

1 **Developmental dynamics are a proxy for selective pressures on**  
2 **alternatively polyadenylated isoforms**

3

4

5 Michal Levin<sup>1\*</sup>, Harel Zalts<sup>2\*</sup>, Natalia Mostov<sup>2\*</sup>, Tamar Hashimshony<sup>2</sup>, and Itai Yanai<sup>3,4</sup>

6

7 <sup>1</sup> Quantitative Proteomics, Institute of Molecular Biology, Mainz, Germany

8 <sup>2</sup> Faculty of Biology, Technion – Israel Institute of Technology, Haifa, Israel

9 <sup>3</sup> Institute for Computational Medicine, NYU School of Medicine, USA

10 <sup>4</sup> Corresponding author: Itai.Yanai@nyumc.org

11 \* These authors contributed equally to this work.

12

13 Running title: RNA-seq characterization of alternatively polyadenylated isoforms in individual  
14 embryos

15 Keywords: Alternative polyadenylation, APA-seq, *C. elegans* embryogenesis

## 16 **Abstract**

17

18 Alternative polyadenylation (APA) leads to multiple transcripts from the same gene, yet their  
19 distinct functional attributes remain largely unknown. Here, we introduce APA-seq to detect the  
20 expression levels of APA isoforms from 3'-end RNA-Seq data by exploiting both paired-end  
21 reads for gene isoform identification and quantification. Applying APA-seq, we detected the  
22 expression levels of APA isoforms from RNA-Seq data of single *C. elegans* embryos, and  
23 studied the patterns of 3' UTR isoform expression throughout embryogenesis. We found that  
24 global changes in APA usage demarcate developmental stages, suggesting a requirement for  
25 distinct 3' UTR isoforms throughout embryogenesis. We distinguished two classes of genes,  
26 depending upon the correlation between the temporal profiles of their isoforms: those with  
27 highly correlated isoforms (HCI) and those with lowly correlated isoforms (LCI) across time.  
28 This led us to hypothesize that variants produced with similar expression profiles may be the  
29 product of biological noise, while the LCI variants may be under tighter selection and  
30 consequently their distinct 3' UTR isoforms are more likely to have functional consequences.  
31 Supporting this notion, we found that LCI genes have significantly more miRNA binding sites,  
32 more correlated expression profiles with those of their targeting miRNAs and a relative lack of  
33 correspondence between their transcription and protein abundances. Collectively, our results  
34 suggest that a lack of coherence among the regulation of 3' UTR isoforms is a proxy for  
35 selective pressures acting upon APA usage and consequently for their functional relevance.

36

37

## 38 **Introduction**

39

40 Alternative polyadenylation (APA) is a crucial regulatory mechanism – widespread and  
41 conserved across all eukaryotes – that diversifies post-transcriptional regulation by selective  
42 mRNA-miRNA and mRNA-protein interactions (Ozsolak et al. 2010; Jan et al. 2011; Derti et al.  
43 2012; Smibert et al. 2012; Ulitsky et al. 2012; Velten et al. 2015; Hu et al. 2017). APA plays an  
44 important role in a vast variety of biological processes, such as maternal to zygotic transition,  
45 cell differentiation, and tissue specification (Wormington 1994; Tadros et al. 2007b; Tadros and

46 Lipshitz 2009) and exerts a tremendous influence over gene expression, the transcript's  
47 cellular localization, stability and translation rate (Mazumder et al. 2003; Keene 2007; Lutz and  
48 Moreira 2011; Elkon et al. 2013; Tian and Manley 2017). Since first discovered in the  
49 immunoglobulin M and the DHFR genes (Alt et al. 1980; Early et al. 1980; Rogers et al. 1980;  
50 Setzer et al. 1980), technological advances and high-throughput sequencing techniques have  
51 made it clear that APA is more widespread than initially thought; 30-70% of the genes undergo  
52 APA in diverse species (Mangone et al. 2010; Oszolak et al. 2010; Jan et al. 2011; Derti et al.  
53 2012; Smibert et al. 2012; Ulitsky et al. 2012; Velten et al. 2015). During the last decade,  
54 widespread APA alterations were detected across different tissues and across distinct stages  
55 of embryogenesis (Tian et al. 2005; Tadros et al. 2007a; Wang et al. 2008; Ji et al. 2009; Li et  
56 al. 2012; Smibert et al. 2012; Ulitsky et al. 2012; Blazie et al. 2015; Blazie et al. 2017; Hu et al.  
57 2017; Khraiweh and Salehi-Ashtiani 2017). However, beyond these classifications and  
58 functional characterization of a short list of single gene APA alterations (Chen et al. 2017) the  
59 global functional significance of APA remains largely elusive. In light of the vast amount of  
60 alternative 3'UTR isoforms that have been detected during the past few years it would be  
61 beneficial to be able to enrich for those alterations which have functional consequences.

62

63 *C. elegans* is a convenient model organism for studying APA and gene expression during  
64 embryogenesis, since the cell lineage is invariant and has been fully traced (Sulston and  
65 Horvitz 1977; Kimble and Hirsh 1979). *C. elegans* was the first multicellular organism to have a  
66 fully sequenced genome (Consortium 1998); and its transcriptome has also been well  
67 characterized (McKay et al. 2003; Shin et al. 2008; Hillier et al. 2009; Ramani et al. 2009;  
68 Lamm et al. 2011; Grün et al. 2014). During embryonic development, the *C. elegans*  
69 transcriptome is highly dynamic; in early stages it is comprised mostly of maternal transcripts,  
70 but as development proceeds, zygotic transcription commences, and maternally supplied  
71 transcripts undergo degradation (Newman-Smith and Rothman 1998; Tadros et al. 2007b;  
72 Tadros and Lipshitz 2009; Walser and Lipshitz 2011). Our previous results detected many  
73 genes with a dynamic overall gene expression profile throughout embryogenesis, as well as  
74 genes with constitutive levels of expression (Levin et al. 2012).

75

76 CEL-Seq is a sensitive multiplexed single-cell RNA-Seq method (Hashimshony et al. 2012;  
77 Hashimshony et al. 2016). One important feature of the CEL-Seq method is that it is restricted  
78 to studying the 3' end of the transcriptome and thus measures overall expression levels;  
79 typically collapsing the various isoforms produced by a gene to one summary profile. While this  
80 has been a useful simplifying criterion, it does ignore possible dynamic profiles across different  
81 splicing isoforms of a particular gene. An important advantage of the CEL-Seq 3' end bias  
82 though is that this information can be used to detect and quantify alternative polyadenylation  
83 patterns *i.e.* 3' UTR isoforms of the same gene, as we propose here.

84

85 Here, we studied APA profiles in individual *C. elegans* throughout embryogenesis using APA-  
86 seq, an approach to detect alternative polyadenylation profiles at a genomic scale. APA-seq is  
87 based on CEL-Seq but further exploits the information from both paired-end reads allows for  
88 gene identification from Read 2 and the exact location of polyadenylation from Read 1.  
89 Combining this information empowers quantitative expression level assessment globally at  
90 both the gene and APA isoform level. Using this approach, we delineated two groups of genes,  
91 those with highly and lowly correlated 3' UTR isoform groups (HCI and LCI), respectively. We  
92 detected unique regulatory features between these groups which supports the notion that  
93 variants across these two groups are under distinct selective biases. Genes with uncorrelated  
94 3' UTR isoform expression (LCI) are predicted to have the highest miRNA regulation compared  
95 to genes with well-correlated 3' UTR isoforms. Integrating extensive previously published  
96 embryonic mRNA and protein expression datasets (Mangone et al. 2008; Grün et al. 2014), we  
97 also found the lowest correlation between total mRNA transcript and protein levels in genes  
98 with dynamic 3' UTR isoform expression (LCI). Extending this analysis to *Drosophila*  
99 *melanogaster* and *Xenopus laevis* datasets we found a consistent relationship (Graveley et al.  
100 2011; Casas-Vila et al. 2017; Sanfilippo et al. 2017; Zhou et al. 2019; Peshkin et al. 2015).  
101 Together, our results suggest that genes of the HCI and LCI groups experience distinct  
102 regulatory pressures upon their alternatively polyadenylated isoforms.

103

104

## 105 Results

### 106 APA-seq identifies expression levels of alternative polyadenylation isoforms

107 Paired-end reads generated by CEL-Seq and CEL-Seq2 contain the sample-specific barcode  
108 on Read 1 and the sequence identifying the transcript on Read 2 (Hashimshony et al. 2012;  
109 Hashimshony et al. 2016) (Fig. 1A). Typically, in CEL-Seq only Read 2 is used for measuring  
110 gene expression levels. However, since CEL-Seq sequences the 3' ends, it can be used in  
111 principle to identify 3' UTR isoforms. Using Read 2 for this purpose is challenging due to the  
112 uneven sizes of the inserts in the sequencing library, producing a smear of mapped reads,  
113 which makes distinguishing different 3' UTR isoforms of the same transcript impossible in  
114 many cases (Fig. 1B, red peak). However, we noted that Read 1 includes the actual 3' end of  
115 the transcript (Fig. 1A), located just upstream of the polyadenylation site (Fig. 1B, black peak).  
116 In CEL-Seq, this region follows 24 Ts used for capturing the polyA tail of the transcript, and the  
117 sequencing quality is relatively poor after this low complexity region rendering conventional  
118 mapping impossible (Supplemental Fig. S1A,B). To overcome this, we found that the  
119 sequence is still of sufficient quality when mapping is performed using relaxed parameters and  
120 restricted to a particular region of the genome. In summary, while permissive Read 1 mapping  
121 matches many genomic loci due to its poor quality, it maps uniquely when restricted to the  
122 sequence of a particular gene, whose identity is detected by using Read 2.

123 We thus devised the APA-seq approach to study 3' UTR isoforms using both Read 1 and Read  
124 2 information and applied it to study the expression of alternative polyadenylated isoforms  
125 during early embryogenesis in the nematode *C. elegans*. For this, we used a dataset  
126 previously published by our lab in which embryos were individually collected and sequenced  
127 throughout embryogenesis (Levin et al. 2016). Pooling together all samples, we detected  
128 distinct 3' UTR isoforms for 5,336 genes using APA-seq, after removing possible artifacts  
129 caused by internal priming (Gruber et al. 2016) (see Methods). For example, the *C. elegans*  
130 genes *gmeb-1* and *rnr-1* (Fig. 1C) show a smear of Read 2 mappings (shown in red), while  
131 Read 1 mappings form distinct peaks (shown in black) identifying previously characterized  
132 polyadenylation sites (shown in green) (Mangone et al. 2008) (see also Supplemental Fig.  
133 S1C).

134 Overall, we detected multiple 3' UTR isoforms for 14% of the expressed genes in our dataset  
135 (746 out of 5336 genes). Of these, less than 1% have more than two isoforms (Supplemental  
136 Fig. S1D). This set is not comprehensive to all possible 3' UTR isoforms produced by the *C.*  
137 *elegans* genome given that we only examined 430 minutes during embryogenesis, and the  
138 stringent thresholds set in our bioinformatics pipeline (in terms of mapping parameters,  
139 mapping level filtering and spurious site removal, see Methods). We assayed the accuracy of  
140 the isoforms detected by comparing our polyadenylation sites with a known repository of 3'  
141 UTR annotations in *C. elegans* (Mangone et al. 2008) and found highly concordant profiles,  
142 with 95% of sites corresponding to well-established annotated sites (Fig. 1D, Supplemental  
143 Table S1). The remaining 5% show significantly lower expression levels than the annotation  
144 overlapping 3' UTR isoforms (Supplemental Fig. S1E) and we therefore excluded these from  
145 further analyses by expression level filtering. To study the temporal dynamics of 3' UTR  
146 isoform expression, we classified the reads according to their embryo of origin (Fig. 1E,  
147 Supplemental Table S1). As an example, two isoforms were identified for *rps-29*, which  
148 encodes a ribosomal protein subunit (Kamath et al. 2003). Interestingly, while the sum of  
149 expression of both *rps-29* isoforms is roughly constant over time, the long 3' UTR isoform is  
150 predominantly expressed in the early stages while the shorter is expressed in later embryos  
151 (Fig. 1E). Correlating the expression levels of all 3' UTR isoforms across stages, we found that  
152 successive stages (near-replicates) show high correlations thus highlighting the reproducibility  
153 of the data (Supplemental Fig. S2). We conclude that while CEL-Seq is typically used to assay  
154 overall expression levels with the mapping location typically ignored, processing the reads  
155 using APA-seq can identify the exact locations of alternative polyadenylation sites and the  
156 expression levels of distinct 3' UTR isoforms.

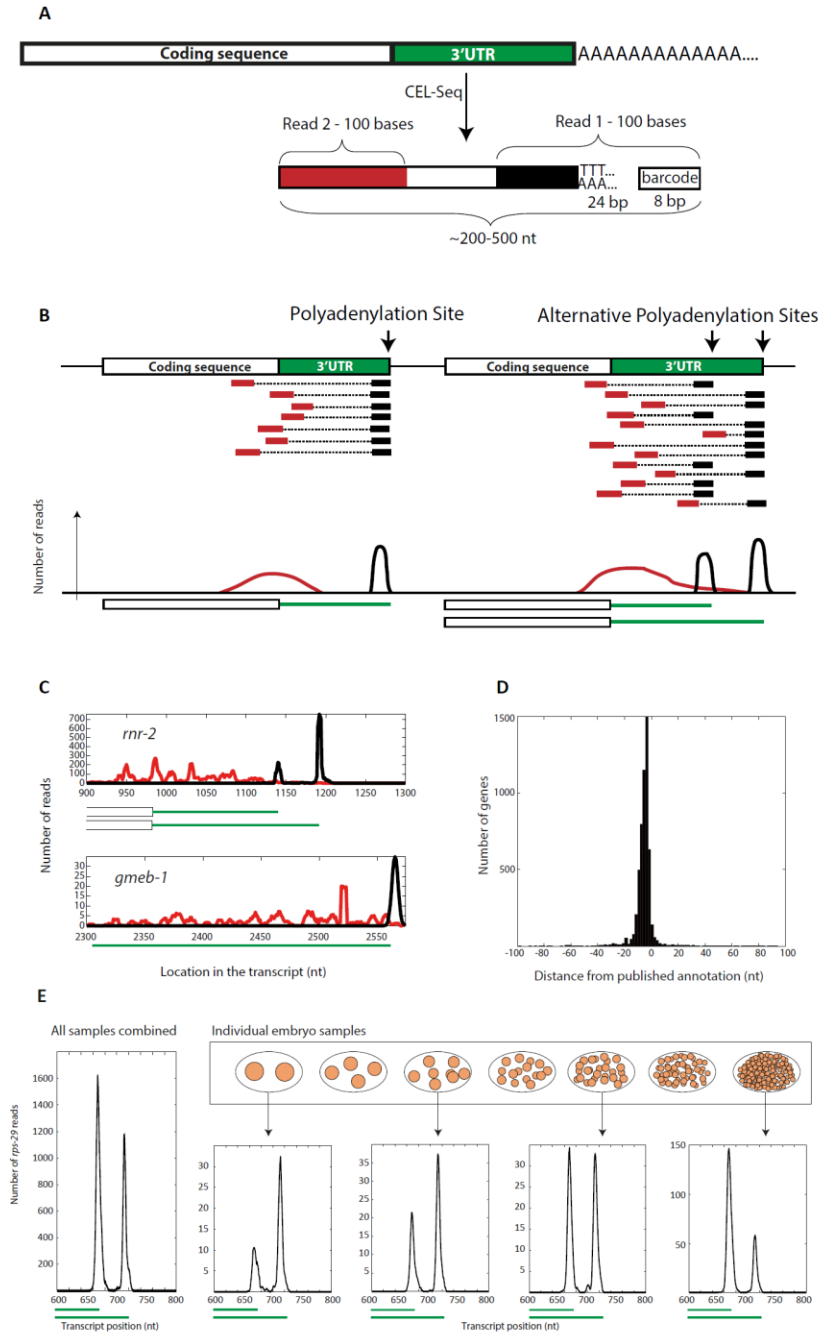


Figure 1

157

158 **Fig. 1.** APA-seq measures expression levels of distinct 3' UTR isoforms in individual *C.*  
 159 *elegans* embryos. (A) APA-seq identifies gene expression levels for distinct alternatively  
 160 polyadenylated isoforms. APA-seq is an adaptation of the CEL-Seq method which utilizes  
 161 paired-end reads: Read 1 contains a sample-specific barcode while Read 2 identifies the  
 162 transcript. If sequenced long enough (100bp in our case), Read 1 also provides information on  
 163 the exact location of the polyadenylation site. APA-seq thus uses Read 2 to identify the  
 164 expressed gene and then maps Read 1 to the gene specific region, thus enabling unique

165 mapping in spite of the low sequencing quality that results from sequencing through the low-  
166 complexity poly-T region. (B) Mapping Read 2 sequences results in a wide distribution within  
167 the gene (red peak) due to the fragmentation step of library preparation. However, Read 1  
168 sequences all map to the site immediately upstream of the poly-A tail thus producing a clear  
169 peak (black peak) when mapped to the gene sequence and revealing exact polyadenylation  
170 sites. The white boxes at the bottom of the distribution plots mark the coding sequence, while  
171 the green lines indicate the determined 3' UTR regions. (C) Using the APA-seq method  
172 enables detection of the exact location of polyadenylation sites in *C. elegans*. The green lines  
173 at the bottom of each plot mark previously annotated 3' UTRs (Mangone et al. 2008) for the  
174 indicated gene, showing good agreement with the APA-seq Read 1 peaks (black). (D) Global  
175 comparison of the detected polyadenylation sites using APA-seq to the *C. elegans* UTRome  
176 annotation (Mangone et al. 2008) shows high consistency between the two datasets. (E) The  
177 expression of two unique 3' UTR isoforms for the *C. elegans rps-29* gene throughout  
178 embryogenesis. Although in total (leftmost panel) the 3' UTR variants show equal expression,  
179 examining expression across developmental stages shows predominant expression of the  
180 shorter 3' UTR variant during early embryogenesis, and the inverse later in development.

181

## 182 **Alternative polyadenylation profiles throughout *C. elegans* embryogenesis**

183 Our dataset enabled us to study overall patterns of 3' UTR isoform usage throughout early  
184 development in individual embryos. For each gene we first computed the ratio of expression  
185 between its 3' UTR isoforms throughout the time-course. Interestingly, the clustering of stages  
186 according to these profiles revealed that adjacent developmental stages show concordant 3'  
187 UTR isoform usage patterns which group into distinct periods with diverging APA dynamics  
188 (Fig. 2A). The four groups correspond to previously characterized developmental periods:  
189 maternal degradation, early period of extensive proliferation and specification, the mid-  
190 embryonic transition period and the subsequent period of morphogenesis (47). The  
191 observation that different periods have different corresponding patterns of 3' UTR isoforms  
192 may reflect that each of these has a distinct functional requirement. Between adjacent periods  
193 we found that the direction and level of change of polyadenylation site usage generally shows  
194 a broad burst of shortening especially between the proliferative and mid-embryonic transition  
195 periods (Fig. 2B). This is consistent with previous work indicating that proliferative states show  
196 an overall trend for 3' UTR shortening (Tian and Manley 2017). Overall, we identified that 89%  
197 of dynamic 3' UTR isoform genes show APA switches in at least one of the period switches  
198 indicated (Fig. 2C).



199 Clustering the 3' UTR isoform ratio profiles throughout the time-course, we identified four main  
200 distinct clusters (Fig. 2D); genes whose 3' UTRs elongate or shorten over time (Cluster 1 and  
201 4, respectively) and genes whose 3' UTRs shorten transiently during the proliferation period  
202 (Cluster 2 and 3). The majority of genes though show continuous shortening of their 3' UTR  
203 regions with progressing development (cluster 4). Interestingly, these genes show enrichments  
204 for MAP kinase cascade, morphogenesis and neuronal differentiation (Supplemental Fig. S3).  
205 All these functions are crucial components of the switch from germ cell biology to proliferative  
206 and differentiation processes. In summary, we show that the 3' UTR dynamics detected by  
207 APA-seq reflect characteristic functionalities of embryonic development and although the  
208 significance of single events is not clear yet the overall contribution of APA events to the  
209 general regulatory states involved in these events might be high or alternatively present a side  
210 effect of the vast changes occurring at the other levels of gene expression regulation.

211

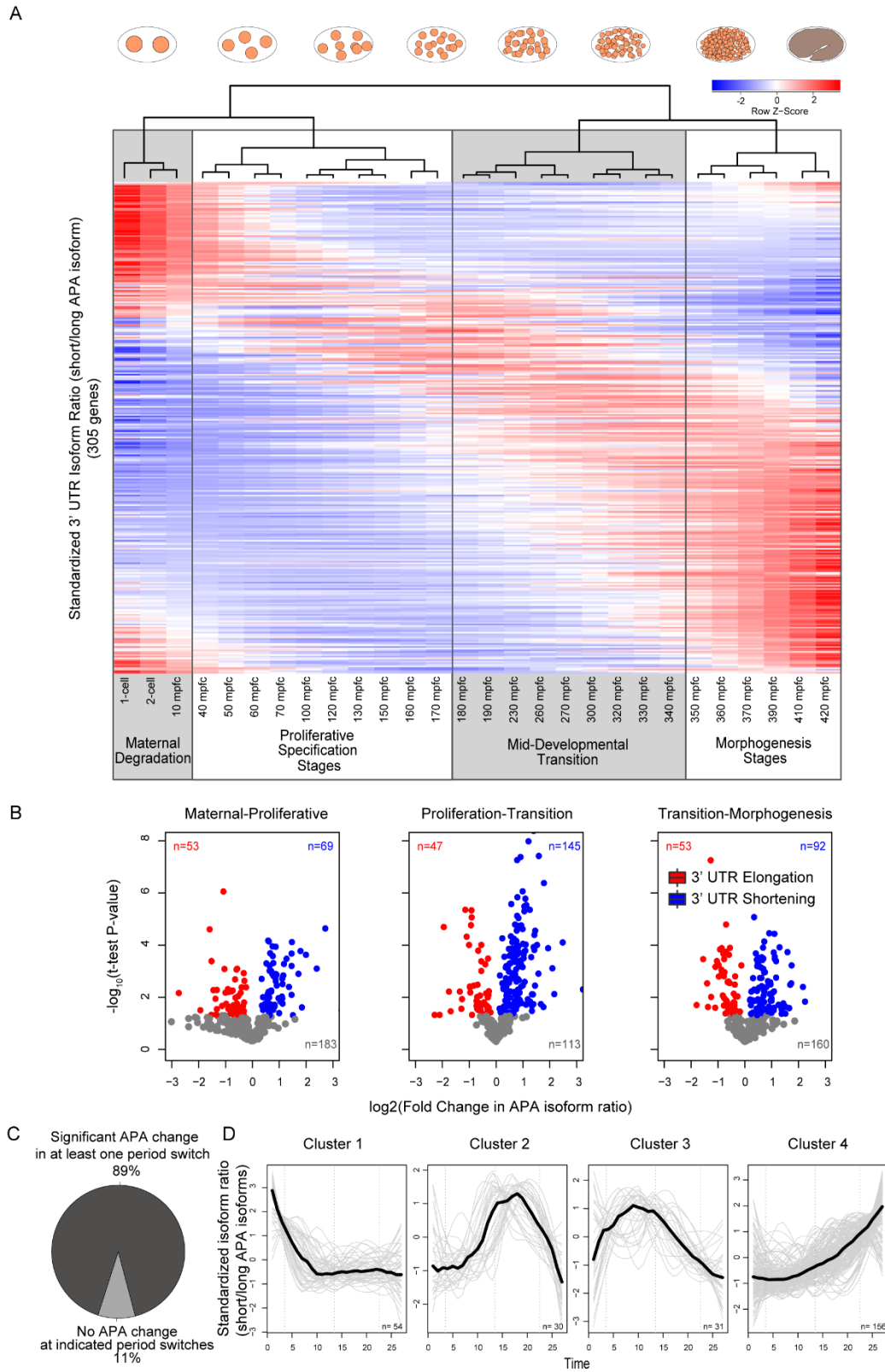


Figure 2

213 **Fig. 2.** 3' UTR isoform expression throughout *C. elegans* embryogenesis. (A) Heatmap  
214 showing the relative expression of 3' UTR isoforms for 305 genes which passed overall  
215 expression and 3' UTR isoform dynamics threshold throughout *C. elegans* embryogenesis.  
216 Each row in the heatmap corresponds to a gene and indicates the ratio between the  
217 expression of its short and long isoforms. Red and blue indicate the maximum and minimum 3'  
218 UTR ratio for each gene, respectively. White and grey shadowed boxes indicate sets of  
219 developmental stages with similar isoform usage (based on Ward clustering, see clustergram  
220 on top of the heatmap) and identify the periods of major isoform switches. mpfc = minutes past  
221 four-cell stage. (B) Studying the level of 3' UTR isoform changes across development. The  
222 plots indicate the fold-change and *P*-value of difference in the isoform ratios for pairs of  
223 successive periods as defined by Ward clustering in A. Red and blue coloring indicates  
224 significant elongation and shortening events, respectively. Grey coloring indicates insignificant  
225 changes. (C) Most of significant 3' UTR isoform switches occur between the identified  
226 developmental switch periods. 89% (271) of the genes show significant changes in 3' UTR  
227 usage at least one of the identified period switches. (D) Clusters of significant 3' UTR isoform  
228 changes indicate four main kinetics – almost constant elongation or shortening over time  
229 (Cluster 1 and 4, respectively) or peaking shortening during mid-developmental transition and  
230 proliferative periods (Cluster 3 and 4). The thick black line indicates the mean 3' UTR isoform  
231 ratio profile of all genes in the cluster and individual genes profiles are shown as thin grey  
232 lines.

233

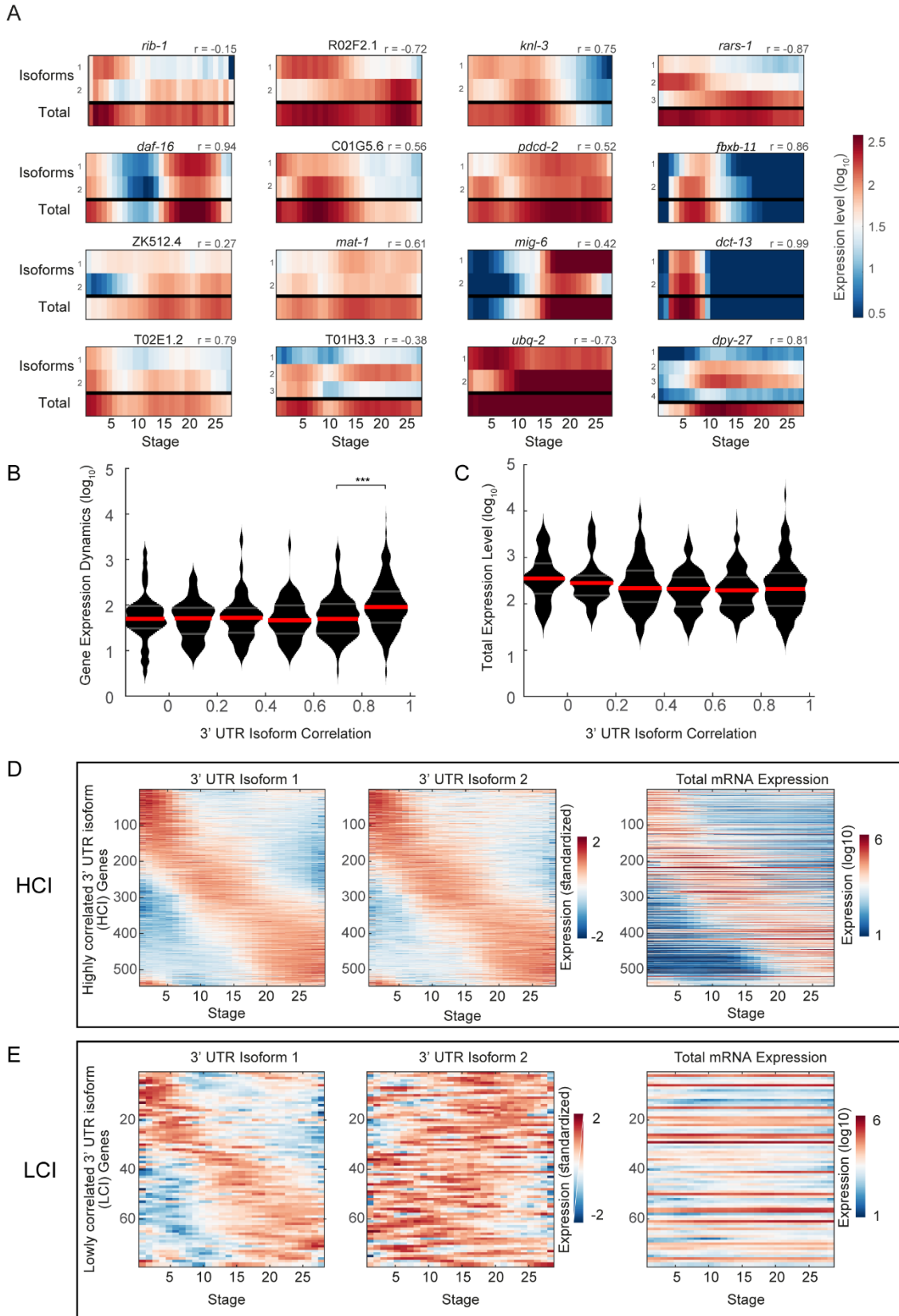
## 234 **Genes with constitutive total expression are enriched for dynamic 3' UTR isoform** 235 **expression**

236 Examining the temporal profiles, we found that 3' UTR isoforms of a particular gene may  
237 exhibit striking dynamics throughout development, while the overall total expression for the  
238 gene may be uniform (Fig. 3A). To study this systematically, we defined the dynamic range of  
239 a gene's overall expression profile as the fold differences between the maximum and minimum  
240 expression values throughout the time-course. We found a positive correlation between the  
241 dynamic overall expression range of a gene and the correlation among its isoforms (Fig. 3B;  
242  $r=0.94$ ,  $P=0.03$ , 2<sup>nd</sup> degree polynomial regression test,  $N=746$ ). Similar results were obtained  
243 using the interquartile range (IQR) of the time-course overall expression levels as a proxy for  
244 expression dynamicity ( $r=0.98$ ,  $P=0.04$ , 3<sup>rd</sup> degree polynomial regression test,  $N=746$ ). This  
245 result provides evidence for the notion that apparently constitutive genes can be highly  
246 dynamic at the level of their individual 3' UTR isoforms.

247 Gene expression levels may explain this correlation, if genes with less dynamic behavior are  
248 also lowly expressed, and in turn may have noisier and uncorrelated 3' UTR isoform profiles.  
249 To control for this possibility, we asked if there is a trend in expression levels according to the  
250 correlation among isoforms (Fig. 3C). Interestingly, genes with non-congruent APA behavior  
251 also show higher expression levels overall (Fig. 3C;  $r=0.99$ ,  $P<0.002$ , 2<sup>nd</sup> degree polynomial  
252 regression test, N=746) eliminating expression noise as a confounding factor.

253 To further examine this phenomenon, we visualized the dynamics of 3' UTR isoform  
254 expression and total expression in genes with highly or lowly correlated isoforms. The 746  
255 genes with multiple 3' UTR isoforms displayed a variety of isoform expression correlations  
256 (Supplemental Fig. S4). More than 70% of the genes (545 genes) have highly correlated 3'  
257 UTR isoform expression ( $r>0.7$ ), while 11% of the profiles (79 genes) are lowly correlated  
258 ( $r<0.3$ ). We henceforth refer to the two gene sets of highly and lowly correlated 3' UTR  
259 isoforms, as HCI and LCI, respectively. As the heat maps in Fig. 3D-E show, we found  
260 dynamic expression for both the HCI and LCI groups at the 3' UTR isoform level. As expected,  
261 the total gene expression for the correlated isoforms is dynamic, recovering the dynamics of  
262 the isoforms. Conversely, the total expression of the LCI genes mostly corresponds to profiles  
263 with constitutive overall expression. Such profiles are frequently attributed as housekeeping  
264 profiles, however as our analysis reveals, at the 3' UTR isoform level they may be very  
265 dynamic. Thus, by de-convolving the total expression into profiles of distinct 3' UTR isoforms  
266 we were able to extract a new layer of information from this dataset.

267 Delineating the LCI and HCI groups, led us to hypothesize that 3' UTR isoforms from the  
268 former group are under stronger selective pressures, and are consequently more likely to be  
269 functionally different. We thus set out to test this hypothesis by studying the post-  
270 transcriptional and translational regulatory characteristics of these two gene sets.



272 **Fig. 3.** Genes with constitutive overall expression have dynamically expressed 3' UTR  
273 isoforms. (A) Expression heatmaps indicating 3' UTR isoform expression for 16 genes with  
274 multiple isoforms displaying varying levels of correlation between the expression of their 3'  
275 UTR isoforms. The Pearson correlation coefficient  $r$  for each of the genes displayed is  
276 indicated at the top of each heatmap. (B) Relationship between 3' UTR isoform expression  
277 correlations and the overall expression dynamics. Genes whose 3' UTR isoform expression  
278 levels are correlated are more dynamic in their overall expression ( $r=0.94$ ,  $P=0.03$ , 2<sup>nd</sup> degree  
279 polynomial regression test). Dynamics for each gene is defined as the fold differences between  
280 its maximum and minimum expression values throughout the time-course. Red and grey  
281 horizontal bars represent the median and the interquartile ranges of the data, respectively. The  
282 last bin with highest 3' UTR isoform correlations exhibit significantly higher overall expression  
283 dynamics than the preceding bin ( $P<10^{-20}$ , Mann-Whitney test). (C) Relationship between 3'  
284 UTR isoform expression correlations and the respective overall total mRNA expression levels.  
285 Genes whose 3' UTR isoform expression levels are uncorrelated show significantly higher  
286 overall expression levels ( $r=0.99$ ,  $P=0.002$ , 2<sup>nd</sup> degree polynomial regression test). Red and  
287 grey horizontal bars represent the median and the interquartile ranges of the data, respectively  
288 (D) Heatmaps in the two leftmost panels show standardized expression of the 3' UTR isoforms  
289 of 545 genes belonging to the group of genes whose 3' UTR isoform's expression show high  
290 correlation (HCI genes). The right panel shows a heatmap of the total mRNA expression (on a  
291 log<sub>10</sub> scale) of the same genes confirming that genes with highly correlated isoform  
292 expression are also dynamically expressed. The time points are the same as indicated in Fig.  
293 2A. (E) Same as D for 79 genes belonging to the group of genes whose 3' UTR isoform's  
294 expression show low correlation (LCI genes). Genes with distinct expression profiles of 3' UTR  
295 isoforms appear constitutively expressed when examined at the total gene expression level.

296

### 297 **LCI genes contain more miRNA binding sites and correlate with miRNA expression**

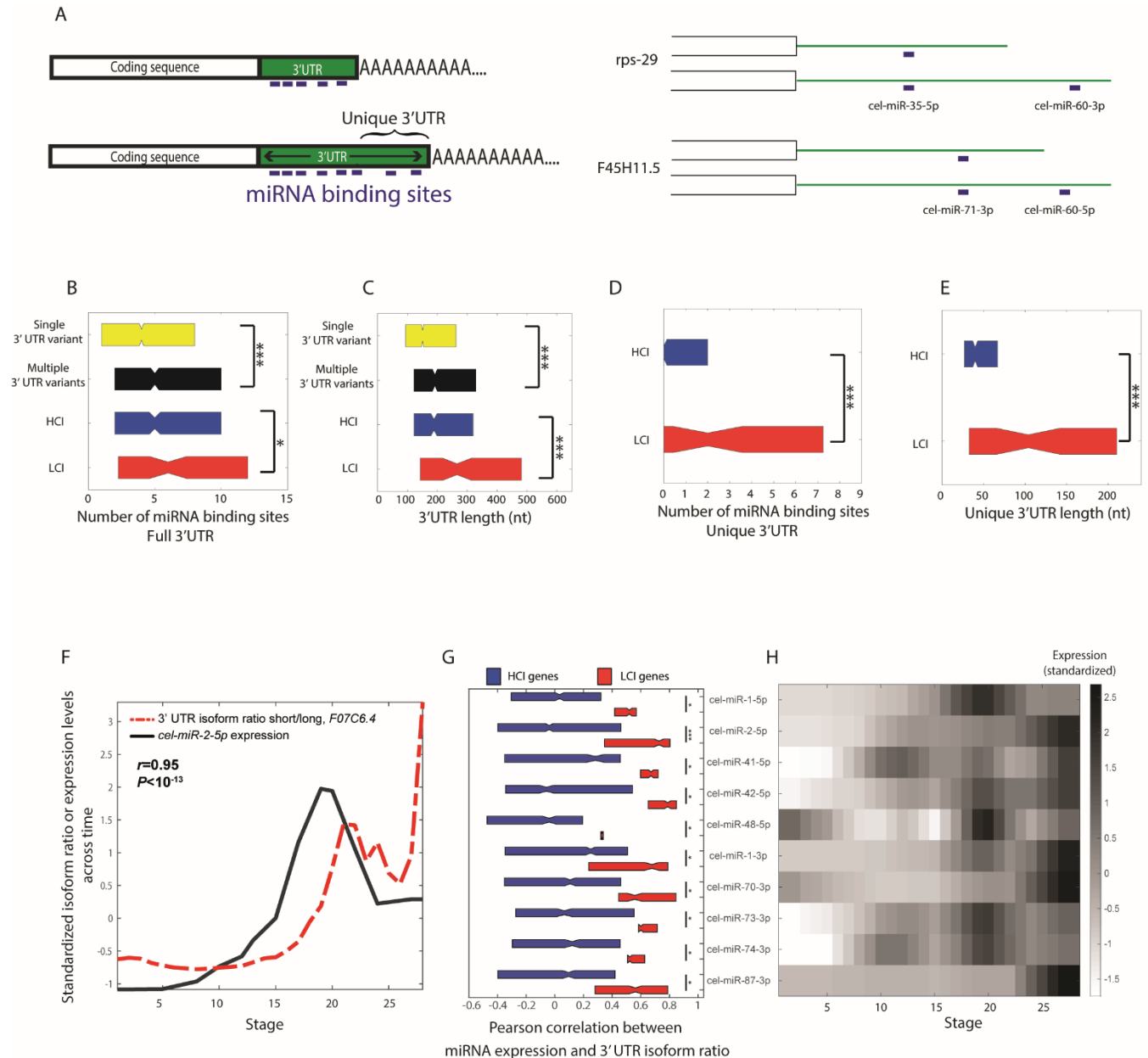
298 The 3' UTR region is known to be a locus of considerable post-transcriptional regulation  
299 (Barrett et al. 2012; Pichon et al. 2012), and the role of miRNAs in this regulation is well  
300 evidenced (Ambros 2004; Bartel 2004). We thus searched for evidence that our detected 3'  
301 UTR isoform expression profiles are regulated by miRNAs. Specifically, we asked if genes with  
302 a different number of 3' UTR variants (single vs. multiple) and different 3' UTR isoform  
303 expression correlations (HCI vs. LCI) have distinguishing sequence properties related to  
304 miRNA regulation (Fig. 4A). We first counted the number of basic miRNA seed matches in the  
305 3' UTR sequence of the different groups (Peterson et al. 2014). We found that genes with  
306 more than one 3' UTR variant have significantly more miRNA binding sites than genes with a  
307 single 3' UTR isoform ( $P<10^{-12}$ , Mann-Whitney test, Fig. 4B). Genes with multiple variants also

308 have significantly longer 3' UTRs ( $P < 10^{-28}$ , Mann-Whitney test, Fig. 4C), though the number of  
309 miRNA binding sites per base is not different across the groups (Supplemental Fig. S5). Thus,  
310 genes with multiple variants are predicted to be regulated more actively by miRNAs.

311 Comparing the number of binding sites between the HCI and LCI groups, we found that the  
312 latter group has significantly more miRNA binding sites in their 3' UTR ( $P < 0.02$ , Mann-Whitney  
313 test, Fig. 4B). We further studied this group of genes by examining the sequence that is unique  
314 to the longer 3' UTR isoform (Fig. 4A). Comparing between the HCI and LCI groups, we found  
315 that the latter have significantly more miRNA binding sites in their unique 3' UTR region ( $P < 10^{-3}$ ,  
316 Mann-Whitney test, Fig. 4D). This is not because LCI genes have denser distribution of  
317 miRNAs, rather their unique 3'UTR regions are significantly longer than in HCI genes ( $P < 10^{-4}$ ,  
318 Mann-Whitney test, Fig. 4E), thus they tend to carry more potential miRNA binding sites.  
319 These results implicate a role for miRNAs in the differences we see between the genes sets of  
320 correlated and uncorrelated isoform genes.

321 To further examine the effect that miRNA regulation exerts on HCI and LCI genes, we  
322 computed correlations between the expression profiles of all dynamically expressed miRNAs  
323 (Avital et al. 2017) and the 3' UTR isoform expression ratio of the genes they regulate. For  
324 example, *cel-miR-2* is conserved in *C. briggsae* as well as in *Drosophila*. Predicted targets of  
325 *mir-2* are enriched for genes involved in neural development (Marco et al. 2012). Expression of  
326 *miR-2* has been detected at all life stages, most abundantly in the L1 larval stage (Marco et al.  
327 2012). Consistently, we detected expression in late embryogenesis in our miRNA expression  
328 data (Fig. 4F and 4H) and interestingly, its profile has a correlation of 0.94 ( $P < 10^{-13}$ ) with the  
329 ratio of the 3' UTR isoform expression of its target gene *F07C6.4*, a gene which is enriched in  
330 the germ line, germline precursor cell, the body wall musculature and in the PVD and OLL  
331 neurons (Smith et al. 2010; Lee et al. 2017). Figure 4G shows the distributions of correlation  
332 coefficients for dynamically expressed miRNAs with the 3' UTR isoform expression ratio of  
333 their targets in the HCI and LCI groups. These ten miRNAs were selected based upon the  
334 most significant difference between the target correlations of the HCI and LCI gene groups (out  
335 of 64 miRNAs with dynamic expression- statistics of all miRNAs can be found in Supplemental  
336 Table S2). We found that in all ten the correlations are significantly higher in LCI than in HCI  
337 genes. For example, *cel-miR-2-5p* shows significantly higher correlations with the expression

338 ratio of the isoforms of all the genes it potentially binds and regulates in the LCI genes ( $P < 10^{-3}$ ,  
 339 Mann-Whitney test, Fig. 4G). This analysis suggests that miRNAs play a major role in  
 340 regulating genes with lowly correlated 3' UTR isoforms (LCI) during *C. elegans*  
 341 embryogenesis.



342 Figure 4

343 **Fig. 4.** Genes with multiple, lowly correlated 3' UTR variants show evidence for increased  
 344 miRNA regulation. (A) The left panel shows a schematic representation of the miRNA analysis.  
 345 The length of the 3' UTR region, the number of basic miRNA seed matches, and the 3' UTR



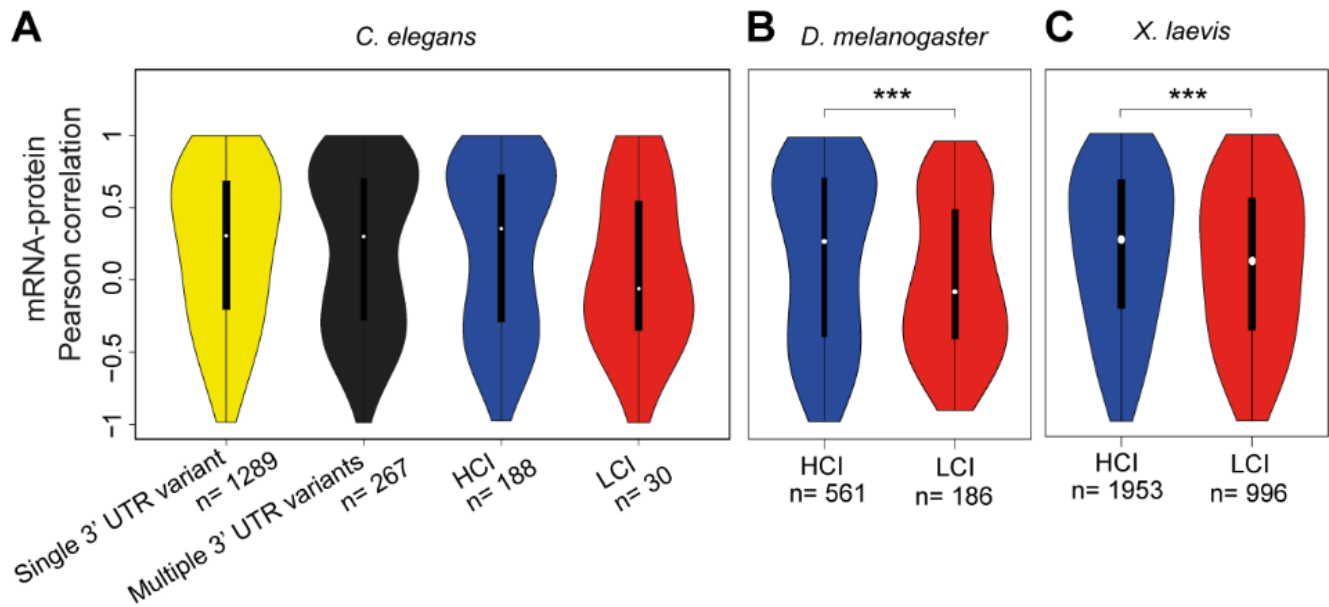
346 region unique to the longer 3' UTR isoform were considered. The right panel shows two  
347 examples of genes and their miRNA binding sites. (B) Boxplots indicating the number of  
348 miRNA binding sites in the full 3' UTR regions of genes with single or multiple 3' UTR isoforms,  
349 highly correlating (HCI) or lowly correlating isoforms (LCI). Genes with multiple 3' UTR variants  
350 have significantly more miRNA targets than genes with a single variant ( $P < 10^{-12}$ , Mann-  
351 Whitney test). Between multiple 3' UTR isoform genes, LCI genes have significantly more  
352 miRNA binding sites than HCI genes ( $P < 0.02$ , Mann-Whitney test). (C) Same as B for the full  
353 3' UTR lengths of genes. Genes with multiple 3' UTR variants have significantly longer 3'  
354 UTRs than single isoform genes ( $P < 10^{-28}$ , Mann-Whitney test). Within multiple isoform genes,  
355 LCI genes have significantly longer 3' UTRs than HCI genes ( $P < 10^{-3}$ , Mann-Whitney test). (D)  
356 Boxplots indicating the number of miRNA binding sites in the unique 3' UTR region of the  
357 longer 3' UTR isoform across HCI and LCI genes. LCI genes have significantly more miRNA  
358 binding sites in their unique 3' UTR region than HCI genes ( $P < 10^{-3}$ , Mann-Whitney test). (E)  
359 Boxplots indicating the length of the unique 3' UTR region across genes HCI and LCI genes.  
360 LCI genes have significantly longer unique 3' UTR regions than HCI genes ( $P < 10^{-4}$ , Mann-  
361 Whitney test). (F) Correlating expression of miRNA expression with expression ratio dynamics  
362 between 3' UTR isoforms. The black line shows the expression profile of *cel-miR-2-5p* miRNA  
363 throughout the developmental time-course. Depicted in red is the ratio between the expression  
364 profiles of the two 3' UTR isoforms (short/long) of the *F07C6.4* gene. The Pearson correlation  
365 coefficients between 3' UTR isoform ratio and miRNA expression of all target genes were used  
366 for the analysis shown in G and H. (G) Shown are the top ten miRNAs showing a significant  
367 difference between the HCI and LCI gene groups (out of 64 miRNAs with dynamic expression-  
368 statistics of all miRNAs can be found in Supplemental Table S2). All ten exhibit a positive  
369 median correlation between miRNA expression and the 3' UTR isoform expression ratio of  
370 genes it is predicted to bind (as in F). Further, this correlation is significantly higher in LCI than  
371 in HCI genes. (H) Heatmap of the standardized expression dynamics of the ten indicated  
372 miRNAs which show differences between HCI and LCI genes across development.

### 373

### 374 **LCI genes exhibit lower mRNA-protein correspondences**

375 Beyond regulation by miRNA, control of translation efficiency constitutes another level of post-  
376 transcriptional regulation. 3' UTR regions are known preferential targets for RNA binding  
377 proteins that regulate translation in terms of localization and efficiency (Szostak and Gebauer  
378 2013; Zhao et al. 2014; Berkovits and Mayr 2015). Indeed, mRNA and protein levels are  
379 notoriously lowly correlated (Grün et al. 2014). We reasoned that one possible explanation for  
380 the low correspondence of some genes may follow from the fact that different 3' UTR isoforms  
381 may exhibit different translation efficiencies. Thus, we predicted that LCI genes would have a

382 worse correspondence – relative to HCI genes – between transcription and protein levels. To  
383 test this, we turned to a previously published mRNA and protein *C. elegans* time-course (Grün  
384 et al. 2014) and examined the distribution of correlations between mRNA and protein  
385 abundances across the HCI and LCI groups. We detected that LCI genes show lower  
386 correlations between mRNA and protein abundances, relative to the HCI genes (Fig. 5A;  
387  $P=0.07$ , Mann-Whitney test; N=188, 30 for HCI and LCI, respectively). This limited significance  
388 may be due the low number of detected LCI genes following the shortness of the time-course,  
389 and the restricted protein data (only about 25% of the RNA-Seq detected transcripts were  
390 detected at the protein level). To further test the prediction we turned to *Drosophila* where an  
391 available high-resolution time-course allowed us to detect more LCI and HCI genes. Coupling  
392 the 3' UTR isoform expression throughout embryogenesis (Sanfilippo et al. 2017) with total  
393 mRNA (Graveley et al. 2011) and protein expression data (Casas-Vila et al. 2017), we  
394 delineated LCI and HCI genes (see Methods) and studied the correlation between their  
395 transcription and protein levels in a *Drosophila melanogaster* embryonic time-course. As in *C.*  
396 *elegans*, we found that LCI genes exhibited reduced mRNA to protein expression correlation  
397 relative to HCI genes (Fig. 5B;  $P=0.00038$ , Mann-Whitney test; N=571, 189 for HCI and LCI,  
398 respectively). We further analyzed three extensive embryonic mRNA, protein and APA  
399 datasets from *Xenopus laevis* (Zhou et al. 2019; Peshkin et al. 2015) and observed similarly  
400 highly significant trends for this vertebrate species (Fig. 5C;  $P<10^{-6}$ , Mann-Whitney test;  
401 N=1953, 996 for HCI and LCI, respectively). These results suggest that weak correlations  
402 between mRNA and protein may be in part explained by the existence of LCI genes with  
403 isoforms with distinct translation efficiencies.



404

405 **Fig. 5.** Genes with lowly correlated 3' UTR isoforms (LCI genes) have a lower correspondence  
406 between total mRNA and protein expression. (A) Violin plots indicate the distribution of  
407 Pearson correlation coefficients between total mRNA and protein expression levels across  
408 developmental stages for *C. elegans* genes with one and multiple 3' UTR isoforms, highly  
409 correlating (HCI), and lowly correlating isoforms (LCI). (B-C) Same as A, for LCI and HCI  
410 genes in *Drosophila melanogaster* (B) and *Xenopus laevis* (C) embryonic development.

411

## 412 Discussion

413 The APA-seq approach allows for the extraction of an additional layer of data from CEL-Seq  
414 data. In addition to quantifying the expression of each gene, APA-seq reveals the alternative  
415 polyadenylation dynamics on a transcriptomic basis. APA-seq uses the CEL-Seq Read 2 to  
416 identify the gene of origin and then maps Read 1 to the respective gene sequence, instead of  
417 to the whole genome, enabling the use of even extremely low-quality sequencing data  
418 resulting from reading through protocol-conditioned poly-T stretches. Hence, by performing a  
419 subtle modification in data analysis without altering the actual CEL-Seq protocol, we were able  
420 to switch from analysis at the gene to the APA level. Although the presented data is based on  
421 CEL-Seq data, our method is in principle applicable to any RNA-Seq method that uses the  
422 polyA tail as anchor and performs paired end sequencing; including inDrop, 10X, and SMART-  
423 Seq (Klein et al. 2015; Picelli et al. 2013; Zheng et al. 2017). Other 3' end sequencing methods

424 have enabled important insights into the biological significance and mechanistic aspects of  
425 APA, though their experimental procedures require either relatively large amount of starting  
426 material or complex protocols combining several amplification steps (IVT and many PCR  
427 cycles) (Shepard et al. 2011; Derti et al. 2012; Yao et al. 2012; Lianoglou et al. 2013;  
428 Wilkening et al. 2013; Gruber et al. 2014; Nam et al. 2014; Li et al. 2015; Velten et al. 2015; Ye  
429 et al. 2018, Gupta et al. 2018).

430

431 In addition to APA-seq, two methods are currently available for assessing APA in low-input  
432 samples: BATSeq (Velten et al. 2015) and ScISOOr-Seq (Gupta et al. 2018). BATSeq achieves  
433 single-cell APA isoform measurements. An important advantage of APA-seq relative to  
434 BATSeq is the formers high sensitivity, deriving from its use of CEL-Seq data; employing far  
435 less amplification and clean-up steps and PCR cycles, and a simple protocol and analysis.  
436 ScISOOr-Seq also constitutes pioneering single-cell isoform work with the advantage of  
437 revealing the complete isoform due to its reliance on long-read sequencing. Relative to  
438 ScISOOr-Seq, APA-Seq has the advantage of higher statistical confidence due to its reliance on  
439 deep Illumina sequencing. We highlight that APA-seq is not aimed towards the identification of  
440 polyA sites from scratch but rather we present it as an efficient method for quantifying the  
441 expression profiles of previously mapped isoforms. Furthermore, while we applied APA-seq  
442 here to single-embryo CEL-Seq data, it could in principle also be applied to single-cell data. As  
443 the protocol is relatively straight-forward omitting unnecessary clean-up and size-selection  
444 steps, the efficiency, complexity and accuracy is as high as that in CEL-Seq (Hashimshony et  
445 al. 2016, 2012). Applying APA-seq at single-cell resolution has many interesting applications,  
446 such as the study of population of cells undergoing cell fate specification, differentiation, and  
447 during tumorigenesis.

448

449 When studying the expression of alternative 3' UTR isoforms throughout development, we  
450 found that the global transitions of isoform usage correspond to distinct developmental periods  
451 (Fig. 2). Our results are reminiscent of those of Lianoglou et al., who studied malignant  
452 transformation across human cell lines, and revealed a change in mRNA abundance levels of  
453 genes with a single 3' UTR isoform, and, in genes with multiple 3' UTR isoforms, a change in  
454 3' UTR isoform ratios (Lianoglou et al. 2013). A similar pattern emerged while comparing

455 embryonic stem cells within differentiated tissues (Lianoglou et al. 2013). Our own results add  
456 an interesting layer to this field, by dissecting the temporal components of normal *C. elegans*  
457 embryonic development, providing insight into APA dynamics across distinct developmental  
458 stages. More generally, results are consistent with the accumulating evidence indicating that  
459 the APA regulatory mechanism is highly conserved across all eukaryotes, regardless of their  
460 morphological complexity (Ara et al. 2006; Wang et al. 2008; Shi 2012; Ulitsky et al. 2012;  
461 Velten et al. 2015; Hu et al. 2017).

462  
463 3' UTR isoforms are ubiquitous and considerable effort has been addressed towards  
464 understanding the distinct functional roles of different isoforms of the same gene (Tian and  
465 Manley 2017). Here, we report two general gene classes with 3' UTR isoforms: those whose 3'  
466 UTR isoforms correlate in their expression profiles across time and those that do not. Most  
467 genes (more than 70%) with alternatively polyadenylated isoforms exhibit a high correlation  
468 among their 3' UTR isoforms. These highly correlated isoform genes (HCI) may be dominantly  
469 regulated at the level of overall transcription. In other words, the main factor influencing the  
470 distribution of the 3' UTR isoforms usage is the intrinsic strength of their polyadenylation sites.  
471 Genes belonging to this class are often referred to as 'dynamic genes', which we previously  
472 showed to be enriched for developmental functions such as specification and differentiation  
473 (Levin et al. 2012; Levin et al. 2016). The HCI genes display a relative paucity of miRNA  
474 binding sites and higher concordance between total mRNA and protein levels (Figs. 4 and 5).

475  
476 A rather small fraction of genes with multiple 3' UTR isoforms (11%) show lowly correlated  
477 isoform expression across time. We have named these LCI genes and provide evidence that  
478 this gene class is under unique regulation. LCI genes show overall less dynamic total  
479 expression profiles; however, at the level of individual 3' UTR isoforms they are highly dynamic  
480 (Fig. 3). LCI genes also exhibit several features which indicate a higher level of post-  
481 transcriptional regulation of 3' UTR isoform usage. Post-transcriptional 3' UTR-isoform  
482 regulation processes include miRNA mediated degradation and RNA-binding protein mediated  
483 stabilization or destabilization of mRNA molecules and control of translation efficiencies. 3'  
484 UTRs of the LCI genes comprise significantly more miRNA binding sites and the 3' UTR  
485 isoform ratio correlates well with the expression of a sub-group of miRNAs, many of which are

486 known regulators of embryogenesis (Fig. 4). Consistently with these findings, the correlation  
487 between total mRNA and protein abundances is lower in LCI genes relative to HCI genes  
488 indicating that the 3' UTR usage of this group of genes is tightly regulated.

489  
490 Our results reveal principles of selective pressures on alternative polyadenylation. By studying  
491 APA dynamics over developmental time, we revealed two classes of genes with alternatively  
492 polyadenylated isoforms, the HCI and LCI. The LCI show the hallmarks of strong regulation on  
493 their 3' UTR isoforms in the form of miRNA. Thus, we revealed here that a powerful litmus test  
494 for a functional distinction among 3' UTR isoforms during a biological process (such as  
495 embryogenesis) is their discordant expression across time. As a corollary, genes with 3' UTR  
496 isoforms showing correlated expression may represent biological noise of no functional  
497 consequence, as is common for other processes such as alternative splicing (Grishkevich and  
498 Yanai 2014). Collectively, our results characterize the regulatory principles of alternative  
499 polyadenylation and provide a context for the incorporation of specific posttranscriptional  
500 regulators such as miRNAs in the modeling of biological pathways.

501

## 502 **Methods**

503 **Detection of polyadenylation site using CEL-Seq reads.** We used our previously published  
504 *C. elegans* time-course data (GSE50548) sequenced using the CEL-Seq protocol and paired-  
505 end 100 bp sequencing mode. The CEL-Seq Read 2 insert was used to identify the gene by  
506 mapping reads to the reference genome (version WS230) using Bowtie 2 version 2.2.3  
507 (Langmead and Salzberg 2012) with default parameters. The htseq-count algorithm (Anders et  
508 al. 2015) coupled with the genomic feature file was used to assign each individual Read 2 to its  
509 gene of origin. For each gene we extracted the whole gene coding sequence as well as the  
510 5000 nucleotides downstream of the stop codon, or fewer than 5000 nucleotides if another  
511 gene was found in closer proximity. We truncated Read 1, removing the barcode sequence  
512 and the polyT stretch, leaving a sequence of approximately 70 nucleotides. We then used  
513 Bowtie 2 (Langmead and Salzberg 2012) in order to map Read 1 exclusively to the coding  
514 sequence of the identified gene. The maximum and minimum mismatch penalty (--mp MX,MN)  
515 parameters for Bowtie 2 were set to 2 and 1, respectively. To summarize the data for each  
516 sample, we counted all 3'-most mapping locations of truncated Read 1 to the respective genes

517 up to a distance of 20 nucleotides upstream of the polyadenylation site. We predicted APA  
518 sites only for those genes passing a threshold of 20 mapped reads in at least two samples. For  
519 these genes we identified the peaks representing the polyadenylation sites by summarizing the  
520 last mapping coordinates of all the truncated Read 1 entities that mapped to a specific gene  
521 using the 'findpeaks' function in MATLAB. Peaks were required to be separated by at least 20  
522 nucleotides in order to be considered as distinct. The height threshold for a peak was 5 reads,  
523 or 1/1000 of the total gene expression in a particular sample. We then filtered out possible  
524 spurious peaks that may have resulted from internal priming, by removing any peak whose  
525 downstream genomic sequence included any of the following nucleotide combinations: AAAA,  
526 AGAA, AAGA or AAAG (Gruber et al. 2016). The exact coordinate of any polyadenylation site  
527 was defined as the most 3' coordinate of the respective peak. To validate the quality of our  
528 data, we compared our polyadenylation site annotation with a previously published *C. elegans*  
529 3' UTR annotation (Mangone et al. 2008). We performed this by measuring the difference  
530 between the mapping positions of polyadenylation sites from both data sets (see Fig. 1D).

531 **3' UTR isoform expression throughout *C. elegans* development.** Expression data for each  
532 sample was obtained by counting all Read 1 sequences whose 3'-most mapping location  
533 mapped up to a distance of 20 nucleotides upstream of the polyadenylation site. The raw  
534 expression data was then converted to transcripts per million (tpm) by dividing by the total  
535 reads and multiplying by one million. We worked with  $\log_{10}$  values unless otherwise noted. The  
536 profiles were further smoothed by computing a running average over 5 time points. To study  
537 dynamics of 3' UTR isoform usage throughout the *C. elegans* embryonic time-course, we first  
538 filtered genes on overall expression, keeping only those in the upper 85 percentile of the sum  
539 of expression throughout the time-course. Ratios between the two most highly expressed 3'  
540 UTR isoforms were calculated for all genes and all stages by dividing the expression levels  
541 (tpm) of the short by the levels of the long 3' UTR isoform. Only genes whose 3' UTR isoform  
542 ratios across time differed by at least a factor of three were kept for further analysis ending up  
543 with 305 genes. Significance of the changes in the ratios during the transition between  
544 successive periods were calculated using Student's t-test on all ratios of one and the ratios of  
545 the successive period. *P*-values and fold changes are shown in Fig. 2B. To generate  
546 temporally sorted expression profiles, we used "ZAVIT" as previously described (Levin et al.  
547 2016; Zalts and Yanai 2017).

548 **Categorization of genes by APA behavior.** To determine whether total gene expression is  
549 considered static or dynamic throughout development, we used the ratio of minimum to  
550 maximum and as validation the interquartile range (IQR) of expression levels across time. To  
551 quantify 3' UTR variant expression deviations, we calculated Pearson correlation for the  
552 expression pattern of the two, or more, 3' UTR isoforms. For genes with more than two  
553 isoforms (<1% of expressed genes), we used the minimal Pearson correlation between any  
554 pair of isoforms. Highly correlated 3' UTR isoforms (HCI) are those with  $r > 0.7$ , while lowly  
555 correlated 3' UTR isoforms (LCI) are those with  $r < 0.3$ .

556 **miRNA target analysis.** 3' UTR sequences were identified according to APA-seq. After  
557 removing gene coding sequences using WormBase annotation, we annotated miRNA binding  
558 sites by searching for the basic seed match sequence within the 3' UTR, i.e. the  
559 complementary sequence for nucleotides 2-7 of the miRNA (Ambros 2004; Bartel 2004). We  
560 disregarded other parameters such as type of seed match or complementarity outside of the  
561 seed region. We then counted the number of miRNA binding sites for each gene subgroup  
562 (LCI and HCI genes), and used Student's t-test to determine the significance of the difference  
563 in the number of miRNA binding sites between the LCI and HCI groups. To determine the  
564 effect specific miRNAs have on their targets, in genes with multiple 3' UTR isoforms, we  
565 calculated Pearson correlation between the isoform expression ratio of the two most highly  
566 expressed isoforms, and the miRNA expression profile (Avital et al. 2017). For this analysis we  
567 examined only dynamically expressed miRNAs, whose expression is above 250 transcripts  
568 overall across the timepoints. To compare the miRNA dataset with our APA-seq dataset, we  
569 examined only matching timepoints to our time-course and used the Matlab 'imresize' function  
570 followed by smoothing to stretch the miRNA expression data.

571 **mRNA-protein correlation analysis.** RNA-Seq and Silac data of different developmental  
572 stages in *C. elegans* was downloaded from Grün et al. (Grün et al. 2014). Replicates were  
573 averaged and data was normalized by division by the sum of all genes for the specific  
574 samples. Correlation between rpkm (mRNA) and Silac (protein) expression values using  
575 Pearson correlation was computed only for genes with rpkm and silac values for at least three  
576 stages. A similar approach was used to calculate the correlation between RNA and protein  
577 levels for the *Drosophila melanogaster* time-course. APA data of an embryonic time-course



578 was downloaded from Sanfilippo et al (Sanfilippo et al. 2017); 3' UTR isoform ratio was  
579 analyzed similarly to our time-course data yielding correlation between 3' UTR isoforms for all  
580 multi-isoform genes. Highly correlated 3' UTR isoforms (HCI) are those with  $r > 0.7$ , while lowly  
581 correlated 3' UTR isoforms (LCI) are those with  $r < 0.3$ , coherent with the thresholds used for  
582 our dataset. Total mRNA and protein expression levels were downloaded from Graveley et al.  
583 (Graveley et al. 2011) and Casas-Vila et al. (Casas-Vila et al. 2017), respectively. The two  
584 datasets were integrated by correlating mRNA (rpkm) and protein (lfq) expression levels using  
585 Pearson correlation. The same procedure was used to analyze the data for *Xenopus laevis*.  
586 mRNA RNA-Seq and protein LFQ data for embryonic time points were downloaded from  
587 (Peshkin et al. 2015) and APA data from (Zhou et al. 2019). For the datasets from all three  
588 species, the Mann-Whitney test was used to determine if LCI genes show different mRNA-  
589 protein expression correlations from HCI genes.

590

## 591 **Acknowledgements**

592 We thank Eitan Winter and Martin Feder for their assistance throughout this project. We also  
593 acknowledge assistance from the Technion Genome Center.

594

## 595 **Author contributions**

596 N.M., T.H., and I.Y. conceived and designed the project. N.M. led the development of the APA-  
597 seq approach. M.L., H.Z., N.M., and I.Y. analyzed the 3' UTR isoform data. H.Z. contributed  
598 the miRNA analysis and comparison with previous annotations. M.L. contributed the  
599 developmental APA dynamics, mRNA-protein correlation and RNA-binding-protein analyses.  
600 M.L., H.Z., N.M. and I.Y. led the interpretation of the data. M.L., H.Z., N.M., and I.Y. drafted the  
601 manuscript. M.L. is supported by a Humboldt-Bayer research fellowship.

602

## 603 **Disclosure declaration**

604 The authors declare that they have no conflicts of interest.

605

## 606 **References**

- 607 Alt FW, Bothwell AL, Knapp M, Siden E, Mather E, Koshland M, Baltimore D. 1980. Synthesis  
608 of secreted and membrane-bound immunoglobulin mu heavy chains is directed by  
609 mRNAs that differ at their 3' ends. *Cell* **20**: 293-301.
- 610 Ambros V. 2004. The functions of animal microRNAs. *Nature* **431**: 350-355.
- 611 Anders S, Pyl PT, Huber W. 2015. HTSeq--a Python framework to work with high-throughput  
612 sequencing data. *Bioinformatics* **31**: 166-169.
- 613 Ara T, Lopez F, Ritchie W, Benech P, Gautheret D. 2006. Conservation of alternative  
614 polyadenylation patterns in mammalian genes. *BMC Genomics* **7**: 189.
- 615 Avital G, Starvaggi França G, Yanai I. 2017. Bimodal evolutionary developmental miRNA  
616 program in animal embryogenesis. *Mol Biol Evol* doi:10.1093/molbev/msx316.
- 617 Barrett LW, Fletcher S, Wilton SD. 2012. Regulation of eukaryotic gene expression by the  
618 untranslated gene regions and other non-coding elements. *Cell Mol Life Sci* **69**: 3613-  
619 3634.
- 620 Bartel DP. 2004. MicroRNAs: genomics, biogenesis, mechanism, and function. *Cell* **116**: 281-  
621 297.
- 622 Berkovits BD, Mayr C. 2015. Alternative 3' UTRs act as scaffolds to regulate membrane  
623 protein localization. *Nature* **522**: 363-367.
- 624 Blazie SM, Babb C, Wilky H, Rawls A, Park JG, Mangone M. 2015. Comparative RNA-Seq  
625 analysis reveals pervasive tissue-specific alternative polyadenylation in *Caenorhabditis*  
626 *elegans* intestine and muscles. *BMC Biol* **13**: 4.
- 627 Blazie SM, Geissel HC, Wilky H, Joshi R, Newbern J, Mangone M. 2017. Alternative  
628 Polyadenylation Directs Tissue-Specific miRNA Targeting in *Caenorhabditis elegans*  
629 Somatic Tissues. *Genetics* **206**: 757-774.
- 630 Casas-Vila N, Bluhm A, Sayols S, Dinges N, Dejung M, Altenhein T, Kappei D, Altenhein B,  
631 Roignant J-Y, Butter F. 2017. The developmental proteome of *Drosophila*  
632 *melanogaster*. *Genome Res* **27**: 1273-1285.
- 633 Chen W, Jia Q, Song Y, Fu H, Wei G, Ni T. 2017. Alternative Polyadenylation: Methods,  
634 Findings, and Impacts. *Genomics Proteomics Bioinformatics* **15**: 287-300.
- 635 Consortium CeS. 1998. Genome sequence of the nematode *C. elegans*: a platform for  
636 investigating biology. *Science* **282**: 2012-2018.
- 637 de Lucas S, Oliveros JC, Chagoyen M, Ortín J. 2014. Functional signature for the recognition  
638 of specific target mRNAs by human Staufen1 protein. *Nucleic Acids Res* **42**: 4516-4526.
- 639 Derti A, Garrett-Engele P, Macisaac KD, Stevens RC, Sriram S, Chen R, Rohl CA, Johnson  
640 JM, Babak T. 2012. A quantitative atlas of polyadenylation in five mammals. *Genome*  
641 *Res* **22**: 1173-1183.
- 642 Early P, Rogers J, Davis M, Calame K, Bond M, Wall R, Hood L. 1980. Two mRNAs can be  
643 produced from a single immunoglobulin mu gene by alternative RNA processing  
644 pathways. *Cell* **20**: 313-319.
- 645 Elkon R, Ugalde AP, Agami R. 2013. Alternative cleavage and polyadenylation: extent,  
646 regulation and function. *Nat Rev Genet* **14**: 496-506.
- 647 Goldstrohm AC, Hall TMT, McKenney KM. 2018. Post-transcriptional Regulatory Functions of  
648 Mammalian Pumilio Proteins. *Trends Genet* **34**: 972-990.

- 649 Graveley BR, Brooks AN, Carlson JW, Duff MO, Landolin JM, Yang L, Artieri CG, van Baren  
650 MJ, Boley N, Booth BW et al. 2011. The developmental transcriptome of *Drosophila*  
651 *melanogaster*. *Nature* **471**: 473-479.
- 652 Grishkevich V, Yanai I. 2014. Gene length and expression level shape genomic novelties.  
653 *Genome Res* **24**: 1497-1503.
- 654 Gruber AJ, Schmidt R, Gruber AR, Martin G, Ghosh S, Belmadani M, Keller W, Zavolan M.  
655 2016. A comprehensive analysis of 3' end sequencing data sets reveals novel  
656 polyadenylation signals and the repressive role of heterogeneous ribonucleoprotein C  
657 on cleavage and polyadenylation. *Genome Res* **26**: 1145-1159.
- 658 Gruber AR, Martin G, Müller P, Schmidt A, Gruber AJ, Gumienny R, Mittal N, Jayachandran R,  
659 Pieters J, Keller W et al. 2014. Global 3' UTR shortening has a limited effect on protein  
660 abundance in proliferating T cells. *Nat Commun* **5**: 5465.
- 661 Grün D, Kirchner M, Thierfelder N, Stoeckius M, Selbach M, Rajewsky N. 2014. Conservation  
662 of mRNA and protein expression during development of *C. elegans*. *Cell Rep* **6**: 565-  
663 577.
- 664 Hashimshony T, Senderovich N, Avital G, Klochendler A, de Leeuw Y, Anavy L, Gennert D, Li  
665 S, Livak KJ, Rozenblatt-Rosen O et al. 2016. CEL-Seq2: sensitive highly-multiplexed  
666 single-cell RNA-Seq. *Genome Biol* **17**: 77.
- 667 Hashimshony T, Wagner F, Sher N, Yanai I. 2012. CEL-Seq: single-cell RNA-Seq by  
668 multiplexed linear amplification. *Cell Rep* **2**: 666-673.
- 669 Hillier LW, Reinke V, Green P, Hirst M, Marra MA, Waterston RH. 2009. Massively parallel  
670 sequencing of the polyadenylated transcriptome of *C. elegans*. *Genome Res* **19**: 657-  
671 666.
- 672 Hu W, Li S, Park JY, Boppana S, Ni T, Li M, Zhu J, Tian B, Xie Z, Xiang M. 2017. Dynamic  
673 landscape of alternative polyadenylation during retinal development. *Cellular and*  
674 *molecular life sciences: CMLS* **74**: 1721-1739.
- 675 Jan CH, Friedman RC, Ruby JG, Bartel DP. 2011. Formation, regulation and evolution of  
676 *Caenorhabditis elegans* 3'UTRs. *Nature* **469**: 97-101.
- 677 Ji Z, Lee JY, Pan Z, Jiang B, Tian B. 2009. Progressive lengthening of 3' untranslated regions  
678 of mRNAs by alternative polyadenylation during mouse embryonic development. *Proc*  
679 *Natl Acad Sci USA* **106**: 7028-7033.
- 680 Kamath RS, Fraser AG, Dong Y, Poulin G, Durbin R, Gotta M, Kanapin A, Le Bot N, Moreno S,  
681 Sohrmann M et al. 2003. Systematic functional analysis of the *Caenorhabditis elegans*  
682 genome using RNAi. *Nature* **421**: 231-237.
- 683 Keene JD. 2007. RNA regulons: coordination of post-transcriptional events. *Nat Rev Genet* **8**:  
684 533-543.
- 685 Khraiweh B, Salehi-Ashtiani K. 2017. Alternative Poly(A) Tails Meet miRNA Targeting in  
686 *Caenorhabditis elegans*. *Genetics* **206**: 755-756.
- 687 Kim YK, Furic L, Desgroseillers L, Maquat LE. 2005. Mammalian Staufen1 recruits Upf1 to  
688 specific mRNA 3'UTRs so as to elicit mRNA decay. *Cell* **120**: 195-208.
- 689 Kimble J, Hirsh D. 1979. The postembryonic cell lineages of the hermaphrodite and male  
690 gonads in *Caenorhabditis elegans*. *Dev Biol* **70**: 396-417.
- 691 Klein AM, Mazutis L, Akartuna I, Tallapragada N, Veres A, Li V, Peshkin L, Weitz DA,  
692 Kirschner MW. 2015. Droplet Barcoding for Single-Cell Transcriptomics Applied to  
693 Embryonic Stem Cells. *Cell* **161**: 1187-1201.
- 694 <http://www.ncbi.nlm.nih.gov/pubmed/26000487> (Accessed May 21, 2015).

- 695 Lamm AT, Stadler MR, Zhang H, Gent JI, Fire AZ. 2011. Multimodal RNA-seq using single-  
696 strand, double-strand, and CircLigase-based capture yields a refined and extended  
697 description of the *C. elegans* transcriptome. *Genome Res* **21**: 265-275.
- 698 Langmead B, Salzberg SL. 2012. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**:  
699 357-359.
- 700 Lee C-YS, Lu T, Seydoux G. 2017. Nanos promotes epigenetic reprogramming of the germline  
701 by down-regulation of the THAP transcription factor LIN-15B. *Elife* **6**.
- 702 Levin M, Anavy L, Cole AG, Winter E, Mostov N, Khair S, Senderovich N, Kovalev E, Silver  
703 DH, Feder M et al. 2016. The mid-developmental transition and the evolution of animal  
704 body plans. *Nature* **531**: 637-641.
- 705 Levin M, Hashimshony T, Wagner F, Yanai I. 2012. Developmental milestones punctuate gene  
706 expression in the *Caenorhabditis* embryo. *Dev Cell* **22**: 1101-1108.
- 707 Li W, You B, Hoque M, Zheng D, Luo W, Ji Z, Park JY, Gunderson SI, Kalsotra A, Manley JL  
708 et al. 2015. Systematic profiling of poly(A)+ transcripts modulated by core 3' end  
709 processing and splicing factors reveals regulatory rules of alternative cleavage and  
710 polyadenylation. *PLoS Genet* **11**: e1005166.
- 711 Li Y, Sun Y, Fu Y, Li M, Huang G, Zhang C, Liang J, Huang S, Shen G, Yuan S et al. 2012.  
712 Dynamic landscape of tandem 3' UTRs during zebrafish development. *Genome Res* **22**:  
713 1899-1906.
- 714 Lianoglou S, Garg V, Yang JL, Leslie CS, Mayr C. 2013. Ubiquitously transcribed genes use  
715 alternative polyadenylation to achieve tissue-specific expression. *Genes Dev* **27**: 2380-  
716 2396.
- 717 Lutz CS, Moreira A. 2011. Alternative mRNA polyadenylation in eukaryotes: an effective  
718 regulator of gene expression. *Wiley Interdiscip Rev RNA* **2**: 22-31.
- 719 Mangone M, Macmenamin P, Zegar C, Piano F, Gunsalus KC. 2008. UTRome.org: a platform  
720 for 3'UTR biology in *C. elegans*. *Nucleic Acids Res* **36**: D57-62.
- 721 Mangone M, Manoharan AP, Thierry-Mieg D, Thierry-Mieg J, Han T, Mackowiak SD, Mis E,  
722 Zegar C, Gutwein MR, Khivansara V et al. 2010. The landscape of *C. elegans* 3'UTRs.  
723 *Science* **329**: 432-435.
- 724 Marco A, Hooks K, Griffiths-Jones S. 2012. Evolution and function of the extended miR-2  
725 microRNA family. *RNA Biol* **9**: 242-248.
- 726 Mazumder B, Seshadri V, Fox PL. 2003. Translational control by the 3'-UTR: the ends specify  
727 the means. *Trends Biochem Sci* **28**: 91-98.
- 728 McKay SJ, Johnsen R, Khattra J, Asano J, Baillie DL, Chan S, Dube N, Fang L, Goszczynski  
729 B, Ha E et al. 2003. Gene expression profiling of cells, tissues, and developmental  
730 stages of the nematode *C. elegans*. *Cold Spring Harb Symp Quant Biol* **68**: 159-169.
- 731 Micklem DR, Adams J, Grünert S, St Johnston D. 2000. Distinct roles of two conserved  
732 Stauf domains in oskar mRNA localization and translation. *EMBO J* **19**: 1366-1377.
- 733 Mitchell SF, Parker R. 2014. Principles and properties of eukaryotic mRNPs. *Mol Cell* **54**: 547-  
734 558.
- 735 Müller-McNicoll M, Neugebauer KM. 2013. How cells get the message: dynamic assembly and  
736 function of mRNA-protein complexes. *Nat Rev Genet* **14**: 275-287.
- 737 Nam J-W, Rissland OS, Koppstein D, Abreu-Goodger C, Jan CH, Agarwal V, Yildirim MA,  
738 Rodriguez A, Bartel DP. 2014. Global analyses of the effect of different cellular contexts  
739 on microRNA targeting. *Mol Cell* **53**: 1031-1043.

- 740 Newman-Smith ED, Rothman JH. 1998. The maternal-to-zygotic transition in embryonic  
741 patterning of *Caenorhabditis elegans*. *Curr Opin Genet Dev* **8**: 472-480.
- 742 Ozsolak F, Kapranov P, Foissac S, Kim SW, Fishilevich E, Monaghan AP, John B, Milos PM.  
743 2010. Comprehensive polyadenylation site maps in yeast and human reveal pervasive  
744 alternative polyadenylation. *Cell* **143**: 1018-1029.
- 745 Pagano JM, Farley BM, Essien KI, Ryder SP. 2009. RNA recognition by the embryonic cell  
746 fate determinant and germline totipotency factor MEX-3. *Proc Natl Acad Sci USA* **106**:  
747 20252-20257.
- 748 Peshkin L, Wühr M, Pearl E, Haas W, Freeman RM, Gerhart JC, Klein AM, Horb M, Gygi SP,  
749 Kirschner MW. 2015. On the Relationship of Protein and mRNA Dynamics in Vertebrate  
750 Embryonic Development. *Dev Cell* **35**: 383–94.  
751 <https://linkinghub.elsevier.com/retrieve/pii/S1534580715006577> (Accessed April 6, 2019).
- 752 Peterson SM, Thompson JA, Ufkin ML, Sathyanarayana P, Liaw L, Congdon CB. 2014.  
753 Common features of microRNA target prediction tools. *Front Genet* **5**: 23.
- 754 Picelli S, Björklund ÅK, Faridani OR, Sagasser S, Winberg G, Sandberg R. 2013. Smart-seq2  
755 for sensitive full-length transcriptome profiling in single cells. *Nat Methods* **10**: 1096–8.  
756 <http://www.ncbi.nlm.nih.gov/pubmed/24056875> (Accessed July 12, 2014)
- 757 Pichon X, Wilson LA, Stoneley M, Bastide A, King HA, Somers J, Willis AEE. 2012. RNA  
758 binding protein/RNA element interactions and the control of translation. *Curr Protein*  
759 *Pept Sci* **13**: 294-304.
- 760 Prasad A, Porter DF, Kroll-Conner PL, Mohanty I, Ryan AR, Crittenden SL, Wickens M, Kimble  
761 J. 2016. The PUF binding landscape in metazoan germ cells. *RNA* **22**: 1026-1043.
- 762 Ramani AK, Nelson AC, Kapranov P, Bell I, Gingeras TR, Fraser AG. 2009. High resolution  
763 transcriptome maps for wild-type and nonsense-mediated decay-defective  
764 *Caenorhabditis elegans*. *Genome Biol* **10**: R101.
- 765 Rogers J, Early P, Carter C, Calame K, Bond M, Hood L, Wall R. 1980. Two mRNAs with  
766 different 3' ends encode membrane-bound and secreted forms of immunoglobulin mu  
767 chain. *Cell* **20**: 303-312.
- 768 Sanfilippo P, Wen J, Lai EC. 2017. Landscape and evolution of tissue-specific alternative  
769 polyadenylation across *Drosophila* species. *Genome Biol* **18**: 229.
- 770 Setzer DR, McGrogan M, Nunberg JH, Schimke RT. 1980. Size heterogeneity in the 3' end of  
771 dihydrofolate reductase messenger RNAs in mouse cells. *Cell* **22**: 361-370.
- 772 Shepard PJ, Choi E-A, Lu J, Flanagan LA, Hertel KJ, Shi Y. 2011. Complex and dynamic  
773 landscape of RNA polyadenylation revealed by PAS-Seq. *RNA* **17**: 761-772.
- 774 Shi Y. 2012. Alternative polyadenylation: new insights from global analyses. *RNA* **18**: 2105-  
775 2117.
- 776 Shin H, Hirst M, Bainbridge MN, Magrini V, Mardis E, Moerman DG, Marra MA, Baillie DL,  
777 Jones SJM. 2008. Transcriptome analysis for *Caenorhabditis elegans* based on novel  
778 expressed sequence tags. *BMC Biol* **6**: 30.
- 779 Smibert P, Miura P, Westholm JO, Shenker S, May G, Duff MO, Zhang D, Eads BD, Carlson J,  
780 Brown JB et al. 2012. Global patterns of tissue-specific alternative polyadenylation in  
781 *Drosophila*. *Cell Rep* **1**: 277-289.
- 782 Smith CJ, Watson JD, Spencer WC, O'Brien T, Cha B, Albeg A, Treinin M, Miller DM. 2010.  
783 Time-lapse imaging and cell-specific expression profiling reveal dynamic branching and  
784 molecular determinants of a multi-dendritic nociceptor in *C. elegans*. *Dev Biol* **345**: 18-  
785 33.

- 786 Sulston JE, Horvitz HR. 1977. Post-embryonic cell lineages of the nematode, *Caenorhabditis*  
787 *elegans*. *Dev Biol* **56**: 110-156.
- 788 Szostak E, Gebauer F. 2013. Translational control by 3'-UTR-binding proteins. *Brief Funct*  
789 *Genomics* **12**: 58-65.
- 790 Tadros W, Goldman AL, Babak T, Menzies F, Vardy L, Orr-Weaver T, Hughes TR, Westwood  
791 JT, Smibert CA, Lipshitz HD. 2007a. SMAUG is a major regulator of maternal mRNA  
792 destabilization in *Drosophila* and its translation is activated by the PAN GU kinase. *Dev*  
793 *Cell* **12**: 143-155.
- 794 Tadros W, Lipshitz HD. 2009. The maternal-to-zygotic transition: a play in two acts.  
795 *Development* **136**: 3033-3042.
- 796 Tadros W, Westwood JT, Lipshitz HD. 2007b. The mother-to-child transition. *Dev Cell* **12**: 847-  
797 849.
- 798 Tian B, Hu J, Zhang H, Lutz CS. 2005. A large-scale analysis of mRNA polyadenylation of  
799 human and mouse genes. *Nucleic Acids Res* **33**: 201-212.
- 800 Tian B, Manley JL. 2017. Alternative polyadenylation of mRNA precursors. *Nat Rev Mol Cell*  
801 *Biol* **18**: 18-30.
- 802 Ulitsky I, Shkumatava A, Jan CH, Subtelny AO, Koppstein D, Bell GW, Sive H, Bartel DP.  
803 2012. Extensive alternative polyadenylation during zebrafish development. *Genome*  
804 *Res* **22**: 2054-2066.
- 805 Velten L, Anders S, Pekowska A, Järvelin AI, Huber W, Pelechano V, Steinmetz LM. 2015.  
806 Single-cell polyadenylation site mapping reveals 3' isoform choice variability. *Mol Syst*  
807 *Biol* **11**: 812.
- 808 Walser CB, Lipshitz HD. 2011. Transcript clearance during the maternal-to-zygotic transition.  
809 *Curr Opin Genet Dev* **21**: 431-443.
- 810 Wang ET, Sandberg R, Luo S, Khrebtkova I, Zhang L, Mayr C, Kingsmore SF, Schroth GP,  
811 Burge CB. 2008. Alternative isoform regulation in human tissue transcriptomes. *Nature*  
812 **456**: 470-476.
- 813 Wilkening S, Pelechano V, Järvelin AI, Tekkedil MM, Anders S, Benes V, Steinmetz LM. 2013.  
814 An efficient method for genome-wide polyadenylation site mapping and RNA  
815 quantification. *Nucleic Acids Res* **41**: e65.
- 816 Wormington M. 1994. Unmasking the role of the 3' UTR in the cytoplasmic polyadenylation and  
817 translational regulation of maternal mRNAs. *Bioessays* **16**: 533-535.
- 818 Yao C, Biesinger J, Wan J, Weng L, Xing Y, Xie X, Shi Y. 2012. Transcriptome-wide analyses  
819 of CstF64-RNA interactions in global regulation of mRNA alternative polyadenylation.  
820 *Proc Natl Acad Sci USA* **109**: 18773-18778.
- 821 Ye C, Long Y, Ji G, Li QQ, Wu X. 2018. APAtrap: identification and quantification of alternative  
822 polyadenylation sites from RNA-seq data. *Bioinformatics* **34**: 1841-1849.
- 823 Zalts H, Yanai I. 2017. Developmental constraints shape the evolution of the nematode mid-  
824 developmental transition. *Nature Ecology & Evolution* **1**: s41559-41017-40113-41017.
- 825 Zhao W, Pollack JL, Blagev DP, Zaitlen N, McManus MT, Erle DJ. 2014. Massively parallel  
826 functional annotation of 3' untranslated regions. *Nat Biotechnol* **32**: 387-391.
- 827 Zheng GXY, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, Ziraldo SB, Wheeler TD,  
828 McDermott GP, Zhu J, et al. 2017. Massively parallel digital transcriptional profiling of  
829 single cells. *Nat Commun* **8**: 14049.
- 830 Zhou X, Zhang Y, Michal JJ, Qu L, Zhang S, Wildung MR, Du W, Pouchnik DJ, Zhao H, Xia Y,

831 et al. 2019. Alternative polyadenylation coordinates embryonic development, sexual  
832 dimorphism and longitudinal growth in *Xenopus tropicalis*. *Cell Mol Life Sci*.  
833 <http://link.springer.com/10.1007/s00018-019-03036-1> (Accessed April 6, 2019).

834

835

836