

16 **Abstract**

17 Transposable elements (TEs) make up a majority of a typical eukaryote's genome,
18 and contribute to cell heterogeneity and fate in unclear ways. Single cell-sequencing
19 technologies are powerful tools to explore cells, however analysis is typically gene-
20 centric and TE activity has not been addressed. Here, we developed a single-cell TE
21 processing pipeline, scTE, and report the activity of TEs in single cells in a range of
22 biological contexts. Specific TE types were expressed in subpopulations of embryonic
23 stem cells and were dynamically regulated during pluripotency reprogramming,
24 differentiation, and embryogenesis. Unexpectedly, TEs were expressed in somatic
25 cells, including human disease-specific TEs that are undetectable in bulk analyses.
26 Finally, we applied scTE to single cell ATAC-seq data, and demonstrate that scTE can
27 discriminate cell type using chromatin accessibility of TEs alone. Overall, our results
28 reveal the dynamic patterns of TEs in single cells and their contributions to cell fate
29 and heterogeneity.

30

31 **Keywords:** single cell RNA-seq; transposable element; gene expression.

32 **Introduction**

33 Transposable elements (TEs) are a heterogeneous collection of genomic elements
34 that have at various stages invaded and replicated extensively in eukaryotic genomes.
35 The vast majority of TEs are fossils, and can no longer duplicate themselves, but they
36 remain inside the genome and in mammals occupy nearly half the total DNA¹.
37 Intriguingly, it is becoming clear that both the active and remnant TEs are participating
38 in evolutionary innovation and in biological processes²⁻⁵, such as embryonic
39 development⁶⁻⁹, and in human disease and cancer^{10,11}. Additionally, TEs carry cis-
40 regulatory sequences and their duplication and insertion can reshape gene regulatory
41 networks by redistributing transcription factor (TF) binding sites and evolving new
42 enhancer activities¹²⁻¹⁴. TEs transcription also has a key influence upon the
43 transcriptional output of the mammalian genome¹⁵. However, the role of TEs in cell
44 type heterogeneity and biological processes has only recently begun to be explored
45 in depth.

46 Single cell RNA-seq (scRNA-seq) has developed as a powerful tool to observe
47 cell activity¹⁶⁻¹⁸. Many new techniques have been developed to recover or reconstruct
48 missing observations, such as spatial, temporal, and cell lineage information. However,
49 an important source of genomic information has so far been overlooked in single cell
50 studies: the effect of TEs. Despite their importance, we lack quantitative understanding
51 of how those genomic elements are involved in cell fate regulation at the single cell
52 level. As TEs pose unique challenges in quantification, due to their degeneracy and
53 multiple genomic copies, a prerequisite to understand TEs at the single cell level is a
54 tool to quantify the hundreds to millions of copies of repetitive elements within the

55 genome. To this end, we developed scTE, an algorithm that quantifies TE expression
56 in single-cell sequence data.

57 We firstly demonstrate scTE's capabilities through an analysis of mouse
58 embryonic stem cells (mESCs), which is one of the best characterized models for TE
59 expression, as the expression of the endogenous retrovirus (ERV) MERVL marks a
60 small population of cells in embryonic stem cell (ESC) cultures that are totipotent^{19,20},
61 scTE could accurately recover the expected pattern of heterogeneous MERVL
62 expression. Then, we applied our approach to several biological systems including
63 human *in vitro* cardiac differentiation, mouse gastrulation, adult mouse somatic cells,
64 the induced pluripotent reprogramming process and human disease data. Overall, we
65 unveil hitherto unknown insights into complex TE expression patterns in mammalian
66 development and human diseases.

67

68 **Results**

69 **Quantification of TE expression in single cells with scTE**

70 Analysis of TEs pose special challenges as they are present in many hundreds to
71 millions of copies within the genome. A common strategy in regular analyses is to
72 discard multiple mapped reads, however this leads to loss of information from TEs²¹.
73 Assigning these reads to the best alignment location is the simplest way to resolve
74 TE-derived reads, but it is not always correct for individual copies^{21,22}. To solve this
75 problem, we designed an algorithm in which TE reads are allocated to TE metagenes
76 based on the TE type-specific sequence. We built a framework named scTE with this
77 strategy, scTE maps reads to genes/TEs, performs barcode demultiplexing, quality
78 filtering, and generates a matrix of read counts for each cell and gene/TE ([Fig. 1a and](#)

79 [Supplementary Fig. 1a](#)). scTE is easy to use, and its output is designed to be easily
80 integrated into downstream analysis pipelines including, but not limited to, Seurat and
81 SCANPY^{23,24}. The algorithm can in principle be applied to infer TE activities from any
82 type of single-cell sequencing based data, like single-cell ATAC-seq data, DNA
83 methylation, and other single-cell epigenetic data.

84 We first tested scTE's ability by *in silico* mixing two cells lines, MEFs (mouse
85 embryonic fibroblasts) and ESCs in different ratios²⁵. Comparison with the gene-based
86 Cell Ranger pipeline²⁶, scTE shows nearly identical topology in a UMAP (Uniform
87 Manifold Approximation and Projection) plot, and in marker genes expression ([Fig. 1b](#)
88 [and Supplementary Fig.1b](#)). Even when one cell type only contributes a 1% minority
89 in the mixture, scTE identified it correctly ([Fig. 1b](#)), indicating that scTE did not
90 influence the global analysis of gene expression. These results demonstrate the
91 sensitivity of scTE.

92 Next, we sought to explore TE expression, around 12-14% of the reads were
93 derived from TEs ([Fig. 1c](#)). Requiring at least 2-fold change and FDR<0.05, scTE
94 detected 150 significantly differentially expressed TEs between ESCs and MEFs
95 ([Supplementary Fig. 1c](#)), including ERVB7_1-LTR_MM, which is highly expressed in
96 ESCs, and RMER10B in MEFs ([Fig. 1d and Supplementary Fig.1d](#)). Furthermore,
97 UMAP based on single cell TE expression alone could distinguish the cell types with
98 the expected ratio ([Fig. 1e](#)), demonstrating TE expression discerns cell identity.

99

100 **Deciphering TE heterogeneity in mouse ESCs and during human cardiac**
101 **differentiation**

102 It is known that a small subset of ESCs acquire a totipotent state named 2C-like cells
103 and express a MERVL TE which also marks the embryonic 2-cell stage^{19,27,28}. scTE
104 could correctly identify this rare 2C-like subpopulation in UMAP plots, based on the
105 specific marker genes *Zscan4c* and *Tcstv3*, and the expression of MERVL and
106 MT2_Mm TEs (Fig. 2a, b and Supplementary Fig. 2a, b)^{19,29}. If we discarded multiple
107 mapped reads and only considered unique reads, the level of MERVLs was reduced,
108 but it was still specifically expressed in the 2C-like cells (Supplementary Fig. 2c). This
109 confirms that scTE can correctly identify known TE patterns.

110 In humans, HERV-H LTRs are expressed in early embryos and human pluripotent
111 stem cells (hPSCs), and contribute to pluripotency maintenance and somatic
112 reprogramming^{6,30-32}, but little is known about TE expression dynamics during
113 differentiation to somatic cells. Applying scTE to an scRNA-seq time series of hPSCs
114 differentiating to cardiomyocytes³³, we accurately recovered the repression of HERV-
115 H LTRs including LTR7 and HERVH-int during differentiation, concomitant with
116 reduction in the expression of the pluripotency factor *POU5F1* (Fig. 2c, d and
117 Supplementary Fig. 2d). During *in vitro* cardiac differentiation of hPSCs there is a
118 bifurcation towards definitive cardiomyocytes (dCM) and non-contractile cells (Fig. 2c).
119 Between these two branches, marked by *NKX2-5* and *SPARC*, respectively, we found
120 differential expression of TEs such as LTR32, MER57A-int and MER45A in the dCM
121 cells, whilst, MLT1H1, HERVIP10B-int and LTR5A were specifically expressed in the
122 non-contractile cells (Fig. 2e, f and Supplementary Fig. 2e). Independent bulk RNA-
123 seq data³⁴ demonstrated that these TEs were expressed in late cardiac differentiation
124 (Supplementary Fig. 2f), however, as the bulk is a mixture of dCM and non-contractile
125 cells, the restriction of these TEs to divergent fates can only be observed in the

126 scRNA-seq data. This highlights the importance of analyzing TE expression in sc-
127 RNA-seq data, as MLT1H1 is very high in the bulk RNA-seq, but this hides the reality
128 that it is restricted to the non-contractile cells and plays no role in dCMs (Fig. 2e, f and
129 [Supplementary Fig. 2f](#)).

130

131 **Analysis of TEs in mouse gastrulation and early organogenesis reveals the** 132 **widespread cell fate-specific expression of TEs**

133 The previous analysis showed how TE expression contributed to *in vitro* cardiac
134 differentiation, next we explored complex *in vivo* developmental processes. TE
135 expression is dynamic during pre-implantation development⁶, however the expression
136 of TEs in gastrulation has not been described. We took advantage of the single-cell
137 time course of mouse gastrulation¹⁶. Analysis with scTE did not introduce any
138 unexpected sample-bias, and a side-by-side comparison could retrieve similar
139 patterns of marker gene expression in the expected lineages (Fig. 3a and
140 [Supplementary Fig. 3a-f](#)). We found every lineage expressed a series of lineage-
141 specific TEs (Fig. 3a, b, and [Supplementary Fig. 4a-c](#)). In the extraembryonic
142 ectoderm cells, IAP and RLTR45-family TEs were activated (Fig. 3b, c), and in *Apoa2+*
143 extraembryonic endoderm cells, MER46C, RLTR20B3 and LTRIS2 were up-regulated
144 (Fig. 3b, d). The expression of these TEs was validated using bulk RNA-seq from *in*
145 *vitro*³⁵⁻³⁷ mimics of these embryonic stages including ESCs, epiblast stem cells
146 (EpiSCs), extraembryonic endoderm cells (XENs) and trophoblast stem cells (TSCs)
147 (Fig. 3e). Other embryonic lineages, particularly the *Gypa+* erythroid and the *Tnnt2+*
148 cardiomyocyte lineages expressed specific TEs such as L1_Mur and L1ME3D,
149 respectively (Fig. 3b, f).

150 As this dataset provides dynamic trajectories for each lineage, we wondered if
151 TEs were transiently activated during cell fate commitment. To this end, we noticed
152 ETnERV3-int, whose expression coincides with the early development of the cardiac
153 fate from the mesoderm, and is reduced in *Tnnt2*⁺ cells, while L1ME3D was expressed
154 in the *Tnnt2*⁺ cells (Fig. 3g). Consistently, ETnERV3-int was specifically expressed in
155 *in vitro* derived cardiomyocytes, which more closely resemble a fetal state, whilst
156 L1ME3D was expressed only in the mature heart (Fig. 3h)^{38,39}. However, the bulk
157 samples could not capture the complexity of the transient expression of ETnERV3-int
158 which extended from the late epiblast into the endoderm and mesoderm. To expand
159 on this, we reanalyzed an scRNA-seq dataset of the developing mouse embryonic
160 heart⁴⁰ (Fig. 3i and Supplementary Fig. 5a-c), and found that ETnERV3-int was
161 expressed in the myocardium and epicardium, but not in the endocardium, neural crest
162 and embryonic cells (Fig. 3j). L1ME3D was expressed in *Tnnt2*⁺ myocardium,
163 however in an inverse pattern with respect to ETnERV3-int (Fig. 3j, k). Therefore,
164 ETnERV3-int activity is present in an intermediate stage in cardiac lineage
165 development. Intriguingly, there was a close relationship between the expression of
166 ETnERV3-int and *Isl1* gene, which marks multipotent progenitors⁴⁰ (Fig. 3j). These
167 results highlight the complex patterns of TE expression in developmental processes.

168

169 **Widespread tissue-specific expression of TEs in somatic cells**

170 TE activity is considered to be silenced in somatic cells except LINE-1 expression and
171 retrotransposition in the developing brain^{27,28,41}. As we revealed unexpected
172 heterogeneity of TEs in somatic MEFs and during organogenesis, we next measured
173 TE expression in somatic cells using the Tabula Muris large scale scRNA-seq dataset

174 that profiles 20 mouse organs⁴² (Fig. 4a). Surprisingly, our analysis revealed in total
175 130 TEs that were specifically expressed in distinct cell types (Fig. 4b and
176 Supplementary Fig. 6a). These associations include the expected expression of LINE1
177 elements in brain cells, of which many L1 family members like L1MEh, L1M, L1MC4a,
178 L1MA7 and L1P5 elements are specifically expressed in oligodendrocytes or microglia
179 (Fig.4c and Supplementary Fig. 6a). We also found expression of LTR58, MLT1EA-
180 int, MER110 and RLTR46 that specifically in B cells, T cells, type B pancreatic cells
181 and hepatocytes, respectively (Fig.4c). Next, we took advantage of the Tabula Muris
182 dataset to measure overall TE expression heterogeneity, and, in general, the LTRs
183 and DNA transposons are the major source of heterogeneity (Fig. 4d, e).

184 TE expression is regulated by chromatin modification and transcription factors
185 (TFs)³, thus, we wondered if we could infer the regulatory network between TFs and
186 TEs from large scale scRNA-seq data, taking advantage of the improved cell type
187 definitions from the scRNA-seq data. The co-expression relationships often reflect
188 biological processes in which many genes with related functions are coordinately
189 regulated. Therefore, we reasoned that if a TE is regulated by a TF, they should be
190 co-expressed. To identify TF-TE regulatory relationships, we performed co-expression
191 analysis, and revealed the specific co-clustering of neural genes and TEs (*Sox2* and
192 *Olig1*), the immune system (*Cebpe*, *Tcf7*, *Pax5* and *Sall1*), the endoderm/pancreas
193 (*Gfi1b*, *Nkx6-1* and *E2f8*), and other lineages (Fig. 4f and Supplementary Fig. 6b).
194 Motif analysis also showed that the SOX2 motif was significantly enriched within
195 RLTR13F TEs (Supplementary Fig. 6c). These results highlight the deep link between
196 TE and TF activity indicating those TFs may be responsible for activating TEs in the
197 corresponding cell types.

198 We next explored in closer detail neural and immune cell lineages as TE activity
199 is known to regulate neural activity and immune responses⁴³⁻⁴⁵. Subgrouping the cells
200 from microglia and neuron samples identified several distinct cell types
201 (Supplementary Fig. 7a-c), within which cell type-specific expression of TEs was
202 observed (Supplementary Fig. 7d, e). Next, with the pooled immune cells from marrow,
203 spleen and thymus, 12 distinct immune cell subtypes were defined (Supplementary
204 Fig. 7f, g). Intriguingly, besides finding additional cell type-specific TEs in T cells, B
205 cells and granulocytes, a series of TEs were restricted to subtypes of T cells and B
206 cells (Supplementary Fig. 7h, i). These data show different degrees of subtype
207 specific signatures of TEs in the neural and immune system, and highlight the
208 importance of looking beyond only genes when exploring how those systems differ.

209

210 **TEs are activated during somatic cell reprogramming, in a heterogenous and** 211 **cell branch restricted manner**

212 The above analysis has revealed the well-ordered dynamic expression of TEs in
213 developmental processes, we then wondered if TEs undergo similar stage-specific
214 regulation during somatic reprogramming. Somatic cells can be reprogrammed to
215 induced pluripotent stem cells (iPSCs) by various methods, such as ectopic
216 expression of a group of pluripotency transcription factors^{25,46,47}, or cocktails of
217 chemicals^{48,49}. The reprogramming process is highly heterogeneous, with abundant
218 non-reprogramming cells and divergent cell fate transition routes^{25,50}. We took
219 advantage of reprogramming scRNA-seq data to investigate the activity of TEs during
220 these drastic cell fate transitions. Reprogramming induced by *Oct4/Pou5f1*, *Klf4*, *Sox2*
221 and *c-Myc* (OKSM) generates detectable intermediate branches, including iPSCs,

222 trophoblast, stromal and neural-like cells (Fig. 5a and Supplementary Fig. 8a-d)⁵⁰. We
223 identified specifically expressed TEs in each cell branch (Supplementary Fig. 8a-d).
224 For example, the TEs ERVB7_1-LTR_MM, IAPEz-int, RLTR4_Mm, and Lx were
225 specifically expressed in iPSCs, trophoblast, stromal and neural-like branches,
226 respectively (Fig. 5b). ERVB7_1-LTR_MM (MusD) and IAPs are up-regulated during
227 reprogramming⁵¹, however using scRNA-seq data we show that only ERVB7_1-
228 LTR_MM, as well as ETnERV-int and RLTR13G, were up-regulated in the successful
229 reprogramming route, initiating at the mesenchymal-to-epithelial transition (MET) and
230 peaking at the iPSCs stage (Fig. 5b and Supplementary Fig. 8a). In contrast, the
231 trophoblast-branch expressed IAPEz-int and IAPLTR1_Mm (Fig. 5b and
232 Supplementary Fig. 8c), which are also expressed in *in vivo* extra embryonic ectoderm
233 cells (Fig. 3c), suggesting consistent regulation between development and
234 reprogramming.

235 We then analyzed reprogramming induced by *Oct4*, *Klf4*, and *Sox2* (OKS)²⁵ or
236 only chemicals²⁹. There are two validated branches during OKS-mediated
237 reprogramming²⁵ (Fig. 5c), and we found many TEs, such as ERVB7_1-LTR_MM, that
238 were specifically up-regulated in the reprogramming-potential (RP) branch, and were
239 excluded from the non-reprogramming branch (Fig. 5d and Supplementary Fig. 8e).
240 IAPEz-int and IAPLTR1_Mm were expressed in the RP branch but were ultimately
241 silenced in the reprogrammed cells (Fig. 5e, f), suggesting IAPs were only activated
242 in a pre-reprogrammed state and may impede the final step of pluripotency acquisition.
243 We validated the expression of ERVB7_1-LTR_MM and IAPs by qRT-PCR
244 (Supplementary Fig. 8f), demonstrating that IAPs are silenced in ESCs. Similar to
245 OKS-mediated reprogramming, chemical-mediated reprogramming bifurcates into two

246 branches (Fig. 5g and Supplementary Fig. 8g)²⁹, and TEs, marking an intermediate
247 2C-like program, were activated at the root of the successful branch (Supplementary
248 Fig. 8h, i). ERVB7_1-LTR_MM and RLTR13G were specifically up-regulated in the
249 successful branch, whilst IAPEz-int and IAPLTR1_Mm were activated in the pre-
250 branch and failed branch (Fig. 5h and Supplementary Fig. 8j, k).

251 The similar expression pattern of TEs among the three distinct reprogramming
252 systems described above, suggests there are common regulatory mechanisms.
253 Indeed, we found IAPLTR1_Mm TEs are rich in DNA-binding motifs for JUN and IRF2
254 (Supplementary Fig. 8l), whose expression closely matched IAP expression in all
255 three reprogramming systems (Supplementary Fig. 8m) and are known to impair
256 reprogramming^{52,53}. This suggests that their downregulation deactivates the IAPs
257 before the finalization of reprogramming, indicating IAPs may impede the final step of
258 reprogramming. Overall, these results indicate TEs have a deeper unappreciated role
259 in iPSC formation.

260

261 **Inferring TE Associated Accessibility from scATAC-seq Data**

262 Beyond scRNA-seq, many other single-cell sequencing techniques⁵⁴⁻⁵⁶ have shown
263 great potential to explore cell heterogeneity and increased insight could be fueled by
264 the additional information provided by scTE. For instance, we reasoned that scTE
265 would be informative for the analysis of scATAC-seq data and potentially other single-
266 cell epigenetic data because TEs have a wide array of chromatin states³, are widely
267 bound by transcription factors⁵⁷, and can act as enhancers¹⁴ (Fig. 6a). We then applied
268 scTE to a dataset of fluorescence-activated cell sorted (FACS) mouse cells⁵⁸,
269 including cardiac progenitor cells (CPCs), CD4⁺ T cells, ESCs and skin fibroblasts

270 (SFs). Intriguingly, scTE could accurately recover the expected cell types, based on
271 only the reads that mapped to TEs (Fig. 6b). Specific accessibility of RLTR13A,
272 RLTR4_Mm, RLTR13G and RMER19B/C was found in the CPCs, CD4+ T cells, ESCs
273 and SFs, respectively (Fig. 6c, d and Supplementary Fig.9a). And motif enrichment of
274 these cell-type specific TEs revealed known master regulators of these cell types, such
275 as GATA4/HAND1/T for CPCs, ETS1/TCF3 for T cells, SOX2/POU5F1/NR5A2 for
276 ESCs and FOS/MAF for SFs (Supplementary Fig. 9b), indicating these TEs may act
277 as cis-regulatory elements bound by transcription factors. For instance, scTE reveals
278 there is an RLTR13A TE within an intron of *Smyd1*, a gene essential for heart
279 development⁵⁹⁻⁶¹, which was specifically open in CPCs (Fig. 6e), and was specifically
280 expressed in the myocardium of the fetal heart (Fig. 6f). Applying scTE to scATAC-
281 seq data of peripheral blood monocyte cells (PBMC) was also able to recover the
282 major cell types and cell type-specific TEs (Supplementary Fig. 8c-f), which can be
283 validated by independent bulk ATAC-seq data from FACS sorted cells (Supplementary
284 Fig. 8g)⁶². These results indicate that quantifying chromatin accessibility on TE regions
285 is informative for characterizing cell types and may assist the problems posed by
286 scATAC-seq analysis due to its especially sparse nature⁶³.

287

288 **Disease-specific expression of TEs**

289 The unexpected widespread TE heterogeneity amongst embryonic and somatic cell
290 types and cell fate transitions raised the question as to whether there is TE
291 heterogeneity in diseased cells. Alzheimer's disease (AD) is an age-associated
292 neurodegenerative disorder that is characterized by progressive memory loss and
293 cognitive dysfunction for which there is no known cure. TEs have been reported to be

294 highly active during aging and may contribute to age-dependent loss of neuronal
295 function⁶⁴. To explore the expression of TEs in AD, we reanalyzed the scRNA-seq
296 data from a mouse model of AD expressing five human familial AD gene mutations,
297 which contained 13,114 single cells with age and sex-matched wild-type (WT) controls
298 using the MARS-seq platform⁶⁵ (Fig. 7a). Projecting the cells with a UMAP, we
299 recovered the major groups of cells in AD and WT, including the unique disease-
300 associated microglia cluster cells (M2) identified in the original study (Fig. 7b and
301 Supplementary Fig. 10a). Differential expression analysis demonstrated significant
302 changes in gene expression in M2, including previously described AD risk factors such
303 as *ApoE*, *Tyrobp*, *Lpl*, *CstII* and *Trem2* (Fig. 7c and Supplementary Fig. 10b).
304 Intriguingly, we also found many TEs such as ERVB7_2-LTR_MM, RLTR17, RLTR28
305 and Lx4B that were significantly higher and specifically expressed in M2 (Fig. 7c, d
306 and Supplementary Fig. 10c), indicating those TEs may also be involved in AD
307 development.

308 Type 2 diabetes (T2D) is a common human disease caused by a combination of
309 increased insulin resistance and reduced mass or dysfunction of pancreatic beta cells.
310 We reanalyzed scRNA-seq from two independent studies of the human pancreas in
311 healthy and T2D individuals^{66,67}. The major cell types in the pancreas, including alpha,
312 beta, gamma/PP and delta cells clustered without a visible disease-specific pattern,
313 indicating no drastic change in cell type (Fig. 7e and Supplementary Fig. 10d).
314 Contrasting the transcriptome from healthy and T2D in each cell type independently,
315 *CD36* and *DLK1* was up-regulated in T2D alpha and beta cells respectively (Fig. 7f),
316 as reported by the original studies^{66,67}. Notably, many TEs were significantly highly
317 expressed in T2D beta cells, including L1MC, L1MA4A, Tigger3a, MLT2B4. This

318 differential expression pattern was near identical between the two independent
319 datasets (Fig. 7f). Critically, none of these observations could be observed using bulk
320 RNA-seq datasets (Fig. 7g and Supplementary Fig. 10e)^{66,68}, which might be due to
321 the high expression of these TEs in both normal and T2D alpha cells, emphasizing the
322 importance of analysis at single-cell resolution.

323 As a final human disease dataset we reanalyzed a glioblastoma scRNA-seq
324 experiment⁶⁹, and were able to identify TEs specifically expressed in neoplastic cells
325 and that were correlated with the expression of *EGFR* (Supplementary Fig. 10f-h), a
326 gene upregulated in a large percentage of glioblastomas⁶⁹. Above all, these results
327 revealed the dysregulation of TE expression in diseased human cells, which deserves
328 further mechanistic study and may help to identify new diagnostic markers and
329 therapeutic targets.

330

331 Discussion

332 TEs are the most abundant elements in the genome, however, the understanding of
333 their impact on genome evolution, function and disease remains limited. The rise of
334 genomics and large-scale high-throughput sequencing has shed light on the multi-
335 faceted role of TEs. However, many genomic studies exclude TEs due to difficulties
336 in their analysis as a consequence of their repetitive nature²¹. Thus, TE analysis often
337 requires the use of specialized tools to extract meaning^{5,22}. Here, we developed scTE
338 specifically for the analysis of TEs from single-cell sequencing data. By taking
339 advantage of this tool we could recover previously identified phenomena such as
340 MERVL and LTR7/HERVH expression in mouse and human ESCs, respectively. We
341 then revealed widespread heterogeneity of TE expression throughout embryonic

342 development, in mature somatic cells, during the reprogramming process and in
343 human diseases, and discovered a wealth of cell fate-specific TE expression. These
344 associations with cell fate cannot be observed when only considering bulk samples,
345 demonstrating the enormous power of single-cell sequencing, and the importance of
346 analyzing TE expression.

347 One of the key findings of our analysis has revealed the various TEs that are
348 specifically expressed in different cell types. The expression of TEs during the pre-
349 implantation development stage has been demonstrated previously⁶, our findings
350 extend this to gastrulation and early organogenesis. We find a wide array of
351 expression of TEs in the extraembryonic tissues, which may be related to their activity
352 as enhancers⁷⁰. Furthermore, we show the expression of TEs within the specific
353 lineages in the developing fetal heart. In addition, TEs are also heterogeneously
354 expressed between cell types in adult somatic cells, which has not been demonstrated
355 before, as TEs are thought to be primarily silent in adult tissues. Notably, we found a
356 vast of trove of TEs that are expressed in the brain and the immune system, and
357 individual TE types that are specifically expressed in different sub cell types.
358 Considering the close relationship between the evolution of immune system, brain and
359 TEs⁴³⁻⁴⁵, these results hint at further functions for TEs in these two systems.

360 How cells decide their fate is a fundamental question in biology. Stem cell
361 differentiation and somatic cell reprogramming are both powerful *in vitro* models that
362 mimic *in vivo* development and have provided great insight into cell fate decisions.
363 However, how TEs are involved in these processes is still largely unknown. In this
364 study, we have identified the TEs LTR32 and MLT1H1 that were differentially
365 regulated between contractile and non-contractile cell fate decisions during human

366 cardiac differentiation. In addition, we also found ERVB7_1-LTR_Mm and IAP
367 elements divergent expression during reprogramming, whereas ERVB7_1-LTR_Mm
368 may promote iPSC formation, IAP elements need to be silenced at the final stage
369 before iPSCs formation (Fig. 5b, f, h). These mechanisms are shared among the
370 Yamanaka factor based and chemical based reprogramming systems, indicating a
371 tight association between TEs and cell fate decisions.

372 Considering the growing implication that TEs are important contributors to human
373 disease, their study is becoming increasingly important. In addition to the ability of TEs
374 to impact genomic stability as they duplicate⁷¹, which has clear implications for the
375 development of cancer⁷², TEs are also playing more subtle roles in epigenetic control
376 and transcript expression. For example, TEs are spliced into chimeric transcripts that
377 drive the expression of oncogenes¹¹. Similarly, the expression of TEs has been
378 associated with several nervous system-related disorders, including
379 neurodegeneration¹⁰, and L1 LINE expression is important in inflammation during
380 aging⁷³. In our work, we demonstrate that in single cells of the pancreas there is
381 substantial TE expression deregulation in the beta cells, which is suggestive of
382 epigenetic dysfunction and a loss of control over TE expression. Critically, this
383 observation cannot be observed from bulk pancreatic islet samples. Considering the
384 growing importance of exploring human disease using primary patient samples, the
385 analysis of TEs should be included. However, to date the contribution of TE expression
386 to the aging and diseased states remains relatively unexplored. Our approach will be
387 an important tool in understanding the contributions of TEs to cellular heterogeneity in
388 a variety of systems and in human disease.

389

390 **Methods**

391 **Software availability**

392 scTE is available at <https://github.com/jphe/scTE>. The code is freely available and is
393 released under the MIT license. scTE requires Python >3.6, and the python module
394 numpy, scTE supports the Linux and Mac platforms. Software code for the analysis of
395 the data in this manuscript can be found at:
396 <https://github.com/jphe/scTE/tree/master/example>.

397

398 **scTE pipeline**

399 The input data for scTE consists of the annotation files for genes and TEs, and
400 alignment files in either the SAM or BAM format⁷⁴. By default, scTE uses GENCODE⁷⁵
401 and the UCSC genome browser Repeatmasker track⁷⁶ annotations for genes and TEs,
402 respectively. The SAM/BAM file contains the aligned read genome locations. Many
403 alignment programs can distinguish reads that have a unique alignment in the genome
404 (termed unique-reads) or map to multiple genomic loci (termed multimapping reads or
405 non-unique reads). Multimapping reads are critical for TE quantification, as TEs
406 contain many repeated sequences and non-unique reads often map inside the TEs.
407 To get an accurate quantitation of the number of reads mapping to TEs these reads
408 should be preserved. However, in many analysis pipelines these reads are discarded.
409 scTE recommends aligners to keep all of the mapped reads, and we recommend that
410 the best single aligned multimapped read be kept. The reads can be aligned by any
411 genome aligner, but the aligned reads must be against the genome (i.e. not against a
412 set of genes or transcript assembly). scTE is most tuned to STAR-solo⁷⁷ or the Cell
413 Ranger pipeline outputs, and can accept BAM files produced by either of these two

414 programs. For other aligners, the barcode should be stored in the 'CR:Z' tag, and the
415 UMI in the 'UR:Z' tag in the BAM file. If the UMI is missing or not used in the scRNA-
416 seq technology (for example on the Fluidigm C1 platform), it can be disabled with –
417 UMI False (the default is True) switch in scTE. If the barcode is missing it can be
418 disabled with the –CB False (the default is True), and instead the cell barcodes will be
419 taken from the names of the BAM files (multiple BAM files can be provided to scTE
420 with the –i option).

421

422 **scTE gene and TE indices**

423 scTE builds genome indices for the fast alignment of reads to genes and TEs. These
424 indices can be automatically generated using the commands:

```
425 scTE_build -g mm10 # mouse genome
```

```
426 scTE_build -g hg38 # human genome
```

427 These two scripts will automatically download the genome annotations, for mouse:

```
428 ftp://ftp.ebi.ac.uk/pub/databases/gencode/Gencode\_mouse/release\_M21/gencode.vM21.annotation.g  
429 tf.gz
```

```
430 http://hgdownload.soe.ucsc.edu/goldenPath/mm10/database/rmsk.txt.gz
```

431 Or for human:

```
432 ftp://ftp.ebi.ac.uk/pub/databases/gencode/Gencode\_human/release\_30/gencode.v30.annotation.gtf  
433 .gz
```

```
434 http://hgdownload.soe.ucsc.edu/goldenPath/hg38/database/rmsk.txt.gz
```

435 These annotations are then processed and converted into genome indices. The scTE
436 algorithm will allocate reads first to gene exons, and then to TEs, by default. Hence
437 TEs inside exon/UTR regions of genes annotated in GENCODE will only contribute to
438 the gene, and not to the TE score. This feature can be changed by setting '–mode/-m
439 exclusive' in scTE, which will instruct scTE to assign the reads to both TEs and genes
440 if a read comes from a TE inside exon/UTR regions of genes.

441

442 **Analysis of 10x-style data**

443 scRNA-seq data was processed using the scTE 10x pipeline, Briefly, reads were
444 aligned to the genome using STARsolo⁷⁷ with the setting ‘--outSAMattributes NH HI
445 AS nM CR CY UR UY --readFilesCommand zcat --outFilterMultimapNmax 100 --
446 winAnchorMultimapNmax 100 --outMultimapperOrder Random --runRNGseed 777 --
447 outSAMmultNmax 1’. The default scTE parameters for 10x were used to get the
448 molecule count matrix. The count matrix was lightly filtered to exclude cell barcodes
449 with low numbers of counts: Cells with less than 1000 UMIs and less than 500 genes
450 detected were filtered out, and only the top 10,000 cells with the highest gene count
451 were kept (these default setting can be altered with the ‘--expect-cells, --min_count
452 and --min_genes’ switches in scTE, note that the cell counts are further filtered on a
453 case-by-case basis for each experiment, as detailed below). Other downstream
454 analysis was performed by SCANPY²⁴. Specific analysis settings for the individual
455 datasets are described below.

456

457 **Analysis of C1/SMART-seq-style data**

458 scRNA-seq data were processed using the scTE C1/SMART-seq pipeline, Briefly,
459 reads were aligned to the genome using STAR⁷⁷, with the setting ‘--
460 winAnchorMultimapNmax 100 --outSAMmultNmax 1 --outSAMmultNmax 1’. The
461 default scTE parameters for C1/SMART-seq were used to get the molecule count
462 matrix. Cells with less than 10,000 counts and less than 2000 expressed genes were
463 filtered out. Cells with more than 20% fraction of mitochondrial counts were discarded.
464 Downstream analysis was performed the same as for the 10x data pipeline. Fluidigm

465 C1/SMART-seq data comes as a single BAM file per barcode. To analyze this data,
466 the 'barcode' is taken from the input BAM filenames, and both -CB and -UMI should
467 be False:

```
468 scTE -i *.bam -p 4 -o <output_name> --genome mm10 -x mm10.exclusive.idx -CB False -UMI  
469 False
```

470 The resulting matrices can then be integrated into an scRNA-seq analysis pipeline.

471

472 **Analysis of human cardiac differentiation scRNA-seq data**

473 The raw data were download from E-MTAB-6268³³. As this data was generated using
474 the Single Cell 3' Library, Gel Bead and Multiplex kit (version 1, 10x Genomics, Cat.
475 #PN-120233), the cell barcode and UMI sequence are not in the same read. First, we
476 merged the cell barcode and UMI sequence into the same read using a custom script,
477 and then aligned the modified fastq file to the hg38 genome using STARsolo, as
478 described above. Cells with less than 500 expressed genes/TEs and cells that have
479 more than 20% fraction of mitochondrial reads were discarded. Single cell trajectory
480 was analyzed by Harmony⁷⁸ and the top 1000 highly variable genes were used for
481 PCA, and the force directed layout was computed using first 150 PCs (principle
482 components). Differentially expressed genes and TEs were analyzed using the
483 SCANPY rank_genes_groups functions by t-test method, the top 500 specifically
484 expressed TEs and genes with Benjamini-Hochberg corrected p-value <0.01 and
485 log₂(fold-change) > 0.5 are selected for downstream analysis.

486

487 **Analysis of the gastrulation scRNA-seq data**

488 The raw data was download from E-MTAB-6967, and aligned to the mm10 genome
489 using STARsolo⁷⁷, with the parameters ‘--readFilesCommand zcat --
490 outFilterMultimapNmax 100 --winAnchorMultimapNmax 100 --outMultimapperOrder
491 Random --runRNGseed 777 --outSAMmultNmax 1’. Cells with less than 3000
492 expressed genes/TEs, and less than 8000 UMIs were discarded. Genes expressed in
493 less than 50 cells were removed from the analysis. The count matrix was normalized
494 using `normalize_total` function of SCANPY, and the top 2000 most highly variable
495 genes were used for PCA, and the first 20 PCs (principle components) were used, as
496 described in the original publication¹⁶. UMAP plots were generated (`min_dist=0.6`).
497 Data is from E-MTAB-6967¹⁶.

498

499 **Analysis of Tabula Muris scRNA-seq data**

500 The C1/Smart-seq2 scRNA-seq raw data was download from GSE109774⁴², the
501 reads were aligned to the mm10 genome using STAR with the parameters ‘--
502 readFilesCommand zcat --outFilterMultimapNmax 100 --winAnchorMultimapNmax
503 100 --outMultimapperOrder Random --runRNGseed 777 --outSAMmultNmax 1’. The
504 genes/TEs and cell expression matrix was generated using scTE. Cells with less than
505 50000 counts or more than 2^7 counts, less than 1000 expressed genes, or more than
506 20% fraction of mitochondrial counts were removed. The filtered matrix was
507 normalized using `scrn`⁷⁹. The top 4000 most highly variable genes were used for PCA,
508 and the first 50 PCs were used for downstream analysis. The cell cluster specific
509 expressed genes/TEs was calculated using SCANPY `rank_genes_groups` functions
510 by t-test method, the top 500 specifically expressed TEs and genes with Benjamini-

511 Hochberg corrected p-value <0.01 and $\log_2(\text{fold-change}) >0.5$ compare to all other
512 groups of cells were kept.

513

514 **Analysis of the OKSM/Chemical reprogramming data**

515 The raw data were download from GSE115943⁵⁰ and GSE114952²⁹. Cells with less
516 than 10000 UMIs or more than 1000000 UMIs, or expressed less than 1000 expressed
517 genes, or more than 20% fraction of mitochondrial counts were removed. The filtered
518 matrices were normalized using scran⁷⁹. The top 4000 most highly variable genes
519 were used for PCA, and the first 50 PCs were used for downstream analysis. The cell
520 trajectory routes were taken from the original studies. Differentially expressed
521 genes/TEs were calculated using SCANPY rank_genes_groups functions by the t-test
522 method, the TEs and genes with Benjamini-Hochberg corrected p-value <0.01 and
523 $\log_2(\text{fold-change}) >0.5$ compared to all other branches of cells were kept.

524

525 **Analysis of the OKS reprogramming data**

526 The C1/SMART-seq data were taken from GSE103221²⁵. the reads were aligned to
527 the mm10 genome using STAR with the parameters '`--readFilesCommand zcat --
528 outFilterMultimapNmax 100 --winAnchorMultimapNmax 100 --outMultimapperOrder
529 Random --runRNGseed 777 --outSAMmultNmax 1`'. The genes/TEs and cells
530 expression matrix was generated using scTE. Cells with less than 10000 counts or
531 more than 2^7 counts, less than 1000 expressed genes, or more than 20% fraction of
532 mitochondrial counts were removed. The filtered matrix was normalized using scran
533 ⁷⁹. The top 4000 most highly variable genes were used for PCA, and the first 50 PCs
534 were used for downstream analysis. The genes/TEs expression trajectories on

535 pseudotemporal orderings of cells (Fig. 5e) were analyzed by LineagePulse
536 (<https://github.com/YosefLab/LineagePulse>) according to the pseudotime taken from
537 the original study.

538

539 **Analysis of the embryonic heart scRNA-seq data**

540 The raw data was download from GSE126128⁴⁰. This data was aligned to the genome
541 using STARsolo⁷⁷, as described above. Cells with less than 3000 expressed
542 genes/TEs and the cells with less than 8000 UMIs or more than 100000 UMIS were
543 deleted from the analysis. The count matrix was normalized using `normalize_total`
544 function of SCANPY. The top 2000 most highly variable genes were used for PCA,
545 and the first 20 PCs were used for downstream analysis. UMA projections were
546 generated (`min_dist=0.7`).

547

548 **Analysis of Alzheimer's disease scRNA-seq data**

549 The MARS-seq scRNA-seq raw data were download from GSE98969⁶⁵. The raw fastq
550 file were modified using custom scripts to embed the cell barcode and UMI in the same
551 read, as in the 10x scRNA-seq format. The modified reads were aligned to the mm10
552 genome with STARsolo as described above. Cells with less than 5000 UMIs or more
553 than 1000000 UMIs, or expressed less than 500 genes, or more than 20% fraction of
554 mitochondrial counts, were removed. The filtered matrix was normalized using `scrn`⁷⁹.
555 The top 4000 most highly variable genes were used for PCA, and the first 50 PCs
556 were used for downstream analysis. The differentially expressed genes and TEs
557 between M2 and M1/3 were analyzed using SCANPY `rank_genes_groups` functions

558 by t-test method, the genes or TEs with Benjamini-Hochberg corrected p-value <0.01
559 and log₂(fold-change) >0.5 compared to each other were kept.

560

561 **Analysis of the Type 2 diabetes/glioblastoma sc-RNA-seq data**

562 The raw data was download from GSE86473⁶⁶, GSE81608⁶⁷. The data was aligned
563 to the hg38 genome using STAR⁷⁷, as described above for C1 data. Cells with less
564 than 5000 expressed genes/TEs and cells with less than 1*10⁶ counts or more than
565 6*10⁶ or were deleted from the analysis. The count matrix was normalized using the
566 normalize_total function of SCANPY. There was a strong batch effect based on the
567 sex of the donor in the type 2 diabetes datasets, this was removed using the
568 regress_out function of SCANPY²⁴. We did not detect any other batch effect from other
569 confounding variables (age, body-mass index, race). The top 2000 most highly
570 variable genes were used for PCA, and the first 15 PCs (type 2 diabetes) or 25 PCs
571 (glioblastoma) were used. UMAP plots were generated using SCANPY (min_dist=0.7).

572

573 **Bulk RNA-seq analysis**

574 Analysis of bulk RNA-seq was performed essentially as previously described^{3,80}, with
575 some modifications. Briefly, reads were aligned to the mouse or human
576 genome/transcriptome (GENCODE transcript annotations, mouse M21 or human 30)
577 using STAR (v2.7.1a)⁷⁷. Tetranscripts⁸¹ or scTE (with the setting -CB False -UMI False)
578 was used to quantitate reads on TEs. Reads were GC normalized using EDASeq
579 (v2.16.3)⁸², and analyzed using glbase⁸³.

580

581 **Motif enrichment analysis**

582 The TF motif enrichment in TEs ([Supplementary Fig. 6c and 8l](#)) was measured using
583 AME from the MEME suite⁸⁴ with the options “--control --shuffle”.

584

585 **Bulk ATAC-seq analysis**

586 Analysis of bulk RNA-seq was performed essentially as previously described^{3,85}.
587 Briefly, reads were aligned to the mouse or human genome (mm10 or hg38) using
588 bowtie2 (v2.3.5.1), with the options: “-p 6 --mm --very-sensitive --no-unal --no-mixed -
589 -no-discordant -X2000”, and reads mapping to TEs were counted using `te_counter`
590 (https://github.com/oaxiom/te_counter). The counts per million (CPM) reads metric
591 was used for enrichment scores.

592

593 **Analysis of the scATAC-seq data**

594 We downloaded the scATAC-seq data from the 10x Illumina website
595 (https://support.10xgenomics.com/single-cell-atac/datasets/1.1.0/atac_pbmc_10k_v1). The
596 barcode was inserted into the read name, so that the mapping could keep track of the
597 cell ID. This yielded reads names inside the FASTQ, such as: (where
598 CCACGTTGTGGACTGA sequence is the cell barcode)

599

```
600 @CCACGTTGTGGACTGA:A00519:269:H7FM2DRXX:1:1101:1325:1000 1:N:0:AAGCATAA
```

601

602 The data was aligned to the human hg38 genome using bowtie2⁸⁶ with the command
603 options “-p 6 --mm --very-sensitive --no-unal --no-mixed --no-discordant -X2000”. The
604 resulting data was then processed using `scTE` with the command:

605

```
606 scTE_scatacseq -i $<in> -x hg38.te.atac.idx -g hg38 -p 1 -UMI False -CB True -o <out>
```

607

608 The genome indices were prebuilt using:

609 `wget -c -O mm10.te.txt.gz 'http://hgdownload.soe.ucsc.edu/goldenPath/mm10/database/rmsk.txt.gz'`

610 `zcat mm10.te.txt.gz | grep -E 'LINE|SINE|LTR|Retroposon' | cut -f6-8,11 >mm10.te.bed`

611 `python3 /share/apps/genomics/unstable/scTE/bin/scTEATAC_build -g mm10.te.bed -o mm10.te.atac`

612

613 `wget -c -O hg38.te.txt.gz 'http://hgdownload.soe.ucsc.edu/goldenPath/hg38/database/rmsk.txt.gz'`

614 `zcat hg38.te.txt.gz | grep -E 'LINE|SINE|LTR|Retroposon' | cut -f6-8,11 >hg38.te.bed`

615 `python3 /share/apps/genomics/unstable/scTE/bin/scTEATAC_build -g hg38.te.bed -o hg38.te.atac`

616

617 References

618 1 Cordaux, R. & Batzer, M. A. The impact of retrotransposons on human genome evolution. *Nat*
619 *Rev Genet* **10**, 691-703, doi:10.1038/nrg2640 (2009).

620 2 Hutchins, A. P. & Pei, D. Transposable elements at the center of the crossroads between
621 embryogenesis, embryonic stem cells, reprogramming, and long non-coding RNAs. *Sci Bull* **60**,
622 1722-1733, doi:10.1007/s11434-015-0905-x (2015).

623 3 He, J. *et al.* Transposable elements are regulated by context-specific patterns of chromatin
624 marks in mouse embryonic stem cells. *Nat Commun* **10**, 34, doi:10.1038/s41467-018-08006-y
625 (2019).

626 4 Lu, X. *et al.* The retrovirus HERVH is a long noncoding RNA required for human embryonic
627 stem cell identity. *Nat Struct Mol Biol* **21**, 423-425, doi:10.1038/nsmb.2799 (2014).

628 5 Bourque, G. *et al.* Ten things you should know about transposable elements. *Genome Biol* **19**,
629 199, doi:10.1186/s13059-018-1577-z (2018).

630 6 Goke, J. *et al.* Dynamic transcription of distinct classes of endogenous retroviral elements
631 marks specific populations of early human embryonic cells. *Cell Stem Cell* **16**, 135-141,
632 doi:10.1016/j.stem.2015.01.005 (2015).

633 7 Grow, E. J. *et al.* Intrinsic retroviral reactivation in human preimplantation embryos and
634 pluripotent cells. *Nature* **522**, 221-225, doi:10.1038/nature14308 (2015).

635 8 Percharde, M. *et al.* A LINE1-Nucleolin Partnership Regulates Early Development and ESC
636 Identity. *Cell* **174**, 391-405.e319, doi:10.1016/j.cell.2018.05.043 (2018).

637 9 Jachowicz, J. W. *et al.* LINE-1 activation after fertilization regulates global chromatin
638 accessibility in the early mouse embryo. *Nat Genet* **49**, 1502-1510, doi:10.1038/ng.3945 (2017).

639 10 Tam, O. H., Ostrow, L. W. & Gale Hammell, M. Diseases of the nERVous system:
640 retrotransposon activity in neurodegenerative disease. *Mobile DNA* **10**, 32,
641 doi:10.1186/s13100-019-0176-1 (2019).

642 11 Jang, H. S. *et al.* Transposable elements drive widespread expression of oncogenes in human
643 cancers. *Nat Genet* **51**, 611-617, doi:10.1038/s41588-019-0373-3 (2019).

- 644 12 Chuong, E. B., Elde, N. C. & Feschotte, C. Regulatory activities of transposable elements: from
645 conflicts to benefits. *Nat Rev Genet* **18**, 71-86, doi:10.1038/nrg.2016.139 (2017).
- 646 13 Feschotte, C. Transposable elements and the evolution of regulatory networks. *Nat Rev Genet*
647 **9**, 397-405, doi:10.1038/nrg2337 (2008).
- 648 14 Kunarso, G. *et al.* Transposable elements have rewired the core regulatory network of human
649 embryonic stem cells. *Nat Genet* **42**, 631-634, doi:10.1038/ng.600 (2010).
- 650 15 Faulkner, G. J. *et al.* The regulated retrotransposon transcriptome of mammalian cells. *Nat*
651 *Genet* **41**, 563-571, doi:10.1038/ng.368 (2009).
- 652 16 Pijuan-Sala, B. *et al.* A single-cell molecular map of mouse gastrulation and early
653 organogenesis. *Nature* **566**, 490-495, doi:10.1038/s41586-019-0933-9 (2019).
- 654 17 Li, H. *et al.* Dysfunctional CD8 T Cells Form a Proliferative, Dynamically Regulated
655 Compartment within Human Melanoma. *Cell* **176**, 775-789 e718,
656 doi:10.1016/j.cell.2018.11.043 (2019).
- 657 18 Tirosh, I. *et al.* Dissecting the multicellular ecosystem of metastatic melanoma by single-cell
658 RNA-seq. *Science (New York, N.Y.)* **352**, 189-196, doi:10.1126/science.aad0501 (2016).
- 659 19 Macfarlan, T. S. *et al.* Embryonic stem cell potency fluctuates with endogenous retrovirus
660 activity. *Nature* **487**, 57-63, doi:10.1038/nature11244 (2012).
- 661 20 Fu, X., Wu, X., Djekidel, M. N. & Zhang, Y. Myc and Dnmt1 impede the pluripotent to totipotent
662 state transition in embryonic stem cells. *Nat Cell Biol* **21**, 835-844, doi:10.1038/s41556-019-
663 0343-0 (2019).
- 664 21 Treangen, T. J. & Salzberg, S. L. Repetitive DNA and next-generation sequencing:
665 computational challenges and solutions. *Nat Rev Genet* **13**, 36-46, doi:10.1038/nrg3117 (2011).
- 666 22 Goerner-Potvin, P. & Bourque, G. Computational tools to unmask transposable elements. *Nat*
667 *Rev Genet* **19**, 688-704, doi:10.1038/s41576-018-0050-x (2018).
- 668 23 Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-1902 e1821,
669 doi:10.1016/j.cell.2019.05.031 (2019).
- 670 24 Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression data
671 analysis. *Genome Biol* **19**, 15, doi:10.1186/s13059-017-1382-0 (2018).
- 672 25 Guo, L. *et al.* Resolving Cell Fate Decisions during Somatic Cell Reprogramming by Single-
673 Cell RNA-Seq. *Mol Cell* **73**, 815-829 e817, doi:10.1016/j.molcel.2019.01.042 (2019).
- 674 26 Zheng, G. X. *et al.* Massively parallel digital transcriptional profiling of single cells. *Nat Commun*
675 **8**, 14049, doi:10.1038/ncomms14049 (2017).
- 676 27 Garcia-Perez, J. L., Widmann, T. J. & Adams, I. R. The impact of transposable elements on
677 mammalian development. *Development* **143**, 4101-4114, doi:10.1242/dev.132639 (2016).
- 678 28 Rodriguez-Terrones, D. & Torres-Padilla, M. E. Nimble and Ready to Mingle: Transposon
679 Outbursts of Early Development. *Trends in genetics : TIG* **34**, 806-820,
680 doi:10.1016/j.tig.2018.06.006 (2018).
- 681 29 Zhao, T. *et al.* Single-Cell RNA-Seq Reveals Dynamic Early Embryonic-like Programs during
682 Chemical Reprogramming. *Cell Stem Cell* **23**, 31-45 e37, doi:10.1016/j.stem.2018.05.025
683 (2018).

- 684 30 Wang, J. *et al.* Primate-specific endogenous retrovirus-driven transcription defines naive-like
685 stem cells. *Nature* **516**, 405-409, doi:10.1038/nature13804 (2014).
- 686 31 Fort, A. *et al.* Deep transcriptome profiling of mammalian stem cells supports a regulatory role
687 for retrotransposons in pluripotency maintenance. *Nat Genet* **46**, 558-566, doi:10.1038/ng.2965
688 (2014).
- 689 32 Theunissen, T. W. *et al.* Molecular Criteria for Defining the Naive Human Pluripotent State. *Cell*
690 *Stem Cell* **19**, 502-515, doi:10.1016/j.stem.2016.06.011 (2016).
- 691 33 Friedman, C. E. *et al.* Single-Cell Transcriptomic Analysis of Cardiac Differentiation from
692 Human PSCs Reveals HOPX-Dependent Cardiomyocyte Maturation. *Cell Stem Cell* **23**, 586-
693 598 e588, doi:10.1016/j.stem.2018.09.009 (2018).
- 694 34 Liu, Q. *et al.* Genome-Wide Temporal Profiling of Transcriptome and Open Chromatin of Early
695 Cardiomyocyte Differentiation Derived From hiPSCs and hESCs. *Circ Res* **121**, 376-391,
696 doi:10.1161/CIRCRESAHA.116.310456 (2017).
- 697 35 Abed, M. *et al.* The Gag protein PEG10 binds to RNA and regulates trophoblast stem cell
698 lineage specification. *PLoS One* **14**, e0214110, doi:10.1371/journal.pone.0214110 (2019).
- 699 36 Zhong, Y. *et al.* Isolation of primitive mouse extraembryonic endoderm (pXEN) stem cell lines.
700 *Stem Cell Res* **30**, 100-112, doi:10.1016/j.scr.2018.05.008 (2018).
- 701 37 Factor, D. C. *et al.* Epigenomic comparison reveals activation of "seed" enhancers during
702 transition from naive to primed pluripotency. *Cell Stem Cell* **14**, 854-863,
703 doi:10.1016/j.stem.2014.05.005 (2014).
- 704 38 Morey, L. *et al.* Polycomb Regulates Mesoderm Cell Fate-Specification in Embryonic Stem
705 Cells through Activation and Repression Mechanisms. *Cell Stem Cell* **17**, 300-315,
706 doi:10.1016/j.stem.2015.08.009 (2015).
- 707 39 Merkin, J., Russell, C., Chen, P. & Burge, C. B. Evolutionary dynamics of gene and isoform
708 regulation in Mammalian tissues. *Science (New York, N.Y.)* **338**, 1593-1599,
709 doi:10.1126/science.1228186 (2012).
- 710 40 de Soysa, T. Y. *et al.* Single-cell analysis of cardiogenesis reveals basis for organ-level
711 developmental defects. *Nature* **572**, 120-124, doi:10.1038/s41586-019-1414-x (2019).
- 712 41 Richardson, S. R., Morell, S. & Faulkner, G. J. L1 retrotransposons and somatic mosaicism in
713 the brain. *Annu Rev Genet* **48**, 1-27, doi:10.1146/annurev-genet-120213-092412 (2014).
- 714 42 Tabula Muris, C. *et al.* Single-cell transcriptomics of 20 mouse organs creates a Tabula Muris.
715 *Nature* **562**, 367-372, doi:10.1038/s41586-018-0590-4 (2018).
- 716 43 Chuong, E. B., Elde, N. C. & Feschotte, C. Regulatory evolution of innate immunity through co-
717 option of endogenous retroviruses. *Science (New York, N.Y.)* **351**, 1083-1087,
718 doi:10.1126/science.aad5497 (2016).
- 719 44 Koonin, E. V. & Krupovic, M. Evolution of adaptive immunity from transposable elements
720 combined with innate immune systems. *Nat Rev Genet* **16**, 184-192, doi:10.1038/nrg3859
721 (2015).
- 722 45 Erwin, J. A., Marchetto, M. C. & Gage, F. H. Mobile DNA elements in the generation of diversity
723 and complexity in the brain. *Nat Rev Neurosci* **15**, 497-506, doi:10.1038/nrn3730 (2014).

- 724 46 Takahashi, K. & Yamanaka, S. Induction of pluripotent stem cells from mouse embryonic and
725 adult fibroblast cultures by defined factors. *Cell* **126**, 663-676, doi:10.1016/j.cell.2006.07.024
726 (2006).
- 727 47 Wang, B. *et al.* Induction of Pluripotent Stem Cells from Mouse Embryonic Fibroblasts by Jdp2-
728 Jhdm1b-Mkk6-Glis1-Nanog-Essrb-Sall4. *Cell Rep* **27**, 3473-3485 e3475,
729 doi:10.1016/j.celrep.2019.05.068 (2019).
- 730 48 Hou, P. *et al.* Pluripotent stem cells induced from mouse somatic cells by small-molecule
731 compounds. *Science (New York, N.Y.)* **341**, 651-654, doi:10.1126/science.1239278 (2013).
- 732 49 Cao, S. *et al.* Chromatin Accessibility Dynamics during Chemical Induction of Pluripotency. *Cell*
733 *Stem Cell* **22**, 529-542 e525, doi:10.1016/j.stem.2018.03.005 (2018).
- 734 50 Schiebinger, G. *et al.* Optimal-Transport Analysis of Single-Cell Gene Expression Identifies
735 Developmental Trajectories in Reprogramming. *Cell* **176**, 928-943 e922,
736 doi:10.1016/j.cell.2019.01.006 (2019).
- 737 51 Friedli, M. *et al.* Loss of transcriptional control over endogenous retroelements during
738 reprogramming to pluripotency. *Genome Res* **24**, 1251-1259, doi:10.1101/gr.172809.114
739 (2014).
- 740 52 Liu, J. *et al.* The oncogene c-Jun impedes somatic cell reprogramming. *Nat Cell Biol* **17**, 856-
741 867, doi:10.1038/ncb3193 (2015).
- 742 53 Chronis, C. *et al.* Cooperative Binding of Transcription Factors Orchestrates Reprogramming.
743 *Cell* **168**, 442-459 e420, doi:10.1016/j.cell.2016.12.016 (2017).
- 744 54 Buenrostro, J. D. *et al.* Single-cell chromatin accessibility reveals principles of regulatory
745 variation. *Nature* **523**, 486-490, doi:10.1038/nature14590 (2015).
- 746 55 Shema, E., Bernstein, B. E. & Buenrostro, J. D. Single-cell and single-molecule epigenomics
747 to uncover genome regulation at unprecedented resolution. *Nat Genet* **51**, 19-25,
748 doi:10.1038/s41588-018-0290-x (2019).
- 749 56 Cusanovich, D. A. *et al.* A Single-Cell Atlas of In Vivo Mammalian Chromatin Accessibility. *Cell*
750 **174**, 1309-1324 e1318, doi:10.1016/j.cell.2018.06.052 (2018).
- 751 57 Sun, X. *et al.* Transcription factor profiling reveals molecular choreography and key regulators
752 of human retrotransposon expression. *Proc Natl Acad Sci U S A* **115**, E5526-E5535,
753 doi:10.1073/pnas.1722565115 (2018).
- 754 58 Chen, X., Miragaia, R. J., Natarajan, K. N. & Teichmann, S. A. A rapid and robust method for
755 single cell chromatin accessibility profiling. *Nat Commun* **9**, 5345, doi:10.1038/s41467-018-
756 07771-0 (2018).
- 757 59 Warren, J. S. *et al.* Histone methyltransferase Smyd1 regulates mitochondrial energetics in the
758 heart. *Proc Natl Acad Sci U S A* **115**, E7871-E7880, doi:10.1073/pnas.1800680115 (2018).
- 759 60 Rasmussen, T. L. *et al.* Smyd1 facilitates heart development by antagonizing oxidative and ER
760 stress responses. *PLoS One* **10**, e0121765, doi:10.1371/journal.pone.0121765 (2015).
- 761 61 Franklin, S. *et al.* The chromatin-binding protein Smyd1 restricts adult mammalian heart growth.
762 *Am J Physiol Heart Circ Physiol* **311**, H1234-H1247, doi:10.1152/ajpheart.00235.2016 (2016).
- 763 62 Calderon, D. *et al.* Landscape of stimulation-responsive chromatin across diverse human
764 immune cells. *Nat Genet*, doi:10.1038/s41588-019-0505-9 (2019).

- 765 63 Bravo Gonzalez-Blas, C. *et al.* cisTopic: cis-regulatory topic modeling on single-cell ATAC-seq
766 data. *Nat Methods* **16**, 397-400, doi:10.1038/s41592-019-0367-1 (2019).
- 767 64 Li, W. *et al.* Activation of transposable elements during aging and neuronal decline in
768 *Drosophila*. *Nat Neurosci* **16**, 529-531, doi:10.1038/nn.3368 (2013).
- 769 65 Keren-Shaul, H. *et al.* A Unique Microglia Type Associated with Restricting Development of
770 Alzheimer's Disease. *Cell* **169**, 1276-1290 e1217, doi:10.1016/j.cell.2017.05.018 (2017).
- 771 66 Lawlor, N. *et al.* Single-cell transcriptomes identify human islet cell signatures and reveal cell-
772 type-specific expression changes in type 2 diabetes. *Genome Res* **27**, 208-222,
773 doi:10.1101/gr.212720.116 (2017).
- 774 67 Xin, Y. *et al.* RNA Sequencing of Single Human Islet Cells Reveals Type 2 Diabetes Genes.
775 *Cell Metab* **24**, 608-615, doi:10.1016/j.cmet.2016.08.018 (2016).
- 776 68 Fadista, J. *et al.* Global genomic and transcriptomic analysis of human pancreatic islets reveals
777 novel genes influencing glucose metabolism. *Proc Natl Acad Sci U S A* **111**, 13924-13929,
778 doi:10.1073/pnas.1402665111 (2014).
- 779 69 Darmanis, S. *et al.* Single-Cell RNA-Seq Analysis of Infiltrating Neoplastic Cells at the Migrating
780 Front of Human Glioblastoma. *Cell Rep* **21**, 1399-1410, doi:10.1016/j.celrep.2017.10.030
781 (2017).
- 782 70 Chuong, E. B., Rumi, M. A., Soares, M. J. & Baker, J. C. Endogenous retroviruses function as
783 species-specific enhancer elements in the placenta. *Nat Genet* **45**, 325-329,
784 doi:10.1038/ng.2553 (2013).
- 785 71 Payer, L. M. & Burns, K. H. Transposable elements in human genetic disease. *Nat Rev Genet*,
786 doi:10.1038/s41576-019-0165-8 (2019).
- 787 72 Burns, K. H. Transposable elements in cancer. *Nat Rev Cancer* **17**, 415-424,
788 doi:10.1038/nrc.2017.35 (2017).
- 789 73 De Cecco, M. *et al.* L1 drives IFN in senescent cells and promotes age-associated inflammation.
790 *Nature* **566**, 73-78, doi:10.1038/s41586-018-0784-9 (2019).
- 791 74 Li, H. *et al.* The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**, 2078-2079,
792 doi:10.1093/bioinformatics/btp352 (2009).
- 793 75 Frankish, A. *et al.* GENCODE reference annotation for the human and mouse genomes.
794 *Nucleic Acids Res* **47**, D766-D773, doi:10.1093/nar/gky955 (2019).
- 795 76 Jurka, J. Repbase update: a database and an electronic journal of repetitive elements. *Trends*
796 *in genetics : TIG* **16**, 418-420, doi:10.1016/s0168-9525(00)02093-x (2000).
- 797 77 Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21,
798 doi:10.1093/bioinformatics/bts635 (2013).
- 799 78 Nowotschin, S. *et al.* The emergent landscape of the mouse gut endoderm at single-cell
800 resolution. *Nature* **569**, 361-367, doi:10.1038/s41586-019-1127-1 (2019).
- 801 79 Lun, A. T., McCarthy, D. J. & Marioni, J. C. A step-by-step workflow for low-level analysis of
802 single-cell RNA-seq data with Bioconductor. *F1000Res* **5**, 2122,
803 doi:10.12688/f1000research.9501.2 (2016).
- 804 80 Hutchins, A. P. *et al.* Models of global gene expression define major domains of cell type and
805 tissue identity. *Nucleic Acids Res* **45**, 2354-2367, doi:10.1093/nar/gkx054 (2017).

- 806 81 Jin, Y., Tam, O. H., Paniagua, E. & Hammell, M. TEtranscripts: a package for including
807 transposable elements in differential expression analysis of RNA-seq datasets. *Bioinformatics*
808 **31**, 3593-3599, doi:10.1093/bioinformatics/btv422 (2015).
- 809 82 Risso, D., Schwartz, K., Sherlock, G. & Dudoit, S. GC-content normalization for RNA-Seq data.
810 *BMC Bioinformatics* **12**, 480, doi:10.1186/1471-2105-12-480 (2011).
- 811 83 Hutchins, A. P., Jauch, R., Dyla, M. & Miranda-Saavedra, D. glbase: a framework for combining,
812 analyzing and displaying heterogeneous genomic and high-throughput sequencing data. *Cell*
813 *Regen (Lond)* **3**, 1, doi:10.1186/2045-9769-3-1 (2014).
- 814 84 McLeay, R. C. & Bailey, T. L. Motif Enrichment Analysis: a unified framework and an evaluation
815 on ChIP data. *BMC Bioinformatics* **11**, 165, doi:10.1186/1471-2105-11-165 (2010).
- 816 85 Li, D. *et al.* Chromatin Accessibility Dynamics during iPSC Reprogramming. *Cell Stem Cell* **21**,
817 819-833 e816, doi:10.1016/j.stem.2017.10.012 (2017).
- 818 86 Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**,
819 357-359, doi:10.1038/nmeth.1923 (2012).

820

821 **Acknowledgements**

822 We grateful to Rujin Huang for the helpful discussions and advice. We appreciate the
823 assistance of Lihui Lin, Huijian Feng and Yuanbang Mai on the data analysis. We
824 thank the GuangZhou Branch of the Supercomputing Center of Chinese Academy of
825 Science for its support. We thank the support from the Center for Computational
826 Science and Engineering of Southern University of Science and Technology. This
827 work was supported by the National Natural Science Foundation of China (31970589,
828 31801217, 31850410463, 31850410486), Guangdong Science and Technology
829 Commission (2019A050510004), and the Shenzhen Peacock plan (201701090668B).

830

831 **Author contributions**

832 J.H., A.P.H., and J.C. initiated the project and wrote the manuscript. J.H. and A.P.H.
833 performed the bioinformatic analysis with the assistance of all other authors. A.P.H.
834 and J.C. supervised and funded the project.

835

836 **Competing interests:**

837 The authors declare no competing interests

Fig. 1

bioRxiv preprint doi: <https://doi.org/10.1101/2020.07.23.218800>; this version posted July 24, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under a [CC-BY-NC-ND 4.0 International license](https://creativecommons.org/licenses/by-nc-nd/4.0/).

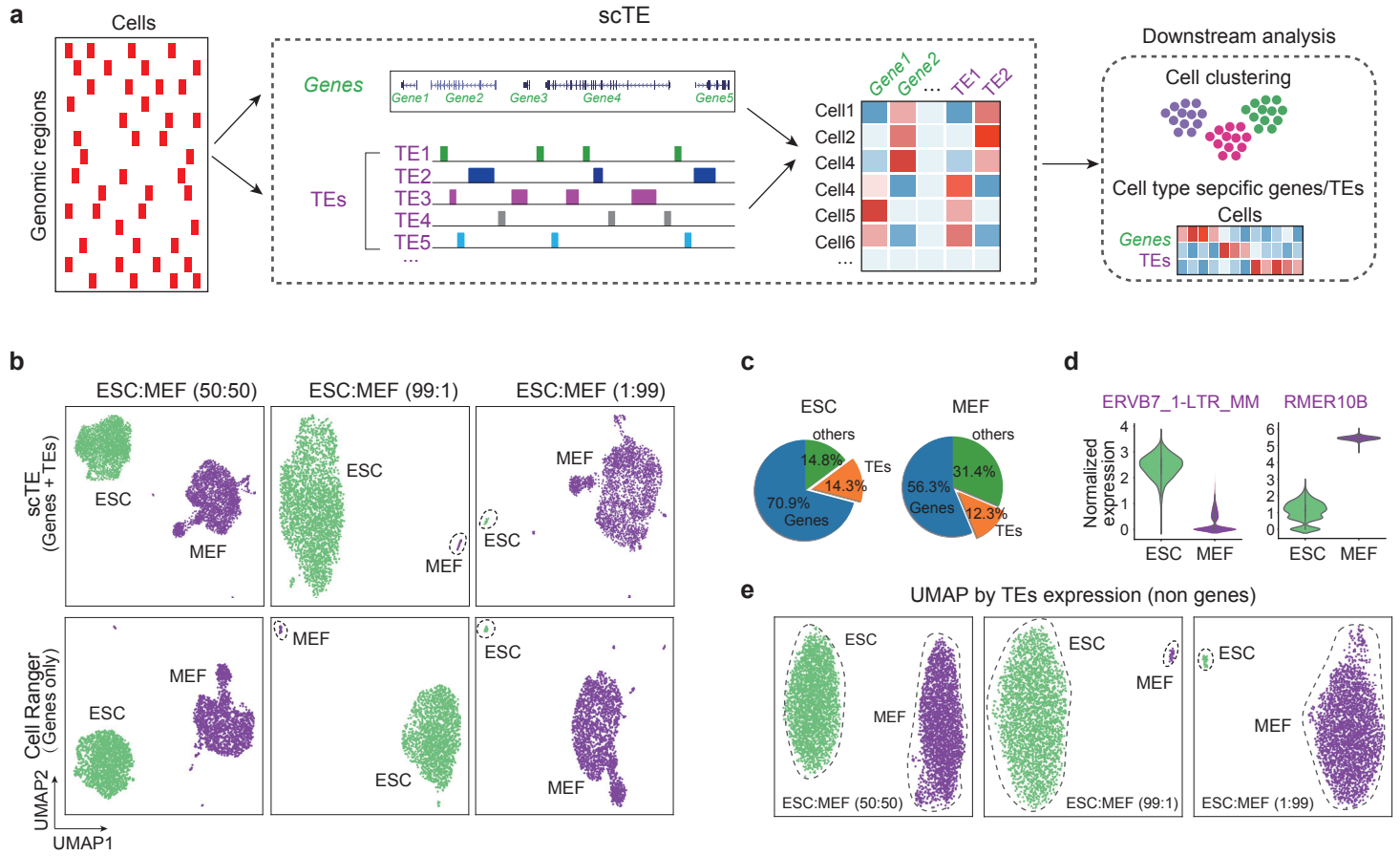


Fig. 1 | scTE workflow and applications. (a) Schematic of the workings of scTE. For scRNA-seq data the reads are mapped to the genome, and assigned to either a gene, or a metagene model of a TE. Multimapping read data will assign the best mapping read to a type of TE. Reads are always mapped to a gene first, and then a TE if no gene is found. The resulting assignments are then collapsed into a matrix of read counts for each cell, versus each gene/TE. This matrix can be used in downstream applications. (b) UMAP plot showing mixtures of MEFs and ESCs in the indicated ratios. The top panels show scTE analysis, the lower panels show Cell Ranger analysis results. Cells are colored by their sample of origin. (c) Percentage of reads mapping to genes, TEs or other regions of the genome in MEFs and ESCs. (d) Violin plot showing the expression of selected TEs in MEFs and ESCs. (e) As in panel b, but only TE expression was used.

Fig. 2

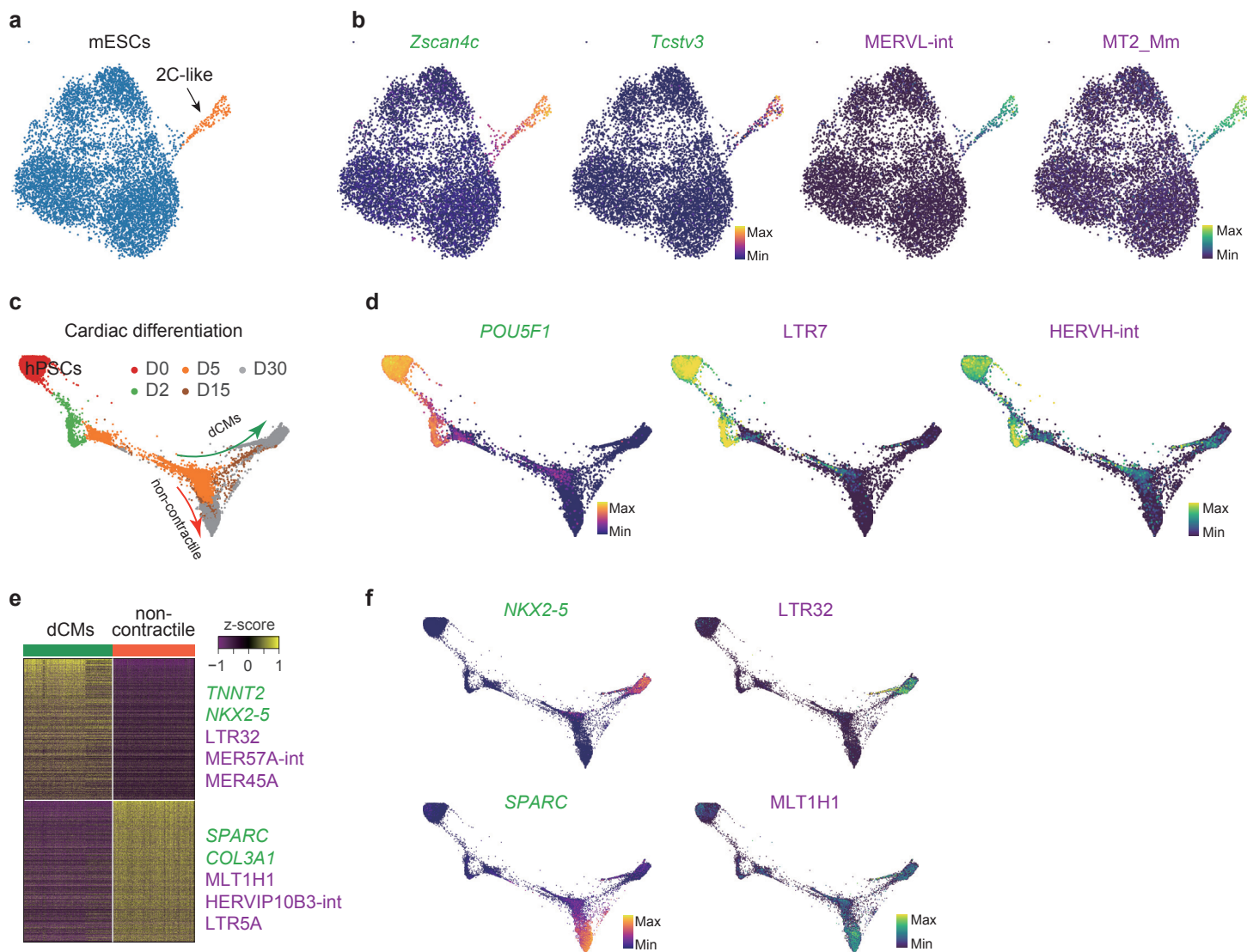


Fig. 2 | Dynamic transcription of TEs in ESCs and during cardiac differentiation. (a) UMAP plot of mouse ESCs. Cells are colored by cell type cluster. (b) Same as panel a, but cells are colored based on the expression of the indicated genes and TEs. *Zscan4c* and *Tcstv3* are marker genes for the 2C-like cells. (c) Trajectory reconstruction of single cells through a cardiac differentiation timecourse showing the definitive cardiomyocytes (dCMs) branch and non-contractile branch. Days of differentiation (D) are labelled. (d) As in panel c, but cells are colored by the expression of the indicated genes and TEs. (e) Heatmap of expression differences between dCM (contractile) branch and non-contractile branch cells, selected differentially expressed genes and TEs are labelled. (f) As in panel d, but cells are colored by the expression level of the indicated genes and TEs.

Fig. 3

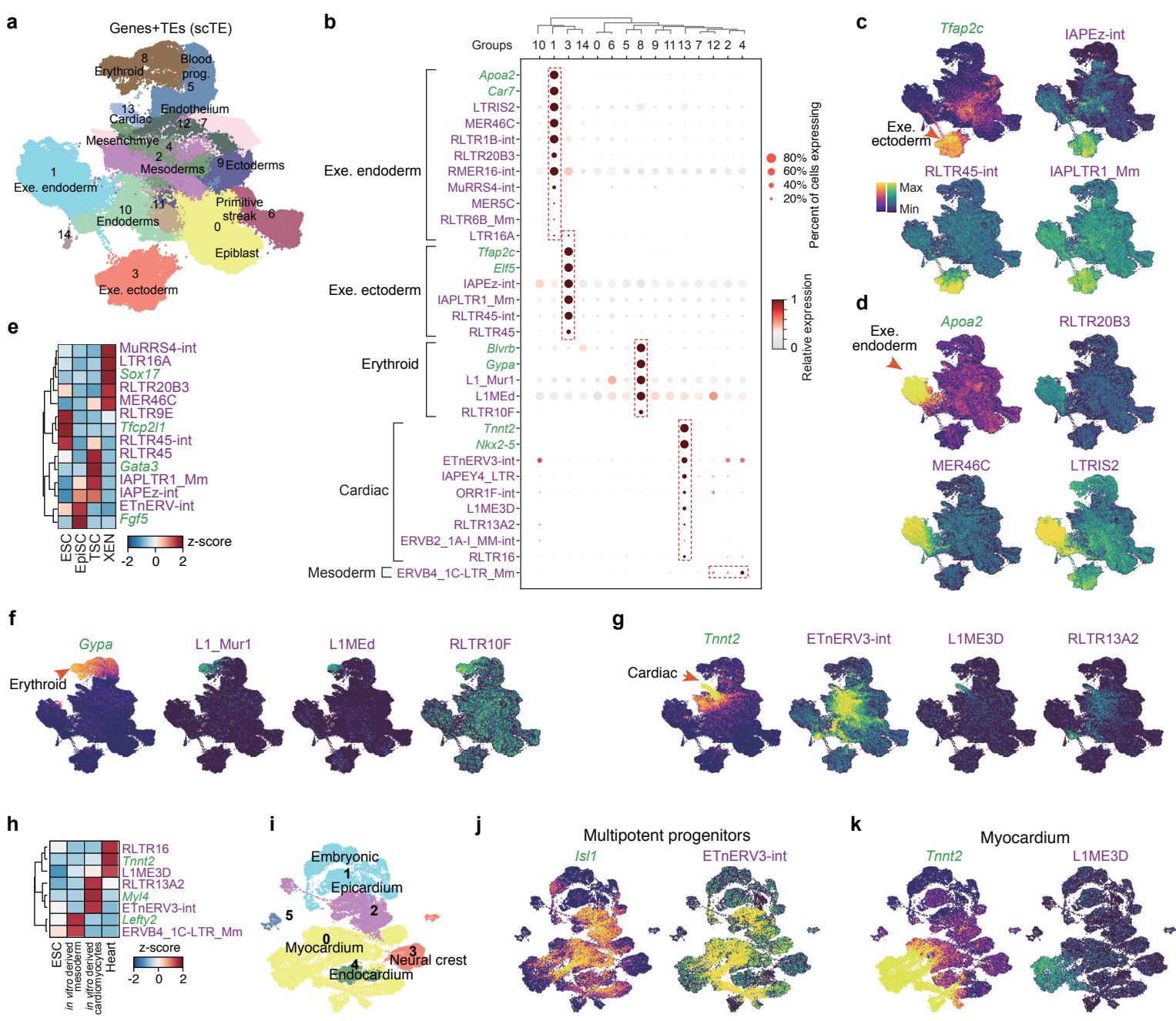


Fig. 3 | Widespread cell type-specific expression of TEs during gastrulation. (a) UMAP plots of the mouse gastrulation data using both genes and TEs. Selected lineages are labelled (Leiden, resolution=0.3). (b) Dot plot showing a selection of marker genes and TEs for the indicated cell lineages. (c) Expression of the indicated extra embryonic ectoderm gene *Tfap2c* and selected TEs. (d) Expression of the extra embryonic endoderm marker gene *Apoa2* and selected TEs. (e) Expression of the indicated TEs and marker genes in bulk RNA-seq data from ESCs, EpiSCs, XEN (extra embryonic endoderm cells) and TSCs (trophoblast stem cells). *Tfcp2l1*, *Fgf5*, *Gata3* and *Sox17* serve as markers for ESCs, EpiSCs, TSCs, and XEN cells, respectively. Data is displayed as a z-score using the variance from all genes. (f) Expression of the erythroid marker gene *Gypa*, and selected TEs. (g) Expression of the cardiac marker gene *Tnnt2* and selected TEs. (h) Expression of the indicated TEs and marker genes from bulk RNA-seq data. (i) UMAP plot of the embryonic mouse heart scRNA-seq data using both TEs and genes. The indicated developmental stages are labelled as in the original study. (j-k) UMAP as panel i, but cells are colored by the expression of indicated genes/TEs.

Fig. 4

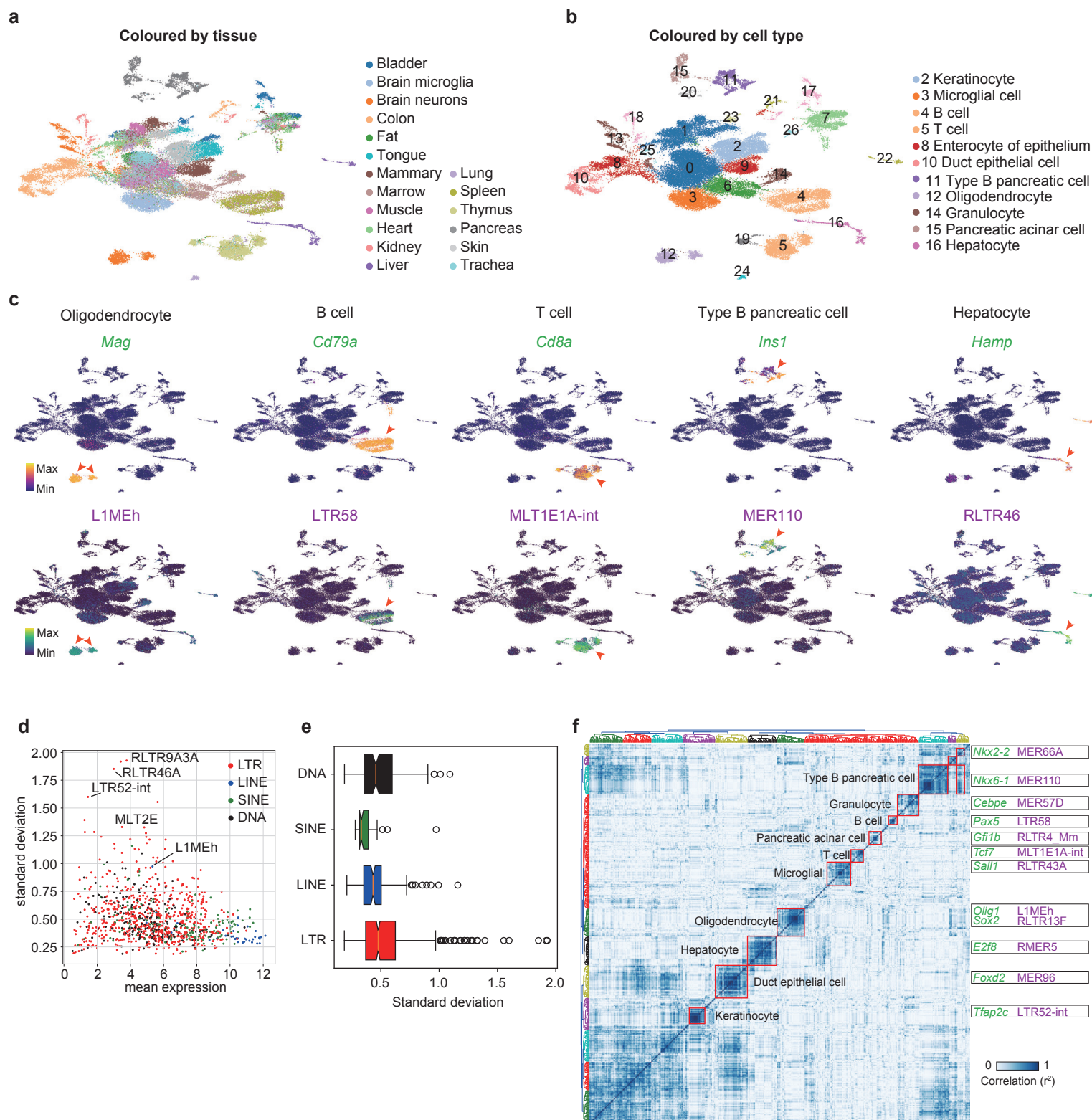


Fig. 4 | Class-specific expression of TEs in somatic cells. (a) UMAP plots of the Tabula Muris data, using both genes and TEs as analyzed with scTE. The tissue sources for the cells is indicated. (b) UMAP plot as in panel a, but clustered into groups (Leiden, resolution=0.5). (c) Same as panel b, but cells are colored by the expression of indicated genes/TEs. (d) Scatter plot showing TE expression heterogeneity. The x-axis is the mean expression for cells from panel b, the y-axis is the standard deviation for each TE type, the higher standard deviation represents higher heterogeneity across cell types. (e) Boxplot for the standard deviations for each class of TEs. (f) Correlation heatmap showing the co-expression of TFs and TEs.

Fig. 5

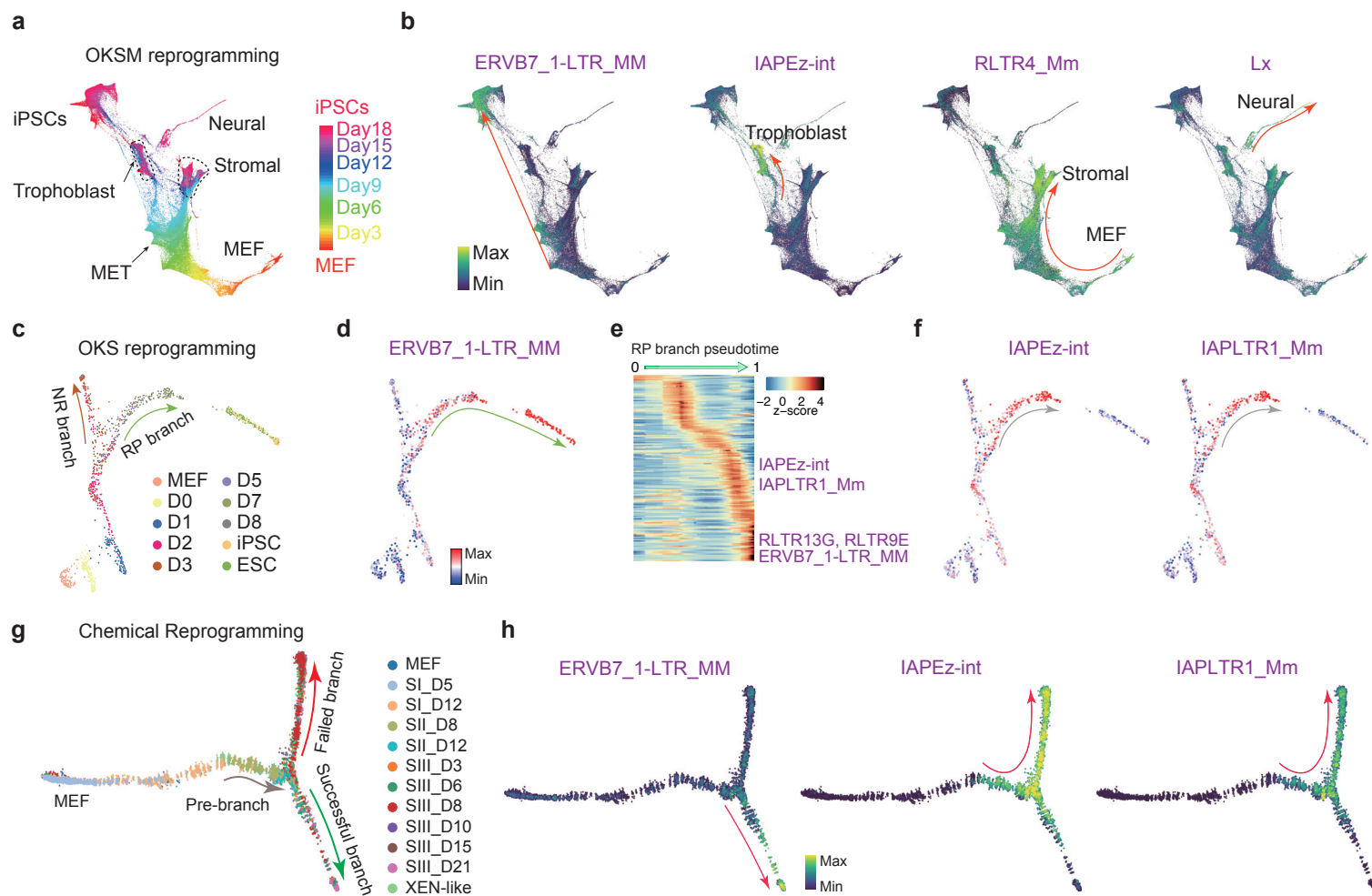


Fig. 5 | Stage-specific expression of TEs in somatic cell reprogramming. (a) Trajectory reconstruction during OKSM reprogramming, cells are colored by time point. (b) As in panel a, but cells are colored by the expression of the indicated TEs. (c) Force-directed (FR) layout of cells during OKS reprogramming, cells are colored by time point. (d) Same with panel c, but cells are colored by the expression change of the ERVB7_1-LTR_MM TE during reprogramming. (e) Expression heatmap of the top 145 dynamically expressed TEs in a pseudotime ordering for the RP branch, selected TEs are indicated. (f) Expression changes of the indicated TEs during reprogramming. (g) Trajectory reconstruction during chemical reprogramming, cells are colored by time point. (h) As in panel g, but showing TE expression specific to the successful or failed branches of reprogramming.

Fig. 6 | Analysis of the Chromatin State of TEs in Single-Cell ATAC-seq data. (a) Schematic plot of scTE for scATAC-seq data analysis. The reads are mapped to the genome, and assigned to a metagene TE, and then the cells were clustered based on the TE matrix. (b) UMAP plot of the TE chromatin state from scATAC-seq data for a selection of FACS-purified mouse cell types. (c) Heatmap of the top 50 cell type-specific opened TEs in the indicated cell types, selected example TEs are indicated. (d) UMAP plot as in panel b, but cells are colored by chromatin-state of the indicated TEs. (e) Genome tracks showing the aggregate scATAC-seq profiles (top panel). Randomly selected 100 single cell profiles are show below the aggregated profiles (bottom panel). With include (unique + multiple) or exclude (unique) multiple mapped reads. (f) UMAP plot of the expression of the myocardium marker gene *Smyd1*, from the cardiogenesis data, see [Fig. 3i](#).

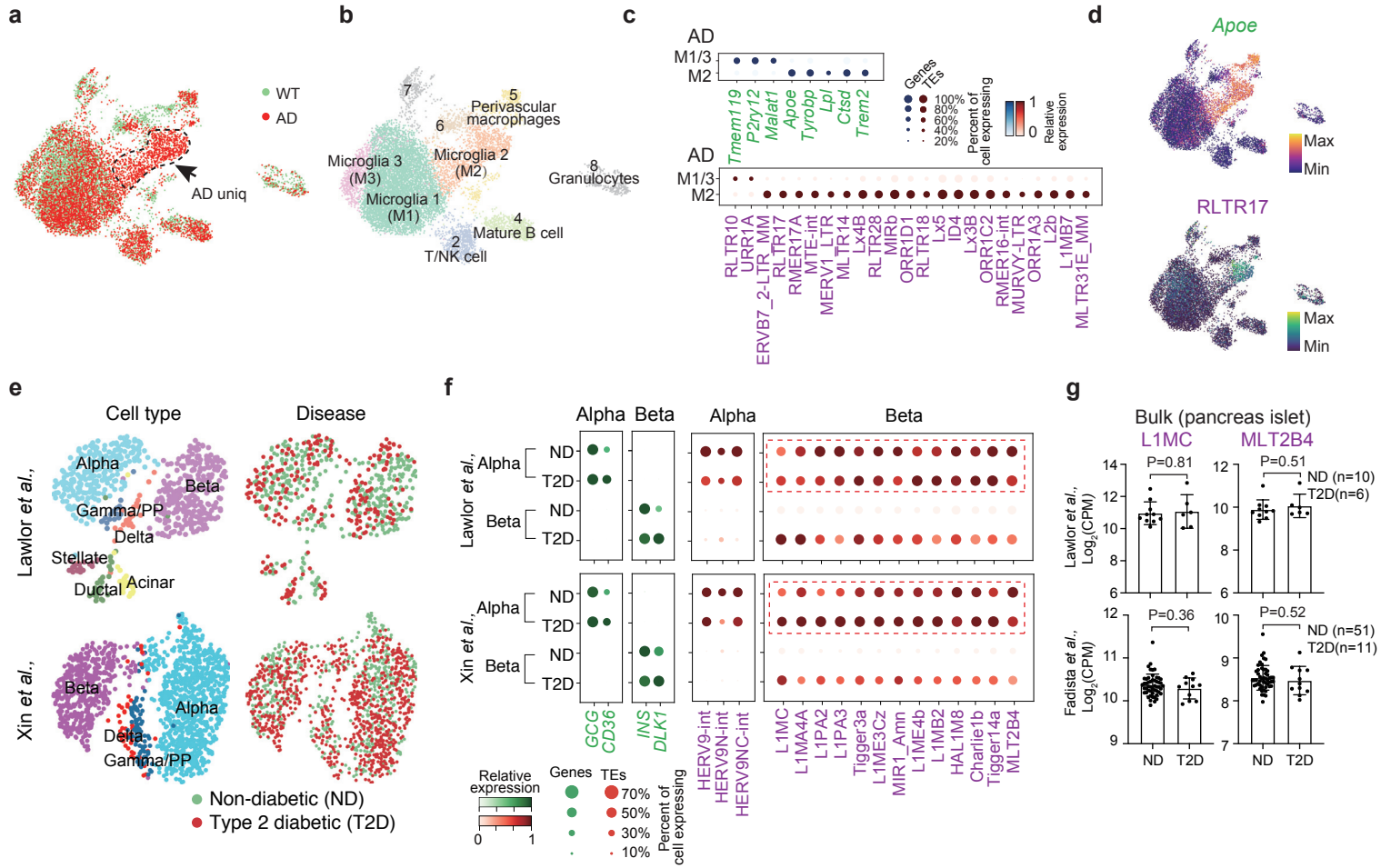


Fig. 7 | TEs are differentially expressed in single cells in the diseased state. (a) UMAP plot of the single cells genes and TE expression, cells are colored by WT (wild-type) and AD (Alzheimer's disease) state. (b) UMAP plot, as in panel a, but clustered into groups (Leiden, resolution=0.5). (c) Dot plot showing the differential expressed genes (top) and TEs (bottom) between disease associated microglia (M2) and homeostatic microglia (M1/3) in AD mice. (d) UMAP plot, as in panel a, but cells are colored by the expression of the indicated *Apoa2* or the TE RLTR17. (e) UMAP plots of pancreatic islet cells. Cells are colored by cell types (left) or disease-state (right). Cell types were annotated according to the metadata from the original study, and matched the expression of known marker genes. (f) Dot plot showing marker gene expression (green) or TEs (red) differentially expressed between healthy and T2D alpha and beta cells (Benjamini-Hochberg corrected Wilcoxon rank-sum test, $P < 0.01$, and at least >2-fold change between groups). (g) Bar charts showing the expression of the indicated TEs from bulk RNA-seq data. P-value was from an unpaired t-test.