

Nonparametric Interrogation of Transcriptional Regulation in Single-Cell RNA and Chromatin Accessibility Multiomic Data

Yuriko Harigaya

University of North Carolina

Zhaojun Zhang

University of Pennsylvania

Hongpan Zhang

University of Virginia

Chongzhi Zang

University of Virginia

Nancy Zhang (✉ nzh@wharton.upenn.edu)

University of Pennsylvania <https://orcid.org/0000-0002-0880-5749>

Yuchao Jiang

University of North Carolina

Article

Keywords: single-cell multiomics, transcriptional regulation, transcription factor, chromatin accessibility

Posted Date: September 28th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-930184/v1>

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

[Read Full License](#)

Nonparametric Interrogation of Transcriptional Regulation in Single-Cell RNA and Chromatin Accessibility Multiomic Data

Yuriko Harigaya¹, Zhaojun Zhang², Hongpan Zhang^{3,4}, Chongzhi Zang^{3,4,5}, Nancy R Zhang^{2,*}, Yuchao Jiang^{6,7,8,*}

- ¹ Curriculum in Bioinformatics and Computational Biology, School of Medicine, University of North Carolina, Chapel Hill, NC 27599, USA.
- ² Department of Statistics, The Wharton School, University of Pennsylvania, Philadelphia, PA 19104, USA.
- ³ Center for Public Health Genomics, University of Virginia, Charlottesville, VA 22908, USA.
- ⁴ Department of Biochemistry and Molecular Genetics, University of Virginia, Charlottesville, VA 22908, USA.
- ⁵ Department of Public Health Sciences, University of Virginia, Charlottesville, VA 22908, USA.
- ⁶ Lineberger Comprehensive Cancer Center, University of North Carolina at Chapel Hill, NC 27599, USA.
- ⁷ Department of Genetics, School of Medicine, University of North Carolina, Chapel Hill, NC 27599, USA.
- ⁸ Department of Biostatistics, Gillings School of Global Public Health, University of North Carolina, Chapel Hill, NC 27599, USA.

* To whom correspondence should be addressed: nzh@wharton.upenn.edu,
yuchaoj@email.unc.edu.

1 **Abstract**

2 Epigenetic control of gene expression is highly cell-type- and context-specific. Yet,
3 despite its complexity, gene regulatory logic can be broken down into modular
4 components consisting of a transcription factor (TF) activating or repressing the
5 expression of a target gene through its binding to a *cis*-regulatory region. Recent
6 advances in joint profiling of transcription and chromatin accessibility with single-cell
7 resolution offer unprecedented opportunities to interrogate such regulatory logic. Here,
8 we propose a nonparametric approach, TRIPOD, to detect and characterize three-way
9 relationships between a TF, its target gene, and the accessibility of the TF's binding site,
10 using single-cell RNA and ATAC multiomic data. We apply TRIPOD to interrogate cell-
11 type-specific regulatory logic in peripheral blood mononuclear cells and contrast our
12 results to detections from enhancer databases, *cis*-eQTL studies, ChIP-seq experiments,
13 and TF knockdown/knockout studies. We then apply TRIPOD to mouse embryonic brain
14 data during neurogenesis and gliogenesis and identified known and novel putative
15 regulatory relationships, validated by ChIP-seq and PLAC-seq. Finally, we demonstrate
16 TRIPOD on SHARE-seq data of differentiating mouse hair follicle cells and identify
17 lineage-specific regulation supported by histone marks for gene activation and super-
18 enhancer annotations.

19

20 **Keywords:** single-cell multiomics, transcriptional regulation, transcription factor,
21 chromatin accessibility.

22 Context-specific regulation of gene transcription is central to cell identity and function in
23 eukaryotes. Precision of transcriptional control is achieved through multitudes of
24 transcription factors (TFs) that bind to the *cis*-regulatory regions of their target genes,
25 dynamically modulating chromatin accessibility and recruiting transcription complexes in
26 response to developmental and environmental cues¹. Dissecting this regulatory logic is
27 fundamental to our understanding of biological systems and our study of diseases. Over
28 the past decades, molecular studies have elucidated the structure of TF complexes and
29 provided mechanistic models into their function². Methods based on high-throughput
30 sequencing have enabled the genome-wide profiling of gene expression³, TF binding⁴,
31 chromatin accessibility⁵, and 3D genome structure⁶. TF knockdown/knockout studies
32 have also identified, *en masse*, their species-, tissue-, and context-specific target genes⁷.
33 Concurrently, novel statistical approaches have allowed for more precise identification
34 and modeling of TF binding sites⁸, and expression quantitative trait loci (eQTLs)
35 databases now include associations that are tissue-specific⁹ and will soon be cell-type
36 specific¹⁰. Yet, despite this tremendous progress, our understanding of gene regulatory
37 logic is still rudimentary.

38 When a TF j activates or represses the expression of a gene g through binding to
39 a regulatory element t in *cis* to the gene, we call such a relationship a *regulatory trio*.
40 Despite its complexity, gene regulatory logic can be broken down into modular
41 components consisting of such peak-TF-gene trios. In this paper, we focus on the
42 identification of regulatory trios using multiomic experiments that jointly profile gene
43 expression and chromatin accessibility at single-cell resolution.

44 Single-cell RNA sequencing (scRNA-seq) and single-cell assay of transposase-
45 accessible chromatin sequencing (scATAC-seq), performed separately, have already
46 generated detailed cell-type-specific profiles of gene expression and chromatin
47 accessibility. When the two modalities are not measured in the same cells, the cells can
48 be aligned by computational methods¹¹, followed by association analyses of gene
49 expression and peak accessibility. While these methods have been shown to align well-
50 differentiated cell types correctly, they often fail for cell populations consisting of transient
51 states. Additionally, the alignment of cells necessarily assumes a peak-gene association
52 model, which is often learned from other datasets. Then, the post-alignment association

53 analysis is plagued by logical circularity, as it is difficult to disentangle new findings from
54 prior assumptions that underlie the initial cell alignment.

55 Single-cell multiomic experiments that sequence the RNA and ATAC from the
56 same cells directly enable joint modeling of a cell's RNA expression and chromatin state,
57 yet methods for the analysis of such data are still in their infancy. Almost all existing
58 methods for detecting and characterizing regulatory relationships between TF, regulatory
59 region, and target gene rely on marginal relationships. For example, Signac¹² and Ma *et*
60 *al.*¹³ use marginal associations between peaks and genes to identify putative enhancer
61 regions, while Signac¹² and Seurat V4¹⁴ link differentially expressed TFs to differentially
62 accessible motifs across cell types. Such pairwise marginal associations are sometimes
63 examined manually, in tandem, using low-dimensional embedding. One exception is in
64 PECA¹⁵, which uses a parametric model to characterize the joint distribution of TF
65 expression, regulatory site accessibility, chromatin remodeler expression, and target
66 gene expression. Although PECA was designed to be applied to matched bulk
67 transcriptomic and epigenomic data, such joint modeling concepts could potentially be
68 very powerful for single-cell multiomic data. Yet, PECA relies heavily on parametric
69 assumptions, is computationally intensive to fit to large datasets, and is difficult to
70 diagnose.

71 Context-specific gene regulation may be masked in marginal associations, as we
72 will show in examples later. We explore in this paper the use of higher-order models that
73 interrogate conditional and three-way interaction relationships for the identification of
74 regulatory trios. First, as proof of principle, we show that a simple model that integrates
75 TF expression with *cis*-peak accessibility significantly improves gene expression
76 prediction, as compared to a comparable model that utilizes peak accessibility alone. We
77 present TRIPOD, a computational framework for transcription regulation interrogation
78 through nonparametric partial association analysis of single-cell multiomic sequencing
79 data. TRIPOD detects two types of trio relationships, which we call conditional level 1 and
80 conditional level 2, through robust nonparametric tests that are easy to diagnose. A novel
81 influence measure allows the detection and visualization of cell states driving these
82 regulatory relationships, applicable both to discrete cell types, as well as continuous cell
83 trajectories.

84 We first apply TRIPOD to single-cell multiomic data of human peripheral blood
85 mononuclear cells (PBMCs) and compare the regulatory trios detected to relationships
86 detected through marginal associations. We show that the detections are coherent with
87 the vast amounts of existing knowledge from enhancer databases, bulk cell-type-specific
88 chromatin immunoprecipitation followed by sequencing (ChIP-seq) experiments, tissue-
89 specific TF knockdown/knockout studies, and *cis*-eQTL studies, but that conditional and
90 marginal models identify different sets of relationships. We next apply TRIPOD to the
91 interrogation of lineage-specific regulation in the developing mouse brain, where
92 relationships detected by TRIPOD are compared against those derived from existing
93 ChIP-seq and proximity ligation-assisted ChIP-seq (PLAC-seq) data. Here, TRIPOD
94 identifies known trio relationships, as well as putative novel regulatory crosstalk between
95 neuronal TFs and glial-lineage genes. We also apply TRIPOD to SHARE-seq data on
96 mouse hair follicle cell differentiation to illustrate trio detection and influence analysis in
97 data collected from different protocols. Through these analyses, we demonstrate how to
98 harness single-cell multiomic technologies in the study of gene regulation, and how the
99 data from these technologies corroborate and complement existing data.

100

101 **Results**

102 **A simple interaction model between TF expression and peak accessibility improves**
103 **RNA prediction.** We started our analyses by making gene expression predictions, a
104 standard procedure carried out by existing methods¹¹. We benchmarked against: (i)
105 Signac¹² and Cicero¹⁶, which compute the gene activity matrix by summing the ATAC
106 reads in gene bodies and promoter regions; (ii) MAESTRO¹⁷, which adopts a regulatory
107 potential model by taking the weighted sum of the ATAC reads based on existing gene
108 annotations; and (iii) sci-CAR¹⁸, which performs a regularized regression using gene
109 expression as the outcome and peak accessibilities as the predictors. We also use a
110 regularized regression model (referred to as peak-TF LASSO model), where instead of
111 peak accessibility, interactions between TF expressions and peak accessibilities serve as
112 predictors. The model considers only peaks within a certain range of the gene's
113 transcription start site (TSS) and only interactions between TFs and peaks containing

114 high-scoring binding motifs for the TFs. To avoid overfitting, we perform leave-one-out
115 prediction and adopt independent training and testing sets. See Methods for details.

116 We analyzed single-cell multiomic datasets from different human and mouse
117 tissues generated by different platforms – PBMC by 10X Genomics, embryonic mouse
118 brain by 10X Genomics, mouse skin by SHARE-seq¹³, and adult mouse brain by SNARE-
119 seq¹⁹. Data summaries are included in Supplementary Table 1; reduced dimensions via
120 uniform manifold approximation and projection (UMAP)²⁰ are shown in Fig. 1a and
121 Supplementary Fig. 1,2a. To mitigate the undesirable consequences of sparsity and
122 stochasticity in the single-cell data, we clustered cells to form metacells¹⁴, and pooled
123 gene expression and chromatin accessibility measurements within each metacell.

124 Our results show that, across window sizes, the peak-TF LASSO model
125 significantly improves prediction accuracy across the transcriptome (Fig.1b), with
126 examples shown in (Fig. 1c). This improvement in prediction accuracy holds true when
127 an independent dataset is used for validation (Supplementary Fig. 3). For the SNARE-
128 seq data¹⁹, sequencing depth is substantially shallower (Supplementary Fig. 4), thus the
129 improvement of the peak-TF LASSO model is diminished but still evident (Supplementary
130 Fig. 2b). This demonstrates that the product of TF expression and peak accessibility
131 significantly improves RNA prediction accuracy beyond simply using peak accessibility,
132 offering strong empirical evidence of three-way interaction relationships between TF
133 expression, peak accessibility, and target gene expression that can be extracted from
134 such multiomic experiments. However, we will not rely on coefficients from the LASSO
135 model to screen for such trios, as their significance is difficult to compute due to the
136 hazards of post-selection inference²¹. Additionally, accessibility of peaks and expression
137 of TF affecting the same gene are often highly correlated, in which case LASSO tends to
138 select the few with the highest associations and ignore the rest. In such cases, we believe
139 it is more desirable to report all trios.

140
141 **TRIPOD for the detection of peak-TF-gene trio regulatory relationships by single-**
142 **cell multiomic data.** We propose TRIPOD, a nonparametric method that screens single-
143 cell RNA and ATAC multiomic data for conditional associations and three-way
144 interactions between the expression of a TF j , the accessibility of a peak region

145 t containing the TF's motif, and the expression of a putative target gene g within a pre-
146 fixed distance of peak t (Fig. 2a). Existing methods¹²⁻¹⁴ screen for marginal associations
147 either between the TF and the peak or between the peak and the target gene. However,
148 three-way relationships may be complex, and true associations may be masked by
149 population heterogeneity, as we demonstrate later. When marginal associations are
150 masked, evidence for cooperation between the TF and the peak in the regulation of a
151 gene can be inferred from partial associations: (i) with the peak open at a fixed
152 accessibility, whether cells with higher TF expression have higher gene expression; and
153 (ii) with the TF expression fixed at a value above a threshold, whether cells with higher
154 peak accessibility have higher gene expression. To identify such conditional associations
155 without making linearity assumptions on the marginal relationships, TRIPOD matches
156 metacells by either their TF expressions or peak accessibilities (Fig. 2b): For each
157 matched metacell pair, the variable being matched is controlled for, and differences
158 between the pair in the other two variables are computed. Then, across pairs, the
159 nonparametric Spearman's test is used to assess the association between difference in
160 target gene expression ΔY_g and difference in the unmatched variable (i.e., ΔY_j if the cells
161 were matched by X_t , or ΔX_t if the cells were matched by Y_j). We call this the "conditional
162 level 1 test."

163 For illustration, consider the metacell denoted by the black point in Fig. 2b: If we
164 were to match by peak accessibility, this metacell would be matched to the metacell
165 colored in red. We would then compute ΔY_j , the difference between TF j expressions of
166 the matched pair. If we were to match by TF expression, the black dot would be matched
167 to the metacell in green, and we would compute ΔX_t , the difference in peak t
168 accessibilities of the pair. In either case, we would compute ΔY_g , the difference in gene g
169 expressions between the pair. We would then mask those metacell-pairs whose values,
170 for the variable being matched, are too low (i.e., those pairs where the TF is off or the
171 peak is closed). Then, ΔX_t or ΔY_j , together with ΔY_g , would be submitted for level 1 test.

172 Even stronger evidence for a regulatory trio could be claimed if the *degree* of
173 association between the pairwise differences depends on the matched variable. For
174 example, we would tend to believe that TF j binds to peak t to regulate gene g if, in cells
175 with high expression of TF j , an increase in peak t accessibility yields a much larger

176 increase in gene g expression, as compared to in cells with low expression of TF j . One
177 could screen for such interactions by matching by either TF j or peak t accessibility.
178 TRIPOD screens for such interaction effects through a “conditional level 2 test”, which
179 assesses the association between ΔY_g and the product of the matched variable with the
180 difference in the unmatched variable, after taking partial residuals on the difference in the
181 unmatched variable.

182 For significant trios, TRIPOD further carries out a sampling-based influence
183 analysis, where phenotypically contiguous sets of metacells are held out to measure their
184 influence on the estimated coefficients. The corresponding cell types/states that lead to
185 significant deviations from the null upon their removal have high influence scores, which
186 can be used to identify cell types/states that drive a regulatory relationship.

187 To highlight the differences between TRIPOD and existing methods based on
188 marginal associations, we show two canonical examples where the two approaches
189 disagree. Fig. 2c outlines a significant trio detected by TRIPOD’s level 2 testing, yet the
190 marginal peak-gene and TF-gene associations were insignificant. It turns out that a
191 subset of cells with high peak accessibility $\{X_t\}$ have close-to-zero TF expressions $\{Y_j\}$,
192 and, meanwhile, another subset of cells with high TF expressions $\{Y_j\}$ have close-to-zero
193 peak accessibilities $\{X_t\}$. In these cells, either the peak is closed, or the TF is not
194 expressed, and this leads to the target gene not being expressed, which masks the
195 marginal associations. The high peak accessibility and TF expression in these cells,
196 which act through other regulatory trios, cancel out when we consider the interaction
197 $\{X_t \times Y_j\}$, leading to a significant interaction term detected by TRIPOD. Conversely, Fig.
198 2d outlines another trio, whose marginal associations were significant, yet TRIPOD did
199 not detect significant conditional associations from either level 1 or level 2 testing. In this
200 case, with almost constant TF expression, the large difference in peak accessibility leads
201 to a small difference in target gene expression. Meanwhile, the cells that drive the
202 significantly positive correlation between $\{Y_g\}$ and $\{Y_j\}$ have almost zero values for $\{X_t\}$.
203 Both observations suggest that this peak has little to do with the regulation of the target
204 gene *FGL2* by this specific TF MAFK. Notably, we do not claim that the significantly linked
205 peaks and TFs through marginal association are false positives, but rather this specific
206 trio is insignificant (i.e., the peak and TF may act through other TF and peak, respectively).

207 In summary, TRIPOD puts peak-TF-gene trios into one unified model, complementing
208 existing methods based on marginal associations and allowing for simultaneous
209 identification of all three factors and prioritization of a different set of regulatory
210 relationships.

211

212 **TRIPOD identifies three-way regulatory relationships in PBMCs with orthogonal**
213 **validations.** We first applied TRIPOD to identify regulatory trios in the 10k PBMC dataset.
214 Cell-type labels for this dataset were transferred from a recently released CITE-seq
215 reference of 162,000 PBMC cells measured with 228 antibodies¹⁴. After quality control,
216 we kept 7790 cells from 14 cell types pooled into 80 metacells, 103,755 peaks, 14,508
217 genes, and 342 TFs; the UMAP reduced dimensions are shown in Supplementary Fig.
218 1a. Distribution of the number of peaks 100kb/200kb upstream and downstream of the
219 TSS per gene, as well as distribution of the number of motifs per peak, are shown in
220 Supplementary Fig. 5.

221 Results from TRIPOD and marginal association tests overlap but exhibit
222 substantial differences (Supplementary Fig. 6). The previous section showed example
223 trios where the two frameworks disagree. As a proof of concept, we now illustrate two
224 trios where the frameworks agree, identified by level 1 conditional testing (regulation of
225 *CCR7* by *LEF1*, Fig. 3a) and level 2 interaction testing (regulation of *GNLY* by *TBX21*,
226 Fig. 3b). From the influence analyses, TRIPOD identified B and T cells as the cell types
227 where *LEF1* regulates *CCR7*, and natural killer (NK) cells as where *TBX21* regulates
228 *GNLY*. These cell type-specific regulatory relationships are corroborated by motif's
229 deviation scores using chromVAR²² (Fig. 3a,b) and the enrichment of Tn5 integration
230 events in the flanking regions using DNA footprinting analyses¹² (Supplementary Fig. 7).
231 Unlike chromVar and DNA footprinting analyses, which only give genome-wide average
232 enrichments, TRIPOD significantly enhances the resolution by identifying the specific *cis*-
233 regulatory regions that the TFs bind for the regulation of target genes.

234 To our best knowledge, no experimental technique can directly validate three-way
235 regulatory relationships at high resolution with high throughput. Therefore, we performed
236 validation and benchmarking by harnessing existing databases and orthogonal

237 sequencing experiments that interrogate each pairwise relationship among the three
238 factors (Table 1).

239 First, to validate the *cis*-linkage between peak region and target gene, we used the
240 enhancer databases of blood and non-cancerous cells from FANTOM5²³ (from HACER²⁴),
241 4DGenome²⁵ (from HACER²⁴), and EnhancerAtlas 2.0²⁶, as well as *cis*-eQTLs in the
242 whole blood reported by the GTEx consortium⁹. We collapsed TRIPOD's trio calls into
243 peak-gene relationships and benchmarked against Signac's LinkPeaks¹² on single cells
244 and marginal testing on metacells; for each target gene, we performed a hypergeometric
245 test for enrichment of the peak-gene linkages in the regulatory databases and annotations
246 (see Methods for details). For all four databases, TRIPOD's *p*-values for enrichment are
247 substantially significant (Fig. 3c).

248 Second, to validate the TF-gene edge in the TRIPOD-identified trios, we referred
249 to knockTF⁷, a TF knockdown/knockout gene expression database, and hTFtarget²⁷, a
250 database of known TF regulatory targets. Specifically, in knockTF, we found seven TF
251 knockdown/knockout RNA-seq experiments in the peripheral blood category. For these
252 TFs, we identified significantly linked genes by marginal association and by TRIPOD and
253 found TRIPOD's results to have significantly higher precision and recall (Fig. 3d). For
254 hTFtarget, we obtained, for each highly variable gene, its blood-specific TFs, and
255 calculated the gene-specific precision-recall rates – TRIPOD is more sensitive compared
256 to marginal association testing, although both suffered from inflated “false positives,”
257 which are likely due to the low sensitivity in the “gold-standard” *in silico* calls by hTFtarget
258 (Fig. 3e).

259 Third, to validate the TF-peak edge representing TF binding to peak regions, in
260 addition to the DNA footprinting analysis shown in Supplementary Fig. 7e, we
261 downloaded from the Cistrome portal²⁸ non-cancerous ChIP-seq data from sorted human
262 blood cells (B lymphocyte, T lymphocyte, and monocyte (Supplementary Table 2). The
263 peaks identified by TRIPOD had a substantially higher percentage of overlap with the
264 ChIP-seq peaks (Fig. 3f). In summary, existing databases and public data of different
265 types from a wide range of studies extensively support each of the three pairwise links in
266 the trios reported by TRIPOD, demonstrating its effectiveness in uncovering true
267 regulatory relationships.

268

269 **TRIPOD identifies known and novel putative regulatory relationships during mouse**

270 **embryonic brain development.** We next applied TRIPOD to single-cell multiomic data

271 of 5k mouse embryonic brain cells at day 18 by 10X Genomics. The cell type labels were

272 transferred from an independent scRNA-seq reference²⁹ using both Seurat V3¹¹ and

273 SAVERCAT³⁰. We kept 3,962 cells that had consistent transferred labels from seven

274 major cell types: radial glia, neuroblast, GABAergic neuron, glutamatergic neuron,

275 glioblast, oligodendrocyte, and Cajal–Retzius neuron (Supplementary Fig. 1b). We

276 applied TRIPOD to 633 TFs, 1000 highly variable genes, and ATAC peaks 200kb

277 up/downstream of the genes' TSSs.

278 First, we investigate a known regulatory trio involving target gene *Sox10*, TF *Olig2*

279 (binding motif CAGCTG), and a *cis*-regulatory element known as the U2 enhancer (chr15:

280 79201691-79201880 from mm10) that has been experimentally validated³¹ The

281 expression of *Sox10* and *Olig2* across the seven cell types, and the ATAC-seq profiles of

282 the region containing the U2 enhancer and *Sox10*, are shown in Fig. 4a. This known

283 regulatory trio was found to be conditional level 2 significant by TRIPOD (Fig. 4b).

284 Importantly, the U2 enhancer resides in one of three *cis*-regulatory elements that were

285 identified by TRIPOD to be enhancers for *Sox10* involving TF *Olig2*; all three significantly

286 linked peaks were validated by *Olig2* ChIP-seq data (Fig. 4a).

287 On the genome scale, the union of TRIPOD's level 1 and 2 tests gave a larger

288 number of unique peak-gene pairs and TF-gene pairs than LinkPeaks¹² and marginal

289 metacell association tests (Supplementary Fig. 8a). To evaluate results, we first examined

290 whether the peak-gene links were enriched in previously reported enhancer-promoter

291 chromatin contacts using PLAC-seq data of mouse fetal brain³² (Table 1, Fig. 4c,

292 Supplementary Fig. 8b). We observed that the regulatory links detected by both marginal

293 association and TRIPOD showed significant enrichment in those supported by PLAC-seq

294 (Fig. 4d, Supplementary Fig. 8c). Importantly, TRIPOD detected sets of trio relationships

295 that were overlapping but distinct from the sets obtained by the marginal model, and a

296 substantial fraction of the links identified by TRIPOD but not by the marginal method were

297 validated by PLAC-seq. This suggests that TRIPOD identifies real regulatory relationships

298 that complement those detected by existing methods. To validate the links between TFs

299 and peaks, we used publicly available ChIP-seq data of mouse embryonic brain for
300 Olig2³³, Neurog2³⁴, Eomes³⁴, and Tbr1³⁵, TFs that play key roles in embryonic brain
301 development (Table 1). TF binding peaks identified by TRIPOD were significantly
302 enriched in the TF ChIP-seq peaks across all embryonic brain datasets; Olig2 ChIP-seq
303 data of mature oligodendrocytes (mOL) serves as a negative control and has insignificant
304 enrichment (Fig. 4e).

305 The validations and global benchmarking demonstrate TRIPOD's effectiveness in
306 finding real regulatory relationships. Next, we focused on a set of TFs known to play
307 essential roles during mouse embryonic brain development. Specifically, we chose Pax6,
308 Neurog2, Eomes, Neurod1, and Tbr1, major TFs mediating glutamatergic neurogenesis³⁶,
309 and Olig2, Sox10, Nkx2-2, Sox9, Nfia, and Ascl1, which initiate and mediate gliogenesis³⁷;
310 the known regulatory cascade is shown in Fig. 4f. TRIPOD's level 1 and level 2 testing
311 successfully captured six out of the seven known regulatory links; interestingly, TRIPOD's
312 results also suggest substantial crosstalk between the two cascades, where
313 neurogenesis-specific TFs activate gliogenesis-specific TFs (Fig. 4g). ChIP-seq data of
314 Neurog2, Eomes, and Tbr1 supported four of the crosstalk links: regulation of Sox9 by
315 Neurog2 and regulation of Nfia by Neurog2, Eomes, and Tbr1, respectively
316 (Supplementary Fig. 9). These crosstalk links that were validated by ChIP-seq are highly
317 significant by both marginal and conditional associations. Thus, we think it is highly
318 plausible that neurogenesis TFs activate gliogenesis genes at day 18 of embryonic
319 mouse brain development, which is exactly when the switch is being made from
320 neurogenesis to gliogenesis. To our best knowledge, these possible links between
321 neurogenesis and gliogenesis pathways have not been systematically explored and thus
322 warrant future investigation. Finally, for each of the neurogenesis and gliogenesis TFs,
323 we performed a gene ontology (GO) analysis of their significantly linked target genes
324 using DAVID³⁸; the enriched terms were largely consistent with the regulatory functions
325 of the TFs during neurogenesis and gliogenesis (e.g., negative regulation of neuron
326 differentiation and oligodendrocyte differentiation) (Fig. 4h).

327 So far, we have taken advantage of the cross-cell-type variation to identify the trio
328 regulatory relationships. To dissect cell-type-specific regulation, we next applied the
329 influence analysis framework (see Methods for details) to the significant trios involving

330 neurogenesis and gliogenesis TFs. For a given TF, the number of trios, for which a given
331 cell type was influential (FDR < 0.01), is summarized in Fig. 4i, with details for specific
332 example trios given in Supplementary Fig. 10. The analyses underpinned the cell types
333 in which the transcriptional regulation was active, and, reassuringly, the neurogenesis
334 and gliogenesis TFs have the most regulatory influence in neuroblasts and glioblasts,
335 respectively. Additionally, *Ascl1* is active in GABAergic neurons in addition to neuroblasts
336 and glioblasts, consistent with its role as a GABAergic fate determinant³⁹. Notably, the
337 highly influential cell types that lead to the significant trios involving several neurogenesis-
338 specific TFs include not only neuroblast but also glioblast, supporting our previous
339 findings on the crosstalk between the two cascades. Overall, TRIPOD allows fine
340 characterization of cell-type- and cell-state-specific functions of the TFs during
341 neurogenesis and gliogenesis.

342

343 **TRIPOD infers lineage-specific regulatory relationships in differentiating mouse**
344 **hair follicle cells.** As a last example, we applied TRIPOD to SHARE-seq¹³ data
345 (Supplementary Fig. 1c) of mouse hair follicle cells, consisting of four broadly defined cell
346 types – transit-amplifying cells (TAC), inner root sheath (IRS), hair shaft, and medulla
347 cells – along a differentiation trajectory. The cell-type labels were curated based on
348 marker genes, TF motifs, and ATAC peaks from the original publication¹³; pseudotime
349 was inferred using Palantir⁴⁰ and overlaid on the cisTopic⁴¹ reduced dimensions of the
350 ATAC domain. Cells were partitioned using both the pseudotime and the UMAP
351 coordinates to construct metacells (Fig. 5a). Due to the low RNA coverage
352 (Supplementary Fig. 4), we focused on 222 highly-expressed TFs, 794 highly expressed
353 genes reported to have more than ten linked *cis*-regulatory peaks¹³, and peaks 100kb
354 up/downstream of the genes' TSSs.

355 For validation, we used H3K4me1 and H3K27ac ChIP-seq data from an isolated
356 mouse TAC population⁴² (Table 1). H3K4me1 and H3K27ac are marks for active
357 enhancers and are used to benchmark TRIPOD's linked peaks against previously
358 reported domains of regulatory chromatin (DORCs)¹³, as well as randomly sampled peaks.
359 The linked peaks by TRIPOD have higher scores for both H3K4me1 and H3K27ac, than
360 DORCs, the latter identified through marginal associations (Fig. 5b). To further validate

361 the regulatory effects of the linked peaks, we obtained previously characterized super-
362 enhancers (SEs) in mouse TACs⁴². Target genes of the 381 SEs were assigned based
363 on the gene's proximity to the SE, as well as the correlation between loss of the SE and
364 loss of the gene transcription⁴². TRIPOD was able to successfully recapitulate the SE
365 regions for the genes considered, with four examples shown in Fig. 5c, where significantly
366 linked peaks mostly reside in the SEs.

367 To demonstrate, Fig. 5d shows regulatory trios that are specific to the IRS lineage,
368 the hair shaft lineage, and the medulla lineage. These trios also show significant pairwise
369 marginal associations (Fig. 5e), lending confidence that they are real. The cell types
370 where the regulation happens are identified by influence analysis, for which the p -values
371 are smoothed along the differentiation trajectory and overlaid on the UMAP embedding
372 (Fig. 5f). DNA footprinting analyses surveyed the enrichment of Tn5 integration events
373 surrounding the corresponding motif sites and showed cell-type-specific enrichment (Fig.
374 5g), corroborating TRIPOD's results.

375

376 Discussion

377 We have considered the detection of regulatory trios, consisting of a TF binding to a
378 regulatory region to activate or repress the transcription of a nearby gene, using single-
379 cell RNA and ATAC multiomic sequencing data. The presented method, TRIPOD, is a
380 new nonparametric approach that goes beyond marginal relationships to detect
381 conditional associations and interactions on peak-TF-gene trios. We applied TRIPOD to
382 three single-cell multiomic datasets from different species and protocols with extensive
383 validations and benchmarks. We started our analyses with predicting gene expression
384 from both peak accessibility and TF expression. Supervised frameworks have been
385 proposed to predict gene expression from DNA accessibility⁴³, and vice versa⁴⁴, using
386 matched bulk transcriptomic and epigenomic sequencing data. Blatti *et al.*⁴⁵ showed that
387 joint analysis of DNA accessibility, gene expression, and TF motif binding specificity
388 allows reasonably good prediction of TF binding as measured by ChIP-seq. However,
389 none of these methods incorporate TF expression. By selecting peaks near the genes'
390 TSSs and TFs with high motif scores in the selected peaks, we constructed biologically

391 meaningful peak-TF pairs as predictors and showed that such a mechanistic model
392 significantly boosts the prediction accuracy of gene expression.

393 We next considered the detection and significance assessment for individual peak-
394 TF-gene trios, comprehensively comparing our detections with those made by tissue- and
395 cell-type-matched PLAC-seq and CHIP-seq experiments, by *cis*-eQTL and TF
396 knockdown/knockout studies, and by those recorded in the main enhancer databases.
397 Our current study is limited in several ways. A study in *Drosophila*⁴⁶ modeled motif binding
398 specificities and chromatin accessibilities in bulk RNA and ATAC sequencing data to
399 predict the cooperative binding of pairs of TFs, using *in vitro* protein-protein binding
400 experiments for validation. The detection of synergies between multiple TFs and peaks
401 on the genome-wide scale and in a cell-type-specific manner needs further investigation.
402 Additionally, while we have not differentiated between positive and negative regulation,
403 TRIPOD reports both types of relationships and categorizes them by sign. While we
404 describe the trios with a positive sign to be enhancers, it is not clear how to interpret the
405 trios with negative signs, the latter having lower overlap with other benchmarking datasets.
406 Transcription activation and repression have been active research areas in biology, with
407 a lot yet unknown⁴⁷. TRIPOD's results provide potential targets for experimental follow-
408 up and detailed characterization.

409 Our analysis focused on three datasets where the RNA and ATAC modalities have
410 sufficient depths of coverage. For the SHARE-seq data, the sequencing depth for RNA is
411 very low, and thus we focused only on highly expressed genes and TFs (Fig. 5). For
412 SNARE-seq data, whose coverage in both modalities is even lower, we focused on
413 prediction models and not trio detection, where we saw only marginal improvement
414 beyond existing methods¹⁹ (Supplementary Fig. 2). For data where the coverage is even
415 lower, e.g., PAIRED-seq, cross-modality metacells could not be stably formed, making
416 such analyses impossible (Supplementary Table 1, Supplementary Fig. 4). With rapidly
417 increasing sequencing capacity and technological advancement, TRIPOD, applied to
418 more cells sequenced at higher depth, can uncover novel regulatory relationships at a
419 finer resolution.

420

421 **Methods**

422 **Data input and construction of metacells.** Denote X_{it} as the peak accessibility for peak
423 t ($1 \leq t \leq T$) in cell i ($1 \leq i \leq N$), Y_{ig} as the gene expression for gene g ($1 \leq g \leq G$), and
424 Y_{ij} as the TF expression for TF j ($1 \leq j \leq M$). The TF expression matrix is a subset of the
425 gene expression matrix, and for single-cell multiomic data, the cell entries are matched.
426 To mitigate the effect of ATAC sparsity⁴⁸ and RNA expression stochasticity⁴⁹, as a first
427 step, TRIPOD performs cell-wise smoothing by pooling similar cells into “metacells.” This,
428 by default, is performed using the weighted-nearest neighbor method by Seurat V4¹⁴ to
429 jointly reduce dimension and identify cell clusters/states across different modalities. In
430 practice, the metacells can also be inferred using one modality – for example, RNA may
431 better separate the different cell types²⁹, and in other cases, chromatin accessibility may
432 prime cells for differentiation¹³. To account for peaks overlapping with other genes
433 (Supplementary Fig. 5b), TRIPOD has the option to either remove the overlapped peaks
434 or to adjust the peak accessibilities by the expressions of the overlapped genes, in a
435 similar fashion to MAESTRO¹⁷. Library size is adjusted for both the RNA and ATAC
436 domain by dividing all counts by a metacell-specific size factor (total read counts divided
437 by 10^6).

438

439 **RNA prediction by TF expression and peak accessibility.** To predict RNA from ATAC,
440 Signac¹² and Cicero¹⁶ take the sum of peak accessibilities in gene bodies and promoter
441 regions to construct a pseudo-gene activity matrix: $\hat{Y}_{ig} = \sum_{t \in E_g} X_{it}$, where E_g is the set of
442 peaks within gene bodies and upstream regions of TSSs. Instead of directly taking the
443 sum, MAESTRO¹⁷ adopts a “regulatory potential” model by taking the weighted sum of
444 accessibilities across all nearby peaks: $\hat{Y}_{ig} = \sum_{t \in E_g} w_t^g X_{it}$, with weights $\{w_t^g\}$ pre-
445 calculated based on existing gene annotations. Specifically, the method weighs peaks by
446 exponential decay from TSS, sums all peaks on the given gene exons as if they are on
447 the TSS, normalizes the sum by total exon lengths, and excludes the peaks from
448 promoters and exons of nearby genes. The strategy to take the unweighted/weighted sum
449 of accessibility as a proxy for expression has been adopted to align the RNA and ATAC
450 modalities when scRNA-seq and scATAC-seq are sequenced in parallel from the same
451 cell population but not the same cells¹¹. For single-cell multiomic data, sci-CAR¹⁸ performs

452 feature selection to identify *cis*-linked peaks via a LASSO regression: $Y_{ig} \sim \sum_{t \in E_g} \beta_t^g X_{it}$,
 453 where an L1 regularization is imposed on β_t^g . Compared to MAESTRO, which pre-fixes
 454 the weights $\{w_t^g\}$, $\{\beta_t^g\}$ are estimated from the data by regressing RNA against matched
 455 ATAC data. What we propose is a feature selection model involving both peak
 456 accessibility and TF expression: $Y_{ig} \sim \sum_{t \in E_g} \sum_{j \in f_t} \beta_{tj}^g X_{it} Y_{ij}$, where f_t contains the set of
 457 TFs with high-scoring binding motifs in peak t inferred from the JASPAR database⁵⁰.

458
 459 **TRIPOD model and trio regulatory relationship.** For a given target gene g , a peak t
 460 within a window centered at the gene's TSS, and a TF j whose binding motif is high-
 461 scoring in the peak, TRIPOD infers the relationship between a regulatory trio (t, j, g) .
 462 TRIPOD focuses on one trio at a time and goes beyond the marginal associations to
 463 characterize the function $Y_g = f(X_t, Y_j)$. In what follows, we first describe TRIPOD's
 464 matching-based nonparametric approach and then describe a linear parametric approach,
 465 followed by a discussion on the connections and contrasts between the two approaches.

466 For each cell i whose TF expression is above a threshold δ (we only carry out
 467 testing in cells that express the TF), we carry out a minimum distance pairwise cross-
 468 match based on $\{Y_{ij} | Y_{ij} > \delta\}$. Let $\{(i_p, i_{p^*})\}$ be the optimal matching, after throwing away
 469 those pairs that have $|Y_{i_p j} - Y_{i_{p^*} j}| > e$. For each pair p , i_p and i_{p^*} are two metacells with
 470 matched TF expression, for which we now observe two, possibly different, values
 471 $\{X_{i_p t}, X_{i_{p^*} t}\}$ for peak t , as well as two corresponding values $\{Y_{i_p g}, Y_{i_{p^*} g}\}$ for gene g . We
 472 then compute the following auxiliary differentials within each pair:

$$473 \quad \Delta X_{pt} = X_{i_p t} - X_{i_{p^*} t},$$

$$474 \quad \Delta Y_{pg} = Y_{i_p g} - Y_{i_{p^*} g},$$

475 as well as

$$476 \quad \bar{Y}_{pj} = (Y_{i_p j} + Y_{i_{p^*} j})/2.$$

477 For level 1 testing of conditional association, we estimate $\hat{r}_t^g = \rho(\Delta X_{pt}, \Delta Y_{pg})$, where ρ is
 478 Spearman correlation, and test $H_1: r_t^g = 0$. For level 2 testing of interaction, we perform a
 479 regression $\Delta Y_{pg} = \alpha \Delta X_{pt} + \gamma \bar{Y}_{pj} \times \Delta X_{pt}$, set $\hat{\gamma}_{tj}$ to be the least-squares solution for γ , and
 480 test $H_2: \gamma_{tj} = 0$. For visualization of the model fitting, we take the partial residuals of ΔY_{pg}

481 and $\bar{Y}_{pj} \times \Delta X_{pt}$ on ΔX_{pt} , respectively. Note that even though TF expression is not included
 482 in the model as a main term, it is controlled for (and not just in the linear sense) by the
 483 matching. Similarly, we can also perform this procedure matching by peak accessibility.
 484 As a summary, for level 1 testing of conditional association, we have:

$$485 \quad \text{Match by } Y_j, \alpha = \rho(\Delta Y_g, \Delta X_t),$$

$$486 \quad \text{Match by } X_t, \beta = \rho(\Delta Y_g, \Delta Y_j).$$

487 For level 2 testing of (TF expression) \times (peak accessibility) interaction effects, we have:

$$488 \quad \text{Match by } Y_j, \Delta Y_g = \alpha^* \Delta X_t + \gamma_1 (\bar{Y}_j \times \Delta X_t),$$

$$489 \quad \text{Match by } X_t, \Delta Y_g = \beta^* \Delta Y_j + \gamma_2 (\bar{X}_t \times \Delta Y_j).$$

490 To test for the conditional associations and interactions, we can also use apply a
 491 parametric method, such as multiple linear regression:

$$492 \quad Y_g = \mu + \alpha_L X_t + \beta_L Y_j,$$

$$493 \quad Y_g = \mu + \alpha_L^* X_t + \beta_L^* Y_j + \gamma_L X_t Y_j.$$

494 See Supplementary Fig. 11 for linear testing results for trios shown in Fig. 3 and Fig. 5.
 495 The estimated coefficients from the nonparametric and parametric methods are
 496 correlated on the global scale (Supplementary Fig. 12), and their interpretations are
 497 similar: α and α_L estimate the change in gene expression per change in peak
 498 accessibility, fixing TF expression; β and β_L estimate the change in gene expression per
 499 change in TF expression, fixing peak accessibility; γ_1 and γ_L measure how the change in
 500 gene expression per change in peak accessibility at each fixed TF expression relies on
 501 the TF expression; γ_2 and γ_L measure how the change in gene expression per change in
 502 TF expression at each fixed peak accessibility relies on the peak accessibility. However,
 503 the underlying models and assumptions are different. Matching controls for not just the
 504 linear variation in the matched variable, but also any nonlinear variation. This contrasts
 505 with adding the variable as a covariate in the linear regression, where we simply remove
 506 linear dependence. The main motivation for using the matching model above is our
 507 reluctance to assume the simple linear relationship. Additionally, we use the rank-based
 508 Spearman correlation, which will not be driven by outliers – a “bulk” association between
 509 ranks is needed for significance. Thus, the nonparametric model of TRIPOD is more
 510 stringent (Supplementary Fig. 13) and more robust to outliers.

511

512 **Identifying regulatory cell type(s) and cell state(s).** For the significant trios detected
513 by TRIPOD, we next seek to identify the underlying regulatory cell type(s). Specifically,
514 we carry out a cell-type-specific influence analysis to identify cell types that are highly
515 influential in driving the significance of the trio. Traditional approaches (e.g., the Cook's
516 distance and the DFFITs) delete observations one at a time, refit the model on remaining
517 observations, and measure the difference in the predicted value from the full model and
518 that from when the point is left out. While they can be readily adopted to detect "influential"
519 metacells one at a time (Supplementary Fig. 7a,b), these methods do not adjust for the
520 degree of freedom properly when deleting different numbers of metacells from different
521 cell types. That is, they do not account for the different numbers of observations that are
522 simultaneously deleted. Additionally, both methods adopt a thresholding approach to
523 determine significance, without returning p -values that are necessary for multiple testing
524 correction. We, therefore, develop a sampling-based approach to directly test for the
525 influence of multiple metacells and to return p -values (Supplementary Fig. 7c).

526 Here, we focus on the linear model for its ease of computation: $\hat{Y}_g = \hat{\mu} + \hat{\alpha}X_t +$
527 $\hat{\beta}Y_j + \hat{\gamma}X_tY_j$. Given a set of observations $I = \{i: i\text{th metacell belongs to a cell type}\}$, we
528 remove these metacells, fit the regression model, and make predictions: $\hat{Y}_g^{(I)} = \hat{\mu}^{(I)} +$
529 $\hat{\alpha}^{(I)}X_t + \hat{\beta}^{(I)}Y_j + \hat{\gamma}^{(I)}X_tY_j$. The test statistics are the difference in the fitted gene
530 expressions $|\hat{Y}_g - \hat{Y}_g^{(I)}|$. We generate the null distribution via sampling. Specifically, within
531 each sampling iteration, we sample without replacement the same number of metacells,
532 denoted as a set of I^* , delete these observations, and refit the regression model on the
533 remaining observations: $\hat{Y}_g^{(I^*)} = \hat{\mu}^{(I^*)} + \hat{\alpha}^{(I^*)}X_t + \hat{\beta}^{(I^*)}Y_j + \hat{\gamma}^{(I^*)}X_tY_j$. The p -value is
534 computed across K sampling iterations as $p_{Y_g} = \sum_{I^*} 1\left(\sum |\hat{Y}_g - \hat{Y}_g^{(I)}| \geq \sum |\hat{Y}_g - \hat{Y}_g^{(I^*)}|\right)/K$,
535 where $1()$ is the indicator function. In addition to testing each cell type separately, the
536 framework can be extended to test for the influence of cell-type groups. For example, in
537 Fig. 3a,b, we reconstruct the cell-type hierarchy using expression levels of highly variable
538 genes from the RNA domain and carry out the aforementioned testing scheme at each
539 split for its descendent cell types in the hierarchical structure.

540 For transient cell states, TRIPOD first identifies the neighbors of each metacell
541 along the trajectory and then carries out metacell-specific testing by simultaneously
542 removing each metacell and its neighbors using the framework described above. The
543 resulting p -values are, therefore, smoothed and can be visualized in the UMAP plot, as
544 shown in Fig. 5f and Supplementary Fig. 10, to identify the underlying branches/segments
545 that are key in defining the significant regulatory trio. This approach can be directly applied
546 to cells with branching dynamics without the need to isolate cell subsets or to identify cell
547 types.

548

549 **Validation resources and strategies.** Resources for validating the trio regulatory
550 relationships are summarized in Table 1. To validate the peak-gene relationships, we
551 referred to existing enhancer databases: FANTOM5²³ links enhancers and genes based
552 on enhancer RNA expression; 4DGenome²⁵ links enhancers and genes based on
553 physical interactions using chromatin-looping data including 3C, 4C, 5C, ChIA-PET, and
554 Hi-C; EnhancerAtlas 2.0²⁶ reports enhancers using 12 high-throughput experimental
555 methods including H3K4me1/H3K27ac ChIP-seq, DNase-seq, ATAC-seq, and GRO-seq.
556 We only focused on blood and non-cancerous cells from these databases (Fig. 3c). A list
557 of *cis*-eQTLs within the whole blood mapped in European-American subjects was
558 downloaded from the GTEx consortium⁹ (Fig. 3c). For the mouse embryonic brain dataset,
559 we additionally adopted H3K4me3-mediated PLAC-seq data³², which reported enhancer-
560 promoter chromatin contacts mapped in mouse fetal forebrain (Fig. 4c,d). For the mouse
561 skin dataset, we adopted TAC-specific ChIP-seq data of H3K4me1 and H3K27ac⁴², both
562 of which are histone marks for active enhancers (Fig. 5b); we also obtained previously
563 reported super-enhancers in mouse TACs from *in vivo* studies⁴² (Fig. 5c). Genomic
564 coordinates were lifted over from mm9 to mm10 when necessary.

565 To validate the TF-gene relationships, we utilized the knockTF⁷ and the hTFtarget²⁷
566 databases. knockTF interrogates the changes in gene expression profiles in TF
567 knockdown/knockout experiments to link the TFs to their target genes in a tissue- or cell-
568 type-specific manner. We downloaded 12 experiments, corresponding to 12 TFs
569 (BCL11A, ELK1, GATA3, JUN, MAF, MYB, NFATC3, NFKB1, STAT3, STAT6, TAL1, and
570 ZNF148) in the peripheral blood category, and focused on seven TFs that have at least

571 one linked gene by any model benchmarked (Fig. 3d). hTFtarget computationally predicts
572 TF-gene relationships using ChIP-seq data, and we manually downloaded the TFs
573 associated with each of the top 100 highly variable genes in the blood tissue (Fig. 3e).

574 To validate the peak-TF relationships, we downloaded non-cancerous cell-type-
575 specific ChIP-seq data of human blood (B lymphocyte, T lymphocyte, and monocyte) from
576 the Cistrome²⁸ portal for the PBMC data (Fig. 3f, Supplementary Table 2), and ChIP-seq
577 data of Olig2³³, Neurog2³⁴, Eomes³⁴, and Tbr1³⁵ for the mouse embryonic brain data. The
578 Olig2 ChIP-seq data were generated in three types of rat cells: data from oligodendrocyte
579 precursor cells (OPC) and immature oligodendrocytes (iOL) were used for validation,
580 while data from mature oligodendrocytes (mOL) serve as a negative control³³. Genomic
581 coordinates were converted from rn4 to mm10. The Neurog2 and Eomes ChIP-seq data
582 were generated in mouse embryonic cerebral cortices at day 14.5³⁴; the Tbr1 ChIP-seq
583 data was generated in the whole cortex dissected from embryos at day 15.5³⁵. In addition,
584 DNA footprinting signatures were corrected for Tn5 sequence insertion bias and stratified
585 by cell types using the Signac package and can be used to validate the identified
586 TFs/motifs in a cell-type-specific manner (Fig. 5g, Supplementary Fig. 7e).

587 For peak enrichment analysis compared to the existing enhancers, *cis*-eQTLs, and
588 enhancer-promoter contacts, we carried out a hypergeometric test as follows. Let k be
589 the number of significantly linked peaks, q be the number of significantly linked peaks that
590 overlap with annotations (e.g., annotated enhancers), m be the number of peaks that
591 overlap with the annotations, and n be the number of peaks that do not overlap with
592 annotations. The p -value of enrichment is derived from the hypergeometric distribution
593 using the cumulative distribution function, coded as `phyper(q, m, n, k, lower.tail=F)` in R.
594 We used this hypothesis testing framework to validate and benchmark the reported peak-
595 gene links, with results shown in Fig. 3c and Fig. 4d. A similar analysis was carried out to
596 test for peak enrichment in TF binding sites by ChIP-seq, thus validating the peak-TF
597 relationships (Fig. 4e).

598

599 **Data availability**

600 This study analyzed existing and publicly available single-cell RNA and ATAC multiomic
601 data. 10X Genomics single-cell multiomic datasets of PBMC (10k and 3k) and mouse

602 embryonic brain were downloaded <https://support.10xgenomics.com/single-cell->
603 [multiome-atac-gex/datasets](https://support.10xgenomics.com/single-cell-multiome-atac-gex/datasets). SNARE-seq data of adult mouse brain and SHARE-seq
604 data of mouse skin are available from the Gene Expression Omnibus (GEO) database
605 with accession numbers GSE126074 and GSE140203. A detailed data summary is
606 provided in Supplementary Table 1. Validation resources based on existing databases
607 and high-throughput sequencing data are summarized in Table 1 and Supplementary
608 Table 2.

609

610 **Code availability**

611 TRIPOD is compiled as an open-source R package available at
612 <https://github.com/yharigaya/TRIPOD>. Scripts used for analyses carried out in this paper
613 are deposited in the GitHub repository.

614

615 **Acknowledgments**

616 This work was supported by the National Institutes of Health (NIH) grant R35 GM133712
617 (to C.Z.), R01 HG006137 (to N.R.Z), R01 GM125301 (to N.R.Z.), and R35 GM138342 (to
618 Y.J.). The authors thank Dr. Sai Ma for support and guidance on the SHARE-seq data,
619 Manas Tiwari for help on accessing the hTFtarget database, and Drs. Yun Li, Michael
620 Love, Li Qian, and Jason Stein for helpful discussions and comments.

621

622 **Author contributions**

623 N.R.Z. and Y.J. initiated and envisioned the study. Y.H., N.R.Z., and Y.J. formulated the
624 model, developed the algorithm, and performed data analysis. Z.Z. processed reference
625 datasets and performed cell-type label transfer. H.Z. and C.Z. provided support on
626 validation, offered consultation, and contributed to result interpretation. Y.H., N.R.Z., and
627 Y.J. wrote the manuscript, which was read and approved by all authors.

628

629 **Competing Interests**

630 The authors declare no competing interests.

631

632 **Figure Legends**

633 **Fig. 1 | Interaction between TF expression and peak accessibility improves RNA**

634 **prediction accuracy. a**, UMAP embedding of 10x Genomics PBMC (left), 10x Genomics

635 embryonic mouse brain (center), and SHARE-seq mouse skin (right) cells from single-

636 cell RNA and ATAC multiomic sequencing. Cell-type labels were transferred from existing

637 single-cell references or curated based on marker genes, motifs, and peaks; metacells

638 were constructed to mitigate sparsity and stochasticity. **b**, Genome-wide distributions of

639 Pearson correlations between observed and leave-one-out predicted RNA expression

640 levels, with varying window sizes. **c**, Predicted and observed RNA expression levels for

641 highly variable genes, *CCR7*, *Adamts6*, and *Ano7*, from the three datasets, respectively.

642

643 **Fig. 2 | TRIPOD infers peak-TF-gene trio regulatory relationships using single-cell**

644 **multiomic data. a**, Data input and schematic on a peak-TF-gene trio. **b**, Overview of

645 TRIPOD for inferring regulatory relationships. TRIPOD complements existing methods

646 based on marginal associations by identifying conditional associations through matching

647 by TF expression or peak accessibility. **c**, An example trio identified by TRIPOD, but not

648 by the marginal associations due to heterogeneity of cell-type-specific regulations. **d**, An

649 example trio identified by the marginal associations, but not by TRIPOD. The peak and

650 TF are significantly linked to the gene, yet they act through other TF and peak, and thus

651 the regulatory trio is insignificant. The points represent metacells (left two panels) and

652 pairs of matched metacells (right two panels). Genomic coordinates for the peaks are

653 from hg38.

654

655 **Fig. 3 | TRIPOD identified trio regulatory relationships in PBMC single-cell**

656 **multiomic dataset. a-b**, Example trios identified by TRIPOD. Violin plots show cell-type-

657 specific distributions of gene expression, peak accessibility, and TF expression.

658 Scatterplots show TRIPOD's level 1 and level 2 testing, respectively. Inner and outer

659 circles around the points are color-coded based on the cell types of the matched metacells.

660 Hierarchical clustering is performed on RNA expression levels of highly variable genes.

661 Red/gray circles indicate whether removal of the corresponding branches of metacells

662 significantly changes the model fitting; crosses indicate that removal of the groups of

663 metacells resulted in inestimable coefficients. Genomic coordinates for the peaks are
664 from hg38. **c**, Peak-gene validation based on enhancer databases (FANTOM5,
665 4DGenome, and EnhancerAtlas) and tissue-specific *cis*-eQTL data from the GTEx
666 Consortium. Box plots show distributions of *p*-values from gene-specific hypergeometric
667 tests. **d**, TF-gene validation based on lists of TF-gene pairs from the knockTF database.
668 **e**, Precision and recall rates for TF-gene pairs using ground truths from the hTFtarget
669 database. **f**, Peak-TF validation based on eight cell-type-specific TF ChIP-seq datasets
670 (B lymphocytes, monocytes, and T lymphocytes). Percentages of significantly linked
671 peaks and all peaks that overlap with the ChIP-seq peaks are shown.

672

673 **Fig. 4 | TRIPOD identified known and novel regulatory relationships during mouse**
674 **neurogenesis and gliogenesis.** **a**, TRIPOD identified a previously reported regulatory
675 trio with gene *Sox10*, TF Olig2, and *cis*-regulatory U2 element. TRIPOD identified two
676 additional linked peaks; all three *cis*-regulatory elements were validated by Olig2 ChIP-
677 seq data. **b**, TRIPOD's level 2 testing matching peak accessibility for the *Sox10* gene, the
678 Olig2 TF, and the U2 enhancer. **c**, Venn diagram of the number of peak-gene pairs
679 captured by PLAC-seq, the marginal model, and the union set of TRIPOD's level 1 and
680 level 2 testing matching TF expression. **d**, Enrichment of peak-gene pairs captured by
681 LinkPeaks, marginal association, and TRIPOD in enhancer-promoter contacts by PLAC-
682 seq. **e**, Peak-TF validation by ChIP-seq data. Olig2 ChIP-seq data of precursor/immature
683 oligodendrocytes (OPC/iOL) were used for validation; data from mature oligodendrocytes
684 (mOL) served as a negative control. **f**, A schematic of well-characterized TF regulatory
685 cascades during neurogenesis and gliogenesis. **g**, Trio examples from known regulatory
686 relationships, as well as from crosstalks supported by ChIP-seq data, captured by
687 TRIPOD. **h**, GO analysis of putative target genes of the neurogenesis and gliogenesis
688 TFs. The number of TRIPOD-identified target genes in the GO categories is shown. The
689 background heatmap shows negative log *p*-values (FDR adjusted) from hypergeometric
690 tests examining enrichment of ATAC peaks in ChIP-seq peaks. **i**, Bar plots showing the
691 number of putative cell-type-specific trios mediated by the neurogenesis- and gliogenesis-
692 specific TFs.

693

694 **Fig. 5 | TRIPOD identified regulatory relationships in mouse hair follicles with**
695 **transient cell states. a**, UMAP embedding of hair follicle cells from the mouse skin data.
696 Cells are colored by cell types (TAC, IRS, hair shaft, and medulla) and pseudotime. **b**,
697 H3K4me1 and H3K27ac ChIP-seq scores for linked peaks identified by TRIPOD, DORCs
698 (regulatory domains identified by gene-peak correlations), and randomly sampled peaks.
699 **c**, TRIPOD's linked peaks for four representative genes were significantly enriched in
700 previously annotated super-enhancers in the mouse TAC population. **d**, Trios detected
701 by TRIPOD that were active in IRS (top), medulla (middle), and hair shaft (bottom),
702 respectively. **e**, Dot plots of gene expressions, peak accessibilities, and TF expressions
703 across different cell types. **f**, Influence analyses identified segments along the
704 differentiation trajectory where the regulation took effect. The colors in the UMAP
705 embedding correspond to the smoothed p -values from a sampling-based approach. **g**,
706 DNA footprinting assays showed cell-type-specific enrichments of Tn5 integration events.
707 The findings were consistent with those from the influence analyses.

708

709 **Table 1 | Resources for validating peak-TF-gene regulatory relationship.** While there
710 is no existing experimental approach to validate all three factors in a trio at high resolution
711 with high throughput, we resort to existing databases and orthogonal sequencing data to
712 validate peak-gene, peak-TF, and TF-gene pairs, completing the loop.

713

714 **References**

- 715 1. Gasperini, M., Tome, J.M. & Shendure, J. Towards a comprehensive catalogue of validated and
716 target-linked human enhancers. *Nat Rev Genet* **21**, 292-310 (2020).
- 717 2. Hobert, O. Gene regulation by transcription factors and microRNAs. *Science* **319**, 1785-1786
718 (2008).
- 719 3. Mortazavi, A., Williams, B.A., McCue, K., Schaeffer, L. & Wold, B. Mapping and quantifying
720 mammalian transcriptomes by RNA-Seq. *Nat Methods* **5**, 621-628 (2008).
- 721 4. Johnson, D.S., Mortazavi, A., Myers, R.M. & Wold, B. Genome-wide mapping of in vivo protein-
722 DNA interactions. *Science* **316**, 1497-1502 (2007).
- 723 5. Buenrostro, J.D., Giresi, P.G., Zaba, L.C., Chang, H.Y. & Greenleaf, W.J. Transposition of native
724 chromatin for fast and sensitive epigenomic profiling of open chromatin, DNA-binding proteins
725 and nucleosome position. *Nat Methods* **10**, 1213-1218 (2013).
- 726 6. Lieberman-Aiden, E. et al. Comprehensive mapping of long-range interactions reveals folding
727 principles of the human genome. *Science* **326**, 289-293 (2009).
- 728 7. Feng, C. et al. KnockTF: a comprehensive human gene expression profile database with
729 knockdown/knockout of transcription factors. *Nucleic Acids Res* **48**, D93-D100 (2020).

730 8. Tompa, M. et al. Assessing computational tools for the discovery of transcription factor binding
731 sites. *Nat Biotechnol* **23**, 137-144 (2005).

732 9. Consortium, G. Genetic effects on gene expression across human tissues. *Nature* **550**, 204-213
733 (2017).

734 10. Kim-Hellmuth, S. et al. Cell type-specific genetic regulation of gene expression across human
735 tissues. *Science* **369** (2020).

736 11. Stuart, T. et al. Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-1902 e1821 (2019).

737 12. Stuart, T., Srivastava, A., Lareau, C. & Satija, R. Multimodal single-cell chromatin analysis with
738 Signac. *bioRxiv*, 2020.2011.2009.373613 (2020).

739 13. Ma, S. et al. Chromatin Potential Identified by Shared Single-Cell Profiling of RNA and Chromatin.
740 *Cell* **183**, 1103-1116 e1120 (2020).

741 14. Hao, Y. et al. Integrated analysis of multimodal single-cell data. *Cell* **184**, 3573-3587 e3529 (2021).

742 15. Duren, Z., Chen, X., Jiang, R., Wang, Y. & Wong, W.H. Modeling gene regulation from paired
743 expression and chromatin accessibility data. *Proc Natl Acad Sci U S A* **114**, E4914-E4923 (2017).

744 16. Pliner, H.A. et al. Cicero Predicts cis-Regulatory DNA Interactions from Single-Cell Chromatin
745 Accessibility Data. *Mol Cell* **71**, 858-871 e858 (2018).

746 17. Wang, C. et al. Integrative analyses of single-cell transcriptome and regulome using MAESTRO.
747 *Genome Biol* **21**, 198 (2020).

748 18. Cao, J. et al. Joint profiling of chromatin accessibility and gene expression in thousands of single
749 cells. *Science* **361**, 1380-1385 (2018).

750 19. Chen, S., Lake, B.B. & Zhang, K. High-throughput sequencing of the transcriptome and chromatin
751 accessibility in the same cell. *Nat Biotechnol* **37**, 1452-1457 (2019).

752 20. Becht, E. et al. Dimensionality reduction for visualizing single-cell data using UMAP. *Nat Biotechnol*
753 (2018).

754 21. Taylor, J. & Tibshirani, R.J. Statistical learning and selective inference. *Proc Natl Acad Sci U S A* **112**,
755 7629-7634 (2015).

756 22. Schep, A.N., Wu, B., Buenrostro, J.D. & Greenleaf, W.J. chromVAR: inferring transcription-factor-
757 associated accessibility from single-cell epigenomic data. *Nat Methods* **14**, 975-978 (2017).

758 23. Andersson, R. et al. An atlas of active enhancers across human cell types and tissues. *Nature* **507**,
759 455-461 (2014).

760 24. Wang, J. et al. HACER: an atlas of human active enhancers to interpret regulatory variants. *Nucleic*
761 *Acids Res* **47**, D106-D112 (2019).

762 25. Teng, L., He, B., Wang, J. & Tan, K. 4DGenome: a comprehensive database of chromatin
763 interactions. *Bioinformatics* **32**, 2727 (2016).

764 26. Gao, T. & Qian, J. EnhancerAtlas 2.0: an updated resource with enhancer annotation in 586
765 tissue/cell types across nine species. *Nucleic Acids Res* **48**, D58-D64 (2020).

766 27. Zhang, Q. et al. hTFtarget: A Comprehensive Database for Regulations of Human Transcription
767 Factors and Their Targets. *Genomics Proteomics Bioinformatics* **18**, 120-128 (2020).

768 28. Mei, S. et al. Cistrome Data Browser: a data portal for ChIP-Seq and chromatin accessibility data
769 in human and mouse. *Nucleic Acids Res* **45**, D658-D662 (2017).

770 29. La Manno, G. et al. Molecular architecture of the developing mouse brain. *Nature* **596**, 92-96
771 (2021).

772 30. Huang, M., Zhang, Z. & Zhang, N.R. Dimension reduction and denoising of single-cell RNA
773 sequencing data in the presence of observed confounding variables. *bioRxiv*,
774 2020.2008.2003.234765 (2020).

775 31. Kuspert, M., Hammer, A., Bosl, M.R. & Wegner, M. Olig2 regulates Sox10 expression in
776 oligodendrocyte precursors through an evolutionary conserved distal enhancer. *Nucleic Acids Res*
777 **39**, 1280-1293 (2011).

- 778 32. Zhu, C. et al. An ultra high-throughput method for single-cell joint analysis of open chromatin and
779 transcriptome. *Nat Struct Mol Biol* **26**, 1063-1070 (2019).
- 780 33. Yu, Y. et al. Olig2 targets chromatin remodelers to enhancers to initiate oligodendrocyte
781 differentiation. *Cell* **152**, 248-261 (2013).
- 782 34. Sessa, A. et al. The Tbr2 Molecular Network Controls Cortical Neuronal Differentiation Through
783 Complementary Genetic and Epigenetic Pathways. *Cereb Cortex* **27**, 5715 (2017).
- 784 35. Notwell, J.H. et al. TBR1 regulates autism risk genes in the developing neocortex. *Genome Res* **26**,
785 1013-1022 (2016).
- 786 36. Mira, H. & Morante, J. Neurogenesis From Embryo to Adult - Lessons From Flies and Mice. *Front*
787 *Cell Dev Biol* **8**, 533 (2020).
- 788 37. Emery, B. & Lu, Q.R. Transcriptional and Epigenetic Regulation of Oligodendrocyte Development
789 and Myelination in the Central Nervous System. *Cold Spring Harb Perspect Biol* **7**, a020461 (2015).
- 790 38. Huang da, W., Sherman, B.T. & Lempicki, R.A. Systematic and integrative analysis of large gene
791 lists using DAVID bioinformatics resources. *Nat Protoc* **4**, 44-57 (2009).
- 792 39. Achim, K., Salminen, M. & Partanen, J. Mechanisms regulating GABAergic neuron development.
793 *Cell Mol Life Sci* **71**, 1395-1415 (2014).
- 794 40. Setty, M. et al. Characterization of cell fate probabilities in single-cell data with Palantir. *Nat*
795 *Biotechnol* **37**, 451-460 (2019).
- 796 41. Bravo Gonzalez-Blas, C. et al. cisTopic: cis-regulatory topic modeling on single-cell ATAC-seq data.
797 *Nat Methods* **16**, 397-400 (2019).
- 798 42. Adam, R.C. et al. Pioneer factors govern super-enhancer dynamics in stem cell plasticity and
799 lineage choice. *Nature* **521**, 366-370 (2015).
- 800 43. Natarajan, A., Yardimci, G.G., Sheffield, N.C., Crawford, G.E. & Ohler, U. Predicting cell-type-
801 specific gene expression from regions of open chromatin. *Genome Res* **22**, 1711-1722 (2012).
- 802 44. Zhou, W. et al. Genome-wide prediction of DNase I hypersensitivity using gene expression. *Nat*
803 *Commun* **8**, 1038 (2017).
- 804 45. Blatti, C., Kazemian, M., Wolfe, S., Brodsky, M. & Sinha, S. Integrating motif, DNA accessibility and
805 gene expression data to build regulatory maps in an organism. *Nucleic Acids Res* **43**, 3998-4012
806 (2015).
- 807 46. Kazemian, M., Pham, H., Wolfe, S.A., Brodsky, M.H. & Sinha, S. Widespread evidence of
808 cooperative DNA binding by transcription factors in Drosophila development. *Nucleic Acids Res*
809 **41**, 8237-8252 (2013).
- 810 47. Panigrahi, A. & O'Malley, B.W. Mechanisms of enhancer action: the known and the unknown.
811 *Genome Biol* **22**, 108 (2021).
- 812 48. Urrutia, E., Chen, L., Zhou, H. & Jiang, Y. Destin: toolkit for single-cell analysis of chromatin
813 accessibility. *Bioinformatics* **35**, 3818-3820 (2019).
- 814 49. Jiang, Y., Zhang, N.R. & Li, M. SCALE: modeling allele-specific gene expression by single-cell RNA
815 sequencing. *Genome Biol* **18**, 74 (2017).
- 816 50. Fornes, O. et al. JASPAR 2020: update of the open-access database of transcription factor binding
817 profiles. *Nucleic Acids Res* **48**, D87-D92 (2020).

818

Figures

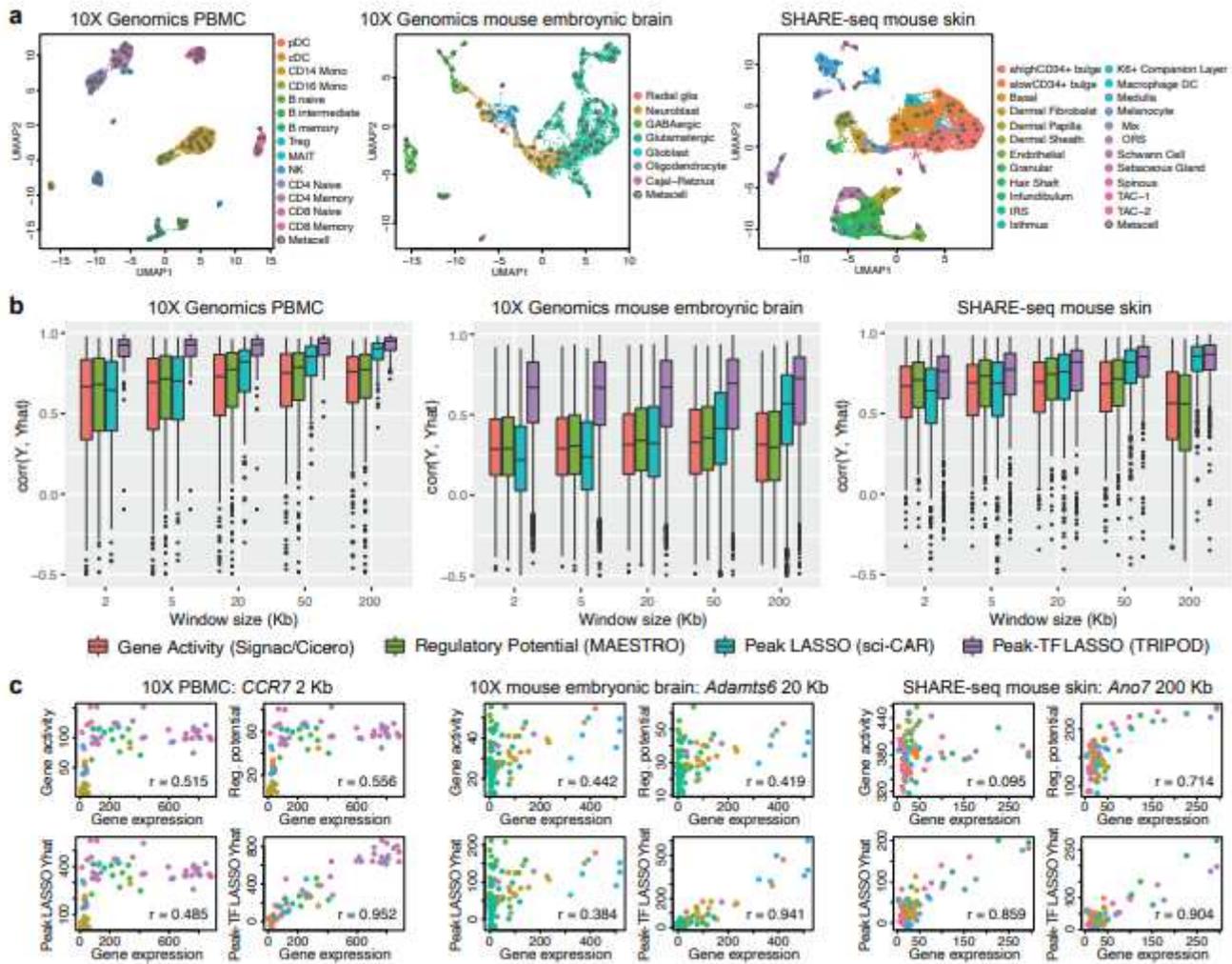


Figure 1

Figure 1

Interaction between TF expression and peak accessibility improves RNA prediction accuracy.

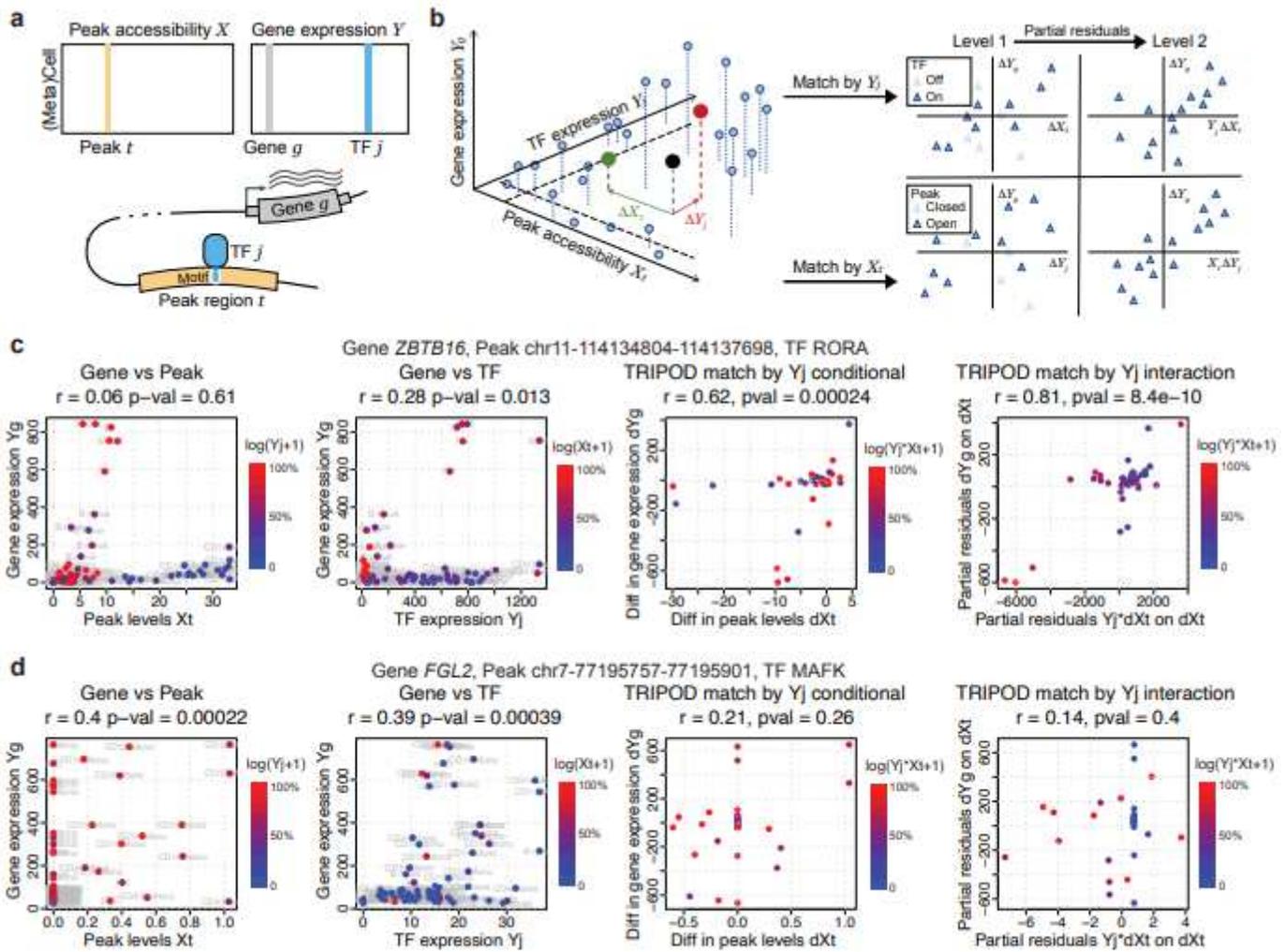


Figure 2

Figure 2

TRIPOD infers peak-TF-gene trio regulatory relationships using single-cell multiomic data.

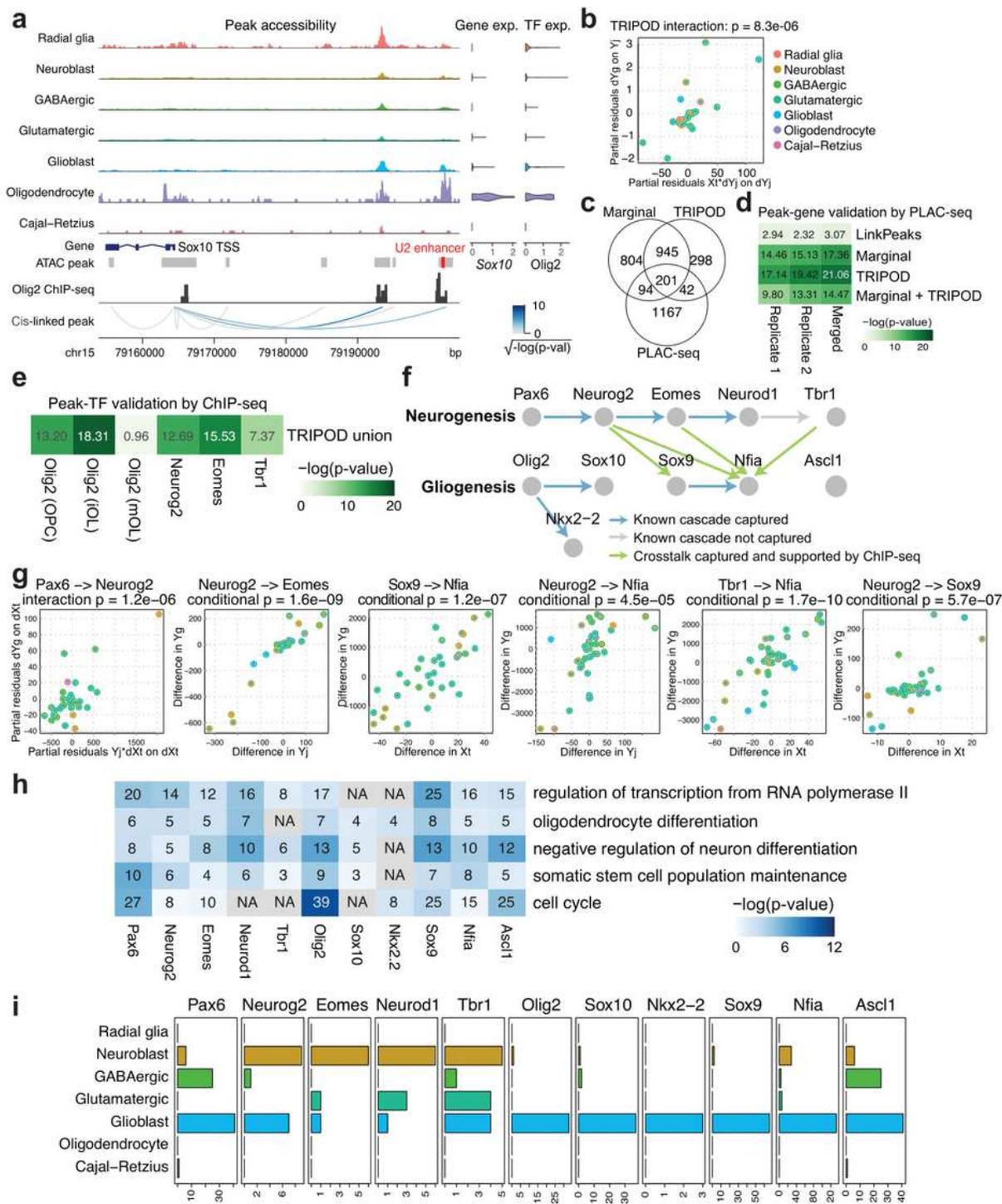


Figure 4

Figure 4

TRIPOD identified known and novel regulatory relationships during mouse neurogenesis and gliogenesis.

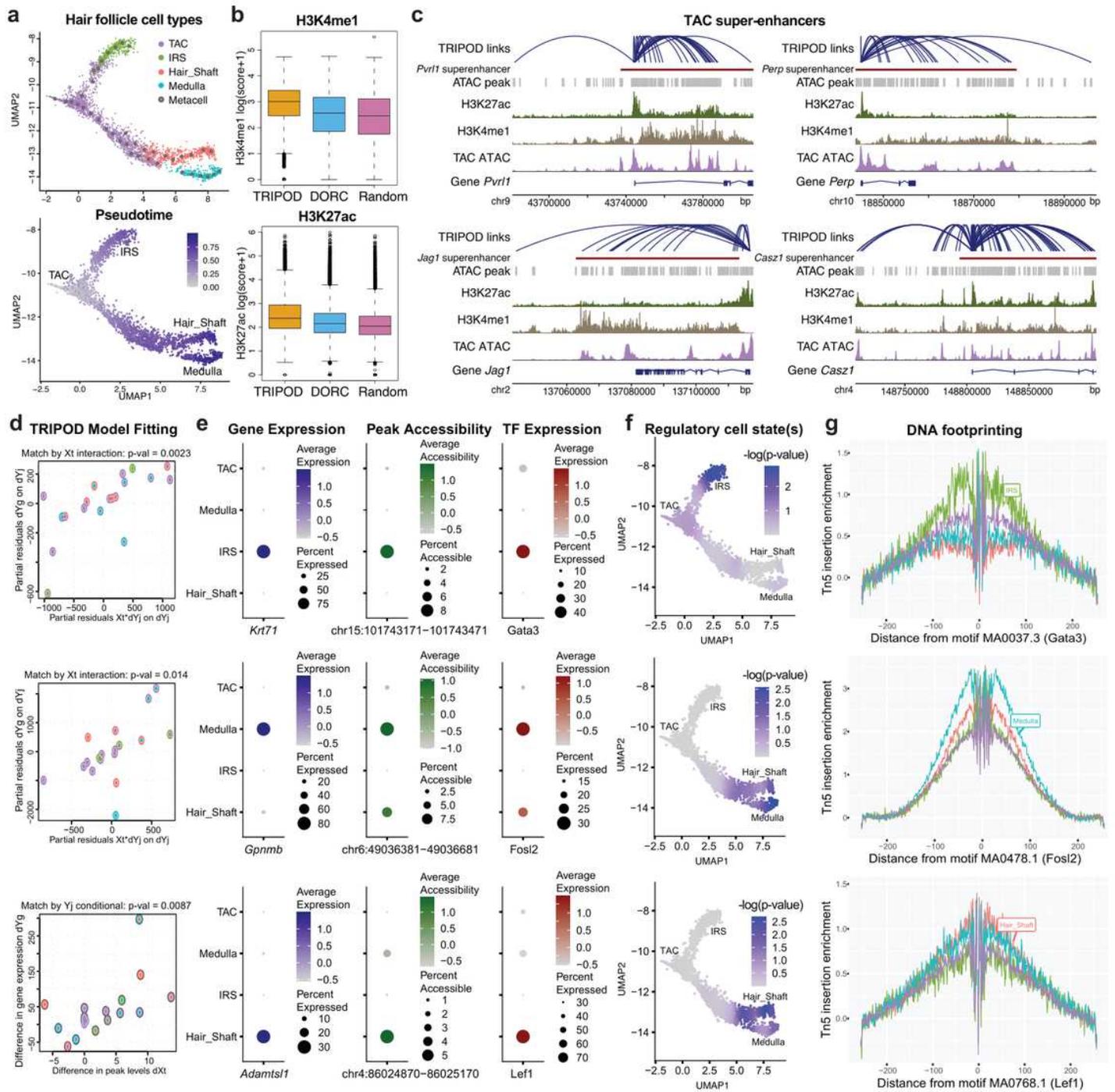


Figure 5

Figure 5

TRIPOD identified regulatory relationships in mouse hair follicles with transient cell states.

Supplementary Files

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

- [Table1validationresources.pdf](#)
- [TRIPODsupplementsv4.pdf](#)