

Single-cell full-length isoform sequencing reveals an abundance of 3' prime partial transcripts in maternally inherited RNAs

Chuan-Le Xiao (✉ xiaochuanle@126.com)

Sun Yat-sen University <https://orcid.org/0000-0002-4680-0682>

Chaoyang Wang

Guangzhou Laboratory

Zhuoxing Shi

Sun Yat-sen University

Qingpei Huang

Guangzhou Laboratory

Kunhua Hu

Proteomics Center, Zhongshan School of Medicine, Sun Yat-sen University,

Rong Liu

Guangzhou Laboratory

Dan Su

Guangzhou Laboratory

Zhuobin Lin

Sun Yat-sen University

Xiaoying Fan

Bioland Laboratory (GuangZhou Regenerative Medicine and Health Guangdong Laboratory)

<https://orcid.org/0000-0002-2599-8319>

Biological Sciences - Article

Keywords:

Posted Date: June 15th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-1752209/v1>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **Single-cell full-length isoform sequencing reveals an abundance of 3'**
2 **prime partial transcripts in maternally inherited RNAs**

3 Chaoyang Wang^{1,2,3,6,7}, Zhuoxing Shi^{1,7}, Qingpei Huang^{2,7}, Kunhua Hu^{4,7}, Rong Liu², Dan
4 Su^{2,3,6}, Zhuobin Lin¹, Xiaoying Fan^{2,3,5,6*}, Chuanle Xiao^{1,*}

5 ¹State Key Laboratory of Ophthalmology, Zhongshan Ophthalmic Center, Sun Yat-sen
6 University, Guangzhou, Guangdong Province, P.R. China

7 ²Guangzhou Laboratory, Guangzhou, Guangdong Province, P.R. China

8 ³The Bioland Laboratory (GuangZhou Regenerative Medicine and Health Guangdong
9 Laboratory), Guangzhou, Guangdong Province, P.R. China

10 ⁴Proteomics Center, Zhongshan School of Medicine, Sun Yat-sen University, Guangzhou,
11 Guangdong Province, P.R. China

12 ⁵The Fifth Affiliated Hospital of Guangzhou Medical University, Guangzhou, Guangdong
13 Province, P.R. China

14 ⁶The Guangzhou Institutes of Biomedicine and Health, Chinese Academy of Sciences,
15 Guangzhou, Guangdong Province, P.R. China

16 ⁷These authors contributed equally to this work.

17 *Corresponding authors: Xiaoying Fan (Fan_xiaoying@gzlab.ac.cn), Chuanle Xiao
18 (xiaochuanle@126.com)

19 **Abstract**

20 Isoform expression in preimplantation embryos has been extensively investigated, and
21 many novel isoforms have been identified. However, the regulation patterns of different
22 types of transcripts along the developing stages remains unexplored. We quantified the
23 expression of full-length isoforms in over one hundred single blastomeres from the mouse
24 oocyte to blastocyst stage and found that the 3' prime partial transcripts that lack stop
25 codons were highly accumulated in oocytes and zygotes. A typical 3' prime partial isoform,
26 *Ncl-S-350*, was demonstrated to be not a transcription by-product but to function in the
27 ZGA process. SRSF4 was identified to be responsible for generating these 3' prime partial
28 transcripts in mouse embryonic stem cells (mESCs), and both *Ncl-S-350* and *Srsf4*
29 overexpression could convert mESCs to a 2-cell (2C)-like state. Our work reveals the

30 important role of isoform switch regulation in early embryonic development and lays the
31 foundation for an alternative way of acquiring totipotent cells.

32 **Introduction**

33 One gene can be transcribed into multiple kinds of transcripts, namely, isoforms, which
34 are then translated into different proteins. Isoform compositions vary between different cell
35 types and states; therefore, isoform switch is important in determining cell identity^{1,2}. In
36 recent years, several third-generation sequencings (TGS) based single-cell RNA-
37 sequencing (scRNA-seq) methods have been developed for direct isoform sequencing<sup>1,3-
38 7</sup>. Among these methods, SCAN-seq has been shown to have high gene detection
39 sensitivity and to enable the discovery of large number of novel transcripts in rare samples,
40 such as preimplantation embryos⁴. Nevertheless, SCAN-seq failed to quantify the absolute
41 abundance of the detected genes and isoforms since it is difficult to specify unique
42 molecular identifiers (UMIs) under the high error rate of Nanopore sequencing^{8,9}. Thereby,
43 HIT-sclISOseq and MAS-ISO-seq have been used to quantify isoform abundance in single
44 cells with improved data throughput by using the PacBio HiFi sequencing platform^{5,7}.

45 Many studies have focused on the molecular regulation of preimplantation embryo
46 development, as it is the basis for reproduction. In particular, the maternal to zygote
47 transition is the foundation of the whole-body development plan¹⁰⁻¹³. Although hundreds of
48 genes have been reported in zygotic genome activation (ZGA), the functional regulators
49 remain largely unclear, not to mention whether the isoform switch participate in the
50 process¹⁴⁻¹⁶.

51 In this study, we modified the HIT-sclISOseq method for low throughput of cells and
52 sequenced the isoforms in single blastomeres of mouse preimplantation embryos. We
53 quantified the isoforms of each gene in each single cell, and different isoform types showed
54 varied proportions among different stages. Specifically, we observed large number of 3'
55 prime partial transcripts (which lack stop codons and generate proteins lacking C-termini)
56 in oocytes and zygotes that quickly degraded from early 2-cell stage (2C) stage. We
57 validated the 3' prime partial nucleolin (*Ncl*) gene transcripts, and surprisingly, the short *Ncl*
58 transcripts may promote ZGA by significantly inducing the expression of *Dux* and other 2C

59 genes. Moreover, we found that the splicing factor *Srsf4* plays an important role in the
60 expression of 3' prime partial transcripts, which could also induce 2C gene expression in
61 mESCs. Therefore, the 3' prime partial transcripts, which are usually regarded as
62 detrimental by-products, may be essential regulators during preimplantation embryo
63 development, and isoform switch may be a specific regulatory mechanism for the ZGA
64 process.

65

66 **Results**

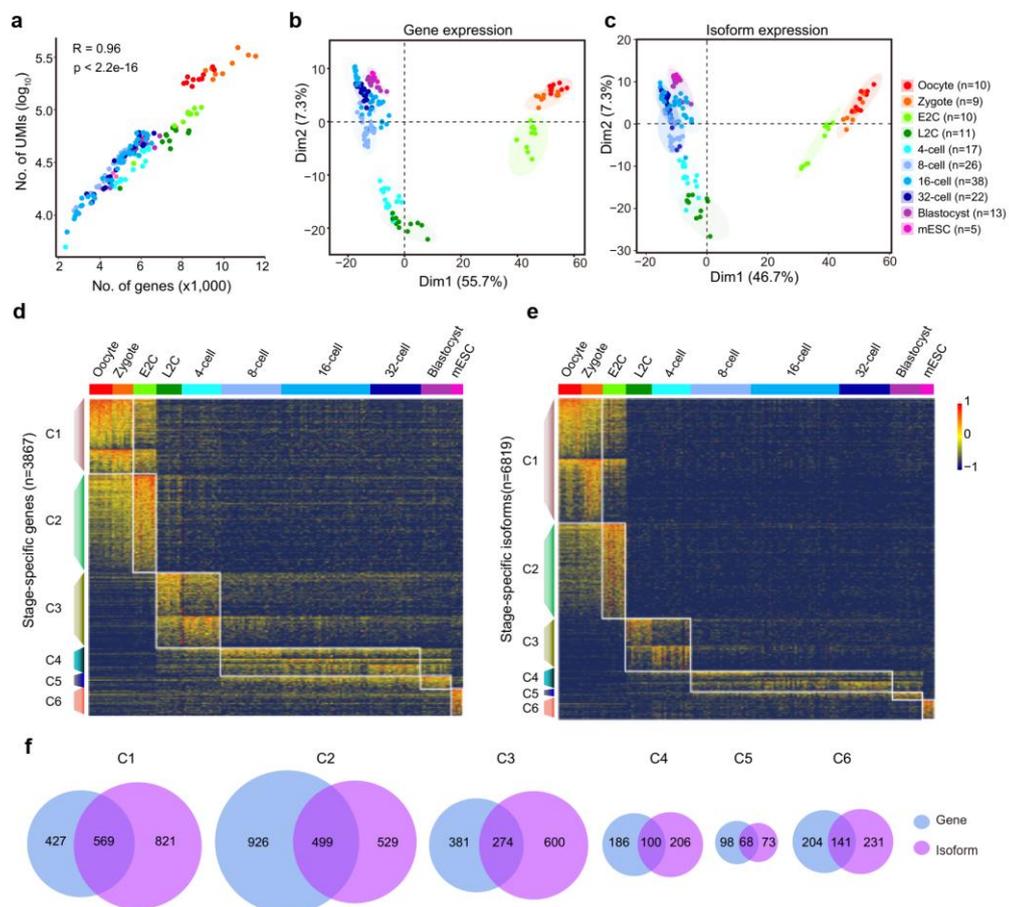
67 **Characterization of the mouse preimplantation embryos by gene and** 68 **isoform expression**

69 To detect gene isoforms in each blastomere of the mouse metaphase II oocytes and
70 preimplantation embryos, we amplified RNAs in each single cell with a 10x gel bead using
71 the Smart-seq2 protocol^{5,17}. Then, the amplified cDNAs of different cells were pooled for
72 further ligation and PacBio library construction following the HIT-sclISOseq method⁵.
73 Meanwhile, the corresponding barcode sequence of each cell was pre-determined through
74 Sanger sequencing of the cDNAs (Extended Data Fig. 1a). Expression data for a total of
75 161 single-cell isoforms were obtained from three batches, comprising the mouse oocyte,
76 zygote, early 2-cell (E2C), late 2-cell (L2C), 4-cell, 8-cell, 16-cell, 32-cell, and blastocyst
77 stages and mESCs (Extended Data Table 1). Each sequencing batch generated about 5
78 million circular consensus sequencing (CCS) reads, with the average length around 4kb,
79 indicating ligation of 2-3 cDNA molecules in most cases. After data splitting and mapping,
80 about 90% of the full-length isoforms could be correctly assigned to cells (Extended Data
81 Table 1). In this way, the samples from each stage could be sequenced with a relatively
82 sufficient depth (Extended Data Fig. 1b). To evaluate the accuracy of measuring absolute
83 numbers of isoforms with our procedure, we also amplified ERCC and SIRV spike-ins¹⁸. At
84 the gene level, we observed high correlation values between the added molecules and the
85 detected UMI counts (Extended Data Fig. 1c). At the isoform level, different isoforms of the
86 same gene could be correctly specified without any false match (Extended Data Fig. 1d).
87 These results indicate that our workflow accurately measures the abundances of the
88 transcripts in each single cell.

89 As expected, the mouse oocyte and zygote contained many more RNA molecules than
90 later stage blastomeres as a result of subsequent maternally inherited RNA degradation.
91 The number of transcript molecules was well correlated ($R=0.96$) with the number of genes
92 expressed in the cells (Fig. 1a). We used gene expression data and isoform expression
93 data to perform principal component analysis (PCA) of all the cells. Both types of data
94 generated similar PCA results and could be used to clearly separate blastomeres of
95 different stages (Fig. 1b, c). In both sets of PCA results, the oocyte and zygote showed
96 similar expression patterns, the L2C and 4-cell stages could be grouped together, the 8-
97 cell, 16-cell and 32-cell stages were close, and the blastocyst cells were analogous to the
98 mESCs. We further extracted stage-specific genes and transcripts with the same criterion
99 and obtained 3867 genes and 6819 isoforms, respectively. Each list could be divided into
100 six corresponding clusters based on their expression patterns across all embryonic stages
101 (Fig. 1d, e). To be specific, Cluster 1 (C1) transcripts were highly abundant in oocytes and
102 zygotes, subsequently degraded from E2C stage. Cluster 2 (C2) included transcripts only
103 upregulated in E2C stage. Cluster 3 (C3), cluster 4 (C4) and cluster 5 (C5) transcripts were
104 highly expressed in the L2C to 4-cell stages, 8-cell to 32-cell stages and blastocyst stages,
105 respectively. The mESC-specific transcripts were in cluster 6 (C6). In addition, we
106 compared the genes and isoforms in each pair of matched clusters. More number of RNAs
107 could be identified in each cluster at the isoform level, and most of the isoforms showed
108 consistency with the genes (Fig. 1f, Extended Data Table 2). Thus, the single-cell isoform
109 expression data could be directly used to illustrate cellular heterogeneity and distinguish
110 different types of cells as single-cell gene expression does.

111 To investigate the relationship between gene and isoform expression during mouse
112 preimplantation development, we divided genes into six groups based on the number of
113 isoform types they expressed (Extended Data Fig. 2a). Although majority of genes
114 expressed only one type of isoform across different stages, more genes expressed multiple
115 types of isoforms in the earlier stages. About 60% of genes in mouse oocytes and zygotes
116 expressed more than one type of isoforms, and nearly 20% of genes were detected with
117 more than five types of isoforms. In comparison, approximately 70% of genes in mESCs
118 expressed only one type of isoform, and few genes expressed over five types of isoforms

119 (Extended Data Fig. 2a). The same isoform expression characteristics could also be
 120 observed with SCAN-seq data (Extended Data Fig. 2b), suggesting rich isoform diversity
 121 in early mouse embryos. In addition, the genes expressing more types of isoforms were
 122 detected with higher expression levels in both our data and SCAN-seq data (Extended
 123 Data Fig. 2c, d). To assess the isoform dominant level in each gene expressing multiple
 124 types of isoforms, we calculated the ratio of the UMI number of the major isoform to the
 125 total UMI number of the corresponding gene. The major isoform ratios increased from early
 126 to late embryonic stages, especially after the ZGA process (Extended Data Fig. 2e). In
 127 comparison, the major isoforms accounted for 90% of most genes in mESCs, indicating a
 128 dominant isoform expression pattern and less isoform diversity in these cells.



129
 130 **Fig. 1 | Gene and isoform expression in the mouse preimplantation embryos**
 131 **a**, Number of genes and UMIs in each cell at different stages. **b-c**, PCA plot of all the
 132 blastomeres and mESCs based on gene expression (**b**) and isoform expression (**c**). **d-e**,
 133 Heatmap of stage-specific genes (**d**) and isoforms (**e**). **f**, Venn plot of pairwise groups of
 134 stage-specific genes and isoform-corresponding genes.

135

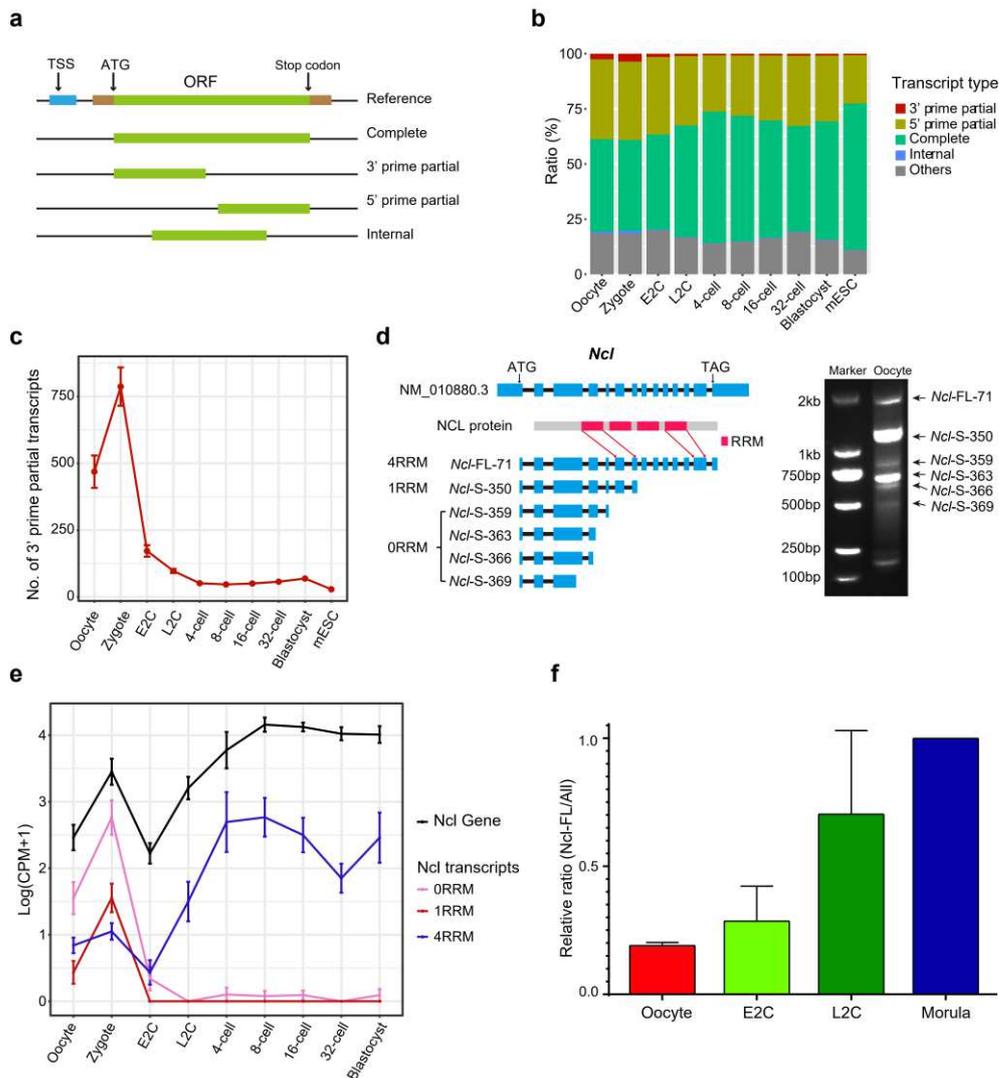
136 **Large abundance of 3' prime partial transcripts was observed in mouse**
137 **oocytes and zygotes**

138 According to the integrity of putative corresponding open reading frames (ORFs), the
139 transcripts were divided into 5 types: complete, representing isoforms coding the complete
140 ORFs of the reference genes; 3' prime partial transcripts and 5' prime partial transcripts,
141 missing the stop codon section and start codon section, respectively; internal, predicted
142 with proteins lacking both ends; and others, including isoforms located outside the
143 reference ORFs¹⁹ (Fig. 2a). As expected, the complete transcripts had the longest lengths,
144 and the internal transcripts were the shortest (Extended Data Fig. 3a). Nevertheless, the
145 predicted protein length was similar for the three incomplete transcript types (Extended
146 Data Fig. 3b). Intriguingly, we found that the 3' prime partial transcripts were highly
147 expressed in oocytes and zygotes and then dramatically decreased from the E2C stage
148 (Fig. 2b, c). The same expression pattern was also observed in the SCAN-seq data
149 (Extended Data Fig. 3c, d). Then, we performed Gene Ontology (GO) analysis on genes
150 detected with 3' prime partial transcripts (Extended Data Table 3). These genes were
151 enriched in the pathways of RNA processing, cell cycle checkpoint, ribonucleoprotein
152 complex biogenesis, DNA metabolic process, chromatin organization, etc. (Extended Data
153 Fig. 3e), all of which have been reported to play essential roles in mouse and human
154 preimplantation embryo development^{10,13,20-23}. In addition, we found that 82% of these 3'
155 prime partial transcripts ended within the exon, while only 44% of the 5' prime partial
156 transcripts ended within the exon (Extended Data Fig. 3f).

157 We further chose some candidates to validate the enrichment of 3' prime partial
158 transcripts in the maternally inherited RNAs. For example, the *Ncl* gene was detected with
159 6 isoform types that could be assigned to three categories based on how many RNA
160 recognition motifs (RRM) they contained (Fig. 2d, left). All 6 isoform types were observed
161 in the reverse transcription and PCR (RT-PCR) products of mouse oocytes (Fig. 2d, right).
162 Meanwhile, the short isoforms appeared to be much more abundant than the full-length
163 isoform (*Ncl*-FL-71) in mouse oocytes. Then, we calculated the abundance of each
164 category of *Ncl* isoform along developing stages. *Ncl* was first downregulated in the E2C

165 stage and then increased from the L2C stage at the gene level (Fig. 2e). According to
166 changes in the abundances of different isoforms, we found that the two categories of short
167 *Ncl* isoforms were highly expressed in oocytes and zygotes and then almost disappeared.
168 In comparison, the full-length *Ncl* isoform only showed low expression in the maternal
169 RNAs and was largely upregulated after ZGA (Fig. 2e). This finding highlights the isoform
170 switch of the *Ncl* gene during the ZGA process; and such regulation approaches are
171 masked in gene-level analysis. To demonstrate this, we performed reverse transcription
172 and real-time quantitative PCR (RT-qPCR) using primers targeting all the *Ncl* isoform
173 types and only the full-length type, respectively, and calculated the relative percentages of
174 the full-length isoform in different stages of mouse preimplantation embryos. As expected,
175 less than 20% of the *Ncl* transcripts were full length in oocytes when we set the full-length
176 relative ratio as 100% at the morula stage (Fig. 2f). The detected 3' prime partial transcripts
177 of some other genes related to RNA processing (*Sf3b2*, *Srpk1*) and protein translation and
178 transporting (*Dnajc3*, *Hsp90aa1*) were also revealed by gel analysis of mouse oocyte RT-
179 PCR products (Extended data Fig. 4a).

180 The 3' prime partial transcripts lack the stop codon by definition. We further sequenced
181 the amplified 3' prime partial transcripts of *Ncl* and genes in Extended data Fig. 4a by
182 Sanger sequencing for validation. Each predicted isoform was confirmed to have a poly(A)
183 tail without a stop codon (Extended data Fig. 4b). Usually, transcripts without stop codons
184 are quickly degraded by a nonstop decay (NSD) mechanism in eukaryotic cells^{24,25}. We
185 analysed the expression of *Pelota*, *Hbs1* and *Abce1*, the three indispensable regulators
186 triggering the NSD process (Extended data Fig. 4c). Both *Pelota* and *Abce1* were
187 expressed at low levels in oocytes and zygotes, indicating the NSD pathway shall be in
188 quite low activity before ZGA process. In this way, the 3' prime partial transcripts can be
189 largely accumulated in the maternal content.



190

191 **Fig. 2| Expression patterns of different types of transcripts during mouse**
 192 **preimplantation embryonic development**

193 **a**, Schematic diagram of different transcript types. **b**, Ratios of each type of transcripts at

194 different stages. **c**, Number of 3' prime partial transcripts detected at each stage. **d**,

195 Schematic diagram of *Ncl* isoforms (left) and gel picture showing these isoforms by RT-

196 PCR of *Ncl* in mouse oocytes (right). **e**, Expression levels of each category of *Ncl* isoforms

197 and the *Ncl* gene at different stages. **f**, The relative ratios of *Ncl* full-length isoforms at

198 different stages detected by RT-qPCR. The relative ratio of NCL-FL/All in morulae was set

199 as 1.0.

200

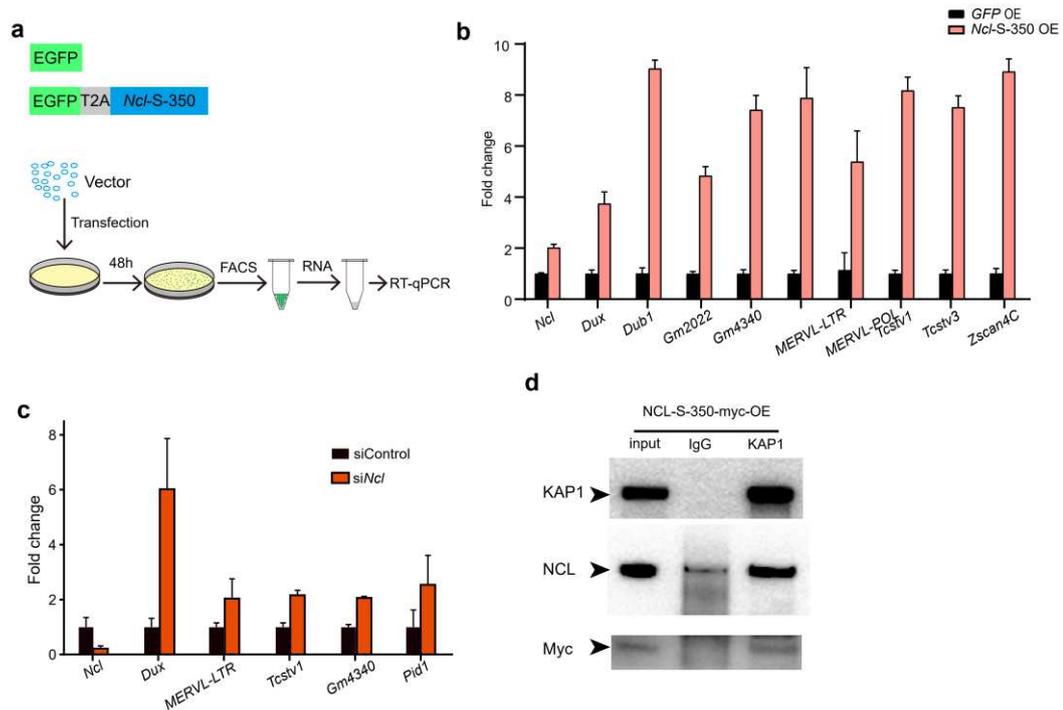
201

202

203 **The 3' prime partial *Ncl* transcript (*Ncl*-S-350) induces the expression of**
204 **2C genes**

205 NCL encoded by the *Ncl* gene has been demonstrated to participate in a wide range
206 of cellular programs, including ribosome biogenesis, chromatin organization and stability,
207 and DNA and RNA metabolism²⁶. This gene is highly expressed in mESCs, and it has been
208 proven that NCL forms a complex with KAP1/TRIM28 and *LINE1* RNA to repress the
209 expression of *Dux*, the master activator of the 2C transcription program²⁷⁻³¹. Knocking
210 down *Ncl* in mESCs by siRNA had been reported to lead to activation of the 2C transcription
211 program characterized by the expression of a series of 2C genes, such as *Dux*, *MERVL*,
212 *Zscan4C* and *Tsctv1*^{27,32}. There were very few 3' prime partial transcripts in mESCs,
213 including *Ncl*-S-350, which was highly detected in the maternal RNA content (Fig. 2b, 2c).
214 To determine whether these 3' prime partial transcripts are functional, we overexpressed
215 *Ncl*-S-350 in mESCs (Fig. 3a). Surprisingly, *Ncl*-S-350 overexpressing cells also
216 significantly upregulated 2C genes, including *Dux*, *MERVL*, and *Zscan4c* (Fig. 3b). This
217 finding was consistent with the results when *Ncl* was knocked down in mESCs (Fig. 3c)^{27,32}.

218 The direct interaction between *LINE1* RNA and NCL has been confirmed *in vitro* and
219 *in vivo*^{33,34}. Moreover, the full integrity of the first and second RRM domains is required to
220 bind RNA³⁵. As NCL inhibits *Dux* expression under the guidance of *LINE1* RNA²⁷, we
221 deduced that the NCL-S-350 lacks the ability to bind *LINE1* RNA compared to the full-
222 length isoform (NCL-FL). Therefore, the NCL-S-350 protein may function as an antagonist
223 of the NCL-FL-71 protein to form a complex with KAP1/TRIM28 without *LINE1* RNA. To
224 verify this assumption, we did co-immunoprecipitation (Co-IP) assay and the result
225 confirmed that NCL-S-350 protein can bind KAP1 as the NCL-FL proteins do (Fig. 3d). The
226 results imply that the 3' prime partial transcripts shall not be transcription by-products and
227 may play an important role in preimplantation embryo development.



228

229 **Fig. 3| Overexpression of the 3' prime partial transcript *Ncl-S-350* induced 2C gene**
 230 **expression in mESCs**

231 **a**, Overview of the overexpression workflow in mESCs. **b-c**, Bloxplot showing increased
 232 expression levels of *Dux* and other 2C genes by RT-qPCR after overexpressing *Ncl-S-350*
 233 (b) or knocking down *Ncl* using siRNA (c) in mESCs (n=3). **d**, Co-IP assay showing that
 234 NCL-S-350 can combine with KAP1 as NCL-FL-71 does.

235

236 ***Srsf4* promotes the accumulation of 3' prime partial transcripts**

237 The 3' prime partial transcripts in the maternally inherited RNAs mostly ended within
 238 the exon region (Extended Data Fig. 3f) and were likely not regulated by spliceosomes,
 239 which function in introns. Exon splicing requires serine/arginine-rich splicing factors
 240 (SRs)^{36,37}. Therefore, we speculated that the generation of 3' prime partial transcripts rely
 241 on SRs. The SR family consists of 12 members in mammals³⁸. We checked the expression
 242 patterns of each SR gene in the mouse preimplantation embryo samples and noticed that
 243 *Srsf4* was highly expressed in oocytes and zygotes and then downregulated at the later
 244 stages (Fig. 4a, Extended data Fig. 5). The same expression patterns of the SR genes
 245 were also observed in previous NGS data (Fig. 4a, Extended data Fig. 6)³⁹. We further
 246 analysed the motifs using the 100 bp sequence flanking the end point of all the 3' prime

247 partial transcripts (Fig. 4b). Interestingly, AGAAAA is consistent with the conserved binding
248 motif of SRSF4, which has been reported to work on *Ncl* mRNA⁴⁰. We checked the
249 distribution of the AGAAAA motif on all the exons of *Ncl*, and its locations matched perfectly
250 to the end sites of those 3' prime partial transcripts of *Ncl* (Extended data Fig. 7a).

251 In order to check the function of SRSF4 in producing the 3' prime partial transcripts,
252 we overexpressed *Srsf4* in mESCs. Subsequently, we amplified the *Ncl* transcripts and
253 examined the isoform types on the gel. Intriguingly, compared to the phenomenon that *Ncl*-
254 FL-71 was dominant in normal mESCs, *Ncl*-S-350 became the most abundant isoform type
255 in *Srsf4*-overexpressing mESCs (Fig. 4c). The other 3' prime partial transcripts of *Ncl* were
256 also increased. Moreover, the expression pattern of *Ncl* isoforms in *Srsf4*-overexpressing
257 mESCs was almost identical to that in mouse oocytes (Fig. 4c). These results suggested
258 that *Srsf4* regulates the generation of the 3' prime partial transcripts of *Ncl*. We further
259 detected the expression of 2C genes in *Srsf4*-overexpressing mESCs by RT-qPCR. As
260 expected, *Dux*, *MERV1*, *Zscan4c* and other 2C genes were all significantly upregulated,
261 which was consistent with the results in mESCs overexpressing *Ncl*-S-350 (Fig. 4d).

262 To further investigate the targets of SRSF4 and whether the cells acquired totipotency
263 when overexpressed with the 3' prime partial transcripts, we performed single-cell full-
264 length isoform sequencing on *Srsf4*-overexpressing mESCs (*Srsf4* OE), *Ncl*-S-350-
265 overexpressing mESCs (*Ncl*-S-350 OE) and GFP-overexpressing control mESCs (*GFP*
266 OE). The *Srsf4* OE and *Ncl*-S-350 OE cells were almost the same with each other but
267 different from the *GFP* OE control mESCs (Extended Data Fig. 7b, c). Overexpression of
268 both *Srsf4* and *Ncl*-S-350 led to elevated expression of totipotent genes in mESCs (Fig.
269 4e, Extended Data Table 4). In addition, most minor and major ZGA genes were also
270 upregulated (Extended Data Fig. 7d, Extended Data Table 4). To check the states of the
271 two groups of cells, we clustered them with the embryo samples and previously reported
272 totipotent-like stem cells (TLSCs)¹⁶. Interestingly, both *Srsf4* OE and *Ncl*-S-350 OE cells
273 were close to the TLSCs, all showing similarity to the 2C embryo (Fig. 4f).

274 Notably, the proportion of 3' prime partial transcripts increased from 7% in the *GFP*
275 OE group to about 19% in the *Srsf4* OE and *Ncl*-S-350 OE groups (Fig. 4g). Genes
276 generating 3' prime partial transcripts under SRSF4 were related to RNA processing,

277 mitotic cell cycle, chromatin organization, etc., which was similar to those enriched in the
278 mouse preimplantation embryos (Extended Data Table 5 and Fig. 7e). Moreover, the
279 isoform diversity increased in the *Srsf4* OE and *Ncl-S-350* OE groups, as 43.7% and 46.2%
280 of genes in the two groups expressed more than 5 isoform types, respectively, which was
281 much higher than the ratio (33.9%) in the *GFP* OE group (Fig. 4h). This finding also
282 resembled the large isoform diversity in oocytes and early embryos (Extended Data Fig.
283 2a, b). Therefore, *Srsf4* may play an essential role in the generation of 3' prime partial
284 transcripts, and both *Srsf4* and *Ncl-S-350* may participate in cellular identity determination
285 in mouse oocytes and early embryos.

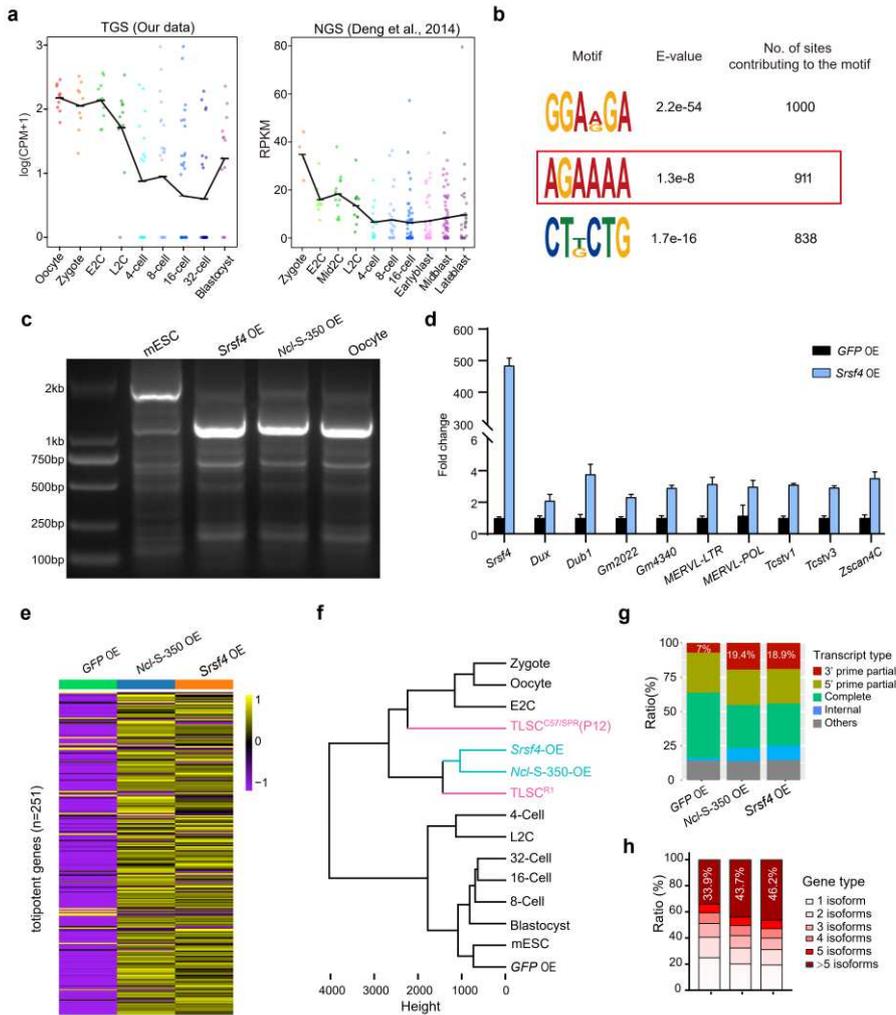
286

287

288

289

290



291

292

Fig. 4| *Srsf4* regulates the production of the 3' prime partial transcripts of *Ncl* and other genes

293

294

a, Expression levels of *Srsf4* at different embryonic stages in our data and in the NGS data

295

by Deng *et al.*, 2014³⁹. **b**, The significant motifs in the 100 bp sequences flanking the end

296

sites of the 3' prime partial transcripts. The red box highlights the one consistent with the

297

SRSF4 binding motif. **c**, Gel analysis of RT-PCR products of normal mESCs, *Srsf4* OE

298

mESCs, *Ncl-S-350* OE mESCs and mouse MII oocytes. **d**, Bloxplot showing increased

299

expression levels of *Ncl*, *Dux* and other 2C genes by RT-qPCR after overexpressing *Srsf4*

300

in mESCs. **e**, Heatmap showing generally elevated expression of totipotent genes in *Srsf4*

301

OE and *Ncl-S-350* OE cells compared with control *GFP* OE cells. **f**, Dendrogram displaying

302

the *Srsf4* OE and *Ncl-S-350* OE cells were similar to previously reported TLSCs and 2-cell

303

embryos in gene expression. **g**, Ratio of each type of transcript in the three groups of cells.

304

h, The ratios of genes containing different numbers of isoform types in the three groups of

305

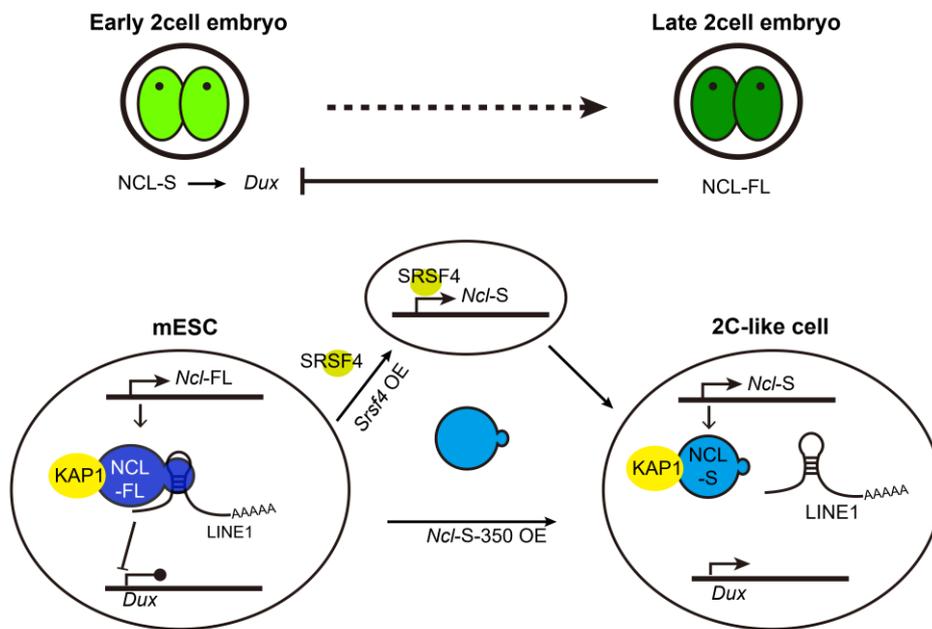
cells.

306

307 **Discussion**

308 An abundance of novel transcripts and splicing events in preimplantation embryos
309 have been annotated by TGS-based full-length isoform sequencing at either the bulk or
310 single-cell level^{4,14,41}. However, it remains unclear how different types of isoforms regulate
311 the developmental process. In the present study, by dividing the transcripts into subtypes
312 according to their coding characteristics, we found large number of 3' prime partial
313 transcripts, which lack stop codons, in mouse MII oocytes and zygotes (Fig. 2b). Defects
314 in the NSD pathway in the maternal contents enable these transcripts to be maintained.
315 This type of transcript has been largely studied in cancers and is considered an oncogenic
316 factor⁴². While in early embryos, these transcripts might be important for promoting the
317 ZGA process, as their host genes were highly enriched in RNA processing, cell cycle
318 checkpoint, ribonucleoprotein complex biogenesis, DNA metabolic process and chromatin
319 organization pathways (Extended Data Fig. 3e).

320 Taking *Ncl* as an example, we proved that the 3' prime partial transcripts of *Ncl*
321 promoted the expression of *Dux* and other 2C genes, possibly by competing with the full-
322 length isoform to form the NCL/Kap1 complex²⁷. However, they lost the capacity to bind
323 *LINE1* RNA to locate the repressive complex to the *Dux* promoter (Fig. 5). After ZGA, these
324 transcripts are degraded rapidly, and the full-length isoform is upregulated to suppress *Dux*
325 again. Such isoform switch regulation is blinded through single-cell gene expression
326 analysis, emphasizing the importance of performing isoform sequencing.



327

328 **Fig. 5| The hypothesis of the regulatory mechanism of *Dux* by *Ncl-S-350***

329 In early 2-cell stage, the numerous existing 3' prime partial types of NCL could enhance
 330 the activation of *Dux*, while in late 2-cell stage, these transcripts degrade rapidly, and the
 331 increased full-length NCL suppresses *Dux* expression. In mESCs, *Dux* is suppressed by
 332 full-length NCL²⁷. Overexpression of 3' prime partial type *Ncl-S-350* directly or by *Srsf4*
 333 overexpression results in competition with the full-length NCL form complex with KAP1 but
 334 lacks LINE1 as the guidance to suppress *Dux*.

335

336 A recent study demonstrated that spliceosomal repression induces totipotency in
 337 mESCs⁴³, indicating that alternative splicing is essential for modulating cell fate. However,
 338 the PlaB they used to inhibit SF3B, which forms spliceosomes to splice the intron region<sup>43-
 339 45</sup>, and the 3' prime partial transcripts in the maternal RNAs primarily ended within the exon
 340 region. We proved that SRSF4 can produce the 3' prime partial transcripts of both *Ncl* and
 341 many other genes, inducing a 2C-like state in mESCs (Fig. 4). We also noticed that the
 342 overexpression of *Ncl-S-350* enhanced the expression of *Srsf4* (data not shown). This
 343 could explain why the *Srsf4* OE and *Ncl-S-350* OE cells were almost the same as each
 344 other and similar to TLSCs from a previous study. This may promote a new way for the *in*
 345 *vitro* culture of totipotent cells, and such an assumption needs further experimental
 346 exploration and verification.

347 We extrapolated that the large number of 3' prime partial transcripts in maternally
348 inherited RNAs are not transcription by-products. These transcripts may play functional
349 roles in the preimplantation embryo development. In current work, we investigated the
350 function of *Ncl*, while there are many other genes worthy of further study. Our work provides
351 new insights into the mechanisms regulating early embryonic development.
352

References

- 354 1 Lebrigand, K., Magnone, V., Barbry, P. & Waldmann, R. High throughput error corrected
355 Nanopore single cell transcriptome sequencing. *Nat Commun* **11**, 4025,
356 doi:10.1038/s41467-020-17800-6 (2020).
- 357 2 Joglekar, A. *et al.* A spatially resolved brain region- and cell type-specific isoform atlas
358 of the postnatal mouse brain. *Nat Commun* **12**, 463, doi:10.1038/s41467-020-20343-5
359 (2021).
- 360 3 Lebrigand, K. *et al.* The spatial landscape of gene expression isoforms in tissue
361 sections. *bioRxiv*, 2020.2008.2024.252296, doi:10.1101/2020.08.24.252296 (2022).
- 362 4 Fan, X. *et al.* Single-cell RNA-seq analysis of mouse preimplantation embryos by third-
363 generation sequencing. *PLoS Biol* **18**, e3001017, doi:10.1371/journal.pbio.3001017
364 (2020).
- 365 5 Zheng, Y.-F. *et al.* HIT-sclISOseq: High-throughput and High-accuracy Single-cell Full-
366 length Isoform Sequencing for Corneal Epithelium. *bioRxiv*, 2020.2007.2027.222349,
367 doi:10.1101/2020.07.27.222349 (2020).
- 368 6 Oguchi, Y., Ozaki, Y., Abdelmoez, M. N. & Shintaku, H. NanoSINC-seq dissects the
369 isoform diversity in subcellular compartments of single cells. *Science Advances* **7**,
370 eabe0317, doi:doi:10.1126/sciadv.abe0317 (2021).
- 371 7 Al'Khafaji, A. M. *et al.* High-throughput RNA isoform sequencing using programmable
372 cDNA concatenation. *bioRxiv*, 2021.2010.2001.462818,
373 doi:10.1101/2021.10.01.462818 (2021).
- 374 8 Liu, Z. *et al.* Towards accurate and reliable resolution of structural variants for clinical
375 diagnosis. *Genome Biol* **23**, 68, doi:10.1186/s13059-022-02636-8 (2022).
- 376 9 Lu, H., Giordano, F. & Ning, Z. Oxford Nanopore MinION Sequencing and Genome
377 Assembly. *Genomics Proteomics Bioinformatics* **14**, 265-279,
378 doi:10.1016/j.gpb.2016.05.004 (2016).
- 379 10 Jukam, D., Shariati, S. A. M. & Skotheim, J. M. Zygotic Genome Activation in
380 Vertebrates. *Dev Cell* **42**, 316-332, doi:10.1016/j.devcel.2017.07.026 (2017).
- 381 11 Tadros, W. & Lipshitz, H. D. The maternal-to-zygotic transition: a play in two acts.
382 *Development* **136**, 3033-3042, doi:10.1242/dev.033183 (2009).
- 383 12 Vastenhouw, N. L., Cao, W. X. & Lipshitz, H. D. The maternal-to-zygotic transition
384 revisited. *Development* **146**, doi:10.1242/dev.161471 (2019).
- 385 13 Zhao, L. W. *et al.* Nuclear poly(A) binding protein 1 (PABPN1) mediates zygotic genome
386 activation-dependent maternal mRNA clearance during mouse early embryonic
387 development. *Nucleic Acids Res* **50**, 458-472, doi:10.1093/nar/gkab1213 (2022).
- 388 14 Qiao, Y. *et al.* High-resolution annotation of the mouse preimplantation embryo
389 transcriptome using long-read sequencing. *Nat Commun* **11**, 2653,
390 doi:10.1038/s41467-020-16444-w (2020).
- 391 15 Macfarlan, T. S. *et al.* Embryonic stem cell potency fluctuates with endogenous
392 retrovirus activity. *Nature* **487**, 57-63, doi:10.1038/nature11244 (2012).
- 393 16 Yang, M. *et al.* Chemical-induced chromatin remodeling reprograms mouse ESCs to
394 totipotent-like stem cells. *Cell Stem Cell* **29**, 400-418 e413,
395 doi:10.1016/j.stem.2022.01.010 (2022).
- 396 17 Picelli, S. *et al.* Smart-seq2 for sensitive full-length transcriptome profiling in single cells.

397 *Nat Methods* **10**, 1096-1098, doi:10.1038/nmeth.2639 (2013).

398 18 Hardwick, S. A. *et al.* Spliced synthetic genes as internal controls in RNA sequencing
399 experiments. *Nat Methods* **13**, 792-798, doi:10.1038/nmeth.3958 (2016).

400 19 Saha, S., Matthews, D. A. & Bessant, C. High throughput discovery of protein variants
401 using proteomics informed by transcriptomics. *Nucleic Acids Res* **46**, 4893-4902,
402 doi:10.1093/nar/gky295 (2018).

403 20 Xue, Z. *et al.* Genetic programs in human and mouse early embryos revealed by single-
404 cell RNA sequencing. *Nature* **500**, 593-597, doi:10.1038/nature12364 (2013).

405 21 Yang, Y. *et al.* RNA 5-Methylcytosine Facilitates the Maternal-to-Zygotic Transition by
406 Preventing Maternal mRNA Decay. *Mol Cell* **75**, 1188-1202 e1111,
407 doi:10.1016/j.molcel.2019.06.033 (2019).

408 22 Sha, Q. Q. *et al.* Characterization of zygotic genome activation-dependent maternal
409 mRNA clearance in mouse. *Nucleic Acids Res* **48**, 879-894, doi:10.1093/nar/gkz1111
410 (2020).

411 23 Fu, X., Zhang, C. & Zhang, Y. Epigenetic regulation of mouse preimplantation embryo
412 development. *Curr Opin Genet Dev* **64**, 13-20, doi:10.1016/j.gde.2020.05.015 (2020).

413 24 Powers, K. T., Szeto, J. A. & Schaffitzel, C. New insights into no-go, non-stop and
414 nonsense-mediated mRNA decay complexes. *Curr Opin Struct Biol* **65**, 110-118,
415 doi:10.1016/j.sbi.2020.06.011 (2020).

416 25 Morris, C., Cluet, D. & Ricci, E. P. Ribosome dynamics and mRNA turnover, a complex
417 relationship under constant cellular scrutiny. *Wiley Interdiscip Rev RNA* **12**, e1658,
418 doi:10.1002/wrna.1658 (2021).

419 26 Jia, W., Yao, Z., Zhao, J., Guan, Q. & Gao, L. New perspectives of physiological and
420 pathological functions of nucleolin (NCL). *Life Sciences* **186**, 1-10,
421 doi:10.1016/j.lfs.2017.07.025 (2017).

422 27 Percharde, M. *et al.* A LINE1-Nucleolin Partnership Regulates Early Development and
423 ESC Identity. *Cell* **174**, 391-405 e319, doi:10.1016/j.cell.2018.05.043 (2018).

424 28 Hendrickson, P. G. *et al.* Conserved roles of mouse DUX and human DUX4 in activating
425 cleavage-stage genes and MERVL/HERVL retrotransposons. *Nat Genet* **49**, 925-934,
426 doi:10.1038/ng.3844 (2017).

427 29 Whiddon, J. L., Langford, A. T., Wong, C. J., Zhong, J. W. & Tapscott, S. J. Conservation
428 and innovation in the DUX4-family gene network. *Nat Genet* **49**, 935-940,
429 doi:10.1038/ng.3846 (2017).

430 30 De Iaco, A. *et al.* DUX-family transcription factors regulate zygotic genome activation
431 in placental mammals. *Nat Genet* **49**, 941-945, doi:10.1038/ng.3858 (2017).

432 31 Ren, W. *et al.* DUX: One Transcription Factor Controls 2-Cell-like Fate. *Int J Mol Sci* **23**,
433 doi:10.3390/ijms23042067 (2022).

434 32 Sun, Z. *et al.* LIN28 coordinately promotes nucleolar/ribosomal functions and represses
435 the 2C-like transcriptional program in pluripotent stem cells. *Protein Cell*,
436 doi:10.1007/s13238-021-00864-5 (2021).

437 33 Peddigari, S., Li, P. W., Rabe, J. L. & Martin, S. L. hnRNPL and nucleolin bind LINE-1
438 RNA and function as host factors to modulate retrotransposition. *Nucleic Acids Res* **41**,
439 575-585, doi:10.1093/nar/gks1075 (2013).

440 34 Moldovan, J. B. & Moran, J. V. The Zinc-Finger Antiviral Protein ZAP Inhibits LINE and

441 Alu Retrotransposition. *PLoS Genet* **11**, e1005121, doi:10.1371/journal.pgen.1005121
442 (2015).

443 35 Serin, G. *et al.* Two RNA-binding domains determine the RNA-binding specificity of
444 nucleolin. *J Biol Chem* **272**, 13109-13116, doi:10.1074/jbc.272.20.13109 (1997).

445 36 Zheng, X. *et al.* Serine/arginine-rich splicing factors: the bridge linking alternative
446 splicing and cancer. *Int J Biol Sci* **16**, 2442-2453, doi:10.7150/ijbs.46751 (2020).

447 37 Jang, H. N. *et al.* Binding of SRSF4 to a novel enhancer modulates splicing of exon 6
448 of Fas pre-mRNA. *Biochem Biophys Res Commun* **506**, 703-708,
449 doi:10.1016/j.bbrc.2018.10.123 (2018).

450 38 Long, J. C. & Caceres, J. F. The SR protein family of splicing factors: master regulators
451 of gene expression. *Biochem J* **417**, 15-27, doi:10.1042/BJ20081501 (2009).

452 39 Deng, Q., Ramskold, D., Reinius, B. & Sandberg, R. Single-cell RNA-seq reveals
453 dynamic, random monoallelic gene expression in mammalian cells. *Science* **343**, 193-
454 196, doi:10.1126/science.1245316 (2014).

455 40 Anko, M. L. *et al.* The RNA-binding landscapes of two SR proteins reveal unique
456 functions and binding to diverse RNA classes. *Genome Biol* **13**, R17, doi:10.1186/gb-
457 2012-13-3-r17 (2012).

458 41 Berrens, R. V. *et al.* Locus-specific expression of transposable elements in single cells
459 with CELLO-seq. *Nat Biotechnol*, doi:10.1038/s41587-021-01093-1 (2021).

460 42 Mohanan, N. K., Shaji, F., Koshre, G. R. & Laishram, R. S. Alternative polyadenylation:
461 An enigma of transcript length variation in health and disease. *WIREs RNA* **13**,
462 doi:10.1002/wrna.1692 (2021).

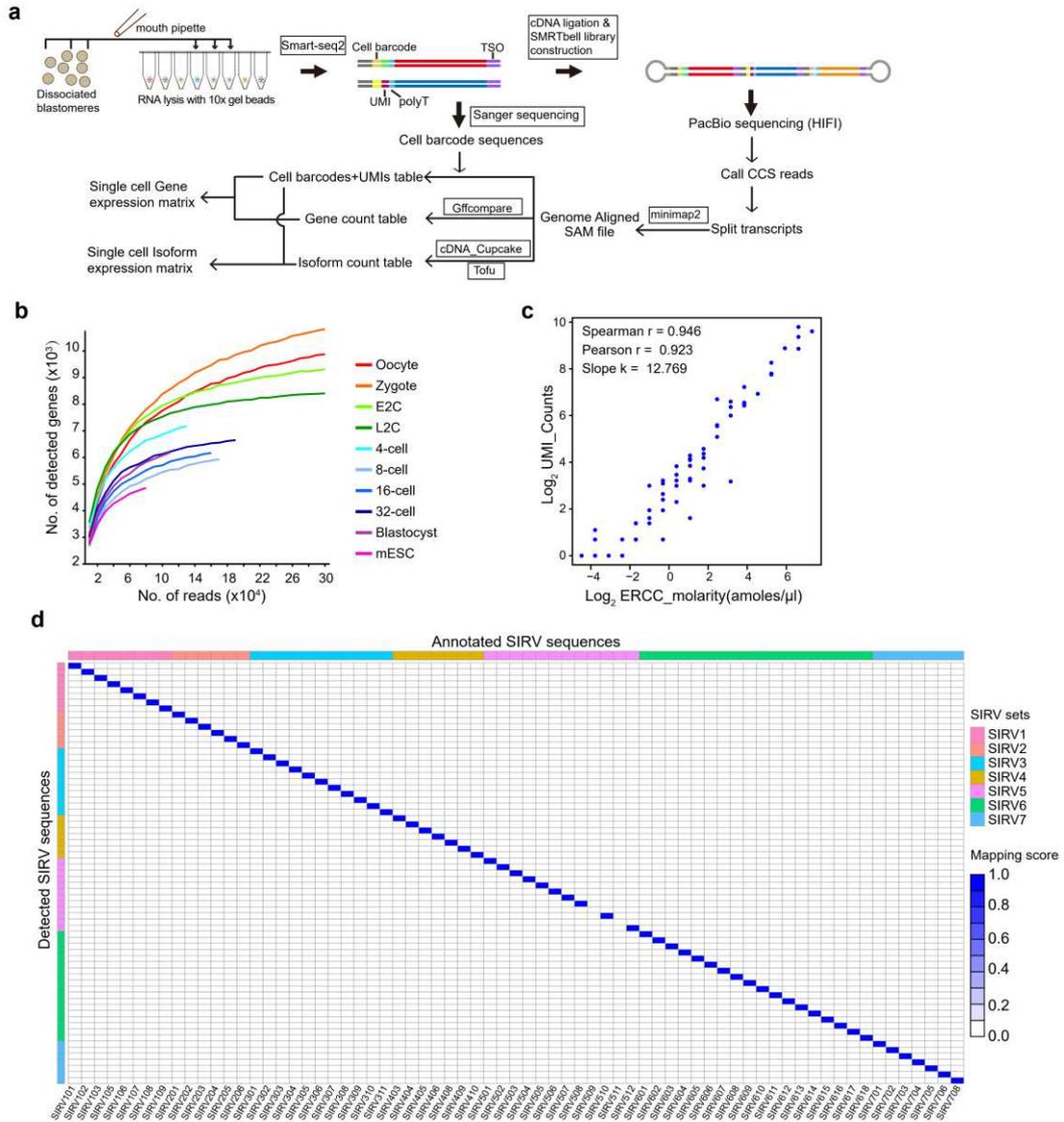
463 43 Shen, H. *et al.* Mouse totipotent stem cells captured and maintained through
464 spliceosomal repression. *Cell* **184**, 2843-2859 e2820, doi:10.1016/j.cell.2021.04.020
465 (2021).

466 44 Kotake, Y. *et al.* Splicing factor SF3b as a target of the antitumor natural product
467 pladienolide. *Nat Chem Biol* **3**, 570-575, doi:10.1038/nchembio.2007.16 (2007).

468 45 Sun, C. The SF3b complex: splicing and beyond. *Cell Mol Life Sci* **77**, 3583-3595,
469 doi:10.1007/s00018-020-03493-z (2020).

470

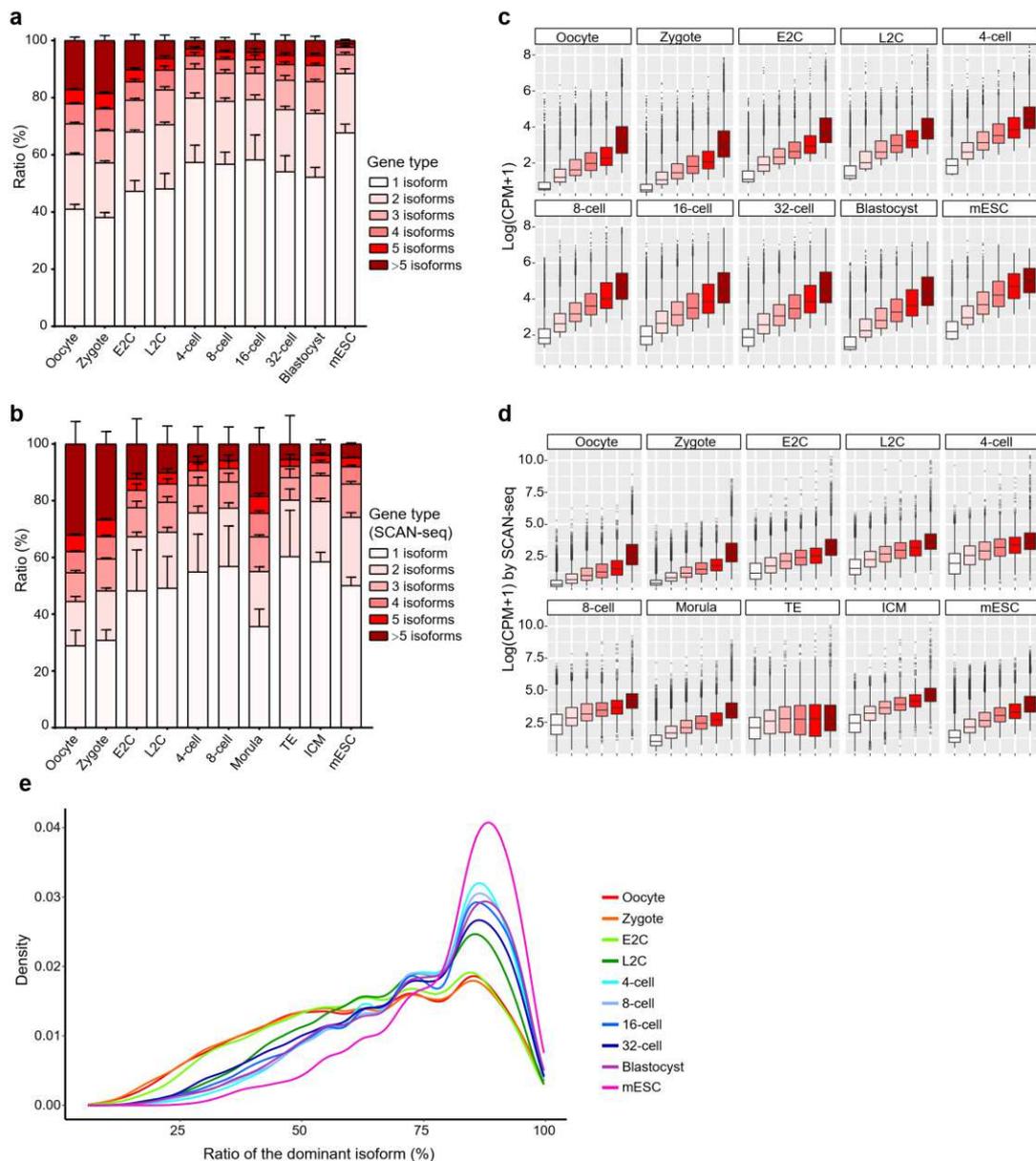
471



474

475 **Extended Data Fig. 1 | Quality evaluation of the single-cell isoform expression data**

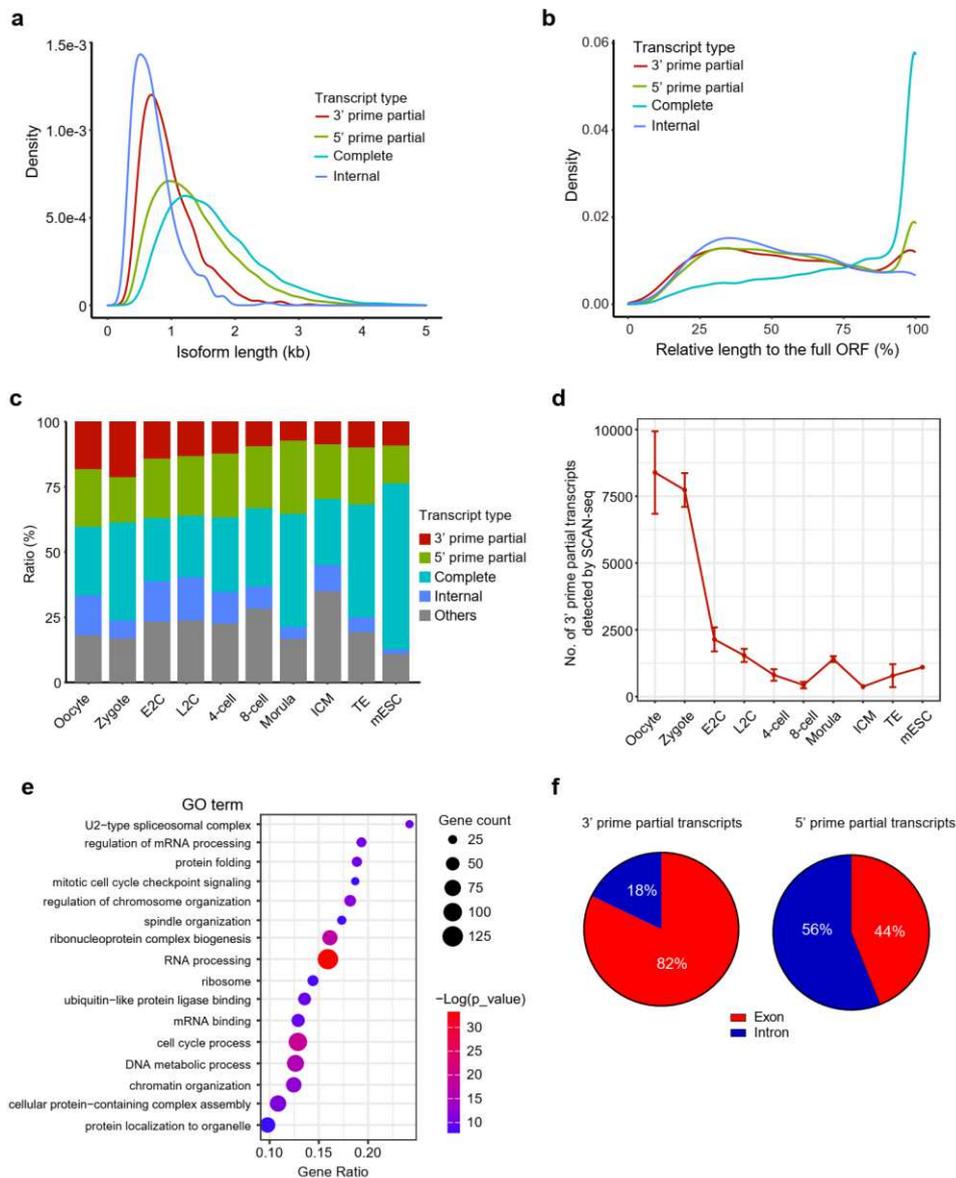
476 **a**, Diagram of the experimental and analysis workflow for single-cell isoform sequencing of
 477 the mouse preimplantation embryos. **b**, Saturation curve of representative cells from each
 478 stage. **c**, Correlation between detected UMI counts and absolute spiked abundances of
 479 each ERCC gene. **d**, Isoform mapping results of the SIRV spike-ins.



480

481 **Extended Data Fig. 2| Relationship between gene and isoform expression**

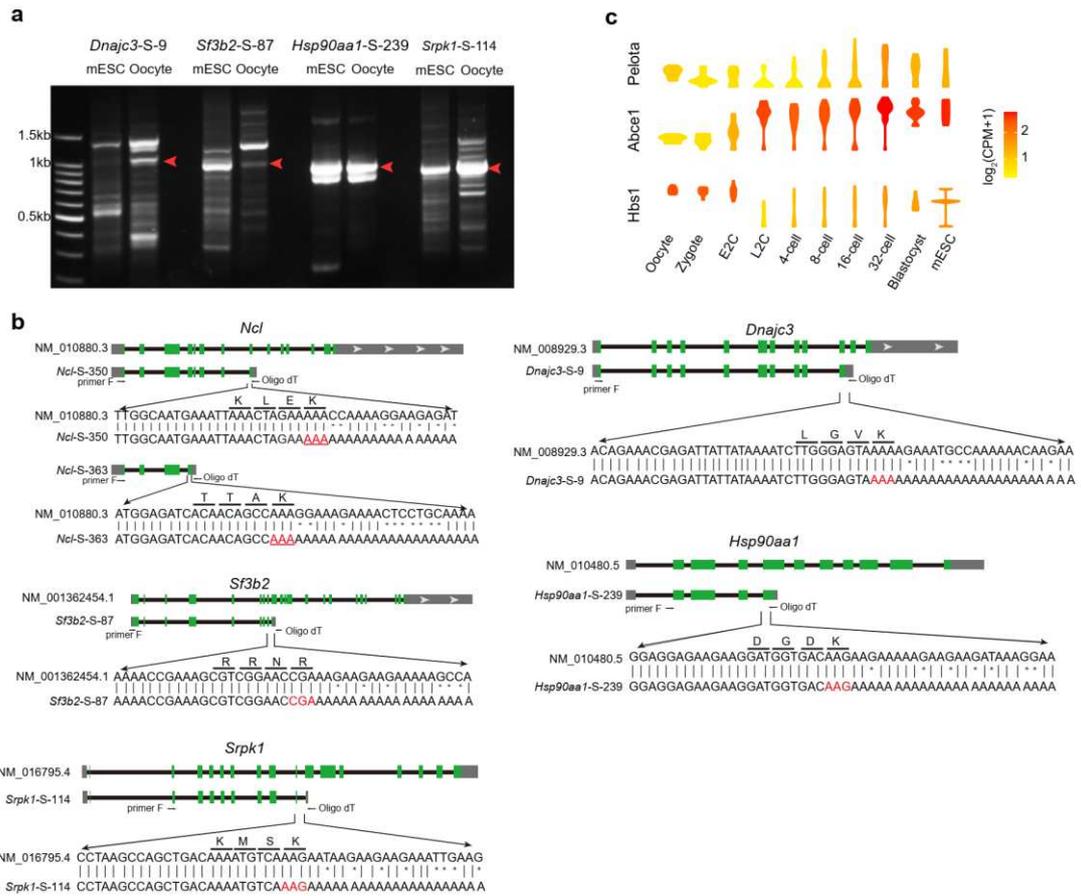
482 **a-b**, The ratios of genes detected with different numbers of isoform types for each stage of
 483 mouse embryos and mESCs in this study (**a**) and SCAN-seq data (**b**). **c-d**, Expression
 484 levels of genes detected with different numbers of isoform types for each stage of mouse
 485 embryos and mESCs in this study (**c**) and SCAN-seq data (**d**). **e**, Density plot showing the
 486 proportion of the major isoforms in genes expressing multiple isoform types. Only genes
 487 detected with UMI counts over 5 were included.



488

489 **Extended Data Fig. 3| The characteristics of different types of transcripts**

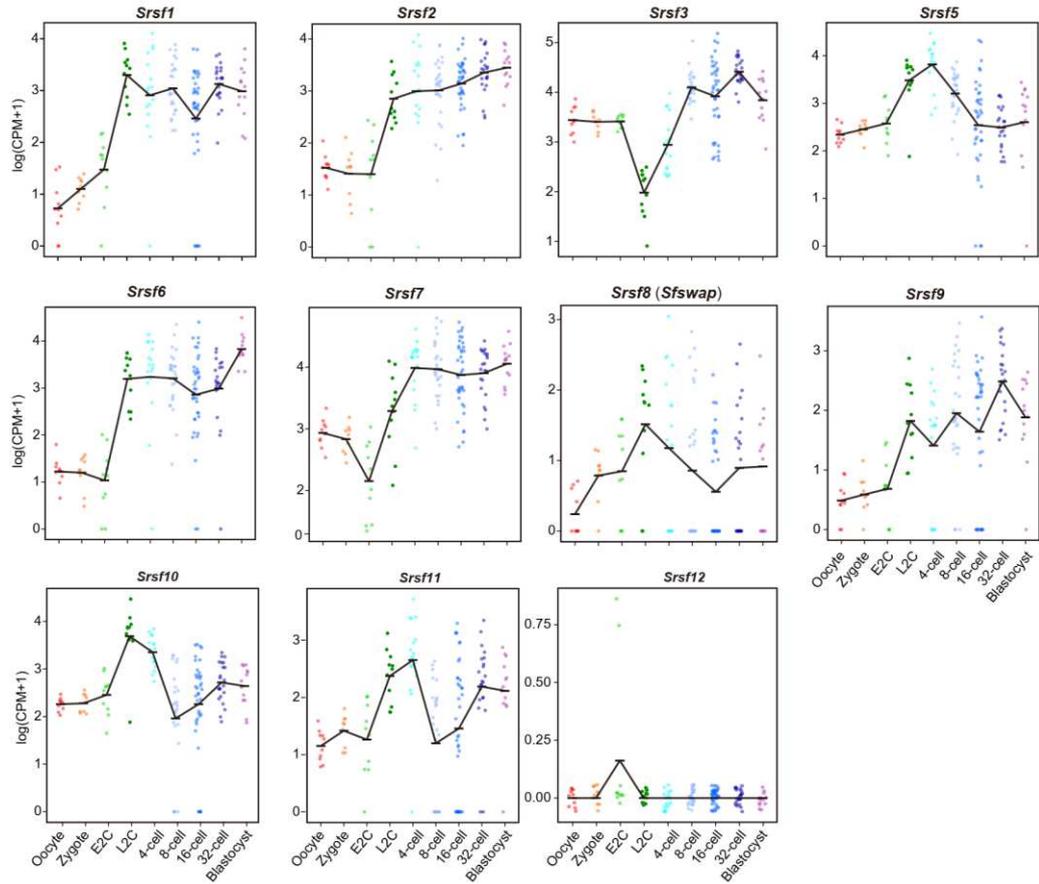
490 **a**, Length distribution of different types of transcripts. **b**, Relative length of the predicted
 491 protein to the complete reference ORF of each type of transcript. **c**, Ratios of each type
 492 of transcript at different stages calculated using SCAN-seq data. **d**, Expression level of
 493 the 3' prime partial transcripts detected at each stage in SCAN-seq data. **e**, The top GO
 494 terms for genes generating 3' prime partial transcripts in the mouse embryo samples. **f**,
 495 The ratios of 3' and 5' prime partial transcripts ended in exons and introns.



496

497 **Extended Data Fig. 4| Validation of the 3' prime partial transcripts**

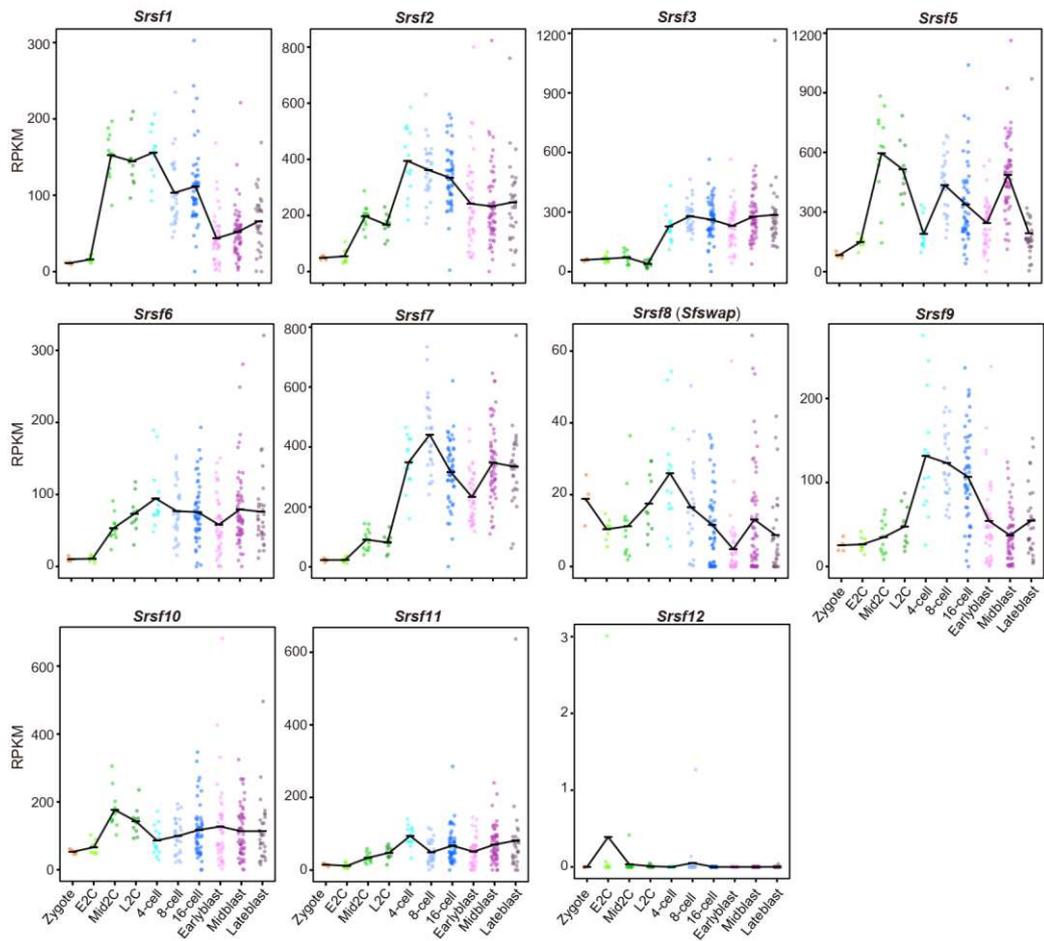
498 **a**, Gel picture showing the isoforms by RT-PCR of the target genes in mESCs and mouse
499 oocytes. The red arrows indicate the candidate 3' prime partial transcripts of the genes. **b**,
500 Sanger sequencing of the candidate 3' prime partial transcripts in panel **a** and the *Ncl* gene.
501 **c**, Violin plot showing the expression levels of the essential NSD genes at different mouse
502 embryonic stages.



503

504 **Extended Data Fig. 5| Expression levels of each SR family gene at different mouse**

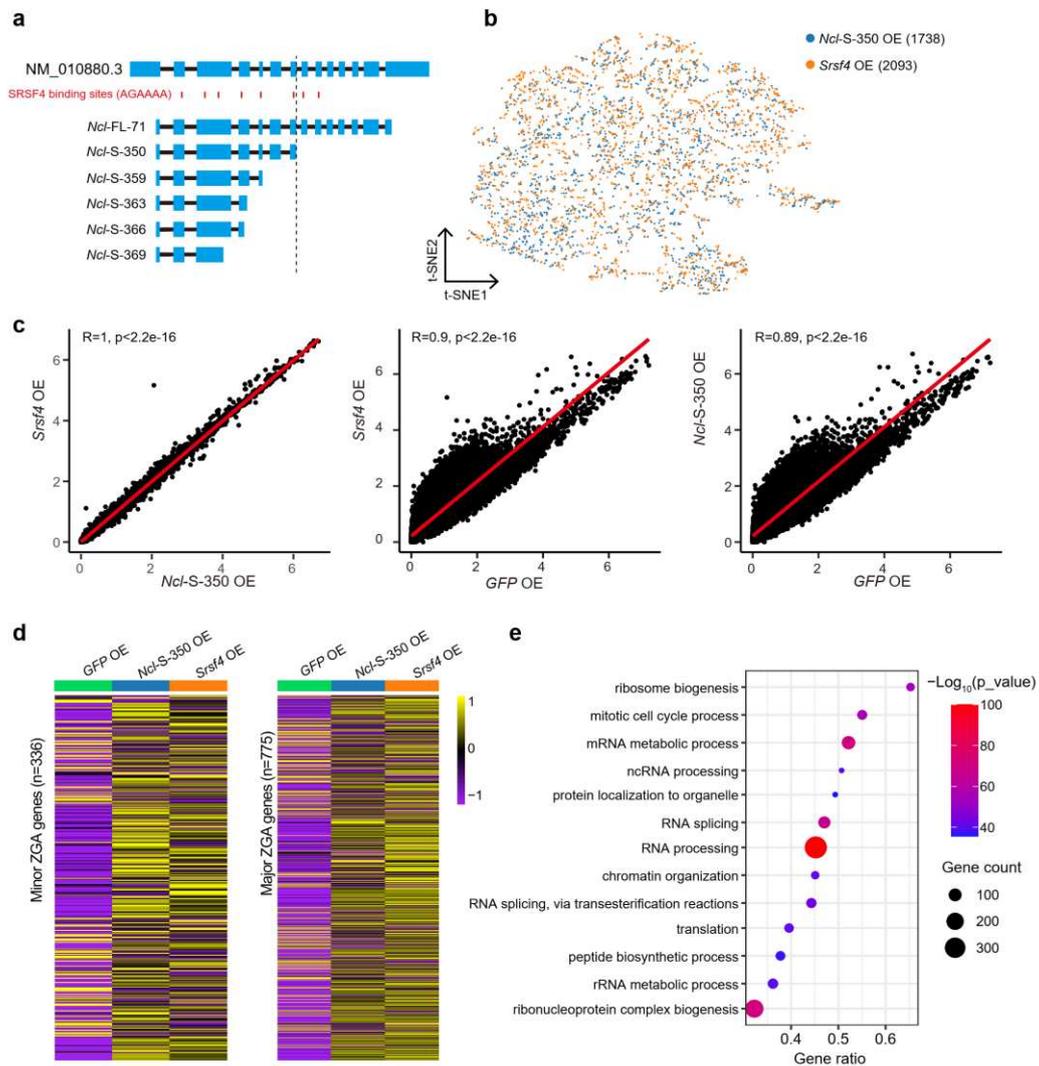
505 **embryonic stages in our data.**



506

507 **Extended Data Fig. 6| Expression level of each SR family gene at different mouse**

508 **embryonic stages in the NGS data by Deng *et al*⁶⁹.**



509

510 **Extended Data Fig. 7 | Transcriptomic changes after overexpression of *Srsf4* or *Ncl-***
 511 ***S-350* in mESCs**

512 **a**, The location distribution of the SRSF4-specific binding motif on the *Ncl* transcript. **b**, The
 513 t-SNE map showing the high consistency of the *Srsf4* OE and *Ncl-S-350* OE mESCs based
 514 on gene expression patterns. **c**, Gene expression correlation between the *Srsf4* OE sample
 515 and the *Ncl-S-350* OE sample, the *GFP* OE sample and the *Srsf4* OE sample, and the
 516 *GFPOE* sample and the *Ncl-S-350* OE sample respectively. **d**, Heatmap showing generally
 517 elevated expression of minor and major ZGA genes in *Srsf4* OE and *Ncl-S-350* OE cells
 518 compared with control *GFP* OE cells. **e**, The top GO terms for genes generating 3' prime
 519 partial transcripts in the *Srsf4* OE mESCs

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [ExtendeddataTable1.xlsx](#)
- [ExtendeddataTable2.xlsx](#)
- [ExtendeddataTable3.xlsx](#)
- [ExtendeddataTable4.xlsx](#)
- [ExtendeddataTable5.xlsx](#)
- [ExtendeddataTable6.xlsx](#)