

13 **Abstract:**

14 Transposable elements (TEs) are an integral part of the host transcriptome. TE-containing
15 noncoding RNAs (ncRNAs) exhibit considerable tissue specificity and play crucial roles during
16 development, including stem cell maintenance and cell differentiation. Recent advances in single
17 cell RNA-seq (scRNA-seq) revolutionized cell-type specific gene expression analysis. However,
18 scRNA-seq quantification tools tailored for TEs are lacking, limiting our ability to dissect TE
19 expression dynamics at single cell resolution. To address this issue, we established a TE expression
20 quantification pipeline that is compatible with scRNA-seq data generated across multiple
21 technology platforms. We constructed TE containing ncRNA references using bulk RNA-seq data
22 and demonstrated that quantifying TE expression at the transcript level effectively reduces noise.
23 As proof of principle, we applied this strategy to mouse embryonic stem cells and successfully
24 captured the expression profile of endogenous retroviruses in single cells. We further expanded
25 our analysis to scRNA-seq data from early stages of mouse embryogenesis. Our results illustrated
26 the dynamic TE expression at pre-implantation stages and revealed 137 TE-containing ncRNA
27 transcripts with substantial tissue specificity during gastrulation and early organogenesis.

28 **Introduction:**

29 Transposable elements (TEs) occupy a large proportion of eukaryotic genomes, representing
30 roughly 50% of the human genome and 40% of the mouse genome. Though once regarded as non-
31 functional parasitic sequences, mounting evidence suggests that TEs play pivotal roles in gene
32 regulation. During evolution, TEs rewire host transcription networks through transposition,
33 resulting in a wide variety of TE-derived regulatory elements, including promoters, enhancers,
34 transcription terminators and chromatin loop anchors (Feschotte and Gilbert 2012; Rebollo et al.

35 2012; Cowley and Oakey 2013; Garcia-Perez et al. 2016; Chuong et al. 2017; Sundaram and
36 Wysocka 2020). In the present-day, despite losing most of their transposition abilities, TEs
37 continue to impact host genomes through transcription, which generates protein-coding-TE
38 chimeric RNAs as well as noncoding RNAs (ncRNAs) that are crucial for normal and cancer
39 development (Gifford et al. 2013; Hadjiargyrou and Delihis 2013; Hutchins and Pei 2015; Anwar
40 et al. 2017; Rodriguez-Terrones and Torres-Padilla 2018).

41 TEs are major contributors of ncRNAs in both human and mouse. More than two thirds of mature
42 long ncRNAs contain at least one TE and almost half of the total base pairs of long ncRNA are
43 derived from TEs (Kelley and Rinn 2012; Kapusta et al. 2013; Veselovska et al. 2015). TE-
44 containing ncRNAs show substantial developmental stage and tissue specificity and play essential
45 roles during embryonic stem cell (ESC) maintenance and early embryogenesis. For instance,
46 endogenous retroviruses (ERVs) are highly expressed in ESCs and are involved in ESC self-
47 renewal and differentiation (Macfarlan et al. 2012; Santoni et al. 2012; Fort et al. 2014; Lu et al.
48 2014; Ohnuki et al. 2014; Wang et al. 2014). During mouse and human embryogenesis, a large
49 number of TEs, including ERVs, long interspersed nuclear element-1 (LINE-1) and short
50 interspersed elements (SINEs), become active and contribute to a significant proportion of total
51 RNAs before blastocyst stage (Kigami et al. 2003; Peaston et al. 2004; Svoboda et al. 2004;
52 Maksakova and Mager 2005; Fadloun et al. 2013; Göke et al. 2015; Grow et al. 2015; De Iaco et
53 al. 2017; Ge 2017; Hendrickson et al. 2017; Jachowicz et al. 2017; Whiddon et al. 2017; Percharde
54 et al. 2018). Moreover, knocking down specific TE families, including LINE-1 and MuERV-L,
55 results in severe developmental defects (Kigami et al. 2003; Huang et al. 2017; Jachowicz et al.
56 2017; Percharde et al. 2018).

57 Despite the importance of TEs, quantifying TE expression using high-throughput sequencing data
58 has been challenging. Due to TEs' repetitive nature, sequencing reads that overlap with TEs are
59 often discarded as a result of ambiguous mapping. Recently, several software tools were developed
60 to address this issue and they enabled TE expression quantification in bulk RNA-seq data (Day et
61 al. 2010; Jin et al. 2015; Srivastava et al. 2016; Lerat et al. 2017; Guffanti et al. 2018; Jeong et al.
62 2018; Valdebenito-Maturana and Riadi 2018; Bendall et al. 2019; Yang et al. 2019). In order to
63 quantify the expression of repetitive elements, these tools either aggregate multi-aligned reads at
64 TE subfamilies/families or redistribute them to individual TEs based on heuristic or statistical rules.
65 Although proven to be successful in a range of biological systems, applications of the current TE
66 quantification strategies were mostly limited to bulk RNA-seq, which lacks the ability to
67 distinguish cell type specific TE expression.

68 Recent developments in single cell RNA-seq (scRNA-seq) provide unprecedented opportunities
69 for examining cell type specific TE expression. However, effective TE quantification tools
70 optimized for scRNA-seq data are lacking. Although the assessment of genome-wide
71 transcriptional activity of TEs in single cells has been attempted by counting signals at TE
72 subfamilies/families (Göke et al. 2015; Ge 2017; Boroviak et al. 2018; Brocks et al. 2018; Yandim
73 and Karakülah 2019; He et al. 2020; Jonsson et al. 2020), such approaches are not optimal.
74 Compared with bulk RNA-seq, scRNA-seq signal is much noisier and often shows 5' or 3' end
75 enrichment along the transcripts. Counting reads at individual TEs or subfamilies/families fails to
76 take into account the structures of the full-length transcripts, which can consist of multiple TEs
77 from different subfamilies/families. Consequently, different expression values will be assigned to
78 individual TEs within the same transcript. This caveat is especially obvious when dealing with
79 scRNA-seq datasets where sequencing reads are enriched at either the 5' or 3' end of the RNA.

80 Counting reads without the knowledge of the full-length transcripts will only capture TEs near the
81 5' end or the polyA signal, resulting in an inaccurate picture of the genome-wide TE expression
82 pattern.

83 In order to address these issues, we have developed an analytical pipeline tailored for TE
84 expression quantification in scRNA-seq datasets. By comparing TE derived RNA-seq reads in
85 bulk and scRNA-seq datasets, we first show that higher percentages of reads are mapped to TEs
86 in scRNA-seq regardless of library preparation protocols. We further demonstrate that counting
87 scRNA-seq signal at individual TEs leads to large amounts of false positives and generate bias due
88 to 3' enrichment of scRNA-seq signal. To overcome these challenges, we quantify TE expression
89 in single cells using TE transcripts assembled from bulk RNA-seq. This approach successfully
90 captures TEs that are exonized into ncRNAs and significantly improves scRNA-seq analysis by
91 enriching for regions with true signal. Furthermore, applying our strategy to mouse early
92 embryogenesis captured dynamic TE expression during pre-implantation stages. Expanding our
93 analysis to mouse gastrulation and early organogenesis revealed 137 TE transcripts with
94 substantial tissue specificity. These TE transcripts are mostly un-annotated and transcripts with
95 different tissue enrichment show distinct TE composition. In summary, our study provides a
96 systematic evaluation of TE derived reads in scRNA-seq datasets and establishes an effective
97 computational approach for quantifying TE expression using scRNA-seq data.

98 **Results:**

99 **A higher percentage of reads are mapped to TEs in scRNA-seq compared with bulk RNA-**
100 **seq**

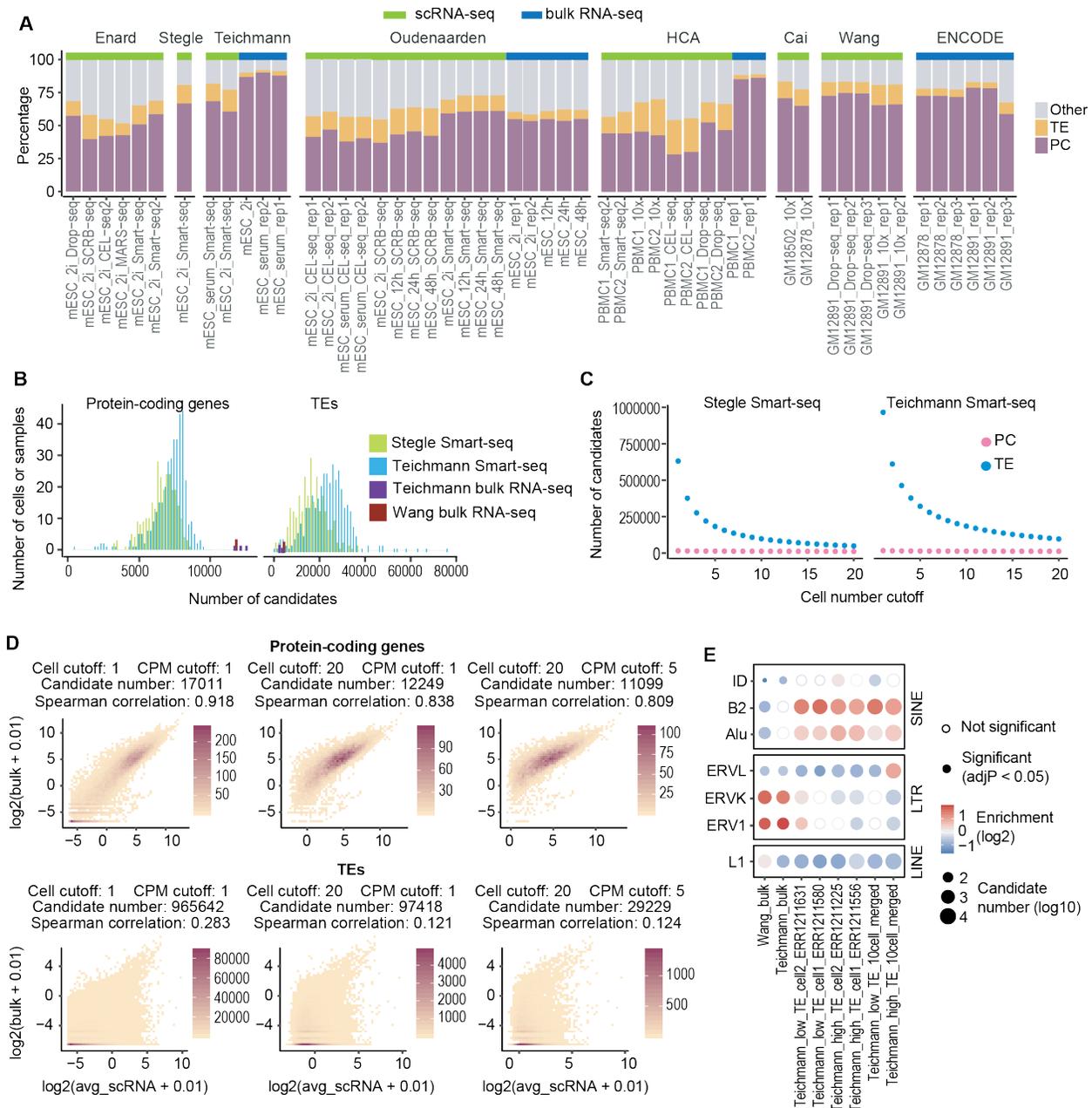
101 Due to the biological significance of TE-containing ncRNAs, we decided to focus our analysis on
102 TEs that are not part of protein-coding transcripts. First, to determine the fraction of reads that can
103 be mapped to these TEs in scRNA-seq data, we processed 36 publicly available single cell datasets
104 ([Supplemental Table S1](#)). These datasets contain both human and mouse samples and were
105 generated using 7 different scRNA-seq protocols. Bulk RNA-seq datasets from the same study or
106 derived from the same cell line were included as controls. To preserve reads that originate from
107 repetitive regions, multiple mapping was implemented during the alignment step. Reads that
108 mapped to multiple locations or overlapped with more than one feature were distributed equally
109 for signal quantification. Calculating the number of mappable reads based on genomic locations
110 revealed that a large proportion of reads overlap with TEs in all tested scRNA-seq datasets,
111 suggesting that TE expression can be captured by scRNA-seq ([Fig. 1A](#)). Notably, we observed a
112 higher percentage of reads mapped to TEs in scRNA-seq compared with bulk RNA-seq, a
113 phenomenon that was consistent across different scRNA-seq platforms even when only uniquely
114 mapped reads were considered ([Fig. 1A, Supplemental Fig. S1A](#)).

115 To evaluate whether the high TE mapping percentage in scRNA-seq was associated with data
116 quality per cell, we further compared scRNA-seq data generated using Smart-seq, 10x Genomics
117 Chromium, SCRB-seq and Drop-seq and examined the relationships between the percentage of
118 TE reads per cell and the two key parameters indicative of scRNA-seq quality: sequencing depth
119 and the percentage of mitochondrial reads (Ilicic et al. 2016). Our analysis revealed that a high TE

120 mapping percentage was observed across individual cells with limited correlation to sequencing
121 depth ([Supplemental Fig. S1B, C](#)). Similarly, no strong correlation was detected between the
122 percentage of TE reads and that of the mitochondrial reads ([Supplemental Fig. S1C](#)). These results
123 suggest that the high TE mapping ratio in scRNA-seq is unlikely an artifact caused by extreme
124 sequencing depth or cell death.

125 To address the concern that genomic DNA contamination may contribute to the majority of TE
126 signal in scRNA-seq, we next quantified the number of total reads mapped to 5 non-overlapping
127 genomic regions: protein-coding exons, TEs within the introns of protein-coding genes, other
128 intronic regions of protein-coding genes, intergenic TEs, and other intergenic regions. A higher
129 percentage of total reads were mapped to the intronic regions of protein-coding genes in scRNA-
130 seq and the majority of TE overlapping reads were located within introns, suggesting that the signal
131 originates from un-processed RNAs ([Supplemental Fig. S1A](#)). Furthermore, consistent with a
132 previous report (La Manno et al. 2018), we observed that regions enriched for un-spliced scRNA-
133 seq reads tend to be flanked by AT-rich sequences ([Supplemental Fig. S1D](#)). Taken together, these
134 observations suggest that scRNA-seq TE signals are unlikely to come from genomic DNA
135 contamination but rather are the products of the polyA priming step during cDNA synthesis.

136



137

138

139 **Figure 1 Counting scRNA-seq signal at individual TEs results in large numbers of false**

140 **positive candidates**

141 (A) Distribution of mappable reads in 16 bulk RNA-seq and 36 scRNA-seq datasets. Compared to

142 bulk RNA-seq, scRNA-seq data have higher percentage of reads mapped to TEs. Samples were

143 arranged by studies. PC: protein-coding exons defined by refSeq. TE: transposable elements that
144 do not overlap with protein-coding exons. Other: other genomic locations. mESC: mouse
145 embryonic stem cell. PBMC: human peripheral blood mononuclear cell. GM12878 and GM12891:
146 human lymphoblastoid cell lines.

147 (B) Number of expressed (Counts Per Million, CPM ≥ 1) protein-coding genes and TEs in mESC
148 bulk RNA-seq and Smart-seq samples. On average, 12,000 protein-coding genes and 6,000 TEs
149 were detected in each bulk RNA-seq sample. In contrast, scRNA-seq captured 7,000 protein-
150 coding genes and 20,000 TEs per cell.

151 (C) Number of candidates as a function of cell number cutoff (the minimum number of cells each
152 candidate is expressed in. Expression cutoff: CPM ≥ 1). Although the majority of protein-coding
153 gene candidates were consistently detected in mESC Smart-seq data, a large number of TE
154 candidates were detected in less than 10 cells.

155 (D) Correlation between bulk RNA-seq and averaged scRNA-seq signal at protein-coding genes
156 and TEs (Teichmann lab, mESC). Low correlation between bulk RNA-seq and averaged Smart-
157 seq signal was observed at TEs regardless of expression cutoff. Cell cutoff: the minimum number
158 of cells each candidate is expressed in. CPM cutoff: the minimum CPM value for one candidate to
159 be considered as expressed.

160 (E) TE-family enrichment analysis using TE candidates identified from mESC bulk RNA-seq and
161 Smart-seq. Enrichment of ERV elements was observed with bulk RNA-seq data, but not in single
162 cells. Smart-seq data of four single cells with different percentage of TE reads and merged Smart-
163 seq data from 10 cells were included.

164

165

166 **Counting scRNA-seq reads at individual TEs leads to large numbers of false positive**
167 **candidates**

168 Current TE expression analyses often quantify RNA-seq signal at individual TEs or TE
169 subfamilies/families (Day et al. 2010; Jin et al. 2015; Lerat et al. 2017; Jeong et al. 2018; Yang et
170 al. 2019; He et al. 2020; Jonsson et al. 2020). Our observation that a large proportion of scRNA-
171 seq reads map to TEs, especially intronic TEs, raises the concern that counting reads at single TEs
172 or TE subfamilies/families will aggregate noise and fail to exclude TEs that are part of protein-
173 coding genes, resulting in high numbers of false positive candidates. To test this, we applied a
174 similar strategy and analyzed bulk and Smart-seq datasets generated using mouse embryonic stem
175 cells (mESCs) cultured in 2i medium (Buettner et al. 2015; Kolodziejczyk et al. 2015). Because
176 mESCs represent a relatively homogeneous population and reads generated by Smart-seq and bulk
177 RNA-seq share a similar distribution along the gene body (Ramsköld et al. 2012), we expected
178 that the expression profiles obtained with scRNA-seq to be similar to those with bulk RNA-seq.

179 We first calculated the number of expressed protein-coding genes and TEs (CPM ≥ 1) in these
180 datasets. On average, 12,000 protein-coding genes and 6,000 TEs were detected in bulk RNA-seq
181 samples. In contrast, scRNA-seq captured an average of 7,000 protein-coding genes and 20,000
182 TEs per cell (Fig. 1B). To evaluate the quality of these TE candidates, we examined the following
183 three parameters: 1) The number of cells each candidate is expressed in, 2) the correlation between
184 the signal in bulk RNA-seq and the average signal across single cells, 3) the over-represented TE
185 families among all candidates. We reasoned that a true candidate should be expressed in a
186 relatively large number of mESCs and exhibit a strong correlation between its bulk RNA-seq and
187 averaged scRNA-seq signal. Strikingly, only protein-coding genes matched this expectation (Fig.

188 [1C, D, Supplemental Fig. S2A](#)). A large proportion of TE candidates were only detected in a small
189 number of cells and exhibited weak correlations between scRNA-seq and bulk RNA-seq signal
190 regardless of the expression cutoff. This observation remained valid after we performed the same
191 analysis by counting signals from individual exons. The exon length distributions were comparable
192 to those of TEs, ruling out the possibility that length discrepancy between TEs and protein-coding
193 genes contributes to false positive TE candidates ([Supplemental Fig. S2B-D](#)). We further
194 compared over-represented TEs within candidates identified from bulk RNA-seq and scRNA-seq
195 by performing a TE-family enrichment analysis ([Supplemental Fig. S3A](#)). Although ERV elements
196 have been shown to be expressed in stem cells (Macfarlan et al. 2012; Santoni et al. 2012; Fort et
197 al. 2014; Lu et al. 2014; Ohnuki et al. 2014; Wang et al. 2014), they were only enriched in bulk
198 RNA-seq in this analysis ([Fig. 1E](#)). scRNA-seq candidates obtained from this analysis were
199 depleted of ERVs and instead enriched for SINEs ([Fig. 1E, Supplemental Fig. S3B](#)), which are
200 often found near protein-coding genes and provide sequences that could act as reverse transcription
201 priming sites (Medstrand et al. 2002).

202 In summary, the extreme high number of TE candidates obtained from scRNA-seq, the weak signal
203 correlation between individual cells, as well as the discordance between bulk and scRNA-seq
204 strongly suggest that counting scRNA-seq reads at individual TEs will result in large numbers of
205 false positive candidates.

206 **Transcript assembly improves TE expression analysis**

207 Transcript annotation serves as the cornerstone for expression quantification. Our ability to
208 accurately assess the expression of protein-coding genes relies on well-annotated gene structures,
209 which help to focus analysis on genomics regions with true signal. Although individual TEs are

210 well annotated, it is usually unclear which TEs are expressed in a biological system and what the
211 underlying transcript structures are. We reason that the large number of false positive candidates
212 in scRNA-seq analysis is due to counting sparse and noisy signal at millions of TE copies, of which
213 only a small proportion are truly expressed ([Supplemental Fig. S3C](#)). Therefore, we hypothesize
214 that incorporating transcript structures of TE-containing ncRNAs into the analysis should help to
215 reduce noise.

216 To obtain ncRNAs with exonized TEs, we performed transcript assembly using mESC bulk RNA-
217 seq data. We identified 692 transcripts that overlap with TEs but not the exons of protein-coding
218 genes ([Fig. 2A, B](#), [Supplemental Fig. S4](#)). These transcripts were termed TE transcripts. To test
219 the accuracy of our assembly, we focused on the promoters of assembled TE transcripts and
220 examined several genomic signatures that are indicative of active transcription. Indeed, the
221 majority of our TE transcript promoters overlap with FANTOM5 CAGE peaks (FANTOM
222 Consortium and the RIKEN PMI and CLST (DGT) et al. 2014) and are enriched for ATAC-seq
223 signal while depleted of CpG methylation ([Fig. 2C](#)).

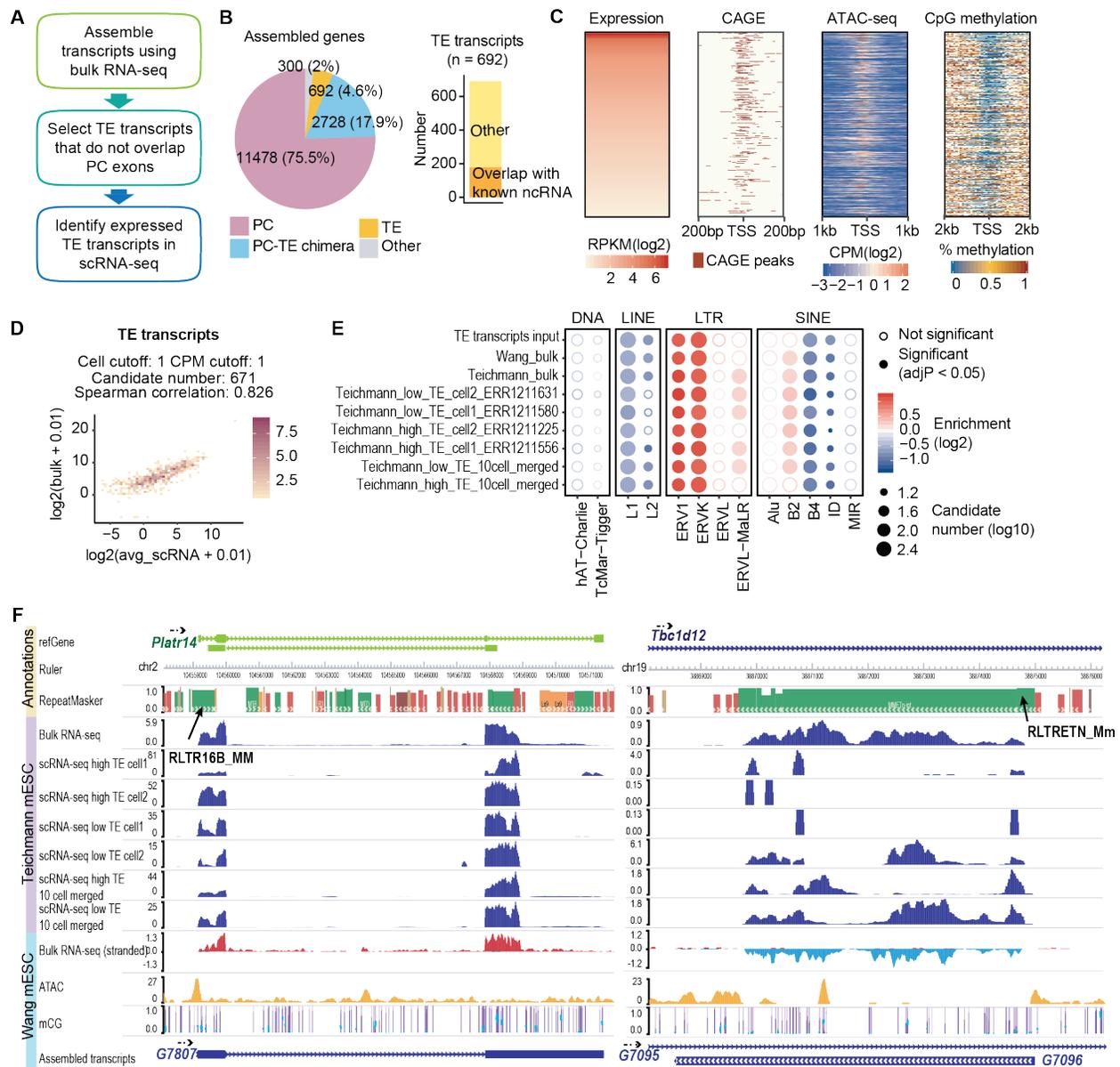
224 Utilizing these newly generated transcript models, we recalculated the expression of TE transcripts
225 and observed a much stronger correlation between mESC bulk and Smart-seq data ([Fig. 2D](#)). More
226 importantly, we obtained much more consistent TE-family enrichment results between bulk and
227 scRNA-seq and were able to identify the expression of ERV elements at single cells ([Fig. 2E, F](#)).

228 While Smart-seq based protocols generate deeper sequencing depth with reads covering the entire
229 gene-body, other popular scRNA-seq strategies such as 10x Genomics Chromium, Drop-seq and
230 SCR-seq produce shallow sequencing with reads biased towards the 5' or 3' end of the RNA
231 ([Ramsköld et al. 2012](#); [Soumillon et al. 2014](#); [Macosko et al. 2015](#); [Zheng et al. 2017](#)). Counting

232 reads at individual TEs using data with 5' or 3' signal enrichment will only capture TEs that are
233 located at either end of the transcripts, thus biasing our understanding about TE expression. We
234 reason that our approach can help to overcome this limitation by utilizing the annotation of full-
235 length TE transcripts.

236 To support our reasoning, we analyzed the dataset from a previously published mESC
237 differentiation study, in which single cell Smart-seq2, single cell SCRB-seq, and bulk RNA
238 processed with SCRB-seq protocol were performed (Semrau et al. 2017). Using individual TEs as
239 reference, we observed a severe discordance of TE expression between SCRB-seq and Smart-seq2,
240 likely due to differences in signal distribution along the transcripts ([Supplemental Fig. S5A](#)).
241 Conversely, using full-length TE transcripts as reference led to significantly improved signal
242 correlations ([Supplemental Fig. S5B](#)). Moreover, quantifying TE expression at transcript level
243 allowed us to recover the enrichment of ERVs from all datasets, whereas only Smart-seq2 showed
244 ERV enrichments when counting at individual TEs ([Supplemental Fig. S5C](#)).

245 Taken together, these results suggest that our analysis approach is applicable not only to scRNA-
246 seq data with high number of reads covering the entire transcript body, but also to other popular
247 scRNA-seq strategies that feature shallow sequencing depth at the 3' end of the transcripts.



248

249 **Figure 2 Transcript assembly improves scRNA-seq TE expression analysis**

250 (A) Flowchart of scRNA-seq TE quantification pipeline. In short, transcript assembly was
 251 performed with bulk RNA-seq data and transcripts that overlap with TEs but not protein-coding
 252 exons were utilized for expression quantification in scRNA-seq data.

253 (B) Transcript assembly using three mESC bulk RNA-seq data (Wang lab) yielded 692 TE

254 transcripts. Among these TE transcripts, 179 overlap with ncRNAs annotated by refSeq.

255 (C) FANTOM5 CAGE peaks, ATAC-seq signals and CpG methylation signals at the promoter
256 region of TE transcripts with RPKM ≥ 1 (Reads Per Kilobase Million).

257 (D) Correlation between mESC bulk RNA-seq and averaged Smart-seq (Teichmann lab) signals
258 at TE transcripts.

259 (E) TE-family enrichment analysis using expressed TE transcripts. Enrichment of ERV elements
260 was observed with both bulk RNA-seq and Smart-seq samples.

261 (F) Examples of TE transcript. Assembled TE transcripts, uniquely mapped reads of mESC bulk
262 RNA-seq, Smart-seq, merged Smart-seq, ATAC-seq and CpG methylation were included. Left: a
263 TE transcript that initiates from RLTR16b_MM. This TE transcript overlaps Platr14, a long
264 ncRNA known to impact the mESC differentiation-associated genes. Right: a TE transcript that
265 initiates from RLRETN_Mm. This transcript is largely composed of TEs and reflect the
266 transcription unit of ERV.

267

268 **Dynamic TE expression in pre-implementation embryos**

269 Encouraged by the results from the mESC data, we decided to apply our strategy to a more complex
270 biological system: mouse embryogenesis. The dynamic regulation of epigenome during
271 development not only fine-tunes protein-coding genes, but also allows specific TE expression at
272 different developmental stages (Rowe and Trono 2011; Gifford et al. 2013; Gerdes et al. 2016;
273 Rodriguez-Terrones and Torres-Padilla 2018; Deniz et al. 2019). Several recent studies utilized
274 scRNA-seq to profile the transcription landscape of mouse embryos from zygote to early
275 organogenesis, providing valuable resources for dissecting the dynamic expression of TEs.

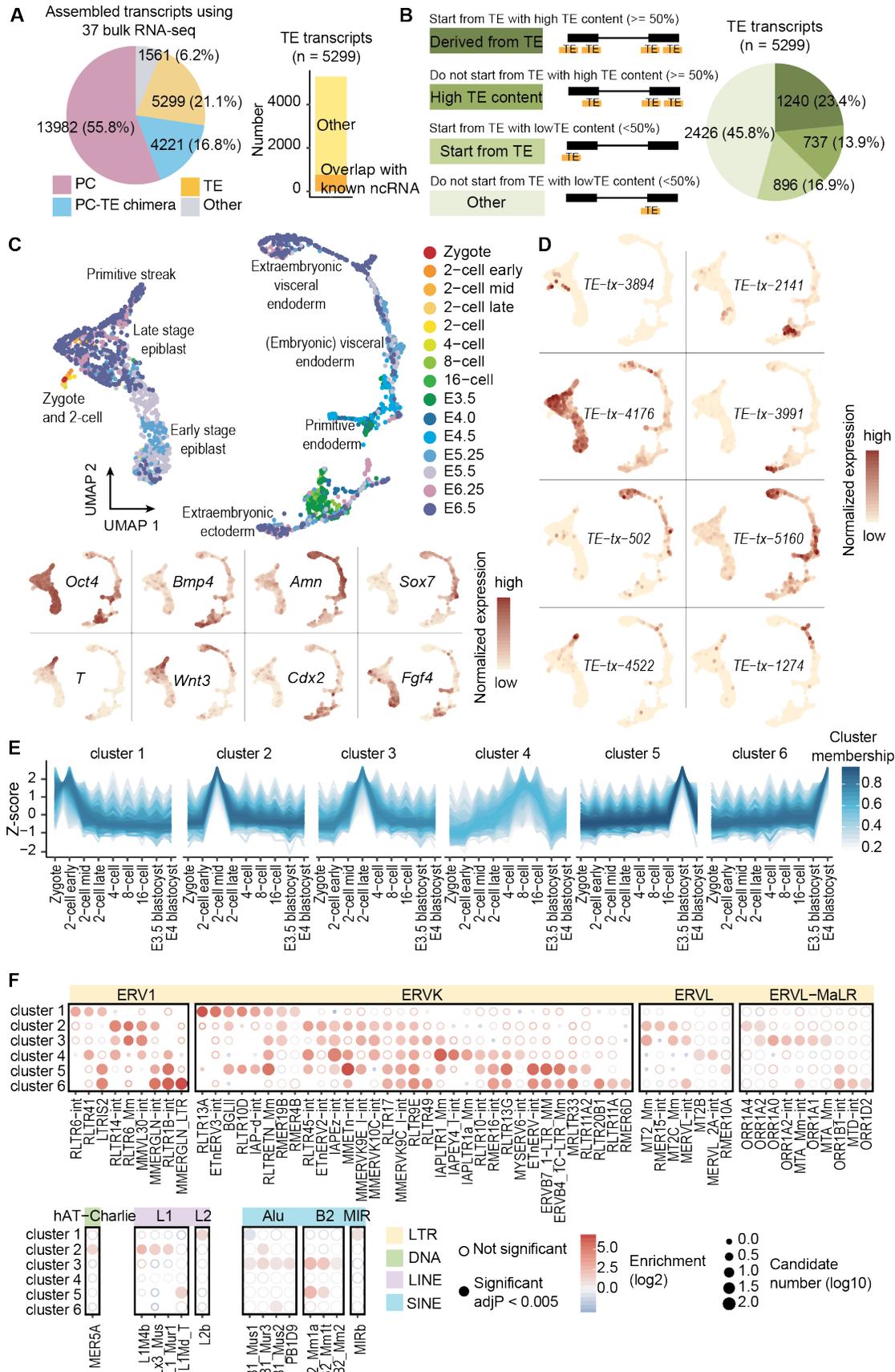
276 To facilitate TE expression quantification, we performed transcript assembly using 37 bulk RNA-
277 seq samples (Supplemental Table S1) that cover a range of tissues and developmental stages and
278 obtained 5299 TE transcripts (Fig. 3A, B, Supplemental Fig. S6A, Supplemental Table S2). 770
279 of these transcripts overlap with known ncRNAs annotated by refSeq (Fig. 3A). Compared with
280 assembled protein-coding transcripts, which show similar length and exon number as refSeq
281 protein-coding gene annotations, assembled TE transcripts are shorter in length and possess fewer
282 exons, a pattern consistent with annotated ncRNAs (Supplemental Fig. S6B, C).

283 Next, we analyzed three publicly available datasets where the transcription landscape of mouse
284 embryos from zygote to gastrulation was profiled using Smart-seq derived protocols (Deng et al.
285 2014; Mohammed et al. 2017; Cheng et al. 2019). In these datasets, a significant number of reads
286 overlap with TEs and about 3% of the total reads are mapped to TE transcripts (Supplemental Fig.
287 S7A, B). After data integration and dimension reduction using the top 4000 variable features, we
288 observed clear clustering patterns that were driven by cell type and developmental stage (Fig. 3C).
289 Among the top 4000 variable features, 377 are TE transcripts, suggesting that the expression of
290 TE transcripts could be cell-type- or developmental-stage-specific (Supplemental Fig. S7C-E).
291 Indeed, we were able to observe TE transcript expression with strong tissue or stage specificity
292 (Fig. 3D).

293 To further investigate the dynamics of TE transcription, we focused on pre-implantation stages,
294 where high TE expression was documented. Due to the limited number of cells, scRNA-seq signals
295 of each TE transcript across all the cells with the same developmental stage were averaged to
296 reduce noise. Grouping TE transcripts based on their expression patterns across pre-implantation
297 stages resulted in the following 6 clusters (Fig. 3E): TE transcripts that are maternally deposited

298 (cluster1), TE transcripts that are up-regulated during minor and major waves of zygotic genome
299 activation (clusters 2 and 3), TE transcripts that are upregulated during zygotic genome activation
300 and keep accumulating till the blastocyst stage (cluster 4), TE transcripts that are up-regulated in
301 the early- and mid-blastocyst stage (clusters 5 and 6).

302 We next performed TE enrichment analysis and observed that TE transcripts with distinct
303 expression profiles tend to be enriched for different TE subfamilies (Fig. 3F). For instance, IAP
304 elements are highly enriched in cluster 4, consistent with previous report that IAP expression
305 initiates from the 2-cell stage, accumulates and then disappears at the blastocyst stage (Pikó et al.
306 1984; Poznanski and Calarco 1991; Svoboda et al. 2004). We also observed the enrichment of
307 ERVL and ERVL-MaLR members in clusters 2 and 3, matching previous studies suggesting that
308 ERVL and ERVL-MaLR members are highly expressed during 2-cell stage, constituting about 5%
309 of the total transcripts (Kigami et al. 2003; Peaston et al. 2004; Svoboda et al. 2004). Furthermore,
310 consistent with functional studies that L1 contributes to the entry and exit of 2-cell stage
311 (Jachowicz et al. 2017; Percharde et al. 2018), we observed that the expression of L1 subfamily
312 members peaks at the 2 cell stage. Moreover, transcription factor binding site analysis using a 500
313 bp window upstream of TE transcripts identified footprints of transcription factors that are crucial
314 for mouse early embryogenesis such as Kruppel-like factors, GABPA and ELF3 (Ristevski et al.
315 2004; Kageyama et al. 2006; Bialkowska et al. 2017), suggesting shared regulatory networks
316 between TE transcripts and protein-coding genes (Supplemental Fig. S8).



318 **Figure 3 Dynamic TE expression in mouse pre-implantation embryos**

319 (A) 5299 TE transcripts were constructed using 37 bulk RNA-seq samples. 770 of these TE
320 transcripts overlap with ncRNAs annotated by refSeq.

321 (B) Over half of all the assembled TE transcripts either initiate from TEs or have more than 50%
322 of their exons composed of TEs.

323 (C) Upper panel: UMAP of scRNA-seq data from mouse zygote to E6.5 embryos. Cells were
324 colored based on developmental stages. Lower panel: expression of cell type specific markers.

325 (D) Examples of developmental stage- and tissue-specific TE transcripts.

326 (E) TE transcripts were grouped into 6 clusters based on their expression pattern across pre-
327 implantation stages.

328 (F) TE-subfamily enrichment analysis using TE transcripts within each of the 6 clusters.

329

330 **Tissue specific TE expression during mouse gastrulation and early organogenesis**

331 Comparing to pre-implantation stages, TE expression during gastrulation and organogenesis is
332 much less well studied and a comprehensive catalog of tissue specific TE transcripts is lacking.

333 To address this, we analyzed a 10x scRNA-seq dataset where more than 100k cells were assayed
334 using mouse E6.5 to E8.5 embryos (Pijuan-Sala et al. 2019). Comparing with mESC or mouse pre-

335 implantation data analyzed in previous sections, this E6.5 to E8.5 10x data contains considerably
336 less TE overlapping reads, with around 1% of the UMI mapping to TE transcripts ([Supplemental](#)

337 [Fig. S9A, B](#)). Even though TE transcripts are in general lowly expressed and lack the extreme

338 standardized variance observed at some protein-coding genes, they still constitute a small

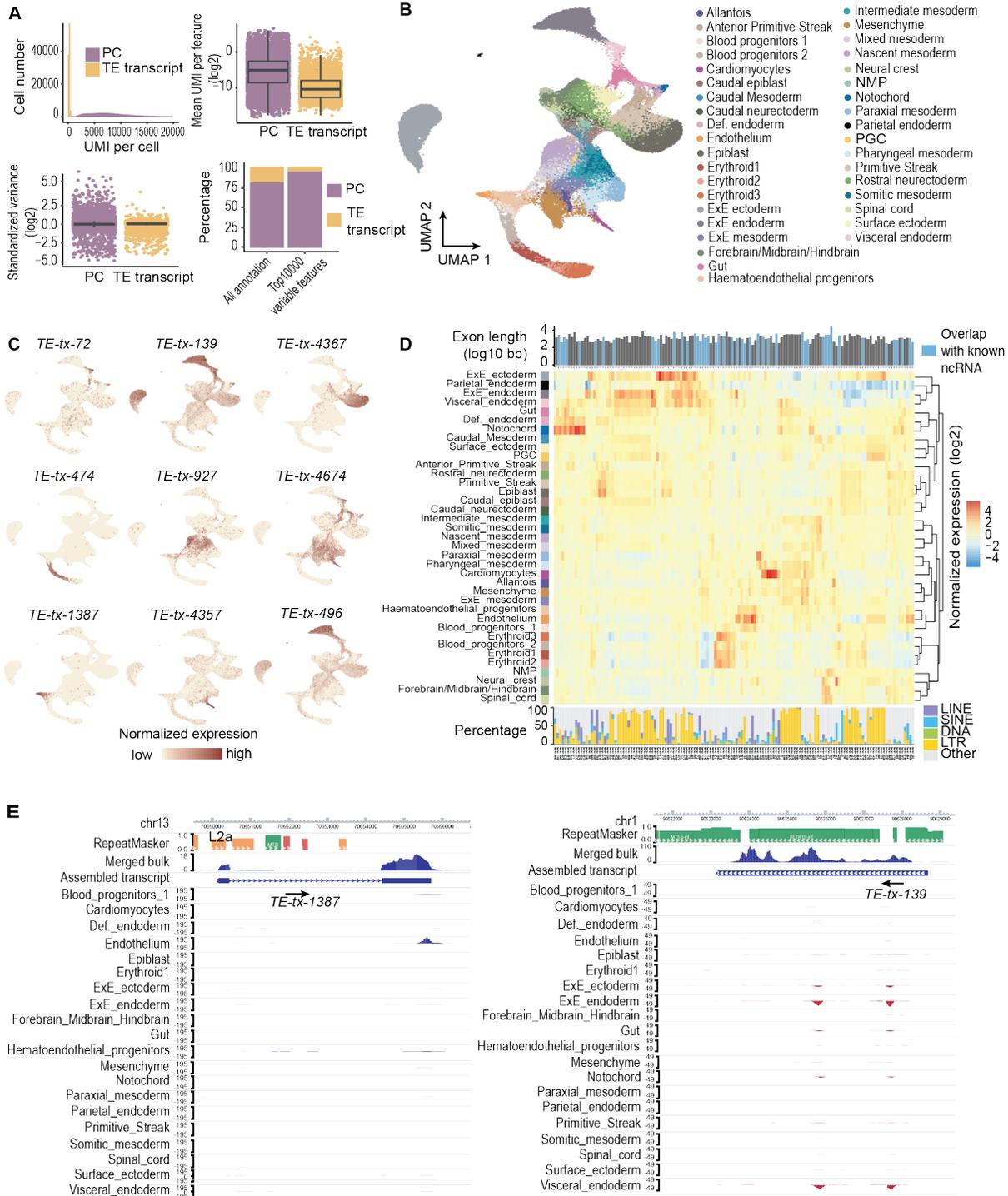
339 proportion of the top 1000 variable features that can be used to recapitulate the clustering pattern

340 in the original study (Fig 4A, B, Supplemental Fig. S9C). More importantly, we were able to
341 observe TE expression patterns that are enriched in small clusters of cells, suggesting that TE
342 transcripts display considerable tissue specificity during these stages (Fig. 4C).

343 To systematically examine the dynamic TE expression, we obtained 137 TE transcripts with
344 substantial tissue enrichment. Hierarchical clustering analysis using the expression of these TE
345 transcripts showed that tissues with similar origins are grouped together (Fig. 4D). For instance,
346 tissues within the hemato-endothelial lineage including hematoendothelial progenitors,
347 endothelium, blood progenitors and erythroids are adjacent to each other, and tissues linked to the
348 neuronal lineage including neuromesodermal progenitor, spinal cord,
349 forebrain/midbrain/hindbrain and neural crest are clustered together.

350 Even though 10x reads were enriched at the 3' end of the transcript, all TEs located along the
351 transcripts were captured using our assembled full-length TE transcripts (Fig. 4E). Among the 137
352 TE transcripts, 82 initiate from TE or have more than 50% of their exonic sequences contributed
353 by TEs (Fig. 4D, Supplemental Fig. S9D). Intriguingly, we observed that TE transcripts enriched
354 in different tissues display distinct TE composition. For example, while TE transcripts enriched in
355 hemato-lineage related tissues contain mostly non-TE sequences, TE transcripts enriched in
356 extraembryonic ectoderm, extraembryonic endoderm, parietal endoderm and visceral endoderm
357 are almost exclusively composed of LTRs. Furthermore, overlapping these TE transcripts with
358 annotated ncRNAs revealed that 47 have been annotated by refSeq. Importantly, we observed that
359 while these known ncRNAs tend to contain a lower percentage of TEs, transcripts that are mostly
360 exclusively composed of TEs are largely un-annotated, demonstrating the value of our approach
361 in capturing transcripts that originated from highly repetitive regions.

362



363

364 **Figure 4 Tissue-specific TE expression during mouse gastrulation and early organogenesis**

365 (A) Upper left: Fewer unique molecular identifiers (UMIs) were mapped to TE transcripts than to
366 protein-coding genes. Upper right: the averaged expression level of TE transcripts across all the
367 cells was lower compared to protein-coding genes. Lower left: TE transcripts lack the extreme
368 standardized variance observed at protein-coding genes. Lower right: TE transcripts account for
369 55 of the top 1000 variable features.

370 (B) UMAP of scRNA-seq data, cells were colored based on tissue information provided by the
371 original study.

372 (C) Examples of tissue specific TE transcripts.

373 (D) Normalized expression pattern (center, heatmap) of 137 TE transcripts (columns) across 37
374 tissues (rows). Transcript length, annotation status (top, bar plot) and TE composition (bottom, bar
375 plot) were shown for each TE transcript.

376 (E) Genome browser view of two TE transcripts with strong tissue enrichment. Assembled TE
377 transcripts, uniquely mapped reads of merged bulk RNA-seq (from 37 samples that were used for
378 transcript assembly) and scRNA-seq signal for selected tissues were shown.

379

380 **Discussion:**

381 Current genome-wide TE expression quantification tools often count signal at individual TEs or
382 TE subfamilies/families (Day et al. 2010; Jin et al. 2015; Lerat et al. 2017; Jeong et al. 2018; Yang
383 et al. 2019; He et al. 2020; Jonsson et al. 2020). This strategy has been widely adopted in bulk
384 RNA-seq and inspired similar analyses with scRNA-seq data (Göke et al. 2015; Ge 2017; Boroviak
385 et al. 2018; Brocks et al. 2018; Yandım and Karakulah 2019). However, we caution that compared
386 with bulk RNA-seq, a higher percentage of scRNA-seq reads are mapped to TEs, making it

387 challenging to identify *bona fide* TE expression. Moreover, quantifying signal at individual TEs
388 or TE subfamilies/families leads to the false impression that transcripts originated from repetitive
389 regions are mostly composed of a single TE or TEs from the same subfamilies/families. While this
390 is true for some well-studied examples such as full length ERVs, in most other cases, TEs only
391 contribute to fragments of the full-length transcript and TEs from different families can be
392 incorporated into the same transcript.

393 A major difference between the expression quantification of protein-coding genes and TEs is that
394 the transcript structures of protein-coding genes are usually well annotated and readily available.
395 Gene annotation guides expression analysis towards genomic regions with true signal and
396 facilitates accurate expression quantification with scRNA-seq data. In this study, we demonstrated
397 that transcripts constructed from bulk RNA-seq can serve as references for TE-containing ncRNAs
398 and improve the accuracy of TE expression analysis in scRNA-seq data generated across multiple
399 sequencing platforms. In comparison to individual TEs or TE subfamilies/families, TE transcripts
400 more accurately reflect the natural transcription units. These transcripts contain not only previously
401 annotated TE transcription units, but also novel ncRNAs that are partially composed of TEs. Out
402 of the 5299 TE transcripts that we assembled, 98 closely resemble the well-studied transcription
403 units of ERVs. These transcripts have more than 80% of their exonic sequences contributed by
404 TEs. They start from 5' LTR, transcribe through internal sequences and end at 3' LTR. Notably,
405 we also obtained another 104 TE transcripts that initiate from TEs and are within 100 bp away
406 from FANTOM5 CAGE peaks (FANTOM Consortium and the RIKEN PMI and CLST (DGT) et
407 al. 2014). 84 of these transcripts were not previously annotated, highlighting the value of our
408 approach in identifying TEs that potentially function as promoters.

409 Using our analytical pipeline, we dissected the expression dynamics of TE transcripts during
410 mouse early embryogenesis and identified 137 TE containing ncRNAs with strong tissue
411 specificity. A close examination of these candidates revealed intricate interaction between TE
412 transcripts and protein-coding genes. For instance, we were able to identify a ChIP-seq peak of
413 Regulatory Factor X3 (RFX3) at the promoter region of the TE transcript *TE-tx-3856*
414 ([Supplemental Fig. 10A](#)). RFX3 is a transcription factor essential for brain development (Baas et
415 al. 2006; Benadiba et al. 2012; Magnani et al. 2015). Our observation that *TE-tx-3856* is highly
416 expressed in mouse neuronal tissues suggests that this TE containing ncRNA is a potential
417 downstream target of RFX3. In another interesting example, the TE transcript *TE-tx-3715* overlaps
418 with *Sonic Hedgehog (Shh)*, a secreted signaling molecule produced by the notochord (Placzek
419 1995; McMahon et al. 2003). Intriguingly, the expression pattern of *TE-tx-3715* strongly resembles
420 that of *Shh*, indicating that they are under the control of a common regulatory circuit ([Supplemental](#)
421 [Fig. S10B](#)). In addition, we also captured TE transcripts *TE-tx-3178* and *TE-tx-2841*, both of which
422 are strongly expressed in the epiblast ([Supplemental Fig. S10C](#)). Both candidates initiate from TEs
423 and were previously annotated as pluripotent associated transcripts (*Platr10* and *Platr14*)
424 (Bergmann et al. 2015). Earlier reports suggested that *Platr10* transcript physically interacts with
425 the promoter of pluripotent transcription factor *Sox2* while the depletion of *Platr14* alters the
426 expression of differentiation- and development-associated genes in stem cells (Bergmann et al.
427 2015; Zhang et al. 2019). Our observation that *Platr10* and *Platr14* are expressed in the epiblast
428 suggests that they may play similar roles during mouse early embryogenesis. Taken together, we
429 dissected the dynamic TE expression during mouse early development and provided a curated list
430 of promising TE candidates for future functional studies.

431 In summary, we established an effective TE quantification pipeline for scRNA-seq data and
432 illustrated the dynamic TE expression during mouse early embryogenesis. In contrast to commonly
433 used bulk RNA-seq tools that evaluate reads at single TEs or TE subfamilies/families, our pipeline
434 emphasizes the importance of full-length TE transcript structures in scRNA-seq TE quantification.
435 Furthermore, our work provides an initial set of TE transcript references during mouse early
436 development through transcript assembly, laying the foundation for future work on constructing a
437 more comprehensive TE transcript database across tissues and developmental stages. Additionally,
438 exploring alternative quantification methods for ambiguously mapped reads, such as using
439 expectation-maximization algorithm or Bayesian-based framework, as well as developing
440 isoform-specific quantification tools for TE-protein-coding-gene chimeras will further expand the
441 TE analysis toolkit for scRNA-seq and greatly advance our knowledge on the expression and the
442 function of TE transcripts.

443 **Methods:**

444 **Data processing and signal quantification of bulk RNA-seq datasets:**

445 Raw sequencing files were downloaded from NCBI Sequence Read Archive and EMBL-EBI
446 ArrayExpress ([Supplemental Table S1](#)) and aligned to the mouse (mm10) or human (hg38)
447 genomes using STAR. To retain reads derived from repetitive regions, “--outFilterMultimapNmax”
448 was set to 500. To facilitate downstream transcript assembly “--outSAMattributes” was set to “NH
449 HI NM MD XS AS”. After alignment, signal quantification at regions of interests was performed
450 using featureCount. Reads aligned to multiple locations or overlapping multiple features were
451 evenly distributed by enabling “-O -M --fraction”.

452 **Data processing and signal quantification of scRNA-seq datasets:**

453 scRNA-seq data generated with Smart-seq derived protocols were processed and quantified using
454 the same procedures as bulk RNA-seq data. scRNA-seq data generated with the other protocols
455 were processed using zUMIs with the following modifications: 1) To retain reads derived from
456 repetitive regions, “--outFilterMultimapNmax” was set to 500 during STAR alignment. 2) To
457 quantify reads that were mapped to multiple locations or features, “allowMultiOverlap”and
458 “countMultiMappingReads” were set to TRUE for function “.runFeatureCount”. Bam files with
459 cell barcode, UMI and the name of overlapping features were reported. 3) A customized R script
460 was used to process the bam file generated in step2. Reads that were mapped to multiple locations
461 or features were evenly distributed. UMIs were then evaluated at each feature.

462

463 Reads of the scRNA-seq datasets from mouse zygote to gastrulation (Smart-seq derived protocols)
464 were quantified at protein-coding genes (refSeq annotation, n= 20779) and TE transcripts
465 (assembled from bulk RNA-seq, n=5299). Only cells that had 200 to 18000 features and less than
466 10% mitochondria reads were kept. To remove batch effect and visualize all three datasets in the
467 same UMAP, data integration was performed using Seurat with the top 4000 variable features.
468 Cell type was determined using the stage information provided by the original studies, the
469 expression patterns of cell type specific markers and Seurat clustering results.

470
471 UMIs of the 10x scRNA-seq dataset from mouse gastrulation to early organogenesis were
472 quantified at protein-coding genes (refSeq annotation, n= 20779) and TE transcripts (assembled
473 from bulk RNA-seq, n=5299). Sample_25 was removed due to higher batch effect. Only cells that
474 have more than 200 features and were annotated by the original study were kept. Cell type
475 information provided by the original study was utilized for identifying tissue specific markers. The
476 137 TE transcripts with strong tissue enrichment were obtained by combining and filtering Seurat-
477 defined markers and customized markers. Seurat-defined markers were obtained by running
478 “FindAllMarkers” with “only.pos = T, min.pct = 0.15” and selecting for TE transcripts with
479 adjusted p-value < 0.05. Customized markers were obtained by identifying TE transcripts with at
480 least 1 UMI in at least 10% of the cells in any tissue and selecting candidates that were expressed
481 in less than 4 tissues. After combining Seurat-defined markers and customized markers, manual
482 curation was performed to remove candidates that were highly expressed in large number of tissues
483 or with sub-optimal transcript structures.

484 **TE transcript clustering in mouse pre-implantation stages:**

485 Due to the limited number of cells, scRNA-seq signals of each TE transcript across all the cells
486 with the same developmental stage were averaged to reduce noise. TE transcripts were then
487 grouped into 6 clusters using soft clustering (R package TCseq) based on their expression patterns
488 across pre-implantation stages.

489

490 **Constructing TE transcripts:**

491 Transcript assembly of each RNA-seq sample was performed using StringTie2. “-j 2 -s 5 -f 0.05 -c
492 2” was used to improve the specificity of the assembly results. To generate the master reference
493 file, assembled transcripts from multiple RNA-seq samples were then merged using TACO with
494 default parameters. Transcripts that overlapped with TEs but not the exons of refSeq protein-
495 coding genes were named as TE transcripts and utilized for TE expression quantification.

496

497 **TE subfamily/family enrichment analysis:**

498 For each TE-family, its enrichment was calculated using the following equation: The observed
499 frequency of TEs belonging to this family in all candidates divided by the expected frequency of
500 TEs belonging to this family in genomic regions that do not overlap with protein-coding genes.
501 The significance for the observed frequency was calculated with Fisher's exact test and corrected
502 for multiple testing with the Benjamini–Hochberg method. Only TE families with more than 20
503 members in the candidates and more than 100 members in the background were included in the
504 figures. TE-subfamily enrichment analysis was performed similarly. Only TE-subfamilies that

505 were significantly enriched in the candidates, had more than 10 members in the candidates and
506 more than 100 members in the background and were plotted.

507 **Publicly available datasets utilized in this study:**

508 Descriptions and accession IDs of all the datasets used in this manuscript are provided in
509 Supplemental Table S1.

510 **Data Access:**

511 All of the datasets analyzed in the paper are currently available in public databases, accession IDs
512 and descriptions are listed in the Supplemental Table S1.

513 **Acknowledgements:**

514 We thank Kara Quaid and Shuonan He for comments on the manuscript, and all the Ting Wang
515 lab members for discussions and critical inputs.

516 **Funding:**

517 This work is supported by NIH grants R01HG007175, U24ES026699, U01CA200060,
518 U01HG009391 and U41HG010972.

519 **Conflicts of interest:**

520 The authors have no conflicts of interest or financial ties to disclose.

521

522

- 523 References:
- 524 Anwar SL, Wulaningsih W, Lehmann U. 2017. Transposable elements in human cancer: causes
525 and consequences of deregulation. *Int J Mol Sci* **18**.
- 526 Baas D, Meiniel A, Benadiba C, Bonnafe E, Meiniel O, Reith W, Durand B. 2006. A deficiency
527 in RFX3 causes hydrocephalus associated with abnormal differentiation of ependymal
528 cells. *Eur J Neurosci* **24**: 1020–1030.
- 529 Benadiba C, Magnani D, Niquille M, Morlé L, Valloton D, Nawabi H, Ait-Lounis A, Otsmane
530 B, Reith W, Theil T, et al. 2012. The ciliogenic transcription factor RFX3 regulates early
531 midline distribution of guidepost neurons required for corpus callosum development.
532 *PLoS Genet* **8**: e1002606.
- 533 Bendall ML, de Mulder M, Iñiguez LP, Lecanda-Sánchez A, Pérez-Losada M, Ostrowski MA,
534 Jones RB, Mulder LCF, Reyes-Terán G, Crandall KA, et al. 2019. Telescope:
535 Characterization of the retrotranscriptome by accurate estimation of transposable element
536 expression. *PLoS Comput Biol* **15**: e1006453.
- 537 Bergmann JH, Li J, Eckersley-Maslin MA, Rigo F, Freier SM, Spector DL. 2015. Regulation of
538 the ESC transcriptome by nuclear long noncoding RNAs. *Genome Res* **25**: 1336–1346.
- 539 Bialkowska AB, Yang VW, Mallipattu SK. 2017. Krüppel-like factors in mammalian stem cells
540 and development. *Development* **144**: 737–754.
- 541 Boroviak T, Stirparo GG, Dietmann S, Hernando-Herraez I, Mohammed H, Reik W, Smith A,
542 Sasaki E, Nichols J, Bertone P. 2018. Single cell transcriptome analysis of human,
543 marmoset and mouse embryos reveals common and divergent features of preimplantation
544 development. *Development* **145**.
- 545 Brocks D, Chomsky E, Mukamel Z, Lifshitz A, Tanay A. 2018. Single cell analysis reveals
546 dynamics of transposable element transcription following epigenetic de-repression.
547 *BioRxiv*.
- 548 Buettner F, Natarajan KN, Casale FP, Proserpio V, Scialdone A, Theis FJ, Teichmann SA,
549 Marioni JC, Stegle O. 2015. Computational analysis of cell-to-cell heterogeneity in
550 single-cell RNA-sequencing data reveals hidden subpopulations of cells. *Nat Biotechnol*
551 **33**: 155–160.
- 552 Cheng S, Pei Y, He L, Peng G, Reinius B, Tam PPL, Jing N, Deng Q. 2019. Single-Cell RNA-
553 Seq Reveals Cellular Heterogeneity of Pluripotency Transition and X Chromosome
554 Dynamics during Early Mouse Development. *Cell Rep* **26**: 2593–2607.e3.
- 555 Chuong EB, Elde NC, Feschotte C. 2017. Regulatory activities of transposable elements: from
556 conflicts to benefits. *Nat Rev Genet* **18**: 71–86.
- 557 Cowley M, Oakey RJ. 2013. Transposable elements re-wire and fine-tune the transcriptome.
558 *PLoS Genet* **9**: e1003234.
- 559 Day DS, Luquette LJ, Park PJ, Kharchenko PV. 2010. Estimating enrichment of repetitive
560 elements from high-throughput sequence data. *Genome Biol* **11**: R69.
- 561 De Iaco A, Planet E, Coluccio A, Verp S, Duc J, Trono D. 2017. DUX-family transcription
562 factors regulate zygotic genome activation in placental mammals. *Nat Genet* **49**: 941–
563 945.
- 564 Deng Q, Ramsköld D, Reinius B, Sandberg R. 2014. Single-cell RNA-seq reveals dynamic,
565 random monoallelic gene expression in mammalian cells. *Science* **343**: 193–196.
- 566 Deniz Ö, Frost JM, Branco MR. 2019. Regulation of transposable elements by DNA
567 modifications. *Nat Rev Genet* **20**: 417–431.
- 568 Fadloun A, Le Gras S, Jost B, Ziegler-Birling C, Takahashi H, Gorab E, Carninci P, Torres-

- 569 Padilla M-E. 2013. Chromatin signatures and retrotransposon profiling in mouse embryos
570 reveal regulation of LINE-1 by RNA. *Nat Struct Mol Biol* **20**: 332–338.
- 571 FANTOM Consortium and the RIKEN PMI and CLST (DGT), Forrest ARR, Kawaji H, Rehli
572 M, Baillie JK, de Hoon MJL, Haberle V, Lassmann T, Kulakovskiy IV, Lizio M, et al.
573 2014. A promoter-level mammalian expression atlas. *Nature* **507**: 462–470.
- 574 Feschotte C, Gilbert C. 2012. Endogenous viruses: insights into viral evolution and impact on
575 host biology. *Nat Rev Genet* **13**: 283–296.
- 576 Fort A, Hashimoto K, Yamada D, Salimullah M, Keya CA, Saxena A, Bonetti A, Voineagu I,
577 Bertin N, Kratz A, et al. 2014. Deep transcriptome profiling of mammalian stem cells
578 supports a regulatory role for retrotransposons in pluripotency maintenance. *Nat Genet*
579 **46**: 558–566.
- 580 Garcia-Perez JL, Widmann TJ, Adams IR. 2016. The impact of transposable elements on
581 mammalian development. *Development* **143**: 4101–4114.
- 582 Ge SX. 2017. Exploratory bioinformatics investigation reveals importance of “junk” DNA in
583 early embryo development. *BMC Genomics* **18**: 200.
- 584 Gerdes P, Richardson SR, Mager DL, Faulkner GJ. 2016. Transposable elements in the
585 mammalian embryo: pioneers surviving through stealth and service. *Genome Biol* **17**:
586 100.
- 587 Gifford WD, Pfaff SL, Macfarlan TS. 2013. Transposable elements as genetic regulatory
588 substrates in early development. *Trends Cell Biol* **23**: 218–226.
- 589 Göke J, Lu X, Chan Y-S, Ng H-H, Ly L-H, Sachs F, Szczerbinska I. 2015. Dynamic
590 transcription of distinct classes of endogenous retroviral elements marks specific
591 populations of early human embryonic cells. *Cell Stem Cell* **16**: 135–141.
- 592 Grow EJ, Flynn RA, Chavez SL, Bayless NL, Wossidlo M, Wesche DJ, Martin L, Ware CB,
593 Blish CA, Chang HY, et al. 2015. Intrinsic retroviral reactivation in human
594 preimplantation embryos and pluripotent cells. *Nature* **522**: 221–225.
- 595 Guffanti G, Bartlett A, Klengel T, Klengel C, Hunter R, Glinsky G, Macciardi F. 2018. Novel
596 Bioinformatics Approach Identifies Transcriptional Profiles of Lineage-Specific
597 Transposable Elements at Distinct Loci in the Human Dorsolateral Prefrontal Cortex. *Mol*
598 *Biol Evol* **35**: 2435–2453.
- 599 Hadjiargyrou M, Delihias N. 2013. The intertwining of transposable elements and non-coding
600 RNAs. *Int J Mol Sci* **14**: 13307–13328.
- 601 He J, Babarinde IA, Sun L, Xu S, Chen R, Wei Y, Li Y, Ma G, Zhuang Q, Hutchins A, et al.
602 2020. Unveiling transposable element expression heterogeneity in cell fate regulation at
603 the single-cell level. *BioRxiv*.
- 604 Hendrickson PG, Doráis JA, Grow EJ, Whiddon JL, Lim J-W, Wike CL, Weaver BD, Pflueger
605 C, Emery BR, Wilcox AL, et al. 2017. Conserved roles of mouse DUX and human
606 DUX4 in activating cleavage-stage genes and MERVL/HERVL retrotransposons. *Nat*
607 *Genet* **49**: 925–934.
- 608 Huang Y, Kim JK, Do DV, Lee C, Penfold CA, Zylicz JJ, Marioni JC, Hackett JA, Surani MA.
609 2017. Stella modulates transcriptional and endogenous retrovirus programs during
610 maternal-to-zygotic transition. *Elife* **6**.
- 611 Hutchins AP, Pei D. 2015. Transposable elements at the center of the crossroads between
612 embryogenesis, embryonic stem cells, reprogramming, and long non-coding RNAs. *Sci*
613 *Bull (Beijing)* **60**: 1722–1733.
- 614 Ilicic T, Kim JK, Kolodziejczyk AA, Bagger FO, McCarthy DJ, Marioni JC, Teichmann SA.

- 615 2016. Classification of low quality cells from single-cell RNA-seq data. *Genome Biol* **17**:
616 29.
- 617 Jachowicz JW, Bing X, Pontabry J, Bošković A, Rando OJ, Torres-Padilla M-E. 2017. LINE-1
618 activation after fertilization regulates global chromatin accessibility in the early mouse
619 embryo. *Nat Genet* **49**: 1502–1510.
- 620 Jeong H-H, Yalamanchili HK, Guo C, Shulman JM, Liu Z. 2018. An ultra-fast and scalable
621 quantification pipeline for transposable elements from next generation sequencing data.
622 *Pac Symp Biocomput* **23**: 168–179.
- 623 Jin Y, Tam OH, Paniagua E, Hammell M. 2015. TETranscripts: a package for including
624 transposable elements in differential expression analysis of RNA-seq datasets.
625 *Bioinformatics* **31**: 3593–3599.
- 626 Jonsson ME, Garza R, Sharma Y, Petri R, Sodersten E, Johansson JG, Johansson PA, Atacho
627 DA, Piracs K, Madsen S, et al. 2020. Activation of endogenous retroviruses during brain
628 development causes neuroinflammation. *BioRxiv*.
- 629 Kageyama S, Liu H, Nagata M, Aoki F. 2006. The role of ETS transcription factors in
630 transcription and development of mouse preimplantation embryos. *Biochem Biophys Res*
631 *Commun* **344**: 675–679.
- 632 Kapusta A, Kronenberg Z, Lynch VJ, Zhuo X, Ramsay L, Bourque G, Yandell M, Feschotte C.
633 2013. Transposable elements are major contributors to the origin, diversification, and
634 regulation of vertebrate long noncoding RNAs. *PLoS Genet* **9**: e1003470.
- 635 Kelley D, Rinn J. 2012. Transposable elements reveal a stem cell-specific class of long
636 noncoding RNAs. *Genome Biol* **13**: R107.
- 637 Kigami D, Minami N, Takayama H, Imai H. 2003. MuERV-L is one of the earliest transcribed
638 genes in mouse one-cell embryos. *Biol Reprod* **68**: 651–654.
- 639 Kolodziejczyk AA, Kim JK, Tsang JCH, Ilicic T, Henriksson J, Natarajan KN, Tuck AC, Gao X,
640 Bühler M, Liu P, et al. 2015. Single Cell RNA-Sequencing of Pluripotent States Unlocks
641 Modular Transcriptional Variation. *Cell Stem Cell* **17**: 471–485.
- 642 La Manno G, Soldatov R, Zeisel A, Braun E, Hochgerner H, Petukhov V, Lidschreiber K,
643 Kastrioti ME, Lönnerberg P, Furlan A, et al. 2018. RNA velocity of single cells. *Nature*
644 **560**: 494–498.
- 645 Lerat E, Fablet M, Modolo L, Lopez-Maestre H, Vieira C. 2017. TETools facilitates big data
646 expression analysis of transposable elements and reveals an antagonism between their
647 activity and that of piRNA genes. *Nucleic Acids Res* **45**: e17.
- 648 Lu X, Sachs F, Ramsay L, Jacques P-É, Göke J, Bourque G, Ng H-H. 2014. The retrovirus
649 HERVH is a long noncoding RNA required for human embryonic stem cell identity. *Nat*
650 *Struct Mol Biol* **21**: 423–425.
- 651 Macfarlan TS, Gifford WD, Driscoll S, Lettieri K, Rowe HM, Bonanomi D, Firth A, Singer O,
652 Trono D, Pfaff SL. 2012. Embryonic stem cell potency fluctuates with endogenous
653 retrovirus activity. *Nature* **487**: 57–63.
- 654 Macosko EZ, Basu A, Satija R, Nemes J, Shekhar K, Goldman M, Tirosh I, Bialas AR,
655 Kamitaki N, Martersteck EM, et al. 2015. Highly Parallel Genome-wide Expression
656 Profiling of Individual Cells Using Nanoliter Droplets. *Cell* **161**: 1202–1214.
- 657 Magnani D, Morlé L, Hasenpusch-Theil K, Paschaki M, Jacoby M, Schurmans S, Durand B,
658 Theil T. 2015. The ciliogenic transcription factor Rfx3 is required for the formation of the
659 thalamocortical tract by regulating the patterning of prethalamus and ventral
660 telencephalon. *Hum Mol Genet* **24**: 2578–2593.

- 661 Maksakova IA, Mager DL. 2005. Transcriptional regulation of early transposon elements, an
662 active family of mouse long terminal repeat retrotransposons. *J Virol* **79**: 13865–13874.
- 663 McMahon AP, Ingham PW, Tabin CJ. 2003. 1 Developmental roles and clinical significance of
664 Hedgehog signaling. In *Current topics in developmental biology volume 53*, Vol. 53 of
665 *Current topics in developmental biology*, pp. 1–114, Elsevier.
- 666 Medstrand P, van de Lagemaat LN, Mager DL. 2002. Retroelement distributions in the human
667 genome: variations associated with age and proximity to genes. *Genome Res* **12**: 1483–
668 1495.
- 669 Mohammed H, Hernando-Herraez I, Savino A, Scialdone A, Macaulay I, Mulas C, Chandra T,
670 Voet T, Dean W, Nichols J, et al. 2017. Single-Cell Landscape of Transcriptional
671 Heterogeneity and Cell Fate Decisions during Mouse Early Gastrulation. *Cell Rep* **20**:
672 1215–1228.
- 673 Ohnuki M, Tanabe K, Sutou K, Teramoto I, Sawamura Y, Narita M, Nakamura M, Tokunaga Y,
674 Nakamura M, Watanabe A, et al. 2014. Dynamic regulation of human endogenous
675 retroviruses mediates factor-induced reprogramming and differentiation potential. *Proc*
676 *Natl Acad Sci USA* **111**: 12426–12431.
- 677 Peaston AE, Evsikov AV, Graber JH, de Vries WN, Holbrook AE, Solter D, Knowles BB. 2004.
678 Retrotransposons regulate host genes in mouse oocytes and preimplantation embryos.
679 *Dev Cell* **7**: 597–606.
- 680 Percharde M, Lin C-J, Yin Y, Guan J, Peixoto GA, Bulut-Karslioglu A, Biechele S, Huang B,
681 Shen X, Ramalho-Santos M. 2018. A LINE1-Nucleolin Partnership Regulates Early
682 Development and ESC Identity. *Cell* **174**: 391–405.e19.
- 683 Pijuan-Sala B, Griffiths JA, Guibentif C, Hiscock TW, Jawaid W, Calero-Nieto FJ, Mulas C,
684 Ibarra-Soria X, Tyser RCV, Ho DLL, et al. 2019. A single-cell molecular map of mouse
685 gastrulation and early organogenesis. *Nature* **566**: 490–495.
- 686 Pikó L, Hammons MD, Taylor KD. 1984. Amounts, synthesis, and some properties of
687 intracisternal A particle-related RNA in early mouse embryos. *Proc Natl Acad Sci USA*
688 **81**: 488–492.
- 689 Placzek M. 1995. The role of the notochord and floor plate in inductive interactions. *Curr Opin*
690 *Genet Dev* **5**: 499–506.
- 691 Poznanski AA, Calarco PG. 1991. The expression of intracisternal A particle genes in the
692 preimplantation mouse embryo. *Dev Biol* **143**: 271–281.
- 693 Ramsköld D, Luo S, Wang Y-C, Li R, Deng Q, Faridani OR, Daniels GA, Khrebtukova I, Loring
694 JF, Laurent LC, et al. 2012. Full-length mRNA-Seq from single-cell levels of RNA and
695 individual circulating tumor cells. *Nat Biotechnol* **30**: 777–782.
- 696 Rebollo R, Romanish MT, Mager DL. 2012. Transposable elements: an abundant and natural
697 source of regulatory sequences for host genes. *Annu Rev Genet* **46**: 21–42.
- 698 Ristevski S, O’Leary DA, Thornell AP, Owen MJ, Kola I, Hertzog PJ. 2004. The ETS
699 transcription factor GABPalpha is essential for early embryogenesis. *Mol Cell Biol* **24**:
700 5844–5849.
- 701 Rodriguez-Terrones D, Torres-Padilla M-E. 2018. Nimble and ready to mingle: transposon
702 outbursts of early development. *Trends Genet* **34**: 806–820.
- 703 Rowe HM, Trono D. 2011. Dynamic control of endogenous retroviruses during development.
704 *Virology* **411**: 273–287.
- 705 Santoni FA, Guerra J, Luban J. 2012. HERV-H RNA is abundant in human embryonic stem cells
706 and a precise marker for pluripotency. *Retrovirology* **9**: 111.

- 707 Semrau S, Goldmann JE, Soumillon M, Mikkelsen TS, Jaenisch R, van Oudenaarden A. 2017.
708 Dynamics of lineage commitment revealed by single-cell transcriptomics of
709 differentiating embryonic stem cells. *Nat Commun* **8**: 1096.
- 710 Soumillon M, Cacchiarelli D, Semrau S, van Oudenaarden A, Mikkelsen TS. 2014.
711 Characterization of directed differentiation by high-throughput single-cell RNA-Seq.
712 *BioRxiv*.
- 713 Srivastava A, Sarkar H, Gupta N, Patro R. 2016. RapMap: a rapid, sensitive and accurate tool for
714 mapping RNA-seq reads to transcriptomes. *Bioinformatics* **32**: i192–i200.
- 715 Sundaram V, Wysocka J. 2020. Transposable elements as a potent source of diverse cis-
716 regulatory sequences in mammalian genomes. *Philos Trans R Soc Lond B, Biol Sci* **375**:
717 20190347.
- 718 Svoboda P, Stein P, Anger M, Bernstein E, Hannon GJ, Schultz RM. 2004. RNAi and expression
719 of retrotransposons MuERV-L and IAP in preimplantation mouse embryos. *Dev Biol*
720 **269**: 276–285.
- 721 Valdebenito-Maturana B, Riadi G. 2018. TEcandidates: prediction of genomic origin of
722 expressed transposable elements using RNA-seq data. *Bioinformatics* **34**: 3915–3916.
- 723 Veselovska L, Smallwood SA, Saadeh H, Stewart KR, Krueger F, Maupetit-Méhouas S, Arnaud
724 P, Tomizawa S-I, Andrews S, Kelsey G. 2015. Deep sequencing and de novo assembly of
725 the mouse oocyte transcriptome define the contribution of transcription to the DNA
726 methylation landscape. *Genome Biol* **16**: 209.
- 727 Wang J, Xie G, Singh M, Ghanbarian AT, Raskó T, Szvetnik A, Cai H, Besser D, Prigione A,
728 Fuchs NV, et al. 2014. Primate-specific endogenous retrovirus-driven transcription
729 defines naive-like stem cells. *Nature* **516**: 405–409.
- 730 Whiddon JL, Langford AT, Wong C-J, Zhong JW, Tapscott SJ. 2017. Conservation and
731 innovation in the DUX4-family gene network. *Nat Genet* **49**: 935–940.
- 732 Yandım C, Karakülah G. 2019. Expression dynamics of repetitive DNA in early human
733 embryonic development. *BMC Genomics* **20**: 439.
- 734 Yang WR, Ardeljan D, Pacyna CN, Payer LM, Burns KH. 2019. SQuIRE reveals locus-specific
735 regulation of interspersed repeat expression. *Nucleic Acids Res* **47**: e27.
- 736 Zhang S, Wang Y, Jia L, Wen X, Du Z, Wang C, Hao Y, Yu D, Zhou L, Chen N, et al. 2019.
737 Profiling the long noncoding RNA interaction network in the regulatory elements of
738 target genes by chromatin in situ reverse transcription sequencing. *Genome Res* **29**:
739 1521–1532.
- 740 Zheng GXY, Terry JM, Belgrader P, Ryvkin P, Bent ZW, Wilson R, Ziraldo SB, Wheeler TD,
741 McDermott GP, Zhu J, et al. 2017. Massively parallel digital transcriptional profiling of
742 single cells. *Nat Commun* **8**: 14049.
- 743