

# 1 **The regulatory landscape of *Arabidopsis thaliana* roots at single-cell resolution**

2

3 **Authors:** Michael W. Dorrity<sup>1</sup>, Cristina M. Alexandre<sup>1</sup>, Morgan Hamm<sup>1</sup>, Anna-Lena  
4 Vigil<sup>3</sup>, Stanley Fields<sup>1,2</sup>, Christine Queitsch<sup>1</sup> and Josh Cuperus<sup>1</sup>

5

## 6 **Affiliations:**

7 <sup>1</sup> Department of Genome Sciences, University of Washington, Seattle, WA 98195

8 <sup>2</sup> Department of Medicine, University of Washington, Seattle, WA 98195

9 <sup>3</sup> School of Life Sciences, University of Nevada, Las Vegas, NV 89154

10

11 † Correspondence to Christine Queitsch (queitsch@uw.edu) and Josh Cuperus  
12 (cuperusj@uw.edu)

13 **Abstract:** In plants, chromatin accessibility – the primary mark of regulatory DNA – is  
14 relatively static across tissues and conditions. This scarcity of accessible sites that are  
15 dynamic or tissue-specific may be due in part to tissue heterogeneity in previous bulk  
16 studies. To assess the effects of tissue heterogeneity, we apply single-cell ATAC-seq  
17 to *A. thaliana* roots and identify thousands of differentially accessible sites, sufficient to  
18 resolve all major cell types of the root. However, even this vast increase relative to bulk  
19 studies in the number of dynamic sites does not resolve the poor correlation at  
20 individual loci between accessibility and expression. Instead, we find that the entirety  
21 of a cell's regulatory landscape and its transcriptome each capture cell type identity  
22 independently. We leverage this shared information on cell identity to integrate  
23 accessibility and transcriptome data in order to characterize developmental  
24 progression, endoreduplication and cell division in the root. We further use the  
25 combined data to characterize cell type-specific motif enrichments of large  
26 transcription factor families and to link the expression of individual family members to  
27 changing accessibility at specific loci, taking the first steps toward resolving the direct  
28 and indirect effects that shape gene expression. Our approach provides an analytical  
29 framework to infer the gene regulatory networks that execute plant development.

30

## 31 **Introduction**

32

33 Single-cell genomics allows an unbiased sampling of cells during development,  
34 with the potential to reveal the order and timing of gene regulatory and gene  
35 expression events that specify cell identity and lineage. An ideal system to test the  
36 ability of single-cell genomics to provide novel insights into development is the  
37 *Arabidopsis thaliana* root: along its longitudinal axis, a single, radially-symmetric root  
38 captures developmental trajectories for several radially-symmetric cell types.  
39 Approaches in this organism have included single-cell RNA-seq to transcriptionally  
40 profile individual root cell types along this developmental axis<sup>1-6</sup> and with respect to  
41 their ploidy.

42

43 Studies of chromatin accessibility in samples enriched for specific plant cell  
44 types have revealed: (i) the existence of cell type-specific regulatory elements; (ii) the  
45 relative scarcity of such elements compared to their prevalence in animals or humans;  
46 (iii) the expected enrichment of transcription factor binding sites within these elements;  
47 and (iv) a higher frequency of dynamic regulatory elements upstream of  
48 environmentally-responsive genes than constitutively expressed genes.<sup>7,8</sup> Although the  
49 correlation between chromatin accessibility and nearby gene expression is generally  
50 weak in both plants and animals,<sup>9</sup> this correlation improves for regulatory elements that  
51 show dynamic changes in chromatin accessibility, for example in response to an  
52 environmental stimulus or developmental signal.<sup>7,9-11</sup> In contrast to animals, however,  
53 the majority of chromatin-accessible sites in plants show little change across tissues,  
54 conditions, or even genetic backgrounds, raising the possibility that cell and tissue  
55 identity is less rigidly engrained in the chromatin landscape in plants than in animals.<sup>7</sup>  
56 Alternatively, cell type-specific regulatory elements and gene expression in plants may  
57 have been obscured by tissue heterogeneity in bulk tissue studies.

58  
59 Cell type-specific chromatin-accessible landscapes are also of interest for  
60 addressing other fundamental biological questions. General transcription decreases  
61 along a cell type's developmental trajectory while expression of cell type-specific  
62 genes increases,<sup>2,12,13</sup> in agreement with Waddington's predictions on epigenetic  
63 landscapes.<sup>14</sup> In the *A. thaliana* root, the increasing maturity of certain cell layers is  
64 accompanied by endoreduplication. The presence of additional gene copies may  
65 contribute to the observed increase in the expression of cell type-specific genes;  
66 alternatively, the initial gene copies may increase their transcription. Although  
67 endoreduplication is a common mechanism to regulate cell size and differentiation in  
68 plants and some human and animal tissues,<sup>15-17</sup> the influence of this phenomenon on  
69 gene regulation and expression has been largely overlooked. In plants,  
70 endoreduplication generally enhances transcription,<sup>17,18</sup> in particular of cell wall-related  
71 genes<sup>19</sup> and genes encoding ribosomal RNA,<sup>20</sup> hinting at a role for this process in  
72 driving increased translation.

73  
74 Here, we provide the first single-cell resolution maps of open chromatin in the *A.*  
75 *thaliana* root to address the issue of tissue heterogeneity and to detect likely  
76 endoreduplication events. We use a droplet-based approach to profile over 5000 nuclei  
77 for chromatin accessibility and identify 8000 regulatory elements that together define  
78 most cell types of the root. We describe an analytical framework that links patterns of  
79 open chromatin with transcriptional states to predict the identity, function and  
80 developmental stage of individual cells in the *A. thaliana* root. We integrate the single-  
81 cell ATAC-seq (scATAC-seq) data with published single-cell RNA-seq (scRNA-seq)  
82 profiles of the same tissue to obtain automated annotations of cells in our scATAC-seq  
83 data. Using the integrated dataset, we link individual cells from our scATAC-seq data  
84 with their nearest neighbors in scRNA space to define relative developmental  
85 progression, level of endoreduplication and the genes differentially expressed in these  
86 nearest neighbors. This approach allows the identification of three distinct

87 developmental states of endodermis cells that had escaped detection using scRNA-  
88 seq alone. Using scATAC-seq data integrated scRNA-seq data, we predict individual  
89 members of large transcription factor families that play a role in epidermis  
90 development, pinpointing individual regulatory events that link peak accessibility and  
91 transcription factor expression in these cells. The combination of binding motifs,  
92 transcription factor expression and chromatin accessibility provides a basis for  
93 predicting the gene regulatory events that underlie development.

94

## 95 **Results**

96

### 97 **scATAC-seq identifies known root cell types**

98

99 We first asked if ATAC-seq profiles at the single-cell level were capable of  
100 capturing known root cell types. We profiled 5283 root nuclei, at a median of 7290  
101 unique ATAC inserts per nucleus. A high fraction of these inserts occurred in one of the  
102 22,749 open chromatin peaks (FRIP score = 0.71) based on pseudo-bulk peak calling  
103 (Cellranger v3.1, 10X Genomics); this fraction is similar to that seen in high-quality bulk  
104 accessibility studies (**Figure S1A, S1B**).<sup>9</sup> Furthermore, the scATAC assay detected  
105 1794 peaks that not been observed at appreciable levels in bulk ATAC-seq. We used  
106 UMAP dimensionality reduction of the peak by cell matrix to build a two-dimensional  
107 representation grouping of cells with similar accessibility profiles (**Figure 1A**).  
108 Subsequent cluster assignment by Louvain community detection identified nine  
109 distinct cell clusters.<sup>21</sup> Across all cells, we identified 4389 peaks (ranging from 307 –  
110 1993 per cell type) with significant differential accessibility, suggesting that around  
111 20% of all accessible sites contain some information on cell type (**Supplementary**  
112 **Table 1**). Though only 16% (707/4,389) of cell type-specific peaks were found to be  
113 distal, or greater than 400 base pair from the nearest gene, this was greater than the  
114 fraction expected by chance. Only 9.4% (2,159/22,749) of all peaks were distal,  
115 suggesting that these distal peaks are slightly (1.7x) enriched for regulatory sites that  
116 define cell identity. To assign cell type annotations to each of these clusters, we  
117 generated “gene activity” scores that sum all ATAC inserts within each gene body and  
118 400 bp upstream of its transcription start site. This approach rests on the assumption  
119 that a chromatin-accessible site in the compact *A. thaliana* genome tends to be  
120 associated with regulation of its most proximal gene.<sup>22</sup> While this assumption may not  
121 hold universally, gene activity scores offer the advantage of allowing a direct  
122 comparison to bulk ATAC-seq and single-cell RNA-seq datasets through a matched  
123 feature set. In this way, we identified genes whose accessibility signal specifically  
124 marks each cell cluster. We visualized peaks with cell type-specific accessibility by  
125 grouping cells of a similar type and “pseudo-bulking” their insert counts at each  
126 position in the genome (**Figure 1B**). Bulk and cell type-specific ATAC signal is similar  
127 to those obtained in prior whole tissue and cell type enrichment-based ATAC-seq  
128 studies for the root (**Figure S1B, S1C, S1D**).<sup>11</sup>

129

130 We used comparisons to tissue-specific genes that were identified from single-  
131 cell RNA-seq studies of the *A. thaliana* root to assign a cell type to each cluster defined  
132 by ATAC markers from “gene activity” scores.<sup>2,5,6</sup> We identified 210 genes with unique  
133 accessibility patterns across all cell types (**Supplementary Table 2**); FRIP scores,  
134 fragment lengths, and total read counts did not vary greatly across cell types (**Figure**  
135 **S1E, S1F, S1G**). For each cell type, the median number of genes with tissue-specific  
136 accessibility was 20 (range 5 to 53) (**Figure 1C**). This small number of genes is  
137 consistent with earlier studies that show few open chromatin sites that define cell type  
138 identity in *A. thaliana*.<sup>7,23</sup> Although thousands of differentially accessible sites have been  
139 found across tissue types,<sup>7</sup> accessibility differences between more closely related cell  
140 types remains largely unexplored, with the exception of root hair vs non-hair, in which  
141 few differences were found.<sup>7,11</sup> These differences, uncovered using a cell-enrichment  
142 based technology,<sup>11</sup> were replicated in the epidermal cells identified in our scATAC  
143 assay (**Figure S1C, S1D**). For three cell clusters (959 cells, or 18% of cells), we could  
144 not identify a coherent set of a markers and therefore could not annotate them (grey  
145 points, **Figure 1A**). However, all other cell clusters were manually annotated and  
146 corresponded to the major cell layers of the root (**Figure S2A**): outer layers including  
147 epidermis cortex, and a precursor of endodermis and cortex (ec pre); endodermal  
148 layers comprised of three distinct types (endo 1, 2, and 3); and the stele comprised of  
149 two main types along with a phloem type (stele phloem). Several traditional marker  
150 genes were used to facilitate annotation of root cell types (**Figure S2B-D**), as were  
151 marker genes identified in previous scRNA-seq studies (**Supplementary Table 3**). In  
152 general, scATAC marker genes did not show a strong overlap with RNA-based marker  
153 genes. Endodermis cells were an exception, as several of their scATAC marker genes  
154 (AT3G32980, AT1G61590, AT1G14580, AT3G22600, AT5G66390) were also found to  
155 be marker genes in single-cell RNA-seq studies.<sup>24</sup> While this lack of overlap makes  
156 annotation more challenging, it is consistent with the reported weak correlation of  
157 chromatin accessibility with gene expression.<sup>23,25</sup> Moreover, the finding that expression  
158 levels are not precisely predicted by nearby accessible sites suggests that accessibility  
159 can add orthogonal information about cell identity to further stratify cell types into  
160 distinct sub-types.

## 161 162 **Sequences motifs of transcription factor families associate with cell type-specific** 163 **sites of open chromatin**

164  
165 Accessibility at regulatory sites is driven by transcription factor binding and  
166 modification of local chromatin.<sup>26</sup> We examined if any of the cell type-specific  
167 accessible sites were associated with the presence of transcription factor binding  
168 motifs. To do so, we used a set of representative motifs for all *A. thaliana* transcription  
169 factor families and nearly every individual transcription factor<sup>27</sup> to tally these motif  
170 counts within all 21,889 peaks in the full scATAC-seq dataset to build a peak-by-motif  
171 matrix. As each peak can be described in terms of its relative accessibility in each of  
172 the identified cell types, we performed a linear regression for each motif to test for

173 significant association of accessibility and motif presence. Relative accessibility values  
174 were calculated by first pseudo-bulking all peak counts by cell type and then  
175 normalizing these cell type-specific peak accessibility scores to a background peak  
176 accessibility of all cells pooled together. By testing the association of motif counts and  
177 cell type-specific accessibility, we identified transcription factor binding motifs whose  
178 presence was correlated with higher accessibility in each cell type. However, because  
179 motif sequence content for individual transcription factors is redundant, we computed  
180 means across each transcription factor family.

181  
182 We found significant associations with motifs from at least one transcription  
183 factor family in all cell types (**Figure 1D**). For example, relative chromatin accessibility  
184 in epidermal cells was strongly associated (q-values ranging from  $1e-24$  to  $1e-133$ )  
185 with the presence of motifs from the WRKY transcription factor family; this family  
186 includes *TTG2*, which, along with *TTG1* and *GL2*, has important roles in atrichoblast  
187 fate in the epidermis.<sup>28</sup> Furthermore, the effects of each motif family on relative  
188 accessibility was sufficient to hierarchically cluster cell types according to broad tissue  
189 classes (**Figure 1D**). Based on similarities in motif associations, hierarchical clustering  
190 grouped all stele clusters (1, 2, and 11), epidermis and cortex (clusters 0 and 3), two  
191 endodermis clusters (4 and 10), and another endodermis cluster with epidermal  
192 precursor cells (clusters 7 and 8). That motif associations alone can distinguish among  
193 clusters and group similar ones together provides independent verification of the cell  
194 type-specific nature of the chromatin-accessible sites detected in the scATAC-seq  
195 data.

## 196 197 **Integration of scATAC-seq and scRNA-seq data improves cell type annotation**

198  
199 Because scATAC-seq data both identified known root cell types and provided  
200 novel cell identity assignments not identifiable through scRNA-seq, we addressed  
201 whether combining these two datasets results in additional insights than what could be  
202 gained from either alone. We first addressed whether both data types could be  
203 embedded in the same low-dimensional space in a manner that maintains the cell  
204 identities defined by both scATAC-seq and scRNA-seq. Such embedding assumes  
205 that the underlying cell identities represented in each dataset are similar. Although the  
206 root tissue sampled for our scATAC-seq experiment was not identical to that used in  
207 previous scRNA-seq experiments, we expected that the same major cell types were  
208 sampled in both experiments. Moreover, the data generated by both methods share  
209 “gene” as a feature, *i.e.* accessibility near or within a given gene; expression of a given  
210 gene.

211  
212 We used the anchor-based multimodal graph alignment tool from the Seurat  
213 package to find nearest-neighbor scRNA-seq matches for each cell in the scATAC-seq  
214 data.<sup>29,30</sup> In short, the tool identifies representative features (shared “anchor” genes in  
215 our case) in each dataset and looks for underlying correlation structure of those

216 features to group similar cells in a co-embedded space. We plotted all cells within the  
217 resulting co-embedded space with cell type labels from each dataset separately. Cells  
218 derived from scRNA-seq and scATAC-seq experiments were well mixed (**Figure 2A**).  
219 Moreover, we found that cells of the same type were co-localized independent of the  
220 source data (**Figure 2B, 2C**), though some separation by data type was apparent, likely  
221 owing to the imputation step of dataset integration.<sup>29</sup> This result suggests that RNA and  
222 ATAC signals, which are only poorly correlated in bulk studies, are capable of grouping  
223 cell identities when determined in individual cells of a complex tissue. We further used  
224 this co-embedded space to refine our earlier manual cell type annotations by  
225 transferring labels of neighboring scRNA cells onto the scATAC cells (**Figure S3A,**  
226 **S3B**); while most of these labels matched, the greatest number of mismatches was  
227 seen in endodermis sub-type 1. The transferred labels matched our manual  
228 annotations, and, in the case of epidermal cells, allowed us to separate a single  
229 scATAC cluster into hair and non-hair cells (**Figure 2A, Figure S3A, S3B**). Furthermore,  
230 this co-embedded space was additionally used to transfer quantitative metrics and  
231 gene expression values derived from scRNA-seq data (**Figure S3C**). The three distinct  
232 scATAC clusters that were assigned an “endodermis” label with this approach are a  
233 striking example of scATAC data yielding, within a single cell type, greater stratification  
234 of “types” than the generally richer scRNA data.

235

### 236 **Epidermal cell layers show increased levels of endoreduplication**

237

238 In contrast to scRNA-seq data, scATAC-seq data can provide insight into DNA  
239 copy number and its impact on gene regulation. DNA copy number is of special  
240 relevance in the *A. thaliana* root, as each cell layer undergoes different rates of  
241 endoreduplication.<sup>19</sup> In a diploid cell, a single accessible locus tends to show 1 or 2  
242 transposition events. In polyploid cells with higher DNA copy number, a single  
243 accessible locus could show 4, 8, or even 16 transpositions. Therefore, cells containing  
244 a large number of peaks with >1 transposition event are likely to represent  
245 endoreduplicated cells. To identify such cells, we classified each cell by the mean  
246 number of cuts it contained per peak and examined the distribution of this metric,  
247 accounting for differences in total UMI counts (see Methods), to draw a threshold  
248 above which cells were classified as likely endoreduplicated (**Figure S4A, S4B**). We  
249 found the expected trend of higher endoreduplication in the outermost cell files, with  
250 reduced prevalence in the stele (**Figure S4C**).

251

252 We then used a second method to identify endoreduplicated cells with a  
253 transcriptional signature. Instead of relying on the number of transpositions in the  
254 accessibility data directly, we instead leveraged the dataset integration described  
255 above (**Figure S3C**) to transfer scRNA-seq-based annotations to the cells in our  
256 scATAC experiment. To identify endoreduplicated cells in scRNA-seq data, we used a  
257 published set of marker genes for ploidy to generate signature scores for 2n, 4n, 8n  
258 and 16n ploidies.<sup>19</sup> With these scores, we predicted endoreduplicated cells by

259 calculating, for each cell, the ratio of the 8n signature relative to the diploid signature.  
260 Similar to the accessibility-based metric, this transcription-based approach identified  
261 endoreduplicated root cells in the expected pattern, with higher fractions in the  
262 epidermis cell layer and diminished levels in the stele (**Figure S4D, S4E**). We found  
263 these two methods of identifying endoreduplicated cells to be concordant (**Figure**  
264 **S4F**), but because the accessibility-based classification was less quantitative, we used  
265 the transcriptionally-based metric in subsequent analyses. This metric captured an  
266 abundance of tetraploid xylem cells in the stele (**Figure S4E**), consistent with previous  
267 findings.<sup>19</sup>

268

### 269 **scATAC-seq captures three distinct endodermis types representing different** 270 **developmental stages**

271

272 We dissected the three endodermis clusters in greater detail using three  
273 approaches: (i) by identifying differentially accessible sites among sub-types; (ii) by  
274 aligning these sub-types to scRNA-seq data that have been annotated for  
275 endoreduplication and developmental progression; and (iii) by determining differentially  
276 expressed genes in the nearest-neighbors to each of these endodermis sub-types in  
277 scRNA-seq space (**Figure 3A**).

278

279 We identified few differentially accessible genes (adjusted p-value < 0.05 and at  
280 least 2-fold change in accessibility) in each endodermis sub-type: 25 for the first sub-  
281 type, 24 for the second, and 17 for the third (**Figure 3A**). The low number of associated  
282 genes precluded gene set enrichment analyses, but genes uniquely accessible in sub-  
283 type 1 included transcription factors *MYB85* (AT4G22680) and *NAC010* (AT1G28470)  
284 as well as genes involved in suberization (*FAR1*, *FAR4*, *FAR5*).<sup>31</sup> Endodermis sub-type  
285 2 showed increased accessibility at *HIPP04* (AT1G2900), encoding a heavy metal-  
286 associated protein, *ANAC038* (AT2G24430), and phenylpropanoid metabolism genes.<sup>32</sup>  
287 Endodermis sub-type 3 showed strong accessibility at the *BLUEJAY* (AT1G14580)  
288 locus encoding a C2H2 transcription factor implicated in endodermis differentiation  
289 (**Figure 3B**),<sup>33</sup> as well as *MYB122* (AT1G74080) and other genes for phenylpropanoid  
290 biosynthesis (*PER22*, *PER32*, *PER72*, *BGLU32*).<sup>32</sup> We addressed whether these  
291 differentially-accessible genes show different expression patterns in endodermis cells  
292 in scRNA-seq space by mapping expression of each gene onto a subclustered set of  
293 endodermis cells combined from several scRNA-seq studies of the *A. thaliana* root.<sup>2-6</sup>  
294 The small set of marker genes identified for each scATAC sub-type showed no  
295 consistent expression pattern in the scRNA-seq data (**Figure S5A**), suggesting that  
296 other features distinguished these three sub-types.

297

298 Structure within two-dimensional embeddings of scRNA-seq and scATAC-seq  
299 data derived from developing tissues is often associated with differences in  
300 developmental progression or other asynchronous processes like the cell cycle.  
301 Furthermore, root tissue has the unique feature of being highly endoreduplicated,

302 which could also account for differences among the sub-types. To assess whether the  
303 endodermal sub-types were associated with these features, we added annotations for,  
304 developmental progression, endoreduplication and cell cycle to the combined root  
305 scRNA-seq data and used data integration (as in **Figure 2**) to test whether cells from  
306 the endodermal sub-types were associated with any of these features (**Figure S3C**).  
307

308 We assessed developmental progression with two orthogonal methods: (i)  
309 correlation with published bulk expression data taken along longitudinal sections of the  
310 root;<sup>1</sup> and (ii) a modified measure of loss in transcriptional diversity (see Methods),  
311 which correlates strongly with developmental progression in a large number of scRNA-  
312 seq datasets, including of the *Arabidopsis* root.<sup>2,34</sup> We found that the developmental  
313 progression metric as measured by loss in transcriptional diversity was strongly  
314 associated with the orthogonal correlation-based classification in cells derived from  
315 scRNA-seq alone (**Figure S5B**).<sup>34</sup> For each cell of the endodermal sub-types, we  
316 calculated the average developmental progression of its 25 nearest neighbors among  
317 root scRNA-seq cells (**Figure S5C, S5D**) and found, assigning this average to each  
318 scATAC endodermis cell, a trend of developmental progression among the endodermis  
319 sub-types (**Figure 3C**). This result was robust to changes in the number of neighbors  
320 used to identify similar cells from scRNA-seq data (**Figure S5E**). This trend was the  
321 same if we calculated the developmental progression metric based on scATAC-seq  
322 data alone (**Figure S5F**), though the correlation to the transcriptional metric was weak  
323 overall (**Figure S5G**).<sup>34</sup> Cells from sub-type 1 were the least developed, while cells from  
324 sub-type 3 tended to co-occur with the most mature endodermal cells in the co-  
325 embedded graph (**Figure 3C**). We conclude that the three endodermal sub-types  
326 primarily represent cells of differing developmental progression and that differences in  
327 chromatin accessibility are able to capture this stratification of endodermis maturity.  
328

329 Developmental progression in the root is associated with increased ploidy  
330 through endoreduplication.<sup>19</sup> Using the transcriptional-based metric for  
331 endoreduplication described above, we examined the predicted ploidy of orthogonally-  
332 classified cells derived from scRNA-seq (**Figure S5H**) and from the nearest RNA-seq  
333 neighbors of each endodermis sub-type (**Figure S5I**). We found that the younger  
334 endodermis sub-type 1 cells had mostly 2n neighbor cells, while the more mature sub-  
335 types 2 and 3 had mostly endoreduplicated neighbor cells, with similar levels in each  
336 (**Figure 3D**).  
337

338 To better understand the differing transcriptional and chromatin accessibility  
339 patterns among endodermis sub-types, we analyzed differentially expressed genes  
340 from each endodermis sub-type. The early endodermis type, which is not yet  
341 endoreduplicated, showed an enrichment of genes (**Supplementary Table 4**) involved  
342 in Casparian strip formation (*CASP3*, *CASP5*) and wax biosynthesis (*HHT1*). The  
343 intermediate sub-type 2 also showed enrichment for genes involved in Casparian strip  
344 formation (*CASP3*, *CASP4*, *CASP5*, *GSO1*), as well as mechanosensitive ion channels

345 (*MSL4*, *MSL6*, *MSL10*) (**Supplementary Table 5**). The most advanced endodermis  
346 sub-type 3 showed enrichment for stress responses and metabolism of toxic  
347 compounds, kinase activity, and aquaporin water channels (**Supplementary Table 6**),  
348 consistent with this mature endodermis cell type modulating water permeability via  
349 aquaporins as well as through suberization.<sup>35</sup> We also identified putative regulators of  
350 these stages by looking for transcription factors among the genes that showed  
351 specificity for each endodermis cluster. The earlier endodermis type showed a single  
352 upregulated transcription factor, *ERF54*, while the intermediate sub-type showed 14  
353 upregulated transcription factors, including *KNAT7*, *SOMNUS*, and *HAT22*. *MYB36*,  
354 which was found expressed in the later endodermis type, activates genes involved in  
355 Casparian strip formation and regulates a crucial transition toward differentiation in the  
356 endodermis.<sup>36</sup> Because *MYB36* regulates early steps of endodermis differentiation,<sup>3,36</sup>  
357 this result suggests that some more mature endodermis types may be absent in these  
358 data, perhaps due to technical differences in their ability to be lysed during nuclear  
359 extraction (see Methods).

360  
361 We used a list of known cell-cycle marker genes (Arabidopsis.org) to generate a  
362 signature score marking proliferating cells. This signature score identified cycling cells  
363 in other cell types, such as early epidermis cells near the quiescent center (**Figure S6A**,  
364 **S6B**) in a meta-analysis of previously published scRNA-seq data. However, when this  
365 signature score was transferred to the scATAC-seq endodermis clusters by the nearest  
366 neighbor procedure described in Figure S3C, we observed no differences  
367 corresponding to each endodermis sub-type (**Figure S6C**). We conclude that cell cycle  
368 does not distinguish the endodermis sub-types.

369  
370 Overall, the combined information gained from transcriptional signatures of  
371 developmental progression and endoreduplication highlights the importance of  
372 integrating both open chromatin and transcriptional profiling to identify cell types or  
373 cell states that may have otherwise been obscured in a single data type.

### 374 375 **Predicting regulatory events using integrated scRNA and scATAC data**

376  
377 We previously identified transcription factor binding motifs that were enriched at  
378 cell type-specific peaks in the root (**Figure 1D**). While individual motifs may be  
379 associated with binding and activation by transcription factors, a sequence-level  
380 analysis cannot distinguish among the many members of plant transcription factor  
381 families that share near-identical sequence preferences. For example, WRKY family  
382 motifs were highly enriched among epidermis and cortex accessible sites, but this  
383 family contains >50 individual genes. In order to narrow down this list of genes to a few  
384 possible candidates, we leveraged our nearest-neighbor annotation approach (**Figure**  
385 **S3C**) to examine expression levels of all WRKY family transcription factors in the  
386 scATAC data (**Figure 4A**). Overall, we found that the majority of WRKY members  
387 showed expression in the epidermis, cortex or epidermal precursor cells (**Figure 4A**),

388 though some members showed stele-specific expression. To identify the most likely  
389 members to bind the abundance of motifs in epidermis-specific peaks, we ranked  
390 these genes by their specificity in the epidermis. The top four most epidermis specific  
391 genes, *WRKY75*, *WRKY9*, *WRK6*, and *TTG2* (**Figure 4A**), have documented roles in  
392 root development.<sup>28,37-39</sup> *TTG2* showed strong specificity for the epidermis, but we also  
393 predict expression in some cortex and precursor cells (**Figure 4B**). Two key interacting  
394 factors of *TTG2* that also contribute to epidermis development, *GL2* and *TTG1*,<sup>40,41</sup>  
395 showed epidermis expression and had correlated patterns (Pearson correlation with  
396 *TTG2* across cells for *GL2* = 0.91, and *TTG1* = 0.47) across all cells (**Figure S7A, S7B**).  
397

398 Given the important role of *TTG2* in specification of atrichoblast fate in the  
399 epidermis, we examined the consequences of its expression on accessibility of  
400 individual peaks. Inference of individual regulatory events, particularly those involving  
401 transcription factors, has long been a goal of studies that profile accessibility at  
402 regulatory sites in bulk tissue. The varied cell states revealed by single-cell profiling  
403 data, even those within a cell type, allow higher-resolution inference of these events.  
404 To identify accessible sites that showed altered accessibility as a function of  
405 transcription factor expression, we used a linear regression approach. We identified  
406 617 peaks that showed significant (q-value < 0.05) associations with *TTG2* expression  
407 levels (**Supplementary Table 7**). To visualize these associations using scATAC data,  
408 we pseudo-bulked epidermis, cortex, and c/e precursor cells into four equal-sized bins  
409 based on their level of *TTG2* expression (**Figure 4C**). We observed peaks whose  
410 accessibility increases (**Figure 4C**, top and lower-left panels) and decreases (**Figure**  
411 **4C**, lower-right panel) in cells with increasing levels *TTG2* expression. Most significant  
412 associations were positive, such that increased *TTG2* expression led to increased peak  
413 accessibility (**Figure 4D**). Using DAP-seq data for *TTG2*, we examined whether peaks  
414 with either positive or negative associations contain *TTG2* binding sites.<sup>27</sup> Positive  
415 associations occurred whether or not a WRKY binding motif was present in the  
416 associated peak (**Figure 4C, 4D**), suggesting that the role of WRKY transcription  
417 factors in specification of the epidermis likely requires both direct and indirect  
418 regulatory events. Of peaks with significant (q-value < 0.05) positive associations with  
419 *TTG2* expression, 80% of these contained a WRKY binding motif, while only 38% of  
420 the peaks with negative associations contained a binding motif (**Figure 4D**). Overall,  
421 this analysis identifies transcription factors and putative target sites that constitute  
422 regulatory events important for specifying cell types; these genes and regulatory sites  
423 are good candidates for further functional studies.  
424

## 425 Discussion

426  
427 By profiling chromatin accessibility in the *A. thaliana* root at single-cell  
428 resolution, we assessed cell types, developmental stages, the transcription factors  
429 likely driving these stages and DNA copy number changes. We assigned over 5,000  
430 root cells to tissues and cell types, demonstrating that these assignments are

431 concordant with single-cell transcriptomic studies. These results answer an unresolved  
432 question in plant gene regulation: does the paucity of dynamic open chromatin sites  
433 seen in bulk profiling experiments represent an accurate reflection of uniform gene  
434 regulation in *A. thaliana* or does it reflect a confounding effect of bulk studies? We  
435 found that distinct root cell types show unique patterns of open chromatin sites, with  
436 approximately 1/3 of all accessible sites showing cell type-specific patterns. This  
437 estimate greatly exceeds the earlier estimates from bulk studies of only 5-10% of  
438 accessible sites showing tissue- or condition-specificity,<sup>9</sup> presumably due in part to  
439 tissue heterogeneity.

440  
441 Although this single-cell ATAC study discovered many more dynamic accessible  
442 sites, the correlation between dynamic accessibility and gene expression in single cells  
443 remained poor, reminiscent of the equally poor correlation seen in bulk studies. These  
444 data types would be integrated more faithfully in a true co-assay experiment.<sup>25,42</sup>  
445 Technical differences in nuclei versus cell-based assays, size selection, developmental  
446 stage, and sequencing depth may also contribute to differences between scRNA and  
447 scATAC datasets. While increasing the depth of our ATAC signal per cell may alleviate  
448 some of this noise, we argue that the poor correlation between chromatin accessibility  
449 and gene expression is not a function of data quality. Instead, we propose that this  
450 weak correlation reflects the complex nature of regulatory processes underlying  
451 development, and the differential aspects of regulation captured in scATAC-seq and  
452 scRNA-seq data, which were notably divergent in the scATAC-specific endodermis  
453 sub-types. Although the correlation of chromatin accessibility and gene expression is  
454 weak at the level of individual loci, either the entirety of a cell's regulatory landscape or  
455 its transcriptome can independently capture its cell identity. It is this feature that allows  
456 joint co-embedding of both data types and the use of scRNA-seq data to annotate  
457 scATAC cells.

458  
459 Thus, while the patterns of both chromatin accessibility and gene expression  
460 contain information on cell identity and development, the relationships between these  
461 patterns are not well-ordered or parsimonious. For the many cells belonging to a  
462 distinct cell type, gene expression results from direct and indirect regulatory events  
463 involving tens or hundreds of transcription factors and chromatin remodelers that do  
464 not necessarily act in concert. For any individual locus, then, the expectation that  
465 average accessibility predicts average expression breaks down. Without a simple one-  
466 to-one model to explain regulatory output, we are left with significant heterogeneity  
467 within and between cell types, and a subset of convergent expression or accessibility  
468 patterns that define cell type specificity. Alternative explanations for the discrepancy in  
469 accessibility and expression include: (1) maintenance of cell identity requires that a  
470 cell's accessibility and expression profile stably reflect the convergent pattern for that  
471 cell type only a fraction of the time; and/or (2) cells have multiple accessibility and  
472 expression patterns that are sufficient to maintain cell identity and together constitute  
473 the convergent patterns we observe. In both scenarios, the heterogeneity in cell type  
474 specification will be buffered by factors outside chromatin accessibility or gene

475 expression, such as spatial location in tissue, metabolic determinants of cell function or  
476 developmental age.

477

478 We posit that scATAC-seq data combined with scRNA-seq data will ultimately  
479 resolve these alternatives by enabling mechanistic models of gene regulatory  
480 networks. scATAC-seq data alone are sufficient to identify the full set of accessible  
481 sites in the *Arabidopsis* genome, and examination of the transcription factor motifs  
482 within these sites can enable predictions of regulatory networks. However, many plant  
483 transcription factor families are large, some containing over fifty members that  
484 recognize near identical motifs. Thus, the accessibility data must be integrated with  
485 single-cell expression data that capture cell type-specific expression of transcription  
486 factors in order to narrow down the most probable transcription factors that are  
487 enacting individual regulatory events. The simple regression framework provided in this  
488 work is only a small step toward more complicated models that capture other relevant  
489 sources of heterogeneity. Building higher resolution models of key regulatory events  
490 will require the expression level of individual transcription factors in a cell type, the  
491 accessibility of individual peaks in this cell type and the presence of binding motifs  
492 corresponding to the relevant transcription factors. Theoretically, a comprehensive  
493 capture of cell states with both open chromatin and transcriptional profiling will allow  
494 the ordering of gene regulatory events and the larger scale ordering of regulatory  
495 programs that underlie development. The ability to take single-cell measurements over  
496 distinct developmental stages will also increase the sampling of key regulatory events.  
497 Ultimately, achieving the goal of building models of gene regulatory events underlying  
498 development will require ever larger datasets to fully capture the range of possible cell  
499 states.

500

501 In the future, single-cell studies of more complex plant tissues in crops and other  
502 species will necessitate larger numbers of profiled cells and higher numbers of cuts per  
503 cell. Deeper coverage in future datasets should enhance our ability to detect rare cell  
504 types and more confidently predict copy number from accessibility data alone. In this  
505 way, approaches that maximize the number of cells profiled at low cost, such as  
506 single-cell combinatorial indexing,<sup>43</sup> will be critical. Annotation in future studies will also  
507 present a substantial challenge if a rich literature and genomic analyses, including  
508 single-cell transcriptome profiles, are not available. Nevertheless, as shown in this  
509 proof-of-principle study of the well-characterized *A. thaliana* root, the knowledge  
510 gained should eventually allow us to manipulate gene expression and organismal  
511 phenotype in a targeted manner.

512

513

## 514 **Methods**

515

### 516 *Plant Material*

517 *Genotype:* *Arabidopsis thaliana* ecotype Col-0 INTACT line *UBQ10::NTF::ACT2::BirA*  
518 (available from ABRC, stock CS68649). *Growth conditions:* LD (16h light/8h dark), 22C,

519 ~100  $\mu$ mol m<sup>2</sup>s, 50% RH. *Sample*: whole roots, harvested 12 days after germination,  
520 from seedlings grown vertically on MS + 1% sucrose, atop filter paper (to facilitate root  
521 harvesting).

522

### 523 *Nuclei Isolation and scATAC-seq*

524 Nuclei were isolated following a modified version of the protocol described in Giuliano  
525 *et al.*, 1988, as follows: 1g of roots was split in two batches of 0.5g, and each batch  
526 chopped with a razor blade in 1 ml of Buffer A (0.8M sucrose, 10mM MgCl<sub>2</sub>, 25mM  
527 Tris-HCl pH 8.0 and 1x Protease Inhibitor Tablet).<sup>44</sup> Extracts were combined, final  
528 volume increased to 5ml with Buffer A, and incubated on ice for 10min, with gentle  
529 swirling. The combined extract was filtered through miracloth, passed through a 26ga  
530 syringe five times and re-filtered through a 40 $\mu$ m cell strainer (BD Falcon). After  
531 centrifugation at 2,000g 5min, the pellet was resuspended in 1ml Buffer B (0.4M  
532 sucrose, 10mM MgCl<sub>2</sub>, 25mM Tris-HCl pH 8.0, 1x Protease Inhibitor Tablet, 1% Triton  
533 X - 100) and loaded atop a 2-step 25/75 Percoll gradient ( 1 volume 25% Percoll in  
534 Buffer B over 1 volume 75% Percoll in Buffer B). After centrifugation at 2,500g for  
535 15min, nuclei were collected either at the 25/75 interface or in the subjacent 75  
536 fraction, washed with 5 vols of Buffer B and recovered by centrifugation at 1,700g for  
537 5min. The nuclei pellet was resuspended in 100 $\mu$ l Buffer B + 1% BSA and any nuclei  
538 clumps broken down by pipetting up and down multiple times. Nuclei yield with this  
539 protocol was ~ 94,000 nuclei per gram of roots (fresh weight).

540 scATAC-seq libraries were built using the 10x Genomics Chromium Single Cell ATAC  
541 Solution platform, following manufacturer's recommendations. Before transposition,  
542 nuclei were spun 5min at 1,500g and resuspended in 10x Genomics Diluted Nuclei  
543 Buffer, at a concentration of 3,200 nuclei/ $\mu$ l. 5 $\mu$ l of nuclei suspension were used for  
544 transposition (16,000 nuclei being the maximum input recommended for 10x  
545 Chromium, and 10,000 nuclei being the expected recovery).

546

### 547 *Combining and processing of root scRNA-seq data*

548 Samples were processed using the CellRanger v1.2.0 pipeline from 10X Genomics,  
549 including updated filtering of "halflet" cells that emerge due to multiply-barcoded  
550 droplets.

551

### 552 *Integration of scRNA and scATAC data*

553 The R package Seurat version 3.1.5 was used to align and co-embed the scATAC-seq  
554 data with scRNA-seq data published by Ryu *et al.* 2019, and to transfer cell type labels  
555 from the scRNA data to the scATAC data.<sup>30,45</sup>

556

557 The standard workflow and default parameters as described in the Seurat vignette  
558 "PBMC scATAC-seq Vignette" ([satijalab.org/seurat/v3.1/atacseq\\_integration\\_vignette](https://satijalab.org/seurat/v3.1/atacseq_integration_vignette))  
559 were used with the exception that all features (genes) were used when identifying  
560 transfer anchors and performing the co-embedding rather than a set of "variable"  
561 features as used in the vignette. Briefly this workflow is as follows:

562 An anchor set was established with the function FindTransferAnchors() linking the two  
563 datasets. Cell type annotations were transferred from the scRNA-seq data to the  
564 scATAC data using the function TransferData(). Imputed RNA-seq count data was  
565 generated for the scATAC cells, again using the TransferData() function. The imputed  
566 RNA data was then merged with the true scRNA-seq dataset and embedded in 2D  
567 UMAP space using Seurat functions.<sup>29</sup>

568  
569 A co-embedding was performed with a super-set of previously published scRNA-seq  
570 data.<sup>2,3,5</sup> In the co-embedded space the scATAC-seq were found to be most closely  
571 co-located with data from root tips.<sup>5</sup> Based on this observation co-embedding was  
572 performed with solely with root tip dataset.<sup>5</sup>

573  
574 *Nearest neighbor analysis for transcriptional characterization of cells identified in*  
575 *scATAC assay*

576  
577 To annotate cells from the scATAC-seq assay with transcriptional features, we used  
578 average feature values from the nearest RNA neighbors in our co-embedded data  
579 (**Figure 2A**). In short, the 'distances' package in R was used to extract cell labels for  
580 the 25 nearest neighbors of each scATAC cell. For a feature of interest (individual gene  
581 expression, cell-cycle signature score, endoreduplication signature score,  
582 developmental progression signature), we calculated the mean expression from the 25  
583 scRNA cells, and assigned that mean score to each ATAC cell (**Figure S3C**).

584  
585 *Endoreduplication signatures*

586  
587 We identified endoreduplicated cells using two different approaches, the first using  
588 scRNA data, and the second using scATAC data. In the first approach (as in **Figure**  
589 **3D, Figure S4D, S4E, Figure S5B, S5I**), validated sets of endoreduplication markers  
590 for 2N, 4N, and 8N cells were used to identify endoreduplicated cells in the scRNA  
591 data.<sup>19</sup> We used the nearest neighbor approach described above to transfer this  
592 transcriptional signature to scATAC cells. The average expression of each gene group  
593 was computed for each individual cell, and subsequently averaged per cluster to  
594 generate cell type-specific levels of each ploidy signature. To identify clusters that  
595 were more likely to be endoreduplicated, rather than typical diploid cells, we examined,  
596 for each cluster, the ratio of the endoreduplicated signatures (4N or 8N) relative to the  
597 diploid (2N) signature. Clusters with a higher ratio are more likely to represent  
598 endoreduplicated cells. In the second approach (as in **Figure S4A-C**), the number of  
599 transposition events derived from scATAC data were used directly to identify  
600 endoreduplicated cells. We assumed that cells containing higher than average cuts per  
601 peak were more likely to be endoreduplicated, as the cut counts for a single peak in a  
602 diploid cell should rarely be above two. A peak with a cut count >2 may indicate an  
603 extra copy of the locus present in that cell. To identify cells more likely to be  
604 endoreduplicated, then, we examined the distribution of cuts per peak for all cells, but

605 found this metric was strongly correlated with total UMIs per cell. To account for  
606 contribution total UMIs per cell, we used the relationship between the cuts/feature and  
607 total UMIs per cell to compute a Loess model fit (**Figure S4B**). We then used residuals  
608 of this model as a metric to identify cells that have higher cuts/feature than would be  
609 expected based on their total UMIs. We set an arbitrary threshold of >1 SD in the  
610 distribution of each cell's deviation from the fit line, and defined endoreduplicated cells  
611 as those beyond the threshold (**Figure S4B**). For each cell, a binary designation of  
612 endoreduplication was applied based on whether the cell crossed this threshold.

613

#### 614 *Transcriptional diversity metric for developmental progression*

615

616 Using the general premise that the number of unique genes expressed (transcriptional  
617 complexity) tends to be reduced across the developmental trajectory of a cell type as it  
618 moves from earlier to later stages,<sup>13</sup> we devised a metric to approximate relative  
619 differences in developmental progression among cells. Measuring the number of  
620 unique genes expressed is distinct from measuring the number of UMIs or transcripts  
621 captured per cell, which can vary across cell types. To account for differential recovery  
622 of UMIs across cells in the transcriptional complexity measure, we modeled as a Loess  
623 fit the relationship between total UMIs captured and the number of unique genes  
624 expressed per cell. With this fit, we identified cells that have many more or fewer  
625 unique genes expressed than would be expected for cells over a range of captured  
626 UMIs. Developmental progression for each cell was defined as the residual of each  
627 point in this fit, allowing separation of earlier cells (more unique genes expressed than  
628 would be expected for a given number of captured UMIs) from later cells (fewer unique  
629 genes expressed than would be expected for a given number of captured UMIs).

630

#### 631 *Motif analysis*

632

633 Position weight matrices from the comprehensive DAP-seq dataset<sup>27</sup> were used as  
634 input into FIMO<sup>46</sup> to search for significant matches for each individual TF motif  
635 (adjusted p-value threshold < 1e-5) in each of the scATAC peaks. With the output of  
636 this motif scan, we generated a matrix that tallied counts of each individual motif within  
637 each peak. Each individual motif in the DAP-seq dataset<sup>27</sup> has an associated TF family,  
638 and the counts per peak were averaged by family. To identify motifs whose counts  
639 were significantly associated with cell type-specific accessibility, we first generated, for  
640 each peak, a relative accessibility score by taking the mean accessibility of that peak in  
641 each cell cluster relative to the overall accessibility of that peak in all clusters. Next, we  
642 used a linear regression framework within Monocle3<sup>47</sup> to identify individual motifs  
643 whose counts showed strong positive or negative correlations with the cell type-  
644 specific accessibility score in each cell cluster. The effect size of each motif's  
645 contribution to cell type-specific accessibility is given as the  $\beta$  of the linear regression,  
646 shown as a mean across all transcription factors in the same family.

647

648

## 649 **Data Availability**

650 An R object containing all accessibility and predicted expression data for each cell has  
651 been deposited to Dryad (accession number pending).

## 652 **Code Availability**

653 We have provided R markdown files with code blocks sufficient to complete the  
654 primary processing of the data, generation of scATAC and scRNA co-embedding,  
655 analysis of motifs, and identification of transcription-factor mediated regulatory events.  
656 (Github repository link pending).

## 657 **Acknowledgements**

658 We thank Dr. Ken Jean-Baptiste and Dr. Kerry Bubb for valuable discussions on ATAC-  
659 seq analysis. We also thank Xavi Guitart for helpful discussions on endoreduplication.  
660 This work was supported by the National Science Foundation (RESEARCH-PGR grant  
661 17488843) to S.F. and C.Q. This work was also supported by NIH grant  
662 1RM1HG010461 to C.Q. and S.F and R01-GM079712 to C.Q. and J.C.

663

## 664 **References**

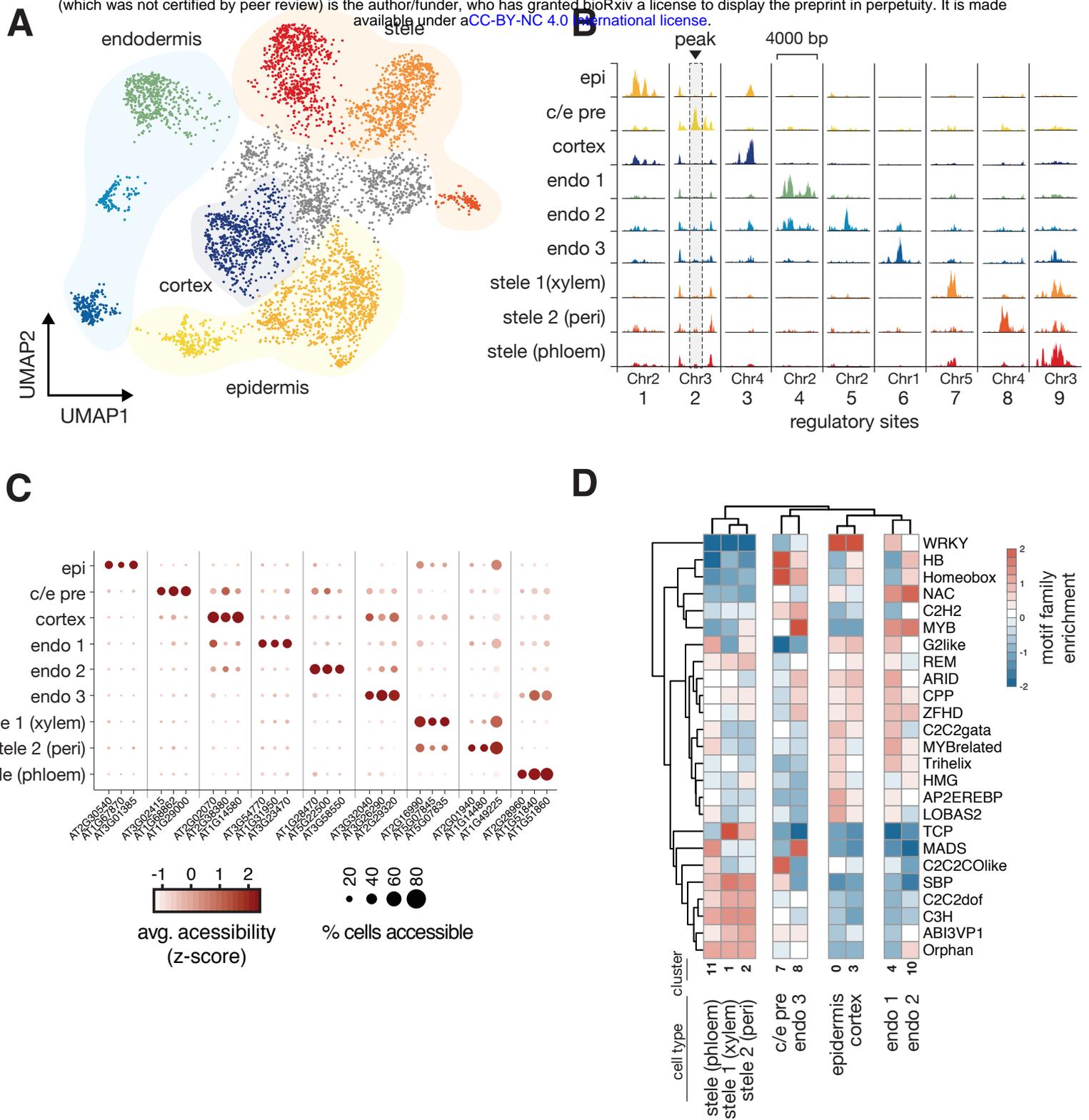
665

- 666 1. Brady, S. M. *et al.* A High-Resolution Root Spatiotemporal Map Reveals  
667 Dominant Expression Patterns. *Science* (80-. ). **318**, 801–806 (2007).
- 668 2. Jean-Baptiste, K. *et al.* Dynamics of gene expression in single root cells of *A.*  
669 *thaliana*. *Plant Cell* **31**, tpc.00785.2018 (2019).
- 670 3. Shulse, C. N. *et al.* High-Throughput Single-Cell Transcriptome Profiling of Plant  
671 Cell Types. *Cell Rep.* **27**, 2241–2247.e4 (2019).
- 672 4. Zhang, T. Q., Xu, Z. G., Shang, G. D. & Wang, J. W. A Single-Cell RNA  
673 Sequencing Profiles the Developmental Landscape of Arabidopsis Root. *Mol.*  
674 *Plant* **12**, 648–660 (2019).
- 675 5. Ryu, K. H., Huang, L., Kang, H. M. & Schiefelbein, J. Single-cell RNA sequencing  
676 resolves molecular relationships among individual plant cells. *Plant Physiol.* **179**,  
677 1444–1456 (2019).
- 678 6. Denyer, T. *et al.* Spatiotemporal Developmental Trajectories in the Arabidopsis  
679 Root Revealed Using High-Throughput Single-Cell RNA Sequencing. *Dev. Cell*  
680 **48**, 840–852.e5 (2019).
- 681 7. Sullivan, A. M. *et al.* Mapping and Dynamics of Regulatory DNA in Maturing  
682 Arabidopsis *thaliana* Siliques. *Front. Plant Sci.* **10**, 1–16 (2019).
- 683 8. Alexandre, C. M. *et al.* Complex relationships between chromatin accessibility,  
684 sequence divergence, and gene expression in *arabidopsis thaliana*. *Mol. Biol.*  
685 *Evol.* **35**, 837–854 (2018).
- 686 9. Sullivan, A. M., Bubb, K. L., Sandstrom, R., Stamatoyannopoulos, J. A. &

- 687 Queitsch, C. DNase I hypersensitivity mapping, genomic footprinting, and  
688 transcription factor networks in plants. *Curr. Plant Biol.* **3–4**, 40–47 (2015).
- 689 10. Reynoso, M. A. *et al.* Evolutionary flexibility in flooding response circuitry in  
690 angiosperms. *Science (80-. )*. **365**, 1291–1295 (2019).
- 691 11. Maher, K. A. *et al.* Profiling of accessible chromatin regions across multiple plant  
692 species and cell types reveals common gene regulatory principles and new  
693 control modules. *Plant Cell* **30**, 15–36 (2018).
- 694 12. Saunders, L. M. *et al.* Thyroid hormone regulates distinct paths to maturation in  
695 pigment cell lineages. *Elife* **8**, (2019).
- 696 13. Gulati, G. S. *et al.* Single-cell transcriptional diversity is a hallmark of  
697 developmental potential. *Science (80-. )*. **367**, 405–411 (2020).
- 698 14. Waddington, C. H. *The Strategy of the Genes.* (1959).
- 699 15. Orr-Weaver, T. L. When bigger is better: The role of polyploidy in organogenesis.  
700 *Trends Genet.* **31**, 307–315 (2015).
- 701 16. Derks, W. & Bergmann, O. Polyploidy in cardiomyocytes: Roadblock to heart  
702 regeneration? *Circ. Res.* **126**, 552–565 (2020).
- 703 17. Lang, L. & Schnittger, A. Endoreplication — a means to an end in cell growth and  
704 stress response. *Curr. Opin. Plant Biol.* **54**, 85–92 (2020).
- 705 18. Pirrello, J. *et al.* Transcriptome profiling of sorted endoreduplicated nuclei from  
706 tomato fruits: how the global shift in expression ascribed to DNA ploidy  
707 influences RNA-Seq data normalization and interpretation. *Plant J.* **93**, 387–398  
708 (2018).
- 709 19. Bhosale, R. *et al.* A spatiotemporal dna endoploidy map of the arabidopsis root  
710 reveals roles for the endocycle in root development and stress adaptation. *Plant*  
711 *Cell* **30**, 2330–2351 (2018).
- 712 20. Robinson, D. O. *et al.* Ploidy and Size at Multiple Scales in the Arabidopsis  
713 Sepal. *Plant Cell* **30**, 2308–2329 (2018).
- 714 21. Blondel, V. D., Guillaume, J.-L., Lambiotte, R. & Lefebvre, E. Fast unfolding of  
715 communities in large networks. *J. Stat. Mech. Theory Exp.* **10**, P10008 (2008).
- 716 22. Shu, H., Wildhaber, T., Siretskiy, A., Gruissem, W. & Hennig, L. Distinct modes of  
717 DNA accessibility in plant chromatin. *Nat. Commun.* **3**, (2012).
- 718 23. Sullivan, A. M. *et al.* Mapping and dynamics of regulatory DNA and transcription  
719 factor networks in *A. thaliana*. *Cell Rep.* **8**, 2015–2030 (2014).
- 720 24. McFaline-Figueroa, J. L., Trapnell, C. & Cuperus, J. T. The promise of single-cell  
721 genomics in plants. *Curr. Opin. Plant Biol.* **54**, 114–121 (2020).
- 722 25. Cao, J. *et al.* Joint profiling of chromatin accessibility and gene expression in  
723 thousands of single cells. *Science (80-. )*. **1385**, 1380–1385 (2018).
- 724 26. Spitz, F. & Furlong, E. E. M. Transcription factors: From enhancer binding to  
725 developmental control. *Nat. Rev. Genet.* **13**, 613–626 (2012).
- 726 27. O'Malley, R. C. *et al.* Cistrome and Epicistrome Features Shape the Regulatory  
727 DNA Landscape. *Cell* **165**, 1280–1292 (2016).
- 728 28. Johnson, Cameron S.; Kolevski Ben, and S. D. R. TRANSPARENT TESTA  
729 GLABRA2 , a Trichome and Seed Coat Development Gene of Arabidopsis,

- 730 Encodes a WRKY Transcription Factor. *Plant Cell* **14**, 1359–1375 (2017).
- 731 29. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888–  
732 1902.e21 (2019).
- 733 30. Butler, A., Hoffman, P., Smibert, P., Papalexi, E. & Satija, R. Integrating single-  
734 cell transcriptomic data across different conditions, technologies, and species.  
735 *Nat. Biotechnol.* **36**, 411–420 (2018).
- 736 31. Domergue, F. *et al.* Three Arabidopsis fatty Acyl-coenzyme a reductases, FAR1,  
737 FAR4, and FAR5, generate primary fatty alcohols associated with suberin  
738 deposition. *Plant Physiol.* **153**, 1539–1554 (2010).
- 739 32. Raudvere, U. *et al.* G:Profiler: A web server for functional enrichment analysis  
740 and conversions of gene lists (2019 update). *Nucleic Acids Res.* **47**, W191–W198  
741 (2019).
- 742 33. Moreno-Risueno, M. A. *et al.* Transcriptional control of tissue formation  
743 throughout root development. *Science (80-. )*. **350**, 426–430 (2015).
- 744 34. Gulati, G. S. *et al.* Single-cell transcriptional diversity is a hallmark of  
745 developmental potential. *Science (80-. )*. **367**, 405–411 (2020).
- 746 35. Kreszies, T., Schreiber, L. & Ranathunge, K. Suberized transport barriers in  
747 Arabidopsis, barley and rice roots: From the model plant to crop species. *J. Plant*  
748 *Physiol.* **227**, 75–83 (2018).
- 749 36. Liberman, L. M., Sparks, E. E., Moreno-Risueno, M. A., Petricka, J. J. & Benfey,  
750 P. N. MYB36 regulates the transition from proliferation to differentiation in the  
751 Arabidopsis root. *Proc. Natl. Acad. Sci. U. S. A.* **112**, 12099–12104 (2015).
- 752 37. Devaiah, B. N., Karthikeyan, A. S. & Raghothama, K. G. WRKY75 transcription  
753 factor is a modulator of phosphate acquisition and root development in  
754 Arabidopsis. *Plant Physiol.* **143**, 1789–1801 (2007).
- 755 38. Chen, Y. F. *et al.* The WRKY6 transcription factor modulates PHOSPHATE1  
756 expression in response to low pi stress in arabidopsis. *Plant Cell* **21**, 3554–3566  
757 (2009).
- 758 39. Zheng, X. *et al.* MdWRKY9 overexpression confers intensive dwarfing in the M26  
759 rootstock of apple by directly inhibiting brassinosteroid synthetase MdDWF4  
760 expression. *New Phytol.* **217**, 1086–1098 (2018).
- 761 40. Long, Y. & Schiefelbein, J. Novel TTG1 Mutants Modify Root-Hair Pattern  
762 Formation in Arabidopsis. *Front. Plant Sci.* **11**, 1–12 (2020).
- 763 41. Schiefelbein, J., Kwak, S. H., Wieckowski, Y., Barron, C. & Bruex, A. The gene  
764 regulatory network for root epidermal cell-type pattern formation in arabidopsis.  
765 *J. Exp. Bot.* **60**, 1515–1521 (2009).
- 766 42. Xing, Q. R. *et al.* Parallel bimodal single-cell sequencing of transcriptome and  
767 chromatin accessibility. *bioRxiv* 1027–1039 (2019). doi:10.1101/829960
- 768 43. Cao, J. *et al.* Comprehensive single-cell transcriptional profiling of a multicellular  
769 organism. *Science (80-. )*. **357**, 661–667 (2017).
- 770 44. Giuliano, G. *et al.* An evolutionarily conserved protein binding sequence  
771 upstream of a plant light-regulated gene. *Proc. Natl. Acad. Sci. U. S. A.* **85**,  
772 7089–7093 (1988).

- 773 45. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data Resource  
774 Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-1902.e21 (2019).  
775 46. Grant, C. E., Bailey, T. L. & Noble, W. S. FIMO: Scanning for occurrences of a  
776 given motif. *Bioinformatics* **27**, 1017–1018 (2011).  
777 47. Cao, J. *et al.* The single-cell transcriptional landscape of mammalian  
778 organogenesis. *Nature* **566**, 496–502 (2019).  
779



**Figure 1. scATAC-seq identifies known root cell types.** (A) UMAP dimensionality reduction plot of root cells using peak-level scATAC data. Cells are colored according to Louvain clusters, and broad tissue types are indicated with transparent shading. (B) Pseudo-bulked peak tracks generated by combining ATAC data from all cells within a cluster. Each column represents a single locus in the genome that shows cell type-specific accessibility; each row represents a cell type, and each column shows an example marker peak for each type. Colors match those in previous panel. A cluster residing between the epidermis and endodermis clusters, with expression of markers from both cell types (**Figure S2B, S2C**) was given the label 'c/e pre' (precursor of cortex/endodermis, second row), and epidermis was shortened to 'epi'. (C) Dotplot showing marker genes for each cell type cluster. Each column represents a single gene's activity score, the summed accessibility of its gene body and promoter sequence (-400bp from transcription start site). The color of each dot indicates the magnitude of accessibility and the size of each dot represents the fraction of cells in each cell type showing accessibility at that gene. (D) Heatmap showing the predicted effect, across all peaks, of motifs from each *Arabidopsis* transcription factor family on cell type-specific accessibility. Darker shades of red indicate that presence of the motif is correlated with increased accessibility in that cell type, whereas shades of blue indicate that the motif is anti-correlated with accessibility. The mean effect all transcription factors within a given family are shown as rows, and each column represents a cell type.

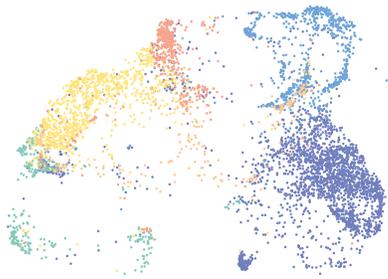
**A**

RNA+ATAC cells  
co-embedded



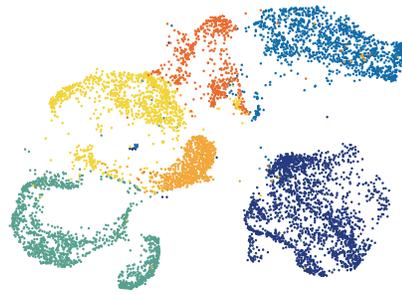
**B**

ATAC cells



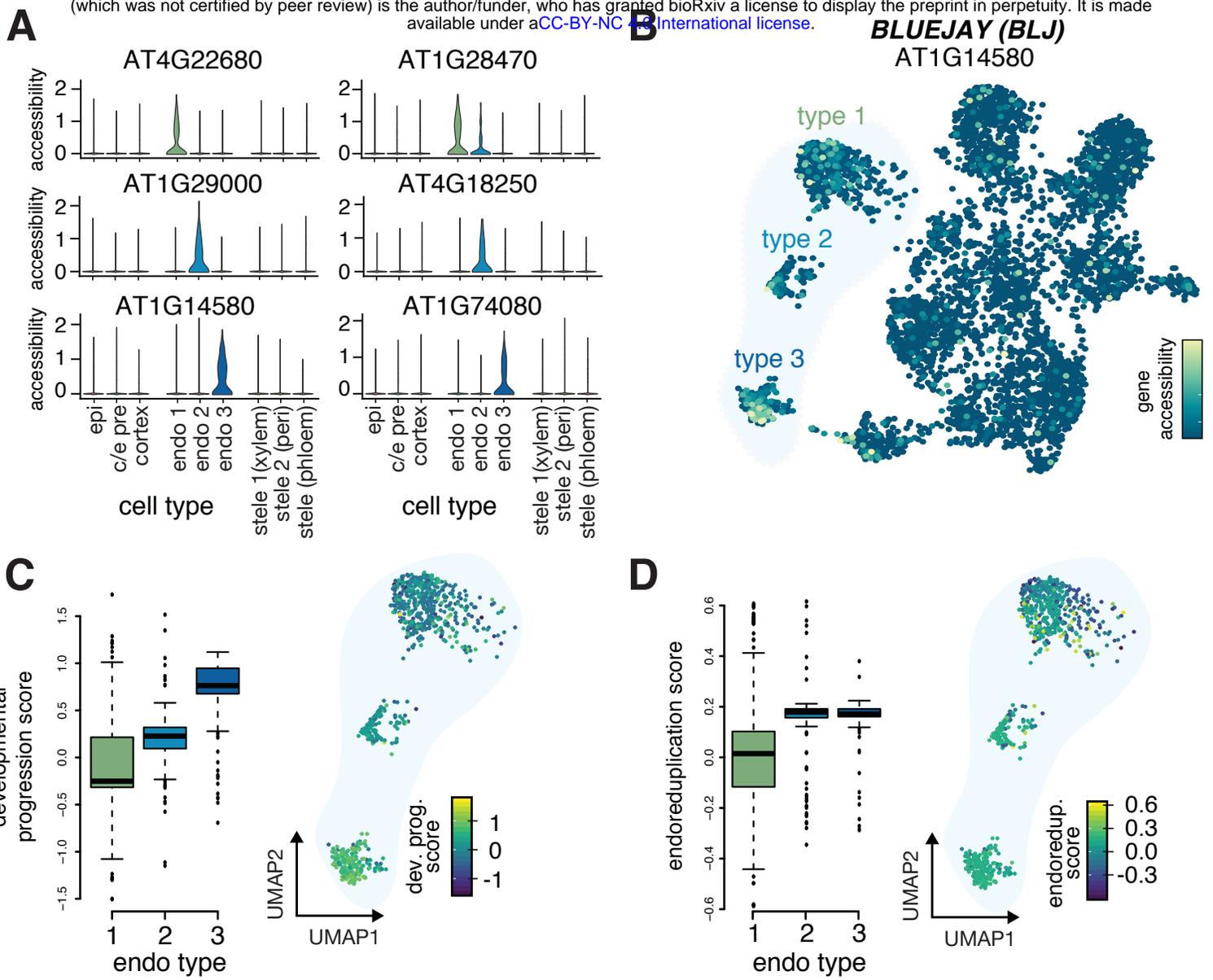
**C**

RNA cells



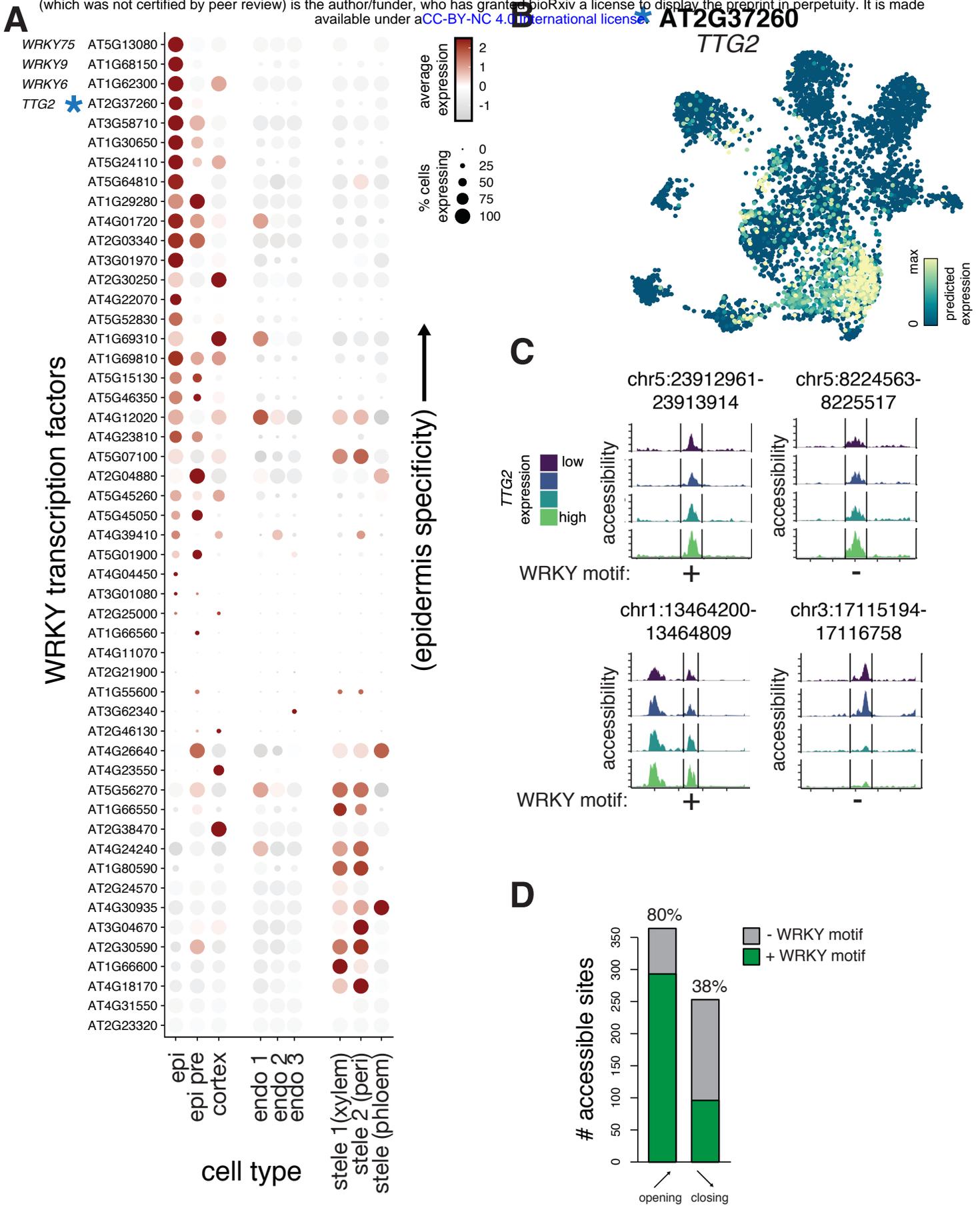
**Figure 2. scATAC-seq data can be integrated with scRNA-seq data to identify cell types.**

**Figure 2. scATAC-seq data can be integrated with scRNA-seq data to identify cell types.** (A) UMAP co-embedding of root scATAC cells alongside root scRNA cells.<sup>5</sup> Cells are colored by broad tissue type, with scATAC cells colored in lighter shades and scRNA cells in darker shades. (B) UMAP from (A), but showing only cells from the scATAC-seq experiment; (C) shows only cells from the scRNA-seq experiment.



**Figure 3. scATAC-seq identifies distinct sub-types of endodermal cells.**

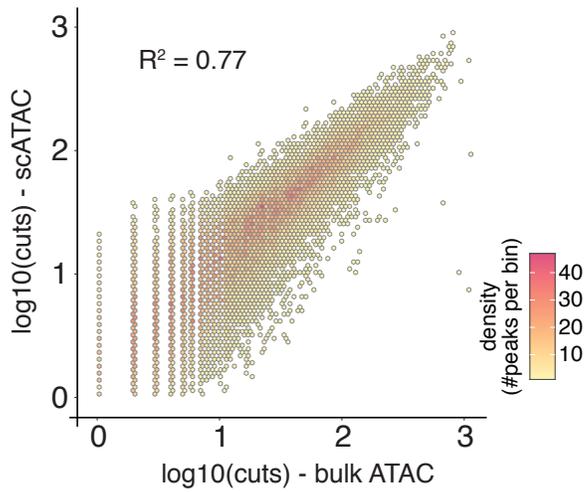
**Figure 3. scATAC-seq identifies distinct sub-types of endodermal cells.** (A) Violin plots showing specific patterns of accessible genes that mark each endodermal sub-type. Two examples are given for each endodermal sub-type, with gene-level accessibility scores indicated for all other cell types. (B) UMAP of all cells colored by accessibility of the *BLUEJAY* gene, which marks endodermal type 3; corresponding violin plot for this gene in lower left panel in (A). (C) Boxplot showing an increase in median developmental progression of each endodermal sub-type, as determined by average transcriptional complexity in the nearest 25 scRNA neighbors of each scATAC cell in the co-embedded representation from Fig. 2A; right inset shows UMAP of endodermal cells with each cell colored by the average developmental progression of its scRNA neighbors, mirroring the gradual increase seen in left panel. (D) Boxplot showing an increase in median levels of endoreduplication across endodermal sub-types, ascertained as in (C), but instead using a gene expression signature of endoreduplication; right inset shows UMAP of endodermal cells with each cell colored by the average endoreduplication score of its scRNA neighbors, with highest levels seen in endodermal sub-types 2 and 3.



**Figure 4. Prediction of candidate regulatory transcription factors from integrated scATAC and scRNA data.**

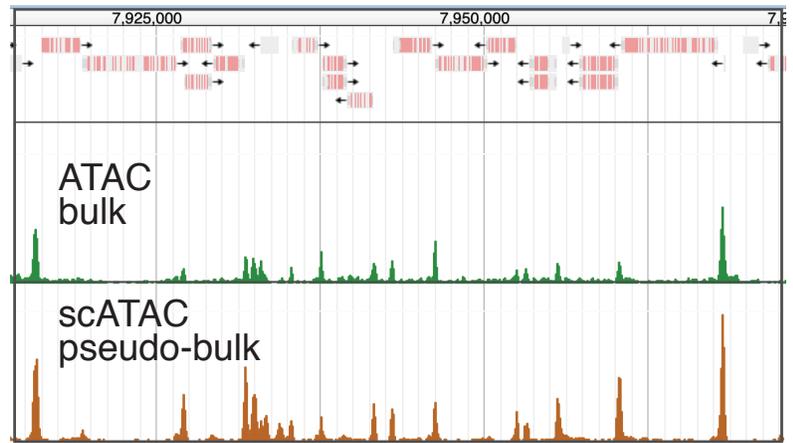
**Figure 4. Prediction of candidate regulatory transcription factors from integrated scATAC and scRNA data.** (A) Dotplot heatmap showing predicted expression of all WRKY family transcription factors across all cells. The color of each dot indicates the magnitude of predicted expression of each gene and the size of each dot represents the fraction of cells in each cell type showing expression at that gene; genes (rows) are ordered by the specificity of their epidermis expression. (B) UMAP plot of cells derived from scATAC experiment, but colored by predicted expression of an epidermis-specific WRKY transcription factor, *TTG2*. (C) Pseudo-bulked accessibility tracks of epidermis peaks whose accessibility showed a significant association with predicted *TTG2* expression. Cells with higher *TTG2* expression are shown in lighter shades. All panels show examples of significant ( $q < 0.05$ ) positive associations of *TTG2* expression with peak accessibility, with exception of the lower right panel. The presence or absence of a WRKY binding motif is indicated below each peak. (D) Barplot showing fraction of WRKY binding motifs in peaks of the epidermis, cortex, and pre-cursor type that showed significant association with *TTG2* expression. Peaks whose accessibility showed positive associations with expression are labelled as “opening”; those with negative associations are labeled as “closing.”

**A**

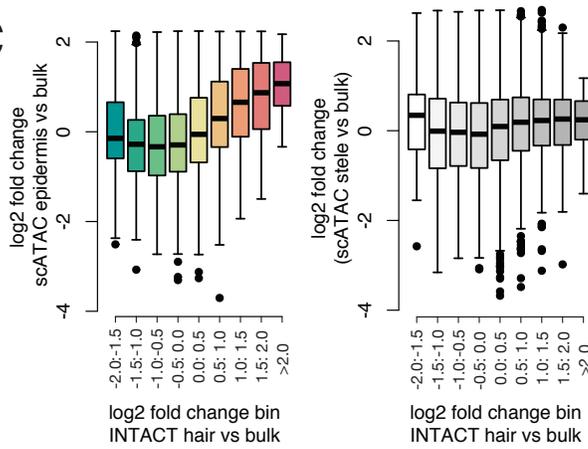


**B**

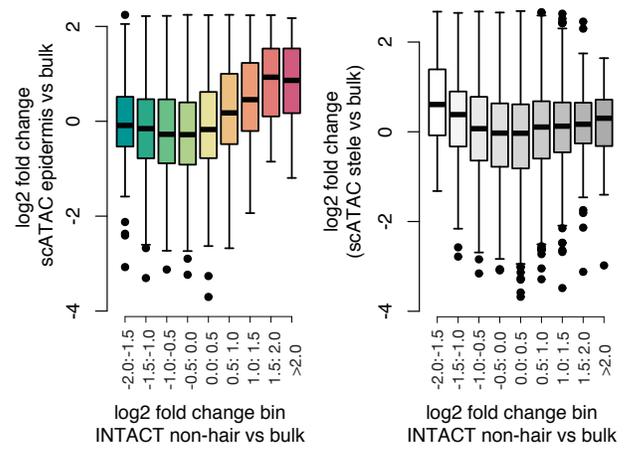
Chr5: 7912194 - 7970893



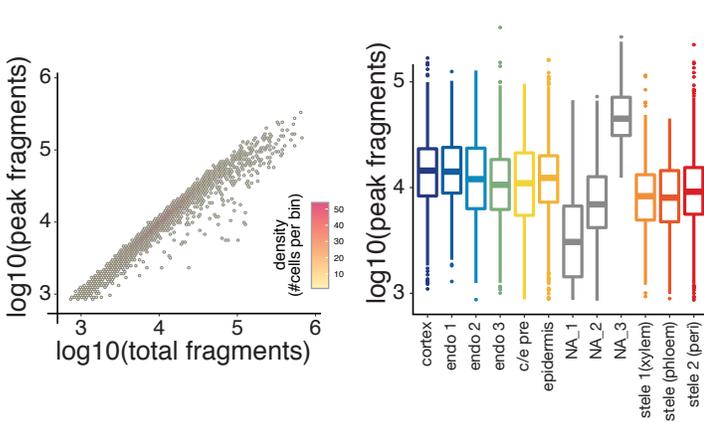
**C**



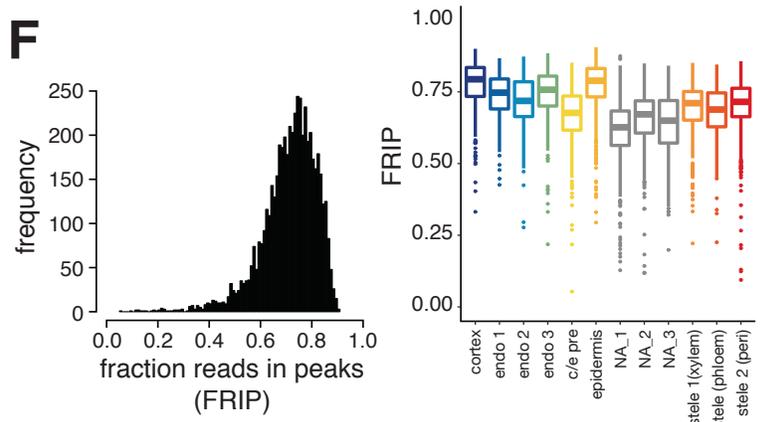
**D**



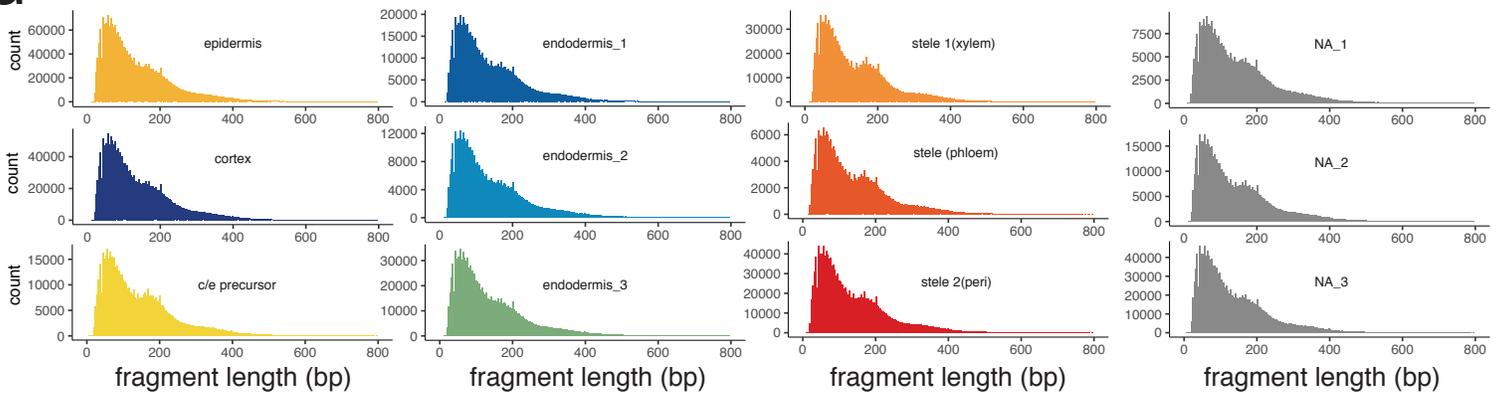
**E**



**F**



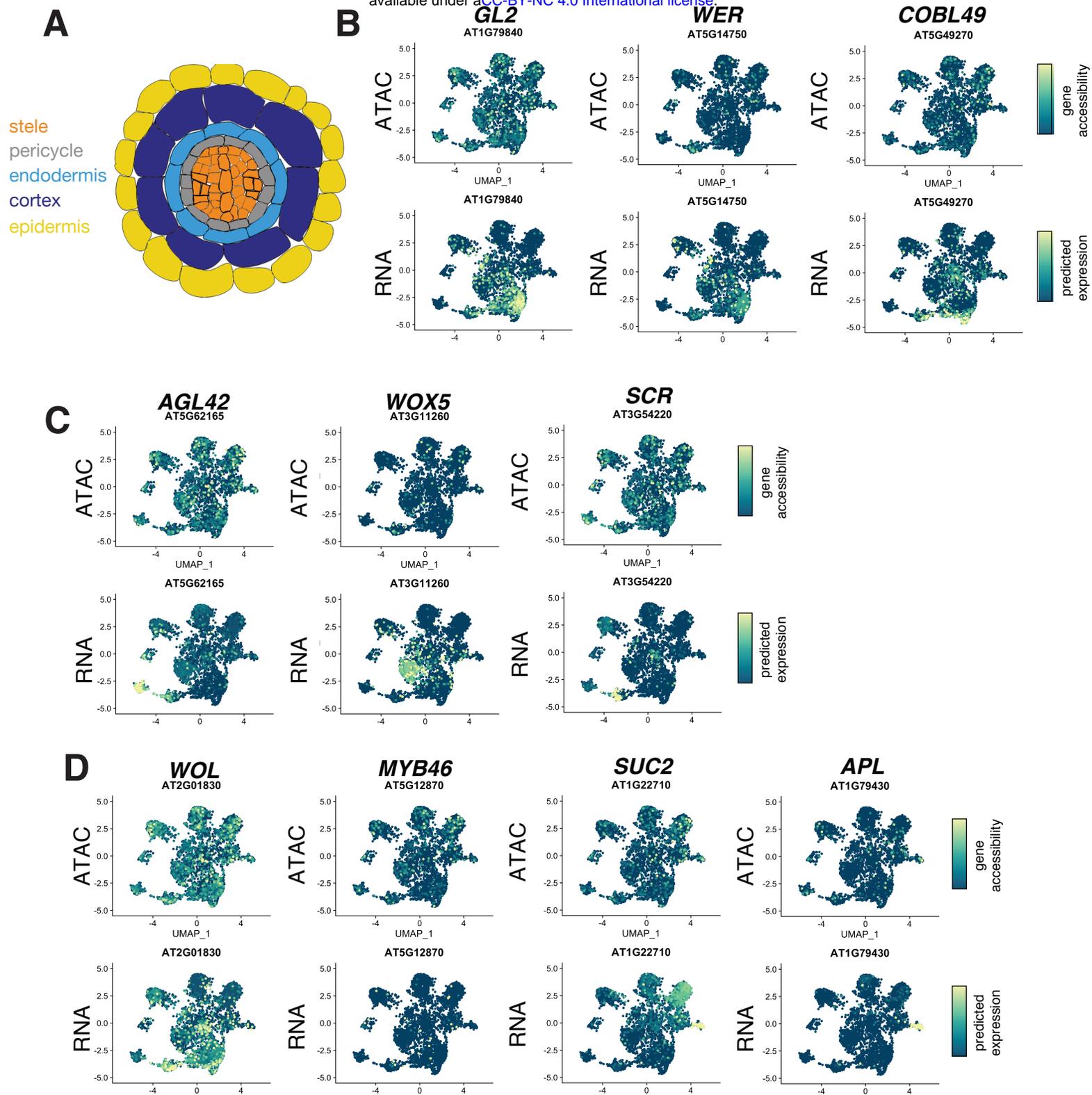
**G**



Supplementary Figure 1. Quality of scATAC-seq data is comparable to bulk ATAC-seq data.

### **Supplementary Figure 1. Quality of scATAC-seq data is comparable to bulk**

**ATAC-seq data.** (A) Scatterplot where each point represents peaks defined in the scATAC data. The x-axis shows the total cutcount within those peaks in bulk ATAC-seq and the y-axis shows the total cutcount within those peaks in scATAC-seq. Point density is indicated by increasing shades of red. (B) Example genomic region showing bulk ATAC accessibility (green) and pseudo-bulked scATAC accessibility (brown). Gene models are indicated above. (C) Boxplots showing peaks from scATAC assay in bins of increasing accessibility from an alternative, cell type-enriched ATAC approach;<sup>11</sup> peaks with low root hair cell-specific accessibility are in the leftmost bin, while those with the greatest root hair cell-specific accessibility are in the rightmost bin ( $n > 300$  for all bins). Root hair-specific accessibility was defined as peak accessibility in INTACT-derived root hair cells relative to a bulk ATAC sample. The y-axis in the left panel denotes epidermis-specific accessibility determined from the scATAC experiment, defined by the accessibility of those peaks in epidermal cells relative to accessibility when all cell types are grouped (simulating a “bulk” sample). The y-axis in the right panel denotes the relative accessibility in stele cells as a control. (D) Identical to (C), except that peaks are grouped by relative accessibility in root non-hair cells, determined by an alternative cell type-enriched ATAC approach.<sup>11</sup>(E) Read recovery per cell: Left panel shows relationship between total reads recovered per cell (x-axis) and reads in peaks (y-axis). Areas with higher point density are shown as in (A). Right panel shows boxplots of total number of reads in peaks recovered for each cell type. (F) ATAC quality per cell: Left panel shows the overall distribution of fraction of reads in peaks (FRIP) across all cells, right panel shows distribution of FRIP scores for each cell type. (G) Read length distributions for all fragments separated by cell type.



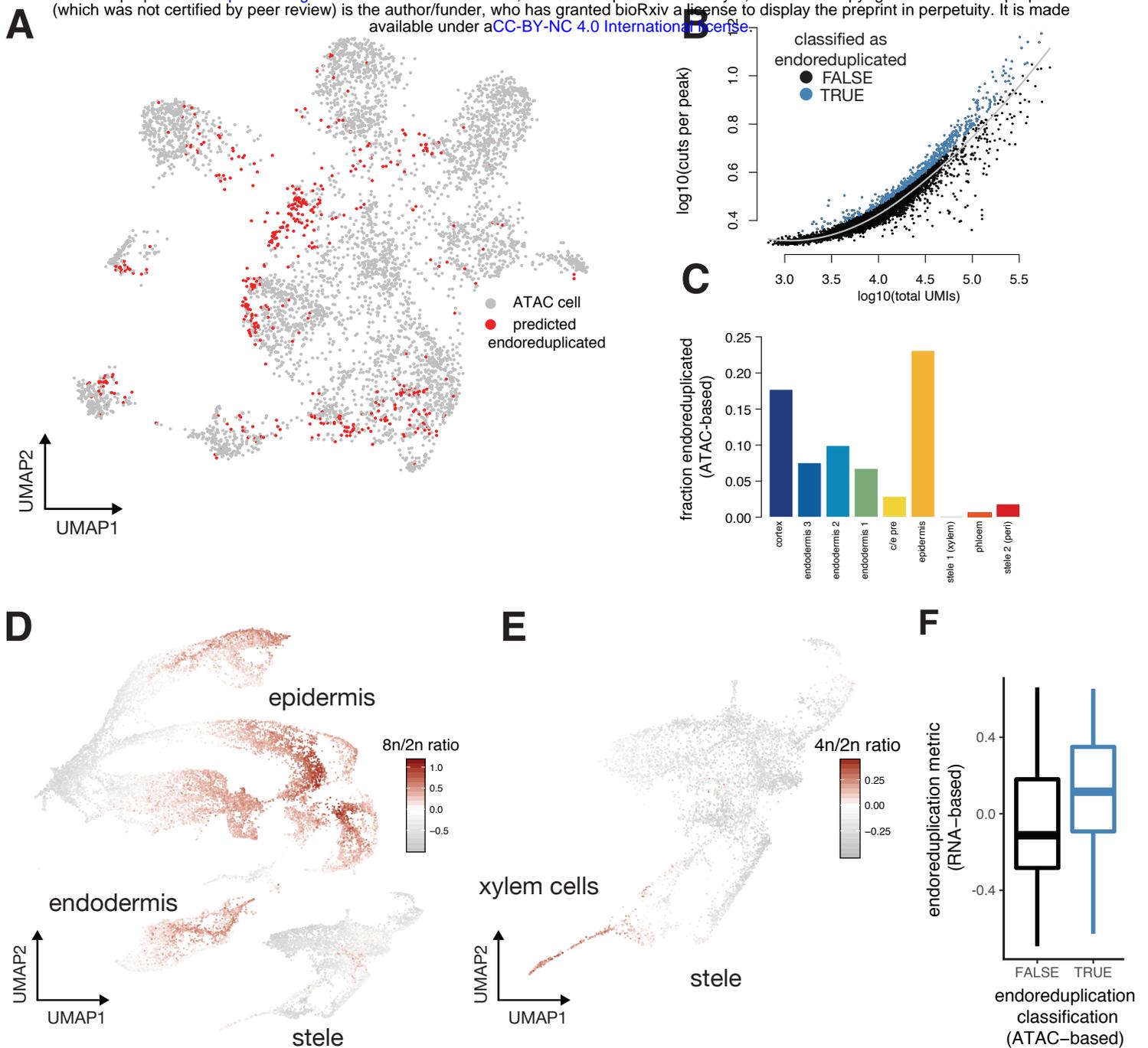
**Supplementary Figure 2. Accessibility and predicted expression levels of traditional marker genes in major cell layers of the root.**

**Supplementary Figure 2. Accessibility and predicted expression levels of traditional marker genes in major cell layers of the root.** (A) Schematic showing the major cell layers of the *Arabidopsis* root, colored as in Figure 1. (B) Marker gene plots for epidermis-specific genes showing accessibility (top) and predicted expression levels (bottom). Common and systematic gene names are indicated above. (C) As in previous panel, showing QC (*AGL42*), cortex (*WOX5*), and endodermis (*SCR*) markers. (D) As in previous panels, showing stele (*WOL* and *MYB46*) and phloem (*SUC2* and *APL*) markers.



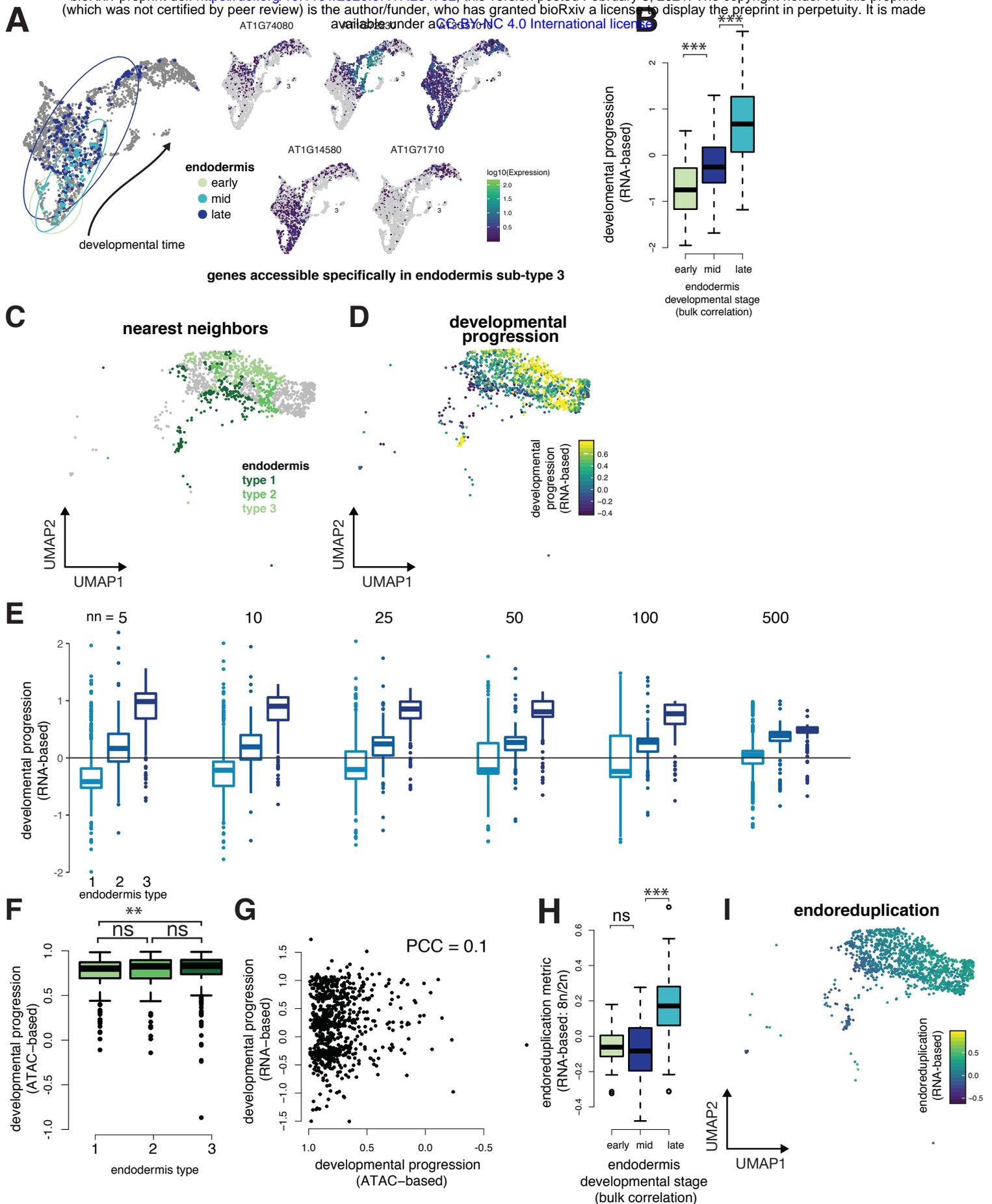
**Supplementary Figure 3. Co-embedding of scATAC and scRNA data allows validation of cell type labels and annotation by scRNA-derived features.** (A)

Confusion matrix showing the correspondence of manual cell annotations (x-axis) with those derived from the label-transfer from RNA to ATAC cells (y-axis). (B) UMAP of scATAC cells as in Fig. 1A, but cells are colored by the cell type label predicted from annotations of scRNA nearest neighbors. These cell type labels broadly match those predicted by manual annotation, and separate the epidermis cluster into hair and non-hair cells. (C) Workflow schematic for annotation of scATAC-cells with transcriptional data. The 25 nearest RNA neighbors from each ATAC cell in the co-embedded graph (**Figure 2A**) were identified, and average expression of individual genes and signatures scores were computed and assigned to each scATAC cell.



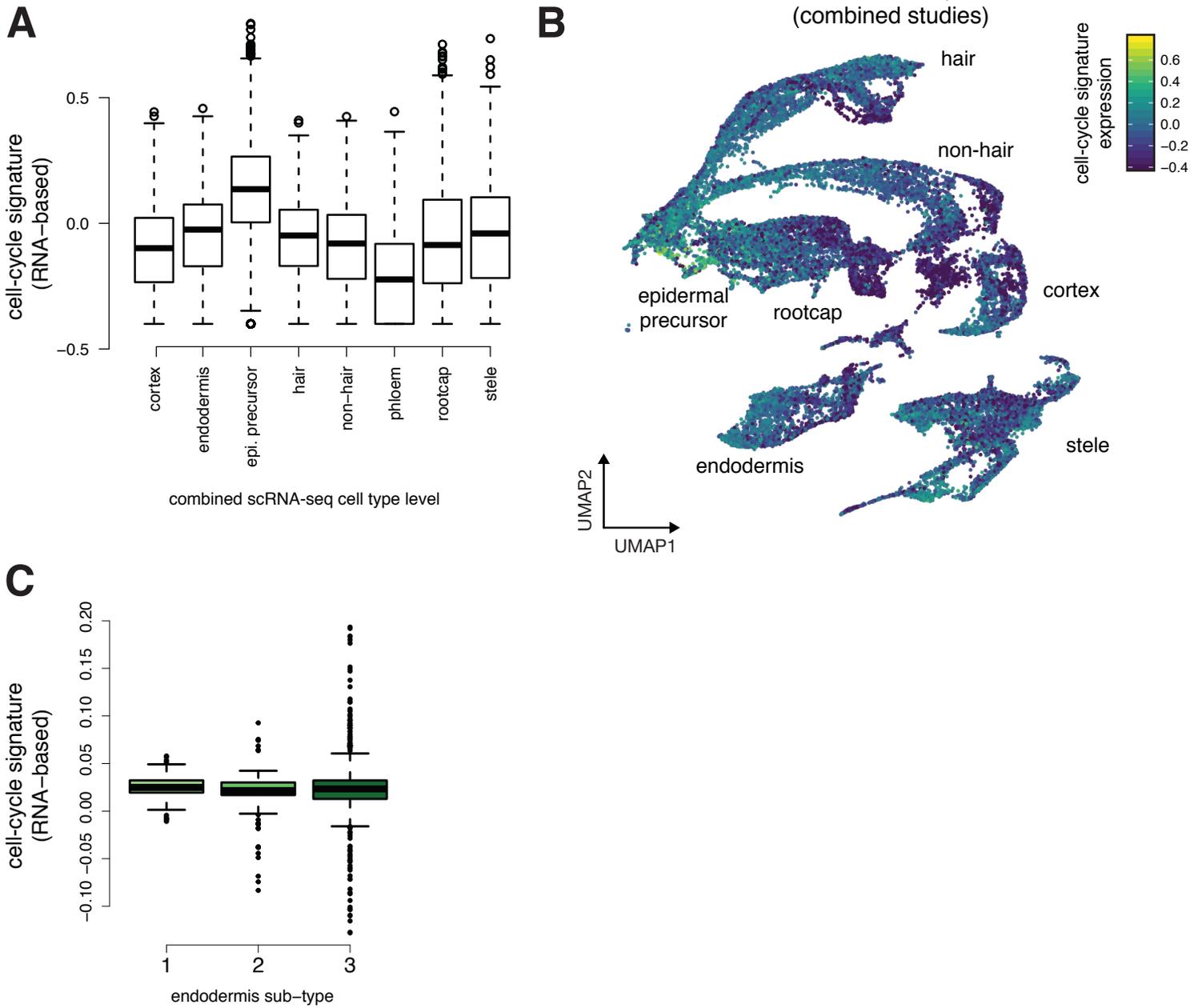
**Supplementary Figure 4. Approaches for identifying endoreduplicated cells in both scATAC and scRNA-seq data.**

**Supplementary Figure 4. Approaches for identifying endoreduplicated cells in both scATAC and scRNA-seq data.** (A) UMAP plot of root scATAC cells, each colored based on whether that cell surpasses a threshold level of cuts per site. Red denotes cells predicted as having undergone endoreduplication. (B) Scatterplot showing the relationship between total UMIs per cell (x-axis) and cuts per peak (y-axis); this relationship was captured in a Loess fit (black line), which was used to determine a threshold for cells with higher cuts per peak than expected based on their total UMIs (cells colored in blue, see Methods for more detail). (C) Barplot showing the fraction of cells in each type that showed putative endoreduplication, as determined by the threshold drawn in (B). In general, outer cell layers showed higher fractions of endoreduplicated cells, while cell layers of the stele showed lower levels. (D) UMAP of root scRNA cells, each colored based on the expression level of a transcriptional signature for endoreduplication, as determined by a ratio of expression levels in genes previously determined as enriched in 8n cells over those enriched in 2n cells.<sup>19</sup> (E) A known instance of endoreduplication in the stele, tetraploid xylem<sup>19</sup>, is identified by a metric similar to (D), except that cells are colored by signature for 4n cells (ratio of 4n-specific genes to 2n-specific genes). (F) Boxplot showing the transcriptional-signature-based endoreduplication metric compared to a binary classification of endoreduplication cells using scATAC data. scATAC cells with high levels of cutcounts at a single locus (suggesting endoreduplication, as in A-C) were analyzed in the co-embedded graph with scRNA-seq cells to calculate the average level of the endoreduplication signature among each scATAC cell's 25 nearest neighbors. The overall trend shows that the cutcount-based classification of endoreduplication is consistent with the transcriptional-signature-based metric (one-sided student's t-test  $p < 1E^{-14}$ ).



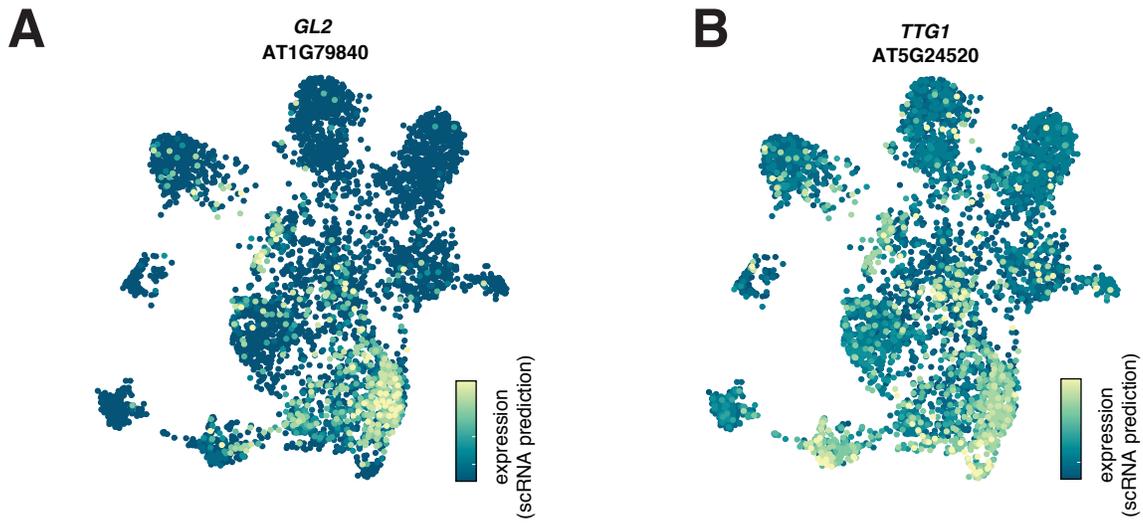
**Supplementary Figure 5. Characterization of endodermal sub-types with combined scATAC and scRNA-seq data.**

**Supplementary Figure 5. Characterization of endodermal sub-types with combined scATAC and scRNA-seq data.** (A) UMAP of endodermal cells from multiple scRNA-seq studies, with previously-determined developmental stages highlighted.<sup>2</sup> Inset shows variable expression patterns of genes with accessibility patterns specific to endodermal sub-type 3 in the scATAC data. (B) Boxplot showing that developmental progression scores are consistent with previously described annotations of developmental progression (early, middle, late) of the endodermis (all comparisons significant in one-sided student's t-test  $p < 1E^{-14}$ ).<sup>2</sup> (C) Subset of co-embedded UMAP from Figure 2A showing only endodermal cells; 25 nearest RNA neighbors for each endodermal type are indicated in shades of green. (D) As in (C), but shows RNA neighbor cells colored by transcription-based developmental progression metric. (E) Boxplots showing data from Figure 3C, with average developmental progression computed with different numbers of nearest neighbors. Above each plot, the number of neighboring cells (nn) from the scRNA-seq data used to predict developmental progression of each scATAC endodermal cell is shown. The relative differences in predicted developmental progression is insensitive to the number of nearest neighbors used in the procedure. (F) Boxplots showing levels of accessible genes (analogous to transcriptional complexity metric from Fig. 3C, only computed as total number of accessible genes rather than total number of transcribed genes). The overall trend remained the same, with progressive loss of complexity in the later endodermal types (significant for sub-type 1 vs 3, one-sided student's t-test p-value = 0.0032, not significant for other comparisons), but the ATAC-based metric showed less sensitivity than the RNA-based one. (G) Scatterplot showing poor correlation (PCC = Pearson correlation coefficient) of ATAC-based developmental progression score and the RNA-based score. (H) Boxplot as in (B), showing transcription-based endoreduplication scores (y-axis) for cells annotated for endodermal developmental stages by a previous scRNA-seq experiment (early and middle comparison not significant [ns], middle and late comparison, one-sided student's t-test  $p < 1E^{-14}$ ). (I) As in (D), with RNA neighbor cells colored by transcription-based endoreduplication metric.



**Supplementary Figure 6. Dividing cells are present in the root, but do not distinguish endodermis sub-types**

**Supplementary Figure 6. Dividing cells are present in the root, but do not distinguish endodermis types.** (A) Boxplots showing levels of a cell-cycle signature in each scRNA-seq root cell type. (B) UMAP plot of combined root scRNA-seq studies with each cell colored by its expression the cell-cycle signature.<sup>32</sup> (C) Cell-cycle signature predicted from nearest neighbors of endodermis types (as in **Figure 3C, 3D**) shows that proliferation is not a strongly distinguishing feature between the sub-types.



**Supplementary Figure 7. Identifying transcription factors involved in tissue specification.**

**Supplementary Figure 7. Identifying transcription factors involved in tissue specification.** (A) UMAP of scATAC cells colored by predicted expression level of epidermal specification factor *GL2*. (B) UMAP of scATAC cells colored by predicted expression level of epidermal specification factor *TTG1*.