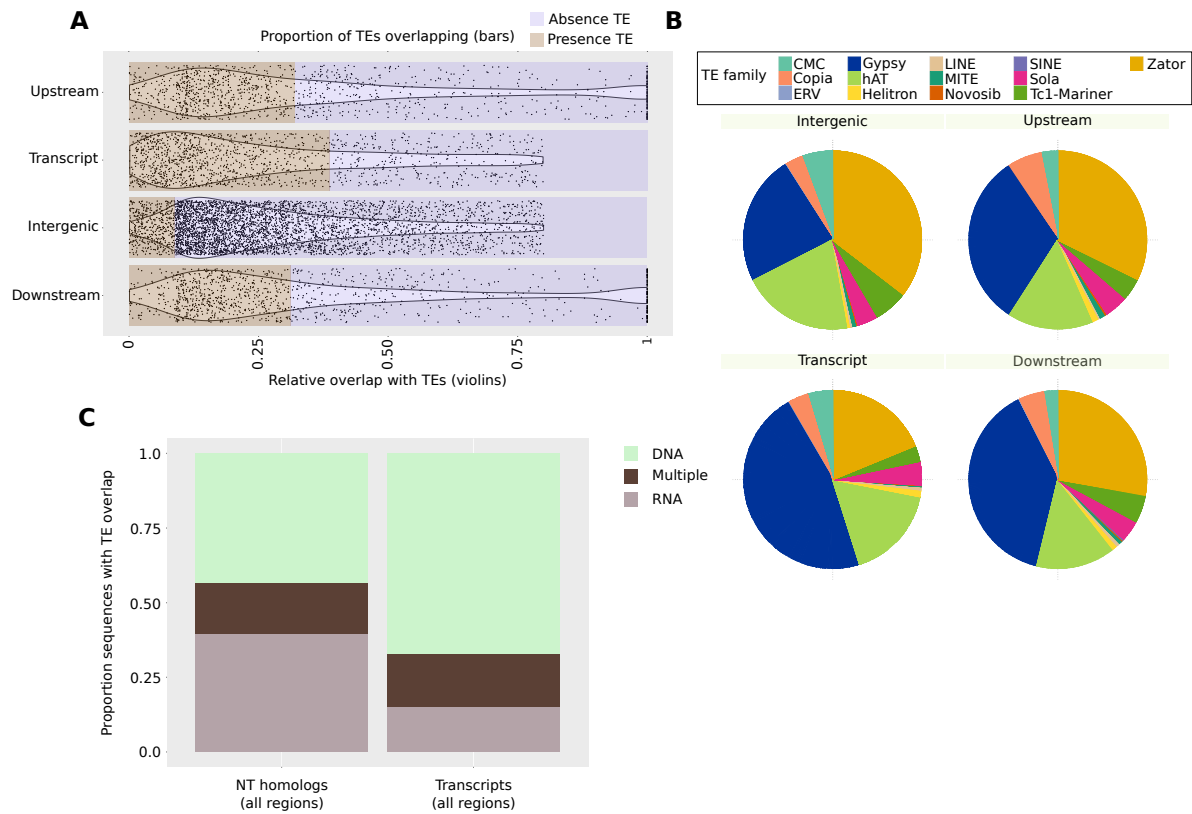


Figure 3: **TEs overlap.** (A) Relative sequence overlap with TEs and number of sequences overlapping with TEs into four datasets: Intergenic sequences, upstream sequences of *de novo* transcripts, downstream sequences of *de novo* transcripts, *de novo* transcripts. (B) Percentage of TEs overlapping with the four datasets according to their families. (C) Major classes of TEs overlapping with *de novo* transcripts and their non-transcribed homologs.



## 137 Motifs enrichment

138 A major factor influencing gene expression is the presence of specific DNA motifs enabling the  
139 transcription machinery to bind to the DNA region. We therefore investigated the role of DNA binding  
140 motifs for the gain of transcription. We compared several measures of motif enrichment (specific to  
141 both tF motifs from enhancers and distal promoters, as well as (core) promoters) upstream of our *de*  
142 *novo* transcripts, as positive controls upstream of genes and expressed TEs, and as negative control  
143 random intergenic regions. Motif enrichments were further divided into two classes according to their  
144 thresholds of similarity to their PSSM matrix : low identity motifs (minimal motifs), with a score of  
145 identity to the matrix > 80%, and high identity motifs, with an ID score of 95% identity as a minimum  
146 (Figure 4). This comparison revealed that TEs and *de novo* transcripts overlapping with TEs have  
147 higher number of low identity tF motifs compared to other sequences (GLMM:  $p < 0.001$ ). Moreover,  
148 *de novo* transcripts that do not overlap with TEs displayed higher numbers of core promoters with high  
149 identity score than TEs and *de novo* transcripts (GLMM:  $p < 0.001$ ; SI-S6,S7). Overall, genes and  
150 intergenic regions displayed a higher enrichment of core promoter motifs (both high and low identity  
151 motifs) and of tF motifs with high identity score, while TEs and *de novo* transcripts displayed an  
152 enrichment of tF motifs of low identity score (GLMM,  $p < 0.001$ , Figure 4, SI-S6,S7).

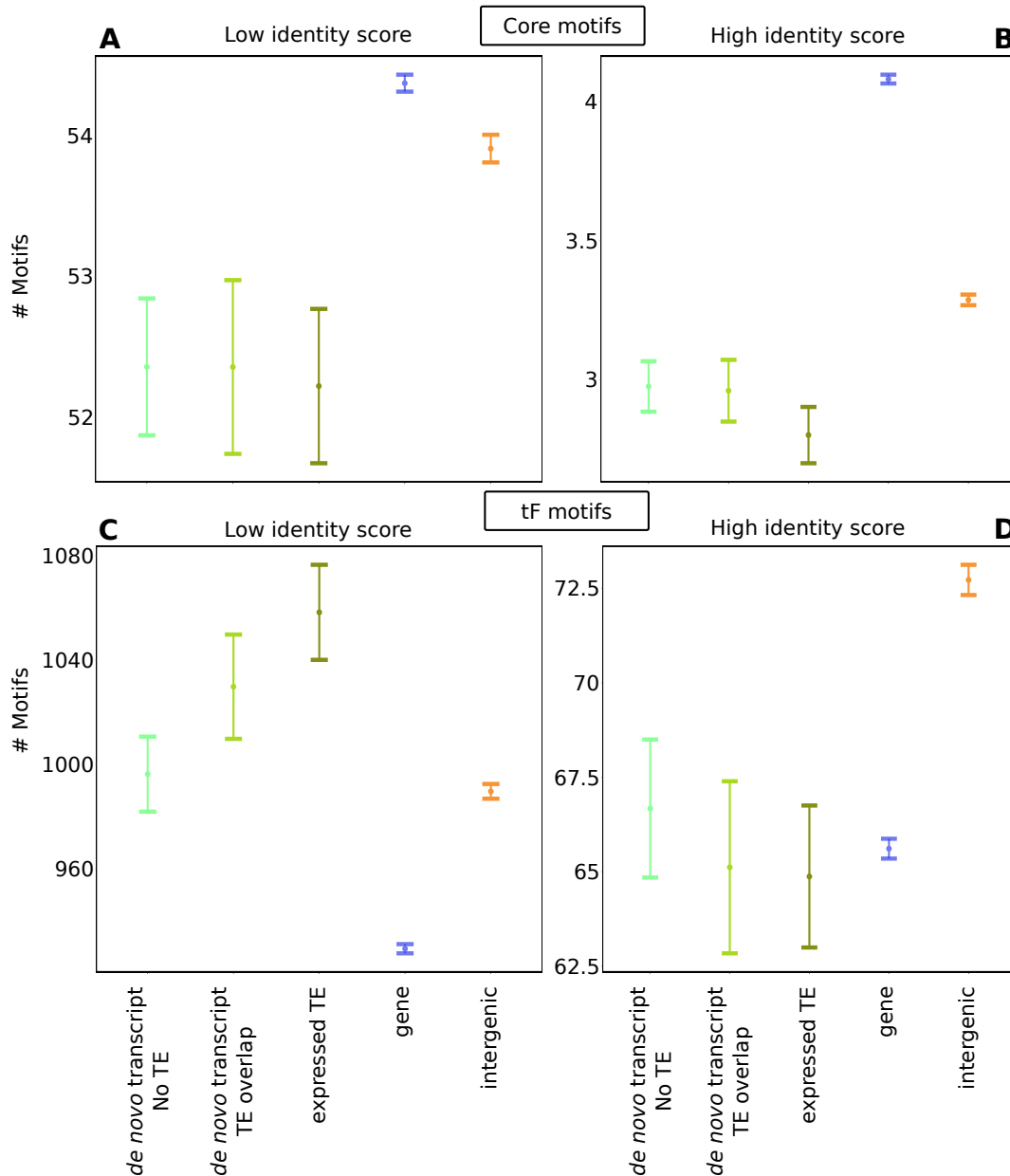


Figure 4: **Number of motifs detected upstream five sequences datasets.** (A) Number of low identity Core promotor (0.8) motifs detected upstream i) *de novo* transcripts overlapping no TE (light green), ii) *de novo* transcripts overlapping with TEs (green), iii) expressed TEs (dark green) iv) genes (blue), v) randomly selected intergenic regions that are not transcribed (orange). (B) Number of high identity Core promotor motifs (0.95) detected upstream the aforementioned dataset of sequences (C) Number of low identity tF motifs (0.80) detected upstream the aforementioned dataset of sequences. (D) Number of high identity tF motifs (0.95) detected upstream the aforementioned dataset of sequences.

153 When studying motifs individually (SI-S6,S7), 13 motifs were enriched upstream *de novo* transcripts  
154 and TEs, compared to intergenic regions, with a high threshold (relative score = 0.95; supp data), four  
155 of them being also significantly enriched in upstream genes. Three of these 13 motifs were significantly  
156 enriched upstream *de novo* transcripts without TE overlap. 11 out of these 13 motifs were specific  
157 for homeo domain factors, with one zinc finger factors (Supplemental deposit). Among the ten most  
158 abundant motifs (ara, mirr, CG4328-RA, lbe, PHDP, H2.0, Deaf.1, caup, C15, lbl), four were enriched  
159 in *de novo* transcribed TEs and in TEs overlapping *de novo* transcripts. We found 78 tF motifs that  
160 were enriched upstream *de novo* transcripts and TEs with a low threshold (relative score = 0.8), 13 of  
161 them being also significantly enriched in genes. Only 18 of them were enriched upstream *de novo*  
162 transcripts that did not overlap any TE. Most of these 73 motifs were specific tF for homeo domain  
163 factors or zinc finger, however they also included one motifs for high mobility group domain factor,  
164 for one heat shock factor, two motifs for leucine zipper factors, two for paired box factors, one fork  
165 head/winged helix factor, for a STAT and TEA domain factor. Out of the ten most frequent motifs  
166 from the dataset using this threshold (CG4328-RA, br, H2.0, PHDP, C15, vvl, Dbx, ct, lbl, ara) seven  
167 were enriched in all *de novo* transcripts, one of them also in genes. Additional two were enriched only  
168 in TEs and *de novo* transcripts overlapping TEs.

169 In addition, we compared directly binding motif enrichment upstream sequences of *de novo* transcripts  
170 and their 'non-transcribed' homologs. We observed no significant difference in motif enrichment  
171 between *de novo* transcripts and their 'non-transcribed' homologs. The best statistical model included  
172 the enrichment of low identity core promoters but it was not significant (GLMM,  $p = 0.136$ , SI-S8).  
173 Furthermore, we implemented the impact of TE insertions along motif enrichment between *de novo*  
174 transcripts and their 'non-transcribed' homologs. *De novo* transcripts exhibit, when TE inserted, a  
175 higher density of tF motifs of low identity, suggesting that TE insertions enable transcription through  
176 low tF motif enrichment (Figure 5). Finally, we accounted for the different TE class (DNA vs. RNA)  
177 inserting among *de novo* transcripts and their non-transcribed homologs. While *de novo* transcripts  
178 have a lower ratio of RNA transposons compared to their 'non-transcribed' homologs, high number  
179 of RNA transposon insertions in *de novo* transcripts is linked with an enrichment low identity core  
180 promoter motifs (GLMM,  $p < 0.001$ , SI-S8).

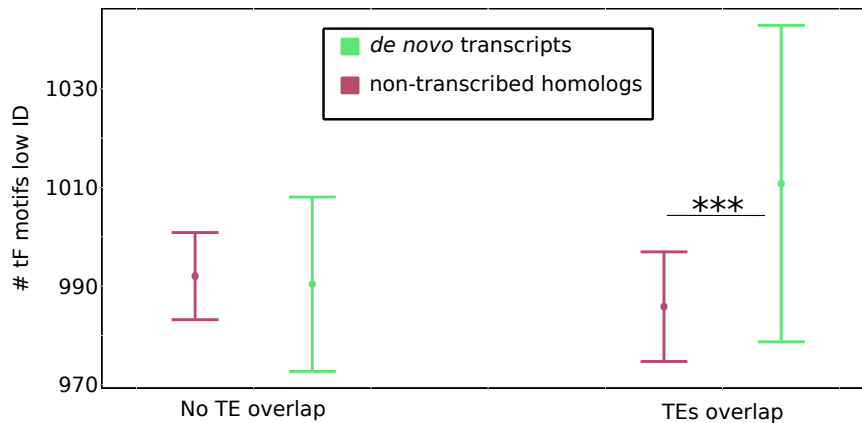


Figure 5: **Enrichment in low tF promoter motifs upstream *de novo* transcripts and their non-transcribed homologs** The green colour represents *de novo* transcripts. The pink colour represents non-transcribed homologs. The bars on the left represent sequences without TE overlap, while the bars on the right represent sequences with TE overlap. The y axis represents the number of low tF motifs.

## 181 Discussion

### 182 Detection of *de novo* transcripts

183 To understand how transcription can be gained in intergenic regions leading to the emergence of *de*  
184 *novo* genes, we searched for *de novo* intergenic transcripts that emerged in seven lines of *Drosophila*  
185 *melanogaster*. Our stringent definition led to the discovery of 3,799 transcripts over 7 *D. melanogaster*  
186 lines, with an average of 504 intergenic *de novo* transcripts per line. This amount of *de novo* transcripts,  
187 while being lower than in a previous study of new transcripts emergence in lines (Everett et al., 2020),  
188 corresponds well to previous estimates (Camilleri-Robles et al., 2022; Huang et al., 2015), if we account  
189 only for intergenic *de novo* transcripts.

190 Moreover, the characteristics of our *de novo* transcripts corresponds well to those of previous studies,  
191 namely a lower expression, lower GC content, lower number of exons, and a shorter sequence than  
192 known genes. Finally, our estimation of *de novo* transcripts could have been minimized by not  
193 accounting for transcripts with low level of expression or tissue- and life-stage specific expression,  
194 resulting in lower detection of *de novo* transcripts (Grandchamp et al., 2022).

### 195 Overlap with transposable elements

196 Among all detected *de novo* transcripts, 34% overlapped fully or by more than 80% with TEs, and  
197 were then considered as "active TEs" rather than *de novo* transcripts. This first outcome suggests

198 that TEs have important mobility inside the species. TEs were massively detected and active in the  
199 telomeric regions of the chromosomes, as previously reported (Kordyukova et al., 2018). While *de*  
200 *novo* transcripts display a higher proportion of TE insertions compared to control sequences, TEs  
201 overlapped mainly with small fractions of the transcripts and of their surrounding regions, rejecting  
202 the hypothesis that such new transcription events correspond to biased transcript activity. However,  
203 such a correlation between TE overlap and new transcription events suggests that TEs insertion could  
204 have contributed to the emergence of the new transcripts that are unrelated to TEs mobility.

205 When comparing *de novo* transcripts with their non-transcribed homologs, they did not differ in  
206 their proportion of TE insertions. Nevertheless, the impact of greater length of TE overlap and  
207 higher number of TE insertions seems detrimental for transcription, since *de novo* transcripts have a  
208 shorter TE overlap and less numerous TE insertions compared to their homologs. Furthermore, not  
209 all TEs seem to impact transcription gain, with RNA transposon being more disruptive than DNA  
210 transposon. Indeed, *de novo* transcripts display a higher proportion of DNA transposons compared  
211 to their homologs. These results suggest first that *de novo* transcripts emerge in regions that are  
212 prone to TE mobility, and are highly variable due to TEs activity. Second, given that DNA TEs  
213 are more associated with new transcription events, the insertion of DNA TEs seem to be the more  
214 likely to initiate novel transcription. Interestingly, the main difference in TE composition of *de novo*  
215 transcripts compared to intergenic sequences, was the higher amount of overlap with retrotransposons  
216 (mainly LTR elements from the gypsy family). In *Drosophila melanogaster* certain TEs, such as  
217 LTR retrotransposons are reported to be more active than others (Kofler et al., 2015; Petrov et al.,  
218 2011). High TE activity can also strongly reshuffle genomes. This could explain why 25% of *de*  
219 *novo* transcripts had no detected transcribed homolog when requiring a high degree (80% identity) of  
220 sequence similarity between transcript and homolog. Finally, most of the *de novo* transcripts show  
221 high CpGoe values which suggest low DNA methylation, indicating that the genomic location of  
222 the transcripts is accessible for the transcription machinery (Roder et al., 2000; Rollins et al., 2006).  
223 Furthermore, the correlation between the length of TE overlap with a *de novo* transcript and CpGoe  
224 values highlights the impact of TEs bringing along their epigenetic marks.

225 Taken together, all these independent outcomes strengthen the hypothesis that TEs are actively  
226 transposing in *D.melanogaster*, and that such activity is noticeable even between lines or individuals.  
227 This lines up with previous studies reporting high activity of several TE families in *Drosophila* (Kofler  
228 et al., 2015; Bourque et al., 2018; Lawlor et al., 2021; Mérel et al., 2020). Moreover, the significant  
229 overlap of active TEs with *de novo* transcripts strongly suggests that TE activity plays a role in  
230 initiating new transcription events in intergenic genome regions.

## 231 Minimal tF motifs enrichment leads to transcription gain

232 Intergenic regions of genomes are known to contain a high proportion of (distal) enhancers which  
233 interact with highly distant promoters (Small and Arnosti, 2020). That was confirmed in our results,  
234 with random intergenic sequences being the most enriched in highly conserved tF motifs. However,  
235 when studying motifs with lower scores of similarity to annotated motifs (80%), *de novo* transcripts  
236 contained the highest amount of such motifs, compared to genes and intergenic sequences. Indeed,  
237 such low tF motifs, also called sub-optimal transcription factor motifs, appears to be a significant  
238 factor for initiating new transcription in genomes. *De novo* transcripts showed lower expression levels  
239 than expressed genes, in line with the finding that transcription is initiated at low levels without the  
240 presence of strong core motifs (Palazzo and Lee, 2015).

241 While *de novo* transcripts showed high motifs enrichment of minimum tF motifs, upstream regions  
242 of transcripts overlapping with TEs showed the highest amount of low TF motifs. Such enrichment  
243 was still lower than in TEs. Most TEs possess a machinery for transcription, which necessitate the  
244 presence of tF motifs in their sequence (Chuong et al., 2017). The enrichment of low tF motifs  
245 upstream of *de novo* transcripts overlapping with TEs opens two hypothesis. First, the insertion of  
246 new TEs in previously untranscribed genomic location could provide sufficient sequence disruption  
247 to mutate into minimum tF motifs. tF motifs are usually shorter than 15 nucleotides, and several  
248 position allow nucleotide variability without affecting the binding. Therefore, the possibility of a motif  
249 emergence caused by mutations due to TE insertions does not seem unlikely. As a second hypothesis,  
250 new transcripts could have benefited from the presence of tF motifs in TEs to initiate new transcription  
251 events. While these two hypothesis could find support in literature (Chuong et al., 2017; Moschetti  
252 et al., 2020), our data seem to give more credit to the second one. Indeed, low tF enrichment was  
253 observed in *de novo* transcript compared to their non-transcribed homologs, only when a TE insertion  
254 within the sequence was present. Furthermore, while *de novo* transcripts and their homologs shared  
255 similar proportion of TE insertions, the TE content of *de novo* transcripts and their homologs diverge.  
256 *De novo* transcripts overlap more with DNA TEs, while non-transcribed homologs overlaps more with  
257 RNA TEs. Therefore, if TE insertions were to disrupt genomics sequences, both TE families would  
258 be expected to generate a similar amount of disruption, and generate a similar amount of motifs.  
259 However, both TEs families do not carry the same tF motifs, as their insertion mechanisms diverge.  
260 Indeed, our results tend to suggest that DNA TE insertion generates more new transcription events,  
261 and that this could be due to the recycling of their tF motifs.

262 Many different regulatory elements were shown to have been gained through a TE insertion, such

263 as enhancers/enhancer-like elements, promoters, splice sites, cis-regulatory elements, poly-A signals  
264 and more (Moschetti et al., 2020). In non-coding regions, transcription can also be initiated through  
265 transposable elements (TE) (Kapusta and Feschotte, 2014). TEs have been shown to have the ability  
266 to induce a regulatory sequence through different mechanisms such as domestication (use of TEs for  
267 a new function), gene duplication, change of gene expression, ectopic recombination (Kapusta and  
268 Feschotte, 2014; Moschetti et al., 2020; Rizzon et al., 2002). About 75% of human and 68% the  
269 mouse lncRNA include at the minimum one (partial) retrotransposon insertion (Kapusta et al., 2013).  
270 In humans TEs provided up to 23 % of non redundant transcription start sites and about 30% of  
271 poly-A sites of lncRNA. (Ganesh and Svoboda, 2016). In *Drosophila*, TE content has been shown to  
272 be high in long noncoding RNA (Ganesh and Svoboda, 2016; Fort et al., 2021), compared to protein  
273 coding genes, which would support our results.

274 Indeed, TEs (and especially DNA families), could have played a (partial) role in the gain of transcription  
275 of new transcripts, e.g. by inserting the motifs enabling the start of transcription. Our outcomes  
276 demonstrate that this gain of transcription through TEs is a frequent event, and can occur independently  
277 in different lines from a same species. Determining how exactly the TEs lead to the transcription  
278 of these regions and which elements (poly-A, promoter, enhancer etc.) they contributed to insert  
279 would need further investigation and more detailed comparisons between the transcript (and up- and  
280 downstream) sequences and their homologous regions in the outgroup lines.

281 60% of *de novo* transcripts emerged without overlapping with TEs. These transcripts showed higher  
282 minimal tF enrichment than control sequences, but the difference was less obvious than for transcripts  
283 overlapping with TEs. Such small enrichment could be explained by the emergence of low identity tF  
284 motifs by other mechanisms than TE insertions, like indels, or other sequence reshuffling that we did not  
285 investigate, e.g. genomic inversions or duplications. Furthermore, the high GpC content in all *de novo*  
286 transcripts could be associated with low methylation, even though genomes methylation is less observed  
287 in invertebrate genomes than vertebrates (Klughammer et al., 2023). Also, we found surprisingly  
288 low amounts of core promoter motifs upstream *de novo* transcripts. If such motif enrichment was  
289 suspected to be lower than upstream genes, it was surprising to find them less enriched than control  
290 intergenic sequences. So far, we have no hypothesis for such an output, but it might also play a role  
291 in new transcripts emergence.

## 292 Conclusion

293 Overall, our study reveals the importance of TEs in transcription gain and loss. At a large scale, a  
294 high TE density seems to enable transcription, most likely through changes of chromatin organization  
295 (Lawson et al., 2023), as TE density was correlated with *de novo* transcripts density within 100  
296 kb windows. At a finer scale, insertions of TEs seems to lead to different outcomes depending on  
297 their insertion patterns. Indeed, a singular insertion of DNA transposon shortly overlapping with the  
298 transcript sequence tends to favour the gain of transcription, most likely through enrichment of the  
299 upstream region with minimal tF motifs. On the contrary, insertions of RNA transposons likely lead  
300 to transcription loss, at the exception of multiple RNA transposon insertions accompanied with an  
301 enrichment of minimal core promoter in the upstream region.



## 302 Methods

### 303 Detection of *de novo* transcripts and their non-transcribed homologs

304 To investigate the molecular mechanisms enabling new transcript emergence, we searched for *de novo*  
305 transcripts and their non-transcribed homologs in the transcriptomes and genomes, respectively, of  
306 seven lines of *D. melanogaster*, six inbred european lines and one from Zambia (NCBI Bioproject  
307 PRJNA929424)(Grandchamp et al., 2022). Transcripts were defined as being *de novo* (i.e. newly  
308 emerged) if they met our four criteria: i) detected in one or several of the seven inbred line transcriptomes  
309 with a TPM value (transcripts per million) above 0.5(Grandchamp et al., 2023a); ii) no homology to  
310 any other annotated transcripts (cRNA and ncRNA) in the *D. melanogaster* reference transcriptome  
311 (Table 1); iii) no homology with annotated transcripts (cRNA and ncRNA) of eleven outgroup  
312 *Drosophila* and five Diptera species (Table 1); iv) no overlap of transcript genome location with TEs  
313 greater than 80%.

Table 1: List of reference species used to build the reference database for the blast search

	Species	Accession number	Assembly
1	<i>Aedes aegypti</i>	GCA_002204515.1	AaegL5
2	<i>Anopheles sinensis</i>	GCA_000472065.2	AsinS2
3	<i>Culex quinquefasciatus</i>	GCA_000209185.1	CpipJ2
4	<i>Drosophila ananassae</i>	GCA_000005115.1	dana_caf1
5	<i>Drosophila erecta</i>	GCA_000005135.1	dana_caf1
6	<i>Drosophila grimshawi</i>	GCA_000005155.1	dgri_caf1
7	<i>Drosophila melanogaster</i>	GCA_000001215.4	BDGP6.32
8	<i>Drosophila mojavensis</i>	GCA_000005175.1	dmoj_caf1
9	<i>Drosophila persimilis</i>	GCA_000005195.1	dper_caf1
10	<i>Drosophila pseudoobscura</i>	GCA_000001765.2	Dpse_3.0
11	<i>Drosophila sechellia</i>	GCA_000005215.1	dsec_caf1
12	<i>Drosophila simulans</i>	GCA_000754195.3	ASM75419v3
13	<i>Drosophila virilis</i>	GCA_000005245.1	dvir_caf1
14	<i>Drosophila williston</i>	GCA_000005925.	dwil_caf1
15	<i>Drosophila yakuba</i>	GCA_000005975.1	dyak_caf1
16	<i>Megaselia scalaris</i>	GCA_000341915.1	Msca1
17	<i>Teleopsis dalmanni</i>	GCA_002237135.2	ASM223713v2

314 Nucleotide BLAST (version 2.12) (Altschul et al., 1990) with the plus option was used to assess  
315 homology between inbred *Drosophila melanogaster* lines and reference transcripts. The lack of  
316 homology was defined if a transcript did not return a BLAST hit (with a threshold E-value of 0.05),  
317 as well as none of its splicing variant.

318 Bedtools (version 2.3, intersect with default parameters) (Quinlan and Hall, 2010) was used to map

319 *de novo* transcripts onto their respective genome. *De novo* transcripts overlapping with a gene in  
320 sense or antisense direction were filtered out, keeping only intergenic *de novo* transcripts.

321 To better understand the frequency of transcription gain and loss, we quantified the amounts of *de*  
322 *novo* transcripts shared across inbred *D. melanogaster* lines. To that end, a BLAST search (plus strand  
323 option, E-value of 0.05) of our *de novo* transcripts were performed against the transcripts of the other  
324 lines. Transcripts were deemed to be homologous if they met those three criteria: i) the transcription  
325 start sites of transcripts match up in a 200 nucleotide window; ii) the transcription termination sites  
326 of transcripts match up in a 200 nucleotide window; iii) transcripts share at least 80% identity.

327 To precisely categorize the mechanisms underlying the gain of transcription, direct comparisons of the  
328 same nucleotide sequences exhibiting different transcription status is mandatory. We, therefore, used  
329 *de novo* transcripts, which were not found across all lines, and their location onto their respective  
330 genome to find their 'non-transcribed homologs'. The unspliced sequences of those *de novo* transcripts  
331 were retrieved using bedtools (get fasta with the -s option)(Quinlan and Hall, 2010). Those unspliced  
332 sequences were then used to identify similar/identical nucleotide sequences in the genome of other  
333 lines, which do not possess this *de novo* transcript, using a nucleotide BLAST search (default settings,  
334 E-value cut-off 0.05) (Altschul et al., 1990). A nucleotide sequence was defined as a 'non-transcribed  
335 homologs', if BLAST hits had 80% query coverage with the *de novo* transcript. If a transcript had  
336 multiple 'non-transcribed homologs' in the same line, only the nucleotide sequence with the lowest  
337 E-value, highest percent identity and highest query coverage, was retained.

338 Non-transcribed homologs were searched per transcript instead of per orthogroup. The original  
339 dataset was reduced to avoid confusion i) Alternative spliceforms were reduced to one spliceform per  
340 orthogroup; ii) Orthogroup containing lines duplication were removed (iii) All orthogroup member and  
341 non-transcribed homologs have their initiation and termination positions in a same window (+/- 200  
342 nt).

### 343 **The role of transposable elements in gain of transcription**

344 To unravel the importance of transposable elements in the emergence of *de novo* transcripts, *de*  
345 *novo* annotations of TEs were performed in each inbred line, using the reasonaTE pipeline from the  
346 TransposonUltimate software (Riehl et al., 2022). This pipeline was chosen as it combines, compiles,  
347 and filters TE annotations from 13 tools with different annotation approaches (Riehl et al., 2022). *De*  
348 *novo* TE annotations of each *D. melanogaster* line genome was used to infer their relative overlap with  
349 *de novo* transcripts, as well as with their upstream and downstream regions, with 'non-transcribed

350 homologs' and their upstream regions, and as a control with random intergenic regions of 1100 bp  
351 length obtained using bedtools (Quinlan and Hall, 2010). Relative overlap was calculated by dividing  
352 the overlap length between a sequence and a TE obtained with bedtools (Quinlan and Hall, 2010)  
353 with the full of length of the sequence. Up- and downstream regions were defined as 1000 bp length  
354 before and after a given sequence, respectively, with a 100 bp overlap with the given sequence (for a  
355 total of 1100 bp length). A given sequence could overlap with more than one TE, in this case relative  
356 overlap was calculated using all overlapping TEs, and the number of TEs as well as their class and  
357 family were calculated.

358 Moreover to evaluate features associated with gain of transcription at the genome scale, the distribution  
359 of *de novo* transcripts and TEs density within a 10kb sliding window, as well as CpGoe (a proxy for DNA  
360 methylation), were plotted along chromosomes for each *D. melanogaster* line, using an R script adapted  
361 from (Ylla et al., 2021) [https://github.com/guillemylla/Crickets\\_Genome\\_Annotation](https://github.com/guillemylla/Crickets_Genome_Annotation).

## 362 Motif enrichment and gain of transcription

### 363 Motif datasets

364 The presence of specific DNA motifs before a gene is a major factor enabling transcription, we therefore  
365 searched for such motif enrichment upstream of *de novo* transcripts and control sequences, using  
366 custom python scripts along the Bio-python motifs (Cock et al., 2009) package. To that end, two  
367 motifs databases were downloaded as position frequency matrices (PFM) from JASPAR: the JASPAR  
368 Core insects (non redundant) database (Castro-Mondragon et al., 2022) and the JASPAR Pol II  
369 database (Fornes et al., 2020), containing 146 tF motifs of *D. melanogaster* and 13 core promoter  
370 motifs, respectively. While the JASPAR Core insects database was used to find general promoter and  
371 proximal enhancer motifs, the JASPAR Pol II database was restricted to the main core promoter motifs.  
372 PFM were used to calculate for each motif a position weight matrix (PWM). The PWM was then  
373 used to determine a position specific scoring matrix (PSSM). The absolute score of the PSSM was  
374 used to calculate the relative score of motif identity (Formula 1), which was then used as a threshold  
375 to determine motifs enrichment. Motifs with a relative score of motif identity superior or equal to  
376 0.8 to the PFM were considered to be enriched in a given sequence. Two types of motifs enrichment  
377 were defined: high similarity motif enrichment when a motif had a score above 0.95 and low similarity  
378 motif enrichment when its score was between 0.8 and 1. Motif enrichment was estimated for upstream  
379 sequences (1000 bp before transcript start and 100 bp after it) of *de novo* transcript, for upstream  
380 sequences of 'non-transcribed' homologs and for random intergenic sequences of 1100 bp (obtained

381 with bedtools (Quinlan and Hall, 2010), N = 53,300) as negative controls, and for upstream sequences  
382 of annotated genes as a positive control. We also restricted our upstream sequences, to 200 bp before  
383 a given sequence start and 100 bp after it, to estimate the core promoter binding motifs enrichment as  
384 those motifs are expected to be closer to the start of a transcript than general promoter and proximal  
385 enhancer (Butler and Kadonaga, 2002).

386 Formula 1:

Absolute score:

$$(pssm.max - pssm.min) * relativescore + pssm.min \quad (1)$$

Relative score:

$$(absolutescore - pssm.min) / (pssm.max - pssm.min) \quad (2)$$

387 All comparisons of transcripts with other sequence types were performed using Generalized Linear  
388 Mixed Models (GLMMs) using the package glmmTMB (Magnusson et al., 2017), retaining best model  
389 after simplifying model with a step-wise factor deletion.

### 390 **Transcripts vs. non-transcribed homologs**

391 To unravel the differences among sequences leading to transcription, four GLMMs were built using a  
392 binomial distribution. The first one assessed the importance of TE overlap, number and presence /  
393 absence in gaining transcription. This model includes as a dependent variable the type of sequence  
394 (transcript or 'non-transcribed' homolog), as fixed factors the relative overlap with TEs, the number of  
395 overlapping TEs, the presence or absence of overlapping TEs, the regions of the sequence (upstream,  
396 sequence, downstream), and their interactions. Moreover to account for pseudo-replication, the  
397 orthogroup ID of the sequence (single ID shared among transcript and non-transcribed homologs) and  
398 *D. melanogaster* line were added as random variables into the GLMM. A second model to account  
399 only for motifs enrichment was built with as fixed factors the number of minimal and optimal tF motifs  
400 and the minimal and optimal number of core promoter. A third model to account simultaneously for  
401 TEs and the different motifs was built by adding as fixed factor the number of the different motifs  
402 (motifs, cores, low and high). Finally, a fourth model was built to disentangle the impact of different  
403 TE classes (DNA vs. RNA transposon) on transcript and non-transcribed homologs, by adding the TE  
404 class as a fixed factor.

## 405 **Transcripts vs. genes and intergenic regions**

406 To understand how transcripts differ from genic and intergenic regions, three GLMMs were built. The  
407 first GLMM compares the relative overlap of transcripts with TEs with the different sequence types,  
408 using a zero-inflated Gamma distribution and as dependent variable: the sequence type, as fixed factor:  
409 the relative of overlap with TEs, and a random variable: the *D. melanogaster* line. The second GLMM  
410 compares the sequence types in term of motif numbers, using a poisson distribution and as dependent  
411 variable: the number of motifs / cores, as fixed factor: the sequence type, and a random variable:  
412 the *D. melanogaster* line. The third GLMM accounts for differences of sequence features among the  
413 different sequence types, using a zero-inflated Gamma distribution and as dependent variable: the  
414 sequence type, as fixed factor: the log TPM, the GC content, spliced length, and exon number, and  
415 a random variable: the *D. melanogaster* line. As the data-sets were of unequal sample size among  
416 the different sequence types and to ensure the robustness of our results, p-values of the best GLMM  
417 was bootstrapped using data-sets with equal sample size, using the package boot (Canty and Ripley,  
418 2017).

419 Furthermore, the density of *de novo* transcript per 100 kb was correlated to its distance to the center  
420 of the chromosome and the density of TEs, using GLMMs with as a dependent variable the number of  
421 *de novo* transcript within a 100kb window, as a random variables the chromosome and population, and  
422 as an explanatory variable the distance from the center of the chromosome (scaled) and the density  
423 of TE per 100kb (scaled), respectively. Furthermore, the levels of CpGoe of *de novo* transcript was  
424 correlated with their relative overlap with TEs, using a GLMM with CpGoe value as a dependent  
425 variable, length of overlap with a TE as explanatory variable and chromosome and population as a  
426 random variable.

## 427 **Visualisation**

428 All graphs and statistics were created with R version > 4.1 (Team, 2022). The packages dplyr  
429 (Wickham et al., 2022), tidyverse (Wickham et al., 2019) and data.table (Dowle and Srinivasan, 2021)  
430 were used for data preparation. The plots were mainly done with ggplot2 (Wickham et al., 2016) and  
431 its extensions ggpubr (Kassambara and Kassambara, 2020).

## 432 **Data access**

433 The files containing processed data is available in the Zenodo archive <https://doi.org/10.5281/zenodo.8403184>,  
434 and is referred in the main text as "Supplemental Deposit". Supplemental figures, information, analyses  
435 and models are found in the Supplementary Information (SI). All programs are stored on GitHub  
436 (<https://github.com/MarieLebh>).

## 437 **Competing interest statement**

438 The authors declare no competing interests.

## 439 **Author contributions**

440 MKL contribution : Data Curation, Formal Analysis, Investigation, Methodology, Writing and Re-  
441 viewing original draft BF contribution : Formal Analysis, Investigation, Methodology, Writing and  
442 Reviewing original draft JS contribution : Data Curation, Formal Analysis EBB contribution : Funding  
443 Acquisition, Reviewing original draft AG contribution : Conceptualization, Funding Acquisition, Project  
444 Administration, Supervision, Validation, Writing and Reviewing original draft.

## 445 **Funding Disclosure**

446 AG and EBB acknowledge funding from the Deutsche Forschungsgemeinschaft priority program  
447 "Genomic Basis of Evolutionary Innovations" (SPP 2349), project BO 2544/20-1 awarded to EBB and  
448 AG.

## 449 References

- 450 Altschul, S. F., Gish, W., Miller, W., Myers, E. W., and Lipman, D. J. Basic local alignment search tool.  
451 *Journal of molecular biology*, 215(3):403–410, 1990.
- 452 Baalsrud, H. T., Tørresen, O. K., Solbakken, M. H., Salzburger, W., Hanel, R., Jakobsen, K. S., and Jentoft,  
453 S. De novo gene evolution of antifreeze glycoproteins in codfishes revealed by whole genome sequence data.  
454 *Molecular biology and evolution*, 35(3):593–606, 2018.
- 455 Boeva, V. Analysis of genomic sequence motifs for deciphering transcription factor binding and transcriptional  
456 regulation in eukaryotic cells. *Frontiers in genetics*, 7:24, 2016.
- 457 Bornberg-Bauer, E., Schmitz, J., and Heberlein, M. Emergence of de novo proteins from 'dark genomic  
458 matter' by 'grow slow and moult'. *Biochemical Society Transactions*, 43(5):867–873, 2015.
- 459 Bornberg-Bauer, E., Hlouchova, K., and Lange, A. Structure and function of naturally evolved de novo proteins.  
460 *Current Opinion in Structural Biology*, 68:175–183, 2021.
- 461 Bourque, G., Burns, K. H., Gehring, M., Gorbunova, V., Seluanov, A., Hammell, M., Imbeault, M., Izsvák, Z.,  
462 Levin, H. L., Macfarlan, T. S., et al. Ten things you should know about transposable elements. *Genome*  
463 *biology*, 19:1–12, 2018.
- 464 Butler, J. E. and Kadonaga, J. T. The rna polymerase ii core promoter: a key component in the regulation of  
465 gene expression. *Genes & development*, 16(20):2583–2592, 2002.
- 466 Camilleri-Robles, C., Amador, R., Klein, C. C., Guigó, R., Corominas, M., and Ruiz-Romero, M. Genomic and  
467 functional conservation of lncrnas: lessons from flies. *Mammalian Genome*, pages 1–15, 2022.
- 468 Canty, A. and Ripley, B. Package 'boot'. *Bootstrap Functions. CRAN R Proj*, 2017.
- 469 Carvunis, A.-R., Rolland, T., Wapinski, I., Calderwood, M. A., Yildirim, M. A., Simonis, N., Charlotheaux, B.,  
470 Hidalgo, C. A., Barbette, J., Santhanam, B., et al. Proto-genes and de novo gene birth. *Nature*, 487(7407):  
471 370–374, 2012.
- 472 Castro-Mondragon, J. A., Riudavets-Puig, R., Rauluseviciute, I., Berhanu Lemma, R., Turchi, L., Blanc-  
473 Mathieu, R., Lucas, J., Boddie, P., Khan, A., Manosalva Pérez, N., et al. Jasp2022: the 9th release of  
474 the open-access database of transcription factor binding profiles. *Nucleic acids research*, 50(D1):D165–D173,  
475 2022.
- 476 Chuong, E. B., Rumi, M., Soares, M. J., and Baker, J. C. Endogenous retroviruses function as species-specific  
477 enhancer elements in the placenta. *Nature genetics*, 45(3):325–329, 2013.
- 478 Chuong, E. B., Elde, N. C., and Feschotte, C. Regulatory activities of transposable elements: from conflicts to  
479 benefits. *Nature Reviews Genetics*, 18(2):71–86, 2017.

- 480 Cock, P. J., Antao, T., Chang, J. T., Chapman, B. A., Cox, C. J., Dalke, A., Friedberg, I., Hamelryck, T.,  
481 Kauff, F., Wilczynski, B., et al. Biopython: freely available python tools for computational molecular biology  
482 and bioinformatics. *Bioinformatics*, 25(11):1422–1423, 2009.
- 483 Delihias, N. An ancestral genomic sequence that serves as a nucleation site for de novo gene birth. *Plos one*,  
484 17(5):e0267864, 2022.
- 485 Delprat, A., Negre, B., Puig, M., and Ruiz, A. The transposon galileo generates natural chromosomal inversions  
486 in drosophila by ectopic recombination. *PLoS One*, 4(11):e7883, 2009.
- 487 Dixon, G. and Matz, M. Invertebrate dna methylation and gene regulation. *bioRxiv*, pages 2021–11, 2021.
- 488 Dowle, M. and Srinivasan, A. data.table: Extension of “data.frame”[r package data.table version 1.14.2],  
489 2021.
- 490 Dunwell, T. L. and Pfeifer, G. P. Drosophila genomic methylation: new evidence and new questions. *Epigenomics*,  
491 6(5):459–461, 2014.
- 492 Durand, É., Gagnon-Arsenault, I., Hallin, J., Hatin, I., Dubé, A. K., Nielly-Thibault, L., Namy, O., and Landry,  
493 C. R. Turnover of ribosome-associated transcripts from de novo orfs produces gene-like characteristics  
494 available for de novo gene emergence in wild yeast populations. *Genome research*, 29(6):932–943, 2019.
- 495 Everett, L. J., Huang, W., Zhou, S., Carbone, M. A., Lyman, R. F., Arya, G. H., Geisz, M. S., Ma, J., Morgante,  
496 F., Armour, G. S., et al. Gene expression networks in the drosophila genetic reference panel. *Genome*  
497 *research*, 30(3):485–496, 2020.
- 498 Fornes, O., Castro-Mondragon, J. A., Khan, A., Van der Lee, R., Zhang, X., Richmond, P. A., Modi, B. P.,  
499 Correard, S., Gheorghe, M., Baranašić, D., et al. Jaspas 2020: update of the open-access database of  
500 transcription factor binding profiles. *Nucleic acids research*, 48(D1):D87–D92, 2020.
- 501 Fort, V., Khelifi, G., and Hussein, S. M. Long non-coding rnas and transposable elements: A functional  
502 relationship. *Biochimica et Biophysica Acta (BBA)-Molecular Cell Research*, 1868(1):118837, 2021.
- 503 Ganesh, S. and Svoboda, P. Retrotransposon-associated long non-coding rnas in mice and men. *Pflügers*  
504 *Archiv-European Journal of Physiology*, 468(6):1049–1060, 2016.
- 505 Grandchamp, A., Kühl, L., Lebherz, M. K., Brüggemann, K., Parsch, J., and Bornberg-Bauer, E. Population  
506 genomics reveals mechanisms and dynamics of de novo proto-gene emergence in drosophila melanogaster.  
507 *bioRxiv*, 2022.
- 508 Grandchamp, A., Czuppon, P., and Bornberg-Bauer, E. High turnover of de novo transcripts in drosophila  
509 melanogaster. *bioRxiv*, pages 2023–02, 2023a.



- 510 Grandchamp, A., Kühl, L., Lebherz, M., Brüggemann, K., Parsch, J., and Bornberg-Bauer, E. Population  
511 genomics reveals mechanisms and dynamics of de novo expressed open reading frame emergence in *Drosophila*  
512 *melanogaster*. *Genome Research*, 33(6):872–890, 2023b.
- 513 Gubala, A. M., Schmitz, J. F., Kearns, M. J., Vinh, T. T., Bornberg-Bauer, E., Wolfner, M. F., and Findlay,  
514 G. D. The goddard and saturn genes are essential for *Drosophila* male fertility and may have arisen de novo.  
515 *Molecular biology and evolution*, 34(5):1066–1082, 2017.
- 516 Guerzoni, D. and McLysaght, A. De novo origins of human genes. *PLoS genetics*, 7(11):e1002381, 2011.
- 517 Haberle, V. and Stark, A. Eukaryotic core promoters and the functional basis of transcription initiation. *Nature*  
518 *reviews Molecular cell biology*, 19(10):621–637, 2018.
- 519 Huang, W., Carbone, M. A., Magwire, M. M., Peiffer, J. A., Lyman, R. F., Stone, E. A., Anholt, R. R., and  
520 Mackay, T. F. Genetic basis of transcriptome diversity in *Drosophila melanogaster*. *Proceedings of the*  
521 *National Academy of Sciences*, 112(44):E6010–E6019, 2015.
- 522 Iyengar, B. R. and Bornberg-Bauer, E. Neutral models of de novo gene emergence suggest that gene evolution  
523 has a preferred trajectory. *Molecular Biology and Evolution*, 40(4):msad079, 2023.
- 524 Kapusta, A. and Feschotte, C. Volatile evolution of long noncoding rna repertoires: mechanisms and biological  
525 implications. *Trends in Genetics*, 30(10):439–452, 2014.
- 526 Kapusta, A., Kronenberg, Z., Lynch, V. J., Zhuo, X., Ramsay, L., Bourque, G., Yandell, M., and Feschotte, C.  
527 Transposable elements are major contributors to the origin, diversification, and regulation of vertebrate long  
528 noncoding rnas. *PLoS genetics*, 9(4):e1003470, 2013.
- 529 Kassambara, A. and Kassambara, M. Package 'ggpubr'. r package version 0.3. 5, 2020.
- 530 Kim, T.-K. and Shiekhattar, R. Architectural and functional commonalities between enhancers and promoters.  
531 *Cell*, 162(5):948–959, 2015.
- 532 Klughammer, J., Romanovskaia, D., Nemeš, A., Posautz, A., Seid, C. A., Schuster, L. C., Keinath, M. C.,  
533 Lugo Ramos, J. S., Kosack, L., Evankow, A., et al. Comparative analysis of genome-scale, base-resolution  
534 dna methylation profiles across 580 animal species. *Nature communications*, 14(1):232, 2023.
- 535 Kofler, R., Nolte, V., and Schlötterer, C. Tempo and mode of transposable element activity in *Drosophila*.  
536 *PLoS Genetics*, 11(7):e1005406, 2015.
- 537 Kordyukova, M., Olovnikov, I., and Kalmykova, A. Transposon control mechanisms in telomere biology. *Current*  
538 *opinion in genetics & development*, 49:56–62, 2018.
- 539 Kurafeiski, J. D., Pinto, P., and Bornberg-Bauer, E. Evolutionary potential of cis-regulatory mutations to cause  
540 rapid changes in transcription factor binding. *Genome Biology and Evolution*, 11(2):406–414, 2019.

- 541 Law, J. A. and Jacobsen, S. E. Establishing, maintaining and modifying dna methylation patterns in plants  
542 and animals. *Nature Reviews Genetics*, 11(3):204–220, 2010.
- 543 Lawlor, M. A., Cao, W., and Ellison, C. E. A transposon expression burst accompanies the activation of  
544 y-chromosome fertility genes during drosophila spermatogenesis. *Nature communications*, 12(1):6854, 2021.
- 545 Lawson, H. A., Liang, Y., and Wang, T. Transposable elements in mammalian chromatin organization. *Nature*  
546 *Reviews Genetics*, pages 1–12, 2023.
- 547 Lyko, F., Ramsahoye, B. H., and Jaenisch, R. Dna methylation in drosophila melanogaster. *Nature*, 408(6812):  
548 538–540, 2000.
- 549 Magnusson, A., Skaug, H., Nielsen, A., Berg, C., Kristensen, K., Maechler, M., van Bentham, K., Bolker, B.,  
550 Brooks, M., and Brooks, M. M. Package 'glmmtmb'. *R Package Version 0.2. 0*, 25, 2017.
- 551 Majic, P. and Payne, J. L. Enhancers facilitate the birth of de novo genes and gene integration into regulatory  
552 networks. *Molecular biology and evolution*, 37(4):1165–1178, 2020.
- 553 Malik, H. S. Retroviruses push the envelope for mammalian placentation. *Proceedings of the National Academy*  
554 *of Sciences*, 109(7):2184–2185, 2012.
- 555 McCullers, T. J. and Steiniger, M. Transposable elements in drosophila. *Mobile genetic elements*, 7(3):1–18,  
556 2017.
- 557 McLysaght, A. and Hurst, L. D. Open questions in the study of de novo genes: what, how and why. *Nature*  
558 *Reviews Genetics*, 17(9):567–578, 2016.
- 559 Meers, M. P., Adelman, K., Duronio, R. J., Strahl, B. D., McKay, D. J., and Matera, A. G. Transcription  
560 start site profiling uncovers divergent transcription and enhancer-associated rnas in drosophila melanogaster.  
561 *BMC genomics*, 19:1–20, 2018.
- 562 Mérel, V., Boulesteix, M., Fablet, M., and Vieira, C. Transposable elements in drosophila. *Mobile DNA*, 11:  
563 1–20, 2020.
- 564 Moschetti, R., Palazzo, A., Lorusso, P., Viggiano, L., and Massimiliano Marsano, R. “what you need, baby, i  
565 got it”: Transposable elements as suppliers of cis-operating sequences in drosophila. *Biology*, 9(2):25, 2020.
- 566 Neme, R. and Tautz, D. Fast turnover of genome transcription across evolutionary time exposes entire  
567 non-coding dna to de novo gene emergence. *elife*, 5:e09977, 2016.
- 568 Palazzo, A. F. and Lee, E. S. Non-coding rna: what is functional and what is junk? *Frontiers in genetics*, 6:2,  
569 2015.
- 570 Petrov, D. A., Fiston-Lavier, A.-S., Lipatov, M., Lenkov, K., and González, J. Population genomics of  
571 transposable elements in drosophila melanogaster. *Molecular biology and evolution*, 28(5):1633–1644, 2011.

- 572 Quinlan, A. R. and Hall, I. M. Bedtools: a flexible suite of utilities for comparing genomic features. *Bioinformatics*,  
573 26(6):841–842, 2010.
- 574 Riehl, K., Riccio, C., Miska, E. A., and Hemberg, M. Transposonultimate: software for transposon classification,  
575 annotation and detection. *Nucleic Acids Research*, 50(11):e64–e64, 2022.
- 576 Rizzon, C., Marais, G., Gouy, M., and Biémont, C. Recombination rate and the distribution of transposable  
577 elements in the drosophila melanogaster genome. *Genome research*, 12(3):400–407, 2002.
- 578 Rödelsperger, C., Prabh, N., and Sommer, R. J. New gene origin and deep taxon phylogenomics: opportunities  
579 and challenges. *Trends in Genetics*, 35(12):914–922, 2019.
- 580 Roder, K., Hung, M.-S., Lee, T.-L., Lin, T.-Y., Xiao, H., Isobe, K.-I., Juang, J.-L., and Shen, C.-K. J.  
581 Transcriptional repression by drosophila methyl-cpg-binding proteins. *Molecular and Cellular Biology*, 20(19):  
582 7401–7409, 2000.
- 583 Rollins, R. A., Haghghi, F., Edwards, J. R., Das, R., Zhang, M. Q., Ju, J., and Bestor, T. H. Large-scale  
584 structure of genomic methylation patterns. *Genome research*, 16(2):157–163, 2006.
- 585 Schlötterer, C. Genes from scratch—the evolutionary fate of de novo genes. *Trends in Genetics*, 31(4):215–219,  
586 2015.
- 587 Scruggs, B. S., Gilchrist, D. A., Nechaev, S., Muse, G. W., Burkholder, A., Fargo, D. C., and Adelman, K.  
588 Bidirectional transcription arises from two distinct hubs of transcription factor binding and active chromatin.  
589 *Molecular cell*, 58(6):1101–1112, 2015.
- 590 Small, S. and Arnosti, D. N. Transcriptional enhancers in drosophila. *Genetics*, 216(1):1–26, 2020.
- 591 Tautz, D. and Domazet-Lošo, T. The evolutionary origin of orphan genes. *Nature Reviews Genetics*, 12(10):  
592 692–702, 2011.
- 593 Team, R. D. C. A language and environment for statistical computing. <http://www.R-project.org>, 2022.
- 594 Thybert, D., Roller, M., Navarro, F. C., Fiddes, I., Streeter, I., Feig, C., Martin-Galvez, D., Kolmogorov, M.,  
595 Janoušek, V., Akanni, W., et al. Repeat associated mechanisms of genome evolution and function revealed  
596 by the mus caroli and mus pahari genomes. *Genome research*, 28(4):448–459, 2018.
- 597 Van Oss, S. B. and Carvunis, A.-R. De novo gene birth. *PLoS genetics*, 15(5):e1008160, 2019.
- 598 Wang, M., Wang, D., Zhang, K., Ngo, V., Fan, S., and Wang, W. Motto: representing motifs in consensus  
599 sequences with minimum information loss. *Genetics*, 216(2):353–358, 2020a.
- 600 Wang, Y.-W., Hess, J., Slot, J. C., and Pringle, A. De novo gene birth, horizontal gene transfer, and gene  
601 duplication as sources of new gene families associated with the origin of symbiosis in amanita. *Genome  
602 biology and evolution*, 12(11):2168–2182, 2020b.

- 603 Wickham, H., Averick, M., Bryan, J., Chang, W., McGowan, L., François, R., Golemund, G., Hayes, A., Henry,  
604 L., Hester, J., et al. Welcome to the tidyverse. *J open source softw* 4: 1686, 2019.
- 605 Wickham, H., François, R., Henry, L., and Müller, K. Rstudio.(2021). dplyr: A grammar of data manipulation  
606 (1.0. 7), 2022.
- 607 Wickham, H., Chang, W., and Wickham, M. H. Package 'ggplot2'. *Create elegant data visualisations using*  
608 *the grammar of graphics. Version, 2(1):1–189*, 2016.
- 609 Ylla, G., Nakamura, T., Itoh, T., Kajitani, R., Toyoda, A., Tomonari, S., Bando, T., Ishimaru, Y., Watanabe,  
610 T., Fuketa, M., et al. Insights into the genomic evolution of insects from cricket genomes. *Communications*  
611 *biology*, 4(1):733, 2021.
- 612 Zhao, L., Saelao, P., Jones, C. D., and Begun, D. J. Origin and spread of de novo genes in drosophila  
613 melanogaster populations. *Science*, 343(6172):769–772, 2014.
- 614 Zhou, W., Liang, G., Molloy, P. L., and Jones, P. A. Dna methylation enables transposable element-driven  
615 genome expansion. *Proceedings of the National Academy of Sciences*, 117(32):19359–19366, 2020.
- 616 Zhuang, X. and Cheng, C.-H. C. Propagation of a de novo gene under natural selection: Antifreeze glycoprotein  
617 genes and their evolutionary history in codfishes. *Genes*, 12(11):1777, 2021.