

Gene duplication and cellular divergence in crops

Kenneth Birnbaum (✉ ken.birnbaum@nyu.edu)

New York University <https://orcid.org/0000-0002-8423-6859>

Bruno Guillotin

NYU <https://orcid.org/0000-0002-0117-2512>

Ramin Rahni

New York University, Center for Genomics and Systems Biology United States

Michael Passalacqua

Cold Spring Harbor Laboratories <https://orcid.org/0000-0002-6344-1175>

Mohammed Mohammed

New York University Abu Dhabi, Center for Genomics and Systems Biology

Xiaosa Xu

Cold Spring Harbor Laboratories

David Jackson

Cold Spring Harbor Laboratory <https://orcid.org/0000-0002-4269-7649>

Simon Groen

University of California, Riverside

Jesse Gillis

University of Toronto

Biological Sciences - Article

Keywords:

Posted Date: June 17th, 2022

DOI: <https://doi.org/10.21203/rs.3.rs-1739501/v1>

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

[Read Full License](#)

1 Gene duplication and cellular divergence in crops

2 Bruno Guillotin^{1,2}, Ramin Rahni¹, Michael Passalacqua³, Mohammed Ateequr Mohammed²,
3 Xiaosa Xu³, David Jackson³, Simon C. Groen⁴, Jesse Gillis⁵, and Kenneth D. Birnbaum^{1,2,*}

4 ¹New York University, Center for Genomics and Systems Biology

5 ²New York University Abu Dhabi, Center for Genomics and Systems Biology

6 ³Cold Spring Harbor Laboratories

7 ⁴University of California, Riverside

8 ⁵University of Toronto, Physiology Department

9 *Corresponding author: ken.birnbaum@nyu.edu

10 Abstract:

11 **Different plant species within the grasses were parallel targets of domestication, giving**
12 **rise to crops with distinct evolutionary histories and traits. Key traits that distinguish these**
13 **species are mediated by specialized cell types within organs. Here, we compare the**
14 **transcriptomes of all cells within roots in three grasses—*Zea mays* (maize), *Sorghum***
15 ***bicolor* (sorghum), and outgroup *Setaria viridis* (*Setaria*). We first show that single-cell and**
16 **single-nucleus RNA-seq provide complementary readouts of cell identity, warranting a**
17 **combined analysis. Comparative cellular analysis shows that the transcriptomes of some**
18 **cell types diverged more rapidly than others, in part by recruiting gene modules from other**
19 **cell types. Furthermore, examining the whole genome duplication in maize, we detect**
20 **extensive dosage compensation in surviving co-expressed homeologs, reinforcing**
21 **genomic balance¹. Homeolog pairs that underwent subfunctionalization², partitioning their**
22 **expression among cell types, represented a minor pattern but showed the highest rate of**
23 **acquiring a novel (non-ancestral) domain. These results fit a conjecture in which**
24 **mechanisms that maintain stoichiometric balance at the molecular level aid in homeolog**
25 **retention for extended periods to allow new functions to arise. An unexpected synergy**
26 **between spatial sub- and neo-functionalization then contributes to changes in**
27 **transcriptional cell identity.**

28 Single-cell mRNA profiling has opened up new opportunities to study cellular evolution by
29 comparing gene regulation in specialized cells across species^{3,4}. In plants, high-resolution cellular
30 profiling also has the potential to associate cell-level transcriptional regulation to key agricultural
31 traits, many of which are mediated by specialized cells⁵.

32 *Zea mays* (maize) is a staple crop and *Sorghum bicolor* (sorghum) is an important dryland crop
33 and biofuel candidate that is closely related to maize, separated by about 12 million years^{6,7}.
34 However, the two species differ substantially in key traits such as drought and chilling tolerance,
35 and release of root exudates that shape soil interactions⁸⁻¹². The importance of the two crops,
36 their evolutionary proximity, and their functional differences present a novel opportunity for
37 comparative analysis of cellular evolution in plants^{13,14}. In addition, since sharing a common
38 ancestor with sorghum, maize underwent a whole genome duplication 5 to 12 million years ago
39 ⁷, offering an opportunity to analyze changes in fine-scale gene regulation among paralogous

40 genes on duplicated chromosomes (homeologs) in the relatively early stages after a
41 duplication^{7,15}.

42 **Cells Provide Depth While Nuclei Give Breadth**

43 Single-cell analyses in plants have relied on the generation of protoplasts by enzymatic digestion
44 of cell walls^{16–21}. However, certain tissues and even some species like sorghum are quite
45 recalcitrant to digestion. There is also historic concern about the effects of protoplast generation
46 on the cellular transcriptome, leading to growing interest in nuclear profiling^{22–24}. To assess the
47 fidelity of nuclear profiling in depth, we first compared single-cell vs single-nucleus vs whole-root
48 profiles in both *Arabidopsis thaliana* (*Arabidopsis/At*, a dicot model with plentiful resources,
49 15,967 cells and 17,373 nuclei) and maize (*Zm*, a monocot model, 4,235 cells²⁵ and 2,668 nuclei)
50 (Supplementary Table 1).

51 Measures of unique molecular indices (UMIs) per dataset showed an increase of 10x (*At*) and 6x
52 (*Zm*) in cells than in nuclei (Extended Data Fig. 1a), similar to animal studies^{26,27}. Accordingly, the
53 average number of genes detected was 2.7x (*At*) and 1.4x (*Zm*) higher in cells than in nuclei
54 (5,281 vs 1,895 in *At*, 4,198 vs 2,304 in maize) (Extended Data Fig.1b, Supplementary Table 1).
55 However, despite the lower mRNA content, nuclear profiling detected 89% (*At*) and 88% (*Zm*) of
56 total genes present in cells (Supplementary Table 1).

57 Both cell and nuclei “pseudo-bulked” transcriptomes displayed a high correlation to whole-root
58 transcriptomes ($r \sim 0.7-0.8$, Extended Data Fig. 1c), confirming that both sampling methods
59 generally reflected expression patterns of intact tissue.

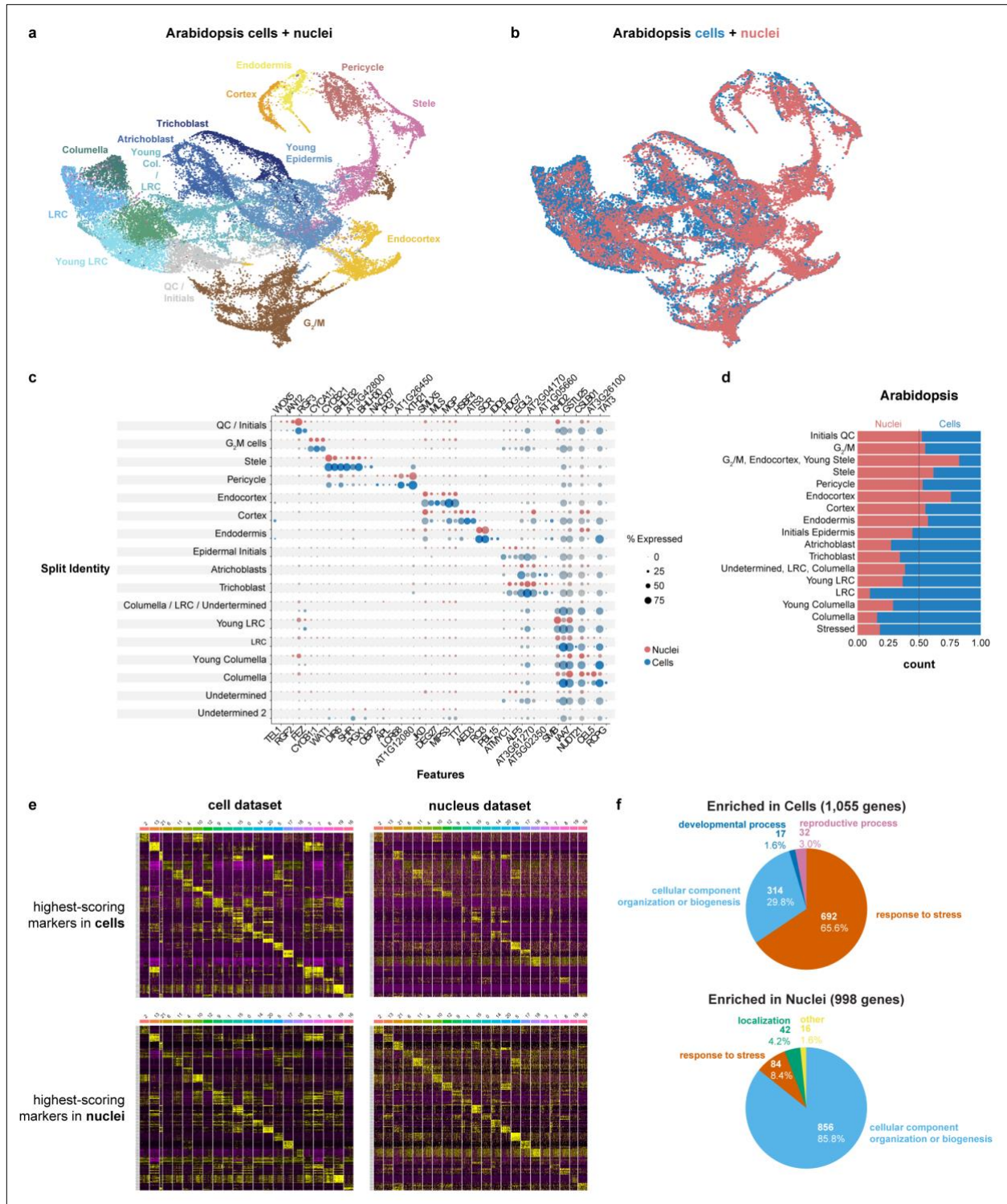
60 In both *Arabidopsis* and maize, cells and nuclei generated UMAP clusters corresponding to all
61 the major cell identities (Extended Data Fig. 2, 5). However, in both species, the nuclear dataset
62 generated fewer distinct clusters, often failing to resolve between closely related or subcellular
63 identities (Extended Data Fig. 2, 5). For example, in maize, stele cells contained a subcluster that
64 we identified as xylem cells, where no such subcluster was apparent in the nuclear UMAP
65 (Extended Data Fig. 5). Using a down-sampling approach on each dataset, a general rule-of-
66 thumb emerged that twice as many nuclei as cells are needed to discover the same number of
67 cell-type clusters as protoplasts (Extended Data Fig. 3a,b). Thus, the shallower depth of nuclear
68 profiles provides less resolution for classification of cell identity—a drawback that down-sampling
69 showed we could rectify, at least in part, by increasing the number of nuclei sampled compared
70 to cells.

71 Either simultaneous or independent analysis of cells and nuclei generated clusters that reflected
72 the same underlying biological patterns (Fig. 1a-c, Extended Data Fig. 3c,d). The highest-scoring
73 markers extracted from nuclei generally matched the highest-scoring ones from cells (Fig. 1c,e,
74 Extended Data Fig. 3d). In addition, the assignment of cells to specific clusters was stable when
75 cells or nuclei were clustered either alone or together (Supplementary Table 2).

76 One advantage of nuclear profiles was their ability to capture a more representative sampling of
77 cells within the tissue (Fig. 1d, Extended Data Fig. 4d). In one example in maize, we detected a
78 unique cluster in single-nucleus profiling not present in single-cell/protoplast profiling, which we
79 confirmed as columella cells using previously published hand-sectioned RNA-seq profiles
80 (Extended Data Fig. 5,²⁵).

81 In *Arabidopsis*, we found that 14% of total genes (3,218) were differentially expressed between
82 cells and nuclei in a cluster-by-cluster analysis (Supplementary Table 3). While cells showed a

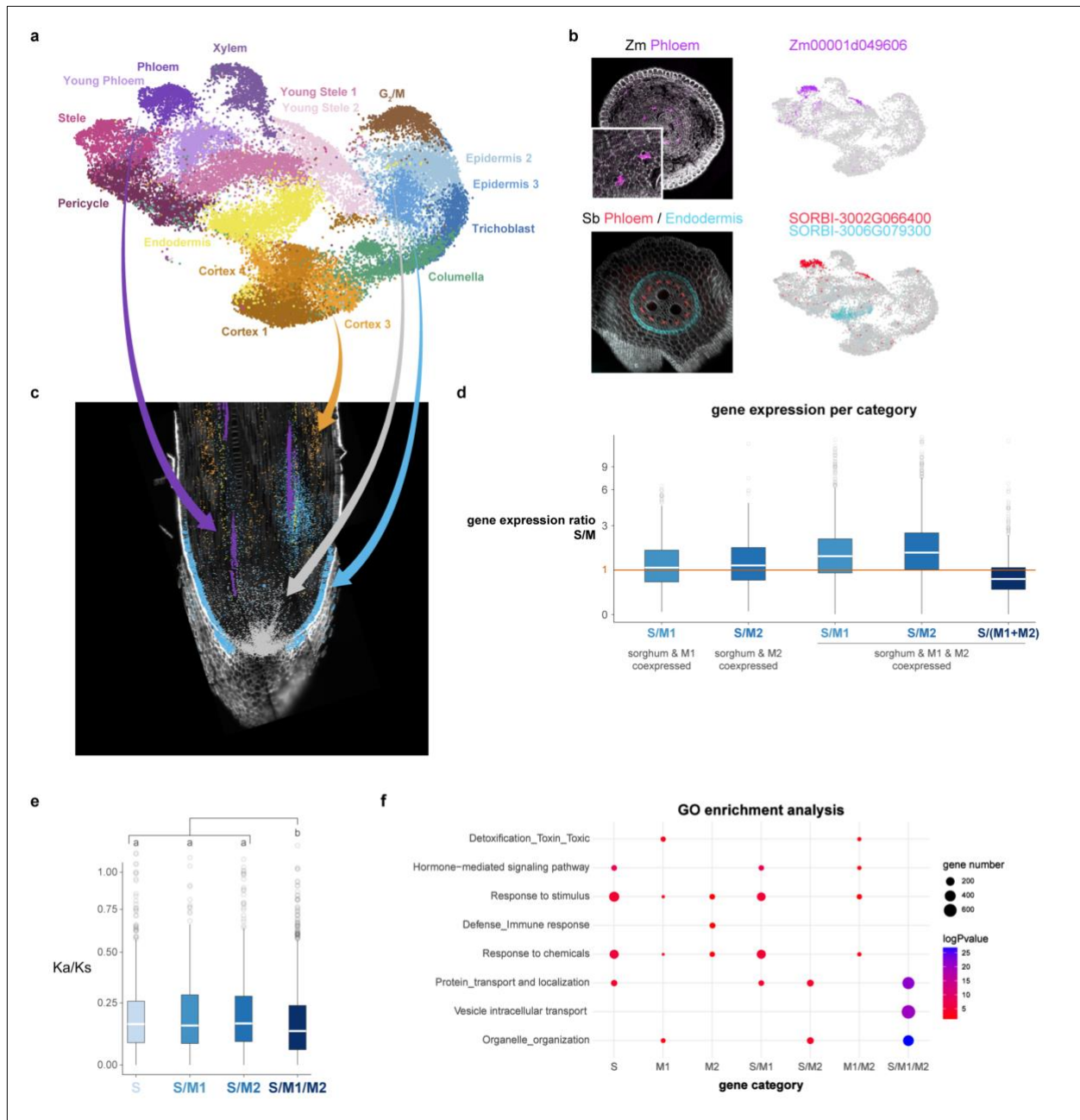
83 higher proportion of stress related genes (Fig. 1f, Extended Data Fig. 4a,b), most of the
 84 differences between cell and nuclear profiling appeared to be related to compartmental RNA
 85 stability, as mRNAs enriched in nuclei significantly overlapped with transcripts shown to have
 86 higher decay rates in the cytoplasm^{28,29} (Extended Data Fig. 4c). This analysis showed that
 87 nuclear profiles did indeed show a lower stress response than protoplasts, but the difference was
 88 subtle and not a dominant component of the difference between cells and nuclei.



90
91 **Fig. 1: Cell and nucleus profiles identify the same markers but show different sensitivities**
92 **and artifacts. a, b** UMAP of combined Arabidopsis cells and nuclei with clusters colored
93 according to assigned cell identity (**a**) or cells vs. nuclei origin (**b**). **c** Dot plots of Arabidopsis
94 marker genes in cells (blue) or in nuclei (red). **d** Proportion of cells vs nuclei present in each cell
95 type cluster. **e** Heatmaps of the 10 highest-scoring marker genes for each cell type found using
96 Seurat in the combined single-cell + single-nucleus dataset. This combined dataset was then
97 divided into a cell dataset and a nucleus dataset. Within each of the four panels, rows correspond
98 to genes, and columns correspond to cell-type clusters. Upper two heatmaps show highest-
99 scoring markers within the cell dataset (left) vs expression of the same markers in the nucleus
100 dataset (right). Lower row shows highest-scoring markers found in nucleus dataset (right) vs
101 expression of the same markers in the cell dataset (left). **f** Pie charts showing the difference in
102 the prevalence of GO terms among differentially expressed genes between cells (top) and nuclei
103 (bottom), aggregated from cluster-by-cluster differences. GO term analysis is for up-regulated
104 genes in either cells or nuclei.

105 **A High-confidence cell-type Map in sorghum and *Setaria* Using a maize Reference.**

106 Given the comprehensive coverage of a combined analysis, we pursued both whole cell and
107 nuclear profiling to investigate cellular evolution in the maize-sorghum clade. We further
108 generated profiles for sorghum (3,510 cells and 7,620 nuclei) and, as an outgroup, *Setaria viridis*
109 (*Setaria*, 974 cells and 12,192 nuclei, Supplementary Table 1). We took advantage of comparative
110 genomic sequence analyses in maize, sorghum, and *Setaria* that mapped orthologs among the
111 three species, including the homeologs created by whole genome duplication in maize^{13,15}
112 (hereafter subgenome M1 and M2). We first used a set of single-copy orthologs in the three
113 species to cluster the maize cells and nuclei and then mapped the sorghum and *Setaria* datasets
114 onto the maize anchor³⁰ (Fig. 2a, Supplementary Table 1). To validate the mapping, we performed
115 an independent MetaNeighbor analysis (Extended Data Fig. 6) as well as whole mount *in situ*
116 hybridizations in maize and sorghum (Fig. 2b, Extended Data Fig. 7,8), and spatial transcriptomics
117 in maize (Fig. 2c, Extended Data Fig. 7), confirming the maize-to-sorghum-to-*Setaria* mapping of
118 cell identities. Thus, we could use the well-annotated maize cell type map to rapidly generate a
119 high confidence single-cell “pan-transcriptome” of these key crop species, including hundreds of
120 new cell-type specific marker genes (Supplementary Table 4).



121

122 **Fig. 2: Mapping cell identities from maize to sorghum reveals global dosage compensation**
 123 **between maize duplicates in cell types.** **a** UMAP of combined maize single-cell and single-
 124 single-nucleus profiles. Clusters are colored and labeled according to cell identity. **b** *in situ* hybridization
 125 in maize (top) and sorghum (bottom). The maize phloem marker in magenta is orthologous to the
 126 sorghum phloem marker in red. Cyan in the lower panel corresponds to a sorghum endodermal
 127 marker. Autofluorescence highlighting anatomy is in grayscale. The minimum/maximum values
 128 for each channel in the fluorescence images have been adjusted to show the localization more
 129 clearly in the merged image. UMAPs next to images show the respective expression in single-
 130 cell / single-nucleus profiles, which were used initially to determine their expression pattern. **c**
 131 Spatial transcriptomics showing simultaneous localization of multiple markers that enable detailed
 132 mapping of single-cell profiles to specific tissues. **d** Expression ratios of sorghum over maize

133 orthologous genes, where the red line indicates equivalent expression levels. The first two
134 boxplots represent cases where a sorghum ortholog is co-expressed with a single maize
135 homeolog (either M1 or M2). The third and fourth boxplots represent cases in which both
136 homeologs are expressed in the same cells. The last boxplot shows the ratio when both of the
137 co-expressed homeologs are summed in the denominator. White bar is median, box limits
138 represent 25th and 75th percentiles, lines are the 95th percentile. **e** Ka/Ks distribution of
139 comparative patterns measured against sorghum. Cell type specific expression is averaged.
140 Letters signify comparative patterns of cell type expression for sorghum and M1 and M2
141 subgenome expression: i.e. S = only the sorghum ortholog is present in the cell type cluster, SM1
142 = sorghum and maize 1 homeologs are co-expressed in the cell type cluster, etc. Statistical
143 analysis was performed using ANOVA followed by the Tukey pairwise test. All “a” designates are
144 statistically significant from “b” at $p < 0.05$. **f** GO terms enriched within each category of genes in
145 each cell types.

146 **Widespread Dosage Compensation between maize homeologs**

147 The comparative cell-type maps provided an opportunity to quantify patterns in preserved gene
148 duplicates – a potential source of innovation particularly following whole genome duplication. One
149 hypothesis, known as the genomic balance model, holds that the viability of a duplicated gene is
150 dependent on the capacity to adjust for dosage effects, such that the stoichiometry of molecular
151 complexes is retained following a whole genome duplication¹.

152 However, we could observe changes in expression patterns in the 5 to 12 million years since
153 whole genome duplication that were likely to disrupt the stoichiometric balance. Prior studies have
154 demonstrated dosage compensation between gene duplicates in aneuploids and in newly
155 generated polyploids^{1,31}. However, it is not clear if this phenomenon has a role after ancient
156 genome duplication events. The cellular resolution of our dataset offered an opportunity to test
157 whether dosage compensation within cells could fine tune genomic balance.

158 To examine this issue, we adapted a ranked-based method to compare gene expression levels
159 across genomes in each cell type³². This enabled us to use sorghum gene expression in each
160 cell type as proxies for the “ancestral” state (expression level and cell-type domains). In addition,
161 ancestral expression domains were supported by at least one maize homeolog. We also made
162 the assumption that both homeologs had identical expression patterns at the time of whole
163 genome duplication³³. We categorized a homeolog as “dominant” if its average expression level
164 per cell type is more than twice the expression of the other homeolog, and as “co-expressed” if
165 both homeologs are detected and neither is dominant.

166 When one of the M1 or M2 homeologs was dominant in a cell type, this dominant homeolog was
167 expressed at the same level as its sorghum ortholog (Fig. 2d). Interestingly, these dominant
168 homeologs were enriched in GO terms for stress adaptation, immunity and response to stimulus
169 (Fig. 2f, Supplementary Table 5). At the same time, the dominant homeologs display a higher
170 cell-type specificity than co-expressed homeolog pairs³⁴ (τ), while the non-dominant ones show a
171 higher nonsynonymous-to-synonymous substitution rate (Ka/Ks) suggesting a more relaxed
172 purifying selection (Extended Data Fig. 9a-c).

173 Alternatively, when maize M1 and M2 homeolog pairs were co-expressed in the same cell types
174 as their sorghum ortholog, each showed an expression level of about half that of the sorghum
175 ortholog (Fig. 2d, Extended Data Fig. 9d,e). These co-expressed homeologs further displayed a
176 marginally lower Ka/Ks ratio than dominant cases, suggesting that they are under more stringent
177 purifying selection (Fig. 2e, Extended Data Fig. 9a). Finally, these homeologs showed GO term

178 enrichments for categories known to be favored for retention after whole genome duplication—
179 cell homeostasis processes, translation, or ribosome biosynthesis^{1,15} (Fig. 2f, Supplementary
180 Table 5). The analysis therefore showed strikingly widespread dosage compensation in cell types
181 where both homeologs are co-expressed likely fine tuning stoichiometric balance over long
182 periods on a cell-type basis³⁵.

183 **Neofunctionalization and subfunctionalization are linked**

184 It has been suggested that mechanisms permitting stoichiometric balance could act as a bridge
185 for the generation of new gene function, promoting the retention of gene duplicates to allow
186 sufficient time for beneficial or complementary mutations to arise¹. Two possible mechanisms for
187 long-term duplicate gene retention are subfunctionalization and neofunctionalization², both of
188 which we consider here at the transcriptional level. Thus, we define subfunctionalization here as
189 a case in which duplicated genes partition the expression of the ancestral² (sorghum) domain
190 (Fig.3a). Neofunctionalization is presumed to occur when one homolog retains ancestral and the
191 second is free to diverge and adopt a new function³⁶. Thus, at the transcriptional level, we define
192 neofunctionalization as the dominance of one homolog in the ancestral domain and the
193 expression of the second homolog in a cell type outside the ancestral domain (Fig.3c). We
194 classified each homeolog pair across cell types on a scale from -1 (full dominance of one
195 homeolog over the other), to 0 (full co-expression of both homeologs) to 1 (full
196 subfunctionalization; Fig. 3a,b). Overall, 70% of the homeolog pairs showed dominance, 11%
197 showed full co-expression, and 19% showed subfunctionalization (Fig. 3a,b, Supplementary
198 Table 6).

199
200 To assess potential evolutionary forces driving each pattern, we assumed both homeologs
201 matched the sorghum ortholog expression pattern at the time of duplication. We then randomly
202 removed gene expression of either homeolog across cell types until their matrix of gene
203 expression matched the overall presence/absence matrix of homeologs in the observed data.

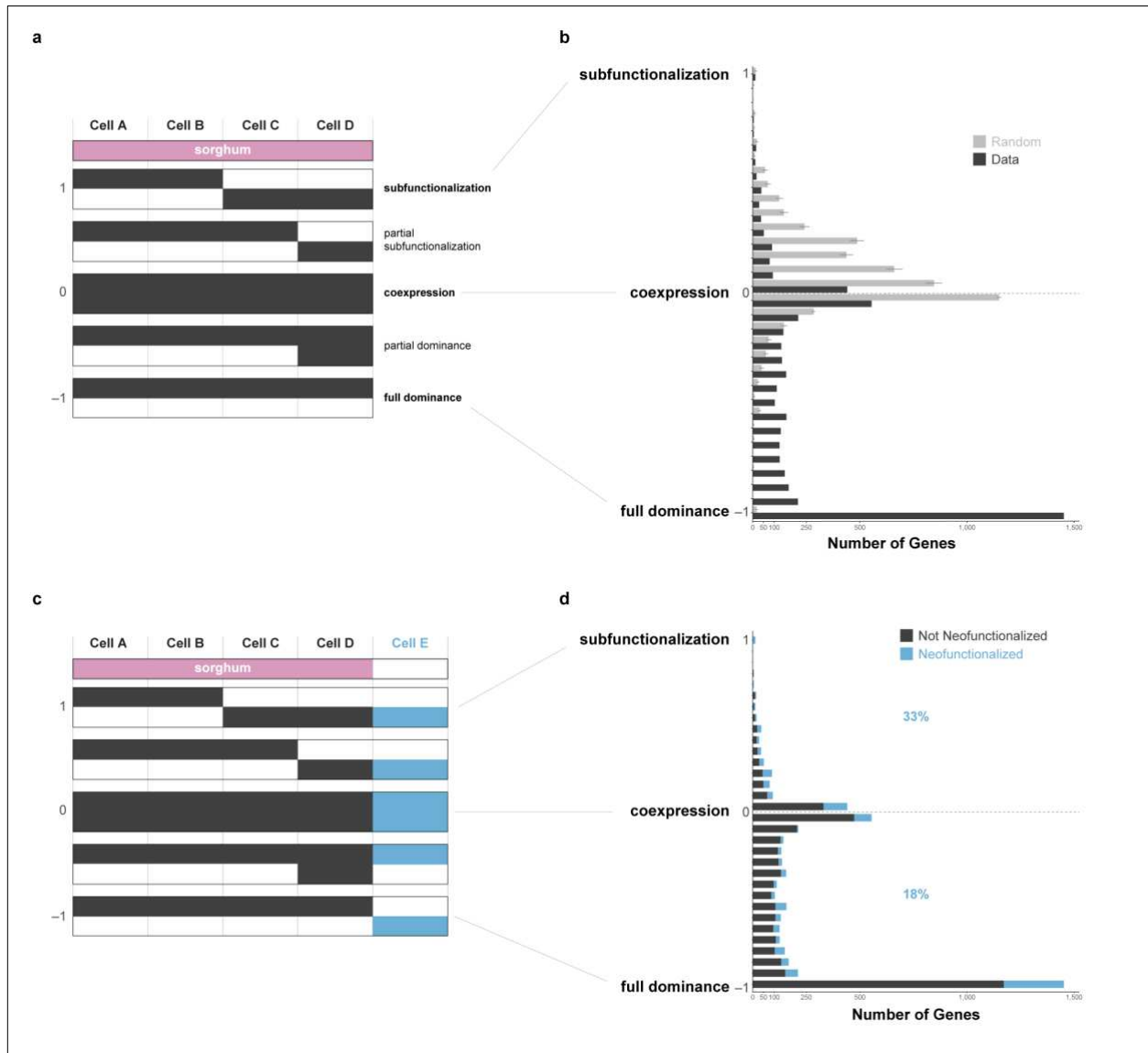
204 The analysis revealed that dominance patterns were highly over represented in the observed data
205 compared to the random model (Fig. 3b), with full dominance being the most abundant category.
206 In instances of dominance, the non-dominant homeolog (in either subgenome) showed slightly
207 relaxed purifying selection on average³³ (Extended Data Fig. 9b). There were many instances of
208 domain expansion/neofunctionalization in the dominance categories (18%; Fig. 3d). However,
209 counter to expectations, 80% of the time it was the dominant homeolog that expanded its domain
210 (Supplementary Table 6). Dominance patterns were largely stable when we examined single cell
211 profiles from the maize inflorescence³⁷, with 72% of the root dominance patterns falling into the
212 same category in the inflorescence cells, swapping the dominant homeolog in a minority of cases,
213 while 26% display a subfunctionalization pattern and 2% co-express (Supplementary Table 6).
214 These observations suggest that the non-dominant homeolog could have a more essential role in
215 at least one context in the plant or may possibly be in the early stages of pseudogenization.

216 While subfunctionalized homeolog patterns were highly underrepresented compared to the
217 random model (Fig. 3b), homeolog pairs in this category showed a significantly higher proportion
218 of neofunctionalization compared to dominance categories (33% vs. 18%, Fisher's exact test, p-
219 value < 0.001; Fig. 3d, Supplementary Table 6). This was unexpected as it suggests that
220 subfunctionalization, rather than dominance, was more likely to be accompanied by
221 neofunctionalization. In addition, in contrast to the dominance categories above, there was no
222 evidence of relaxed purifying selection in either the homeolog that extended or the one that
223 retained the ancestral domain (Extended Data Fig. 10a). The result indicates that transcriptional
224 neofunctionalization and subfunctionalization are somehow linked. We cannot determine which

225 precedes the other, although our model suggests that completely random forces acting on both
 226 homeologs could account for the observed subfunctionalization, which could then serve as a
 227 transition state to neofunctionalization³⁸.

228
 229 Overall, we propose a model in which dosage compensation maintains stoichiometric balance for
 230 an extended period after whole genome duplication. This could be considered a form of
 231 quantitative subfunctionalization in which selection for full, ancestral dosage preserves the activity
 232 of both homeologs². In any case, such a mechanism would prevent pseudogenization and permit
 233 the retention of gene duplicates. The prolonged survival of both homeologs would thus enable
 234 spatial subfunctionalization that then enhances the likelihood that homeologs expand their
 235 expression domains (spatial neofunctionalization)—presumably through changes in cis-
 236 regulatory regions³⁹—leading to new transcriptional states in specialized cells.

237



238
 239

240 **Fig. 3: Gene pairs that subfunctionalize undergo a high rate of neofunctionalization.** a
 241 Conceptual schematic of patterns that characterize different categories of subfunctionalization vs.

242 dominance. **b** Observed distribution of genes pairs by their dominance (all categories <0) vs.
243 subfunctionalization (all categories >0) score, where a score of 1 reflects complete dominance
244 and -1 equal partitioning of the ancestral expression domain (subfunctionalization). Complete co-
245 expression = 0. Dark bars represent the observed data while light gray bars represent the random
246 model where expression domains of homeologs are randomly removed (see Methods), showing
247 an over-representation of dominance patterns and an under-representation of
248 subfunctionalization patterns. **c** Same schematic as in (a) but now showing novel expression
249 domains (neofunctionalization) in blue. **d** The same distribution as shown in (b) with
250 neofunctionalized events (domain expansion) mapped onto the distribution. Domain expansion is
251 more frequent in subfunctionalized homeolog pairs, while dominance patterns show less frequent
252 domain expansion with the vast majority being the dominant homeolog.

253 **Root Cap “Slime” Drives a Case of Rapid Cell-Type Divergence**

254 Using a three-taxa comparative approach, we next asked which cell types diverged most rapidly
255 in maize and sorghum compared to the outgroup *Setaria*. To compare cell identity across species,
256 we adapted MetaNeighbor, which uses neighbor voting to quantify the similarity of cell clusters
257 across datasets using a given marker set of genes and their orthologs⁴⁰ (Fig. 4a).

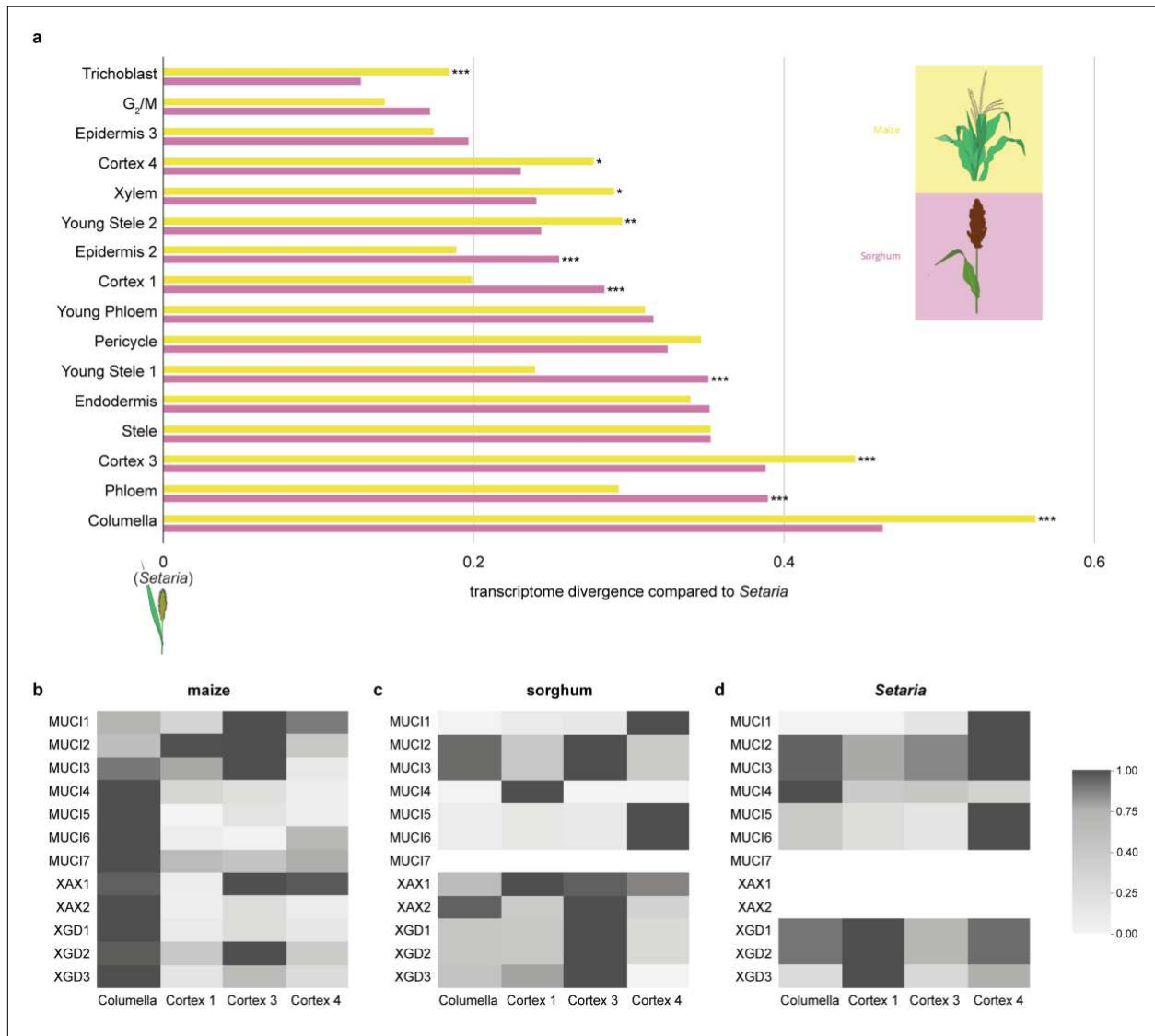
258 The analysis showed that transcriptomes of columella, phloem, cortex subcluster 3, endodermis,
259 pericycle, and stele cell types are the most divergent compared to *Setaria*, suggesting that the
260 function of these tissues diverged from *Setaria* before the maize-sorghum split. In addition, certain
261 cell types—such as cortex subcluster 1, phloem, and young stele subcluster 1—were statistically
262 different between maize and sorghum, implying additional divergence after the maize-sorghum
263 split. Interestingly, in maize, columella was among the most divergent cell types relative to *Setaria*
264 (Fig. 4a).

265 To further investigate the potential functions involved in columella divergence, we used a measure
266 of co-expression conservation to identify transcripts within clusters of interest that showed
267 divergent patterns of expression across species in co-expression networks^{35,41} (Supplementary
268 Table 7). We identified 443 genes displaying high expression divergence across species in
269 Columella. Many of these genes showed dramatic expression changes in spatial domain between
270 species, such as downy mildew resistant 6 (DMR6), which is expressed in columella and
271 epidermis in maize but in cortex and endodermis in sorghum (Extended Data Fig. 10b,c).
272 Furthermore, GO term analysis showed enrichment in enzymes leading to the synthesis of
273 mannose, raffinose, and oligosaccharides (Supplementary Table 7). These sugars and
274 carbohydrates are key components of mucilage, also called slime, whose roles include shaping
275 of the root-associated microbiome and lubricating the root-soil interface^{8,10,42–44}.

276 To further examine the potential change in mucilage-associated genes, we examined all genes
277 implicated in mucilage component synthesis^{10,11,45}. We found that this set of mucilage-annotated
278 genes were mostly expressed in maize columella, while, in sorghum and *Setaria*, they were
279 predominantly expressed in cortical layers (Fig. 4b,c,d). Overall, these results suggest that maize
280 underwent a relatively rapid cellular divergence in columella, in part, by recruiting a mucilage gene
281 expression module from ancestral expression pattern in the cortex. The most parsimonious model
282 is that the recruitment of the mucilage module occurred before the maize whole genome
283 duplication, as both maize homeologs, when conserved, tended to share expression in the
284 columella.

285 Overall, the high-resolution comparative analysis provides evidence for three phenomena—
286 dosage compensation, subfunctionalization, and neofunctionalization—as a driver of cell type

287 divergence following whole genome duplication. The results extend the genomic balance model³⁵
 288 by showing its effects are reinforced after duplication at the transcriptional level. Furthermore, the
 289 data show that both single-cell and single-nucleus profiling produce high-fidelity maps of cell-type
 290 specific expression, with a combination of the two leveraging their complementary strengths. Cell
 291 type specific maps in a well-annotated species can then serve as an anchor for neighboring
 292 taxonomic groups, rapidly generating homologous maps in a related group of plant species.
 293 These maps provide powerful tools for analysis of cell type evolution that can identify genes
 294 associated with specific traits.



295

296 **Fig. 4: Differential evolution between cell types reveals that columella is highly divergent**
 297 **compared to *Setaria*.** **a** MetaNeighbor analysis showing transcriptome conservation scores
 298 between cell types in maize and sorghum compared to the outgroup *Setaria*. High conservation
 299 scores reflect similar transcriptomes across species. Statistical significance between maize and
 300 sorghum was performed using the Hanley McNeil test (see Methods; p-values:
 301 * <0.05 , ** <0.01 , *** <0.001). **b–d** Mucilage gene expression heatmaps for maize (b), sorghum (c),
 302 and *Setaria* (d) columella and cortical layers, showing predominant columella expression in maize

303 compared to cortical layers in sorghum and *Setaria*. Maize genes are identified in (b) and their
304 corresponding sorghum and *Setaria* orthologs are listed on the same row in (c) and (d).

305

306 **Acknowledgements**

307 We thank Michael Purugganan and Gloria Coruzzi for helpful comments. This work was funded
308 by National Science Foundation (IOS-1934388) to K.D.B., D.J, and J.G., the National Institutes
309 of Health (R35GM136362) to K.D.B., and Human Frontiers of Science (LT000972/2018-L) to
310 B.G., startup funds from the University of California Riverside to S.C.G. In addition, M.P. is funded
311 by the William Randolph Hearst Scholarship from the School of Biological Sciences. J.G. is also
312 supported by the National Institutes of Health (R01 LM012736 and R01 MH113005).

313

314 **Contributions**

315

316 B.G. and K.D.B designed the research. B.G. generated all single-cell and single-nucleus RNA-
317 seq data. M.A.M. and B.G. designed the single-nucleus RNA-seq protocol. R.R. and B.G.
318 performed the whole mount in-situ hybridization analysis. X.X. and D.J. performed the tissue
319 preparation and histology for the spatial transcriptomics analysis. S.C.G. and B.G. conceived the
320 analysis strategy and performed the tests for dosage compensation. M.P and, J.G. performed the
321 MetaNeighbor and CoCoCoNet analysis. B.G. analyzed all the data. K.D.B., B.G., and R.R. wrote
322 the manuscript.

323

324 The authors declare no competing interests.

325

326 Supplementary Information is available for this paper.

327

328 Material requests should be addressed to K.D.B.

329 **References**

330 1. Birchler, J. A. & Veitia, R. A. Gene balance hypothesis: Connecting issues of dosage
331 sensitivity across biological disciplines. *Proc. Natl. Acad. Sci. U. S. A.* **109**, 14746–14753
332 (2012).

333 2. Force, A. *et al.* Preservation of duplicate genes by complementary, degenerative
334 mutations. *Genetics* **151**, 1531–1545 (1999).

335 3. Marioni, J. C. & Arendt, D. How Single-Cell Genomics Is Changing Evolutionary and
336 Developmental Biology. *Annu. Rev. Cell Dev. Biol.* **33**, 537–553 (2017).

337 4. Shafer, M. E. R. Cross-Species Analysis of Single-Cell Transcriptomic Data. *Front. Cell*
338 *Dev. Biol.* **7**, 175 (2019).

339 5. Kajala, K. *et al.* Innovation, conservation, and repurposing of gene function in root cell
340 type development. *Cell* **184**, 3333–3348.e19 (2021).

341 6. Swigonova, Z. *et al.* On the tetraploid origin of the maize genome. *Comp. Funct.*
342 *Genomics* **5**, 281–284 (2004).

- 343 7. Swigonova, Z. Close Split of Sorghum and Maize Genome Progenitors. *Genome Res.* **14**,
344 1916–1923 (2004).
- 345 8. Van Deynze, A. *et al.* Nitrogen fixation in a landrace of maize is supported by a mucilage-
346 associated diazotrophic microbiota. *PLoS Biol.* **16**, 1–21 (2018).
- 347 9. Hasan, S. A., Rabei, S. H., Nada, R. M. & Abogadallah, G. M. Water use efficiency in the
348 drought-stressed sorghum and maize in relation to expression of aquaporin genes. *Biol.*
349 *Plant.* **61**, 127–137 (2017).
- 350 10. Kozlova, L. V., Nazipova, A. R., Gorshkov, O. V., Petrova, A. A. & Gorshkova, T. A.
351 Elongating maize root: zone-specific combinations of polysaccharides from type I and
352 type II primary cell walls. *Sci. Rep.* **10**, 1–20 (2020).
- 353 11. Ma, W. *et al.* The mucilage proteome of maize (*Zea mays* L.) primary roots. *J. Proteome*
354 *Res.* **9**, 2968–2976 (2010).
- 355 12. Schittenhelm, S. & Schroetter, S. Comparison of Drought Tolerance of Maize, Sweet
356 Sorghum and Sorghum-Sudangrass Hybrids. *J. Agron. Crop Sci.* **200**, 46–53 (2014).
- 357 13. Zhang, Y. *et al.* Differentially regulated orthologs in sorghum and the subgenomes of
358 maize. *Plant Cell* **29**, 1938–1951 (2017).
- 359 14. Zheng, Z. *et al.* Shared Genetic Control of Root System Architecture between *Zea mays*
360 and *Sorghum bicolor*1[OPEN]. *Plant Physiol.* **182**, 977–991 (2020).
- 361 15. Schnable, J. C., Springer, N. M. & Freeling, M. Differentiation of the maize subgenomes
362 by genome dominance and both ancient and ongoing gene loss. *Proc. Natl. Acad. Sci. U.*
363 *S. A.* **108**, 4069–4074 (2011).
- 364 16. Denyer, T. *et al.* Spatiotemporal Developmental Trajectories in the Arabidopsis Root
365 Revealed Using High-Throughput Single-Cell RNA Sequencing. *Dev. Cell* **48**, 840–852.e5
366 (2019).
- 367 17. Efroni, I. *et al.* Root Regeneration Triggers an Embryo-like Sequence Guided by
368 Hormonal Interactions. *Cell* **165**, 1721–1733 (2016).
- 369 18. Jean-Baptiste, K. *et al.* Dynamics of Gene Expression in Single Root Cells of
370 *Arabidopsis thaliana*. *Plant Cell* **31**, 993–1011 (2019).
- 371 19. Ryu, K. H., Huang, L., Kang, H. M. & Schiefelbein, J. Single-Cell RNA Sequencing
372 Resolves Molecular Relationships Among Individual Plant Cells. *Plant Physiol* **179**, 1444–
373 1456 (2019).
- 374 20. Shulse, C. N. *et al.* High-Throughput Single-Cell Transcriptome Profiling of Plant Cell
375 Types. *Cell Rep.* **27**, 2241–2247.e4 (2019).
- 376 21. Zhang, T.-Q., Xu, Z.-G., Shang, G.-D. & Wang, J.-W. A Single-Cell RNA Sequencing
377 Profiles the Developmental Landscape of Arabidopsis Root. *Mol. Plant* **12**, 648–660
378 (2019).

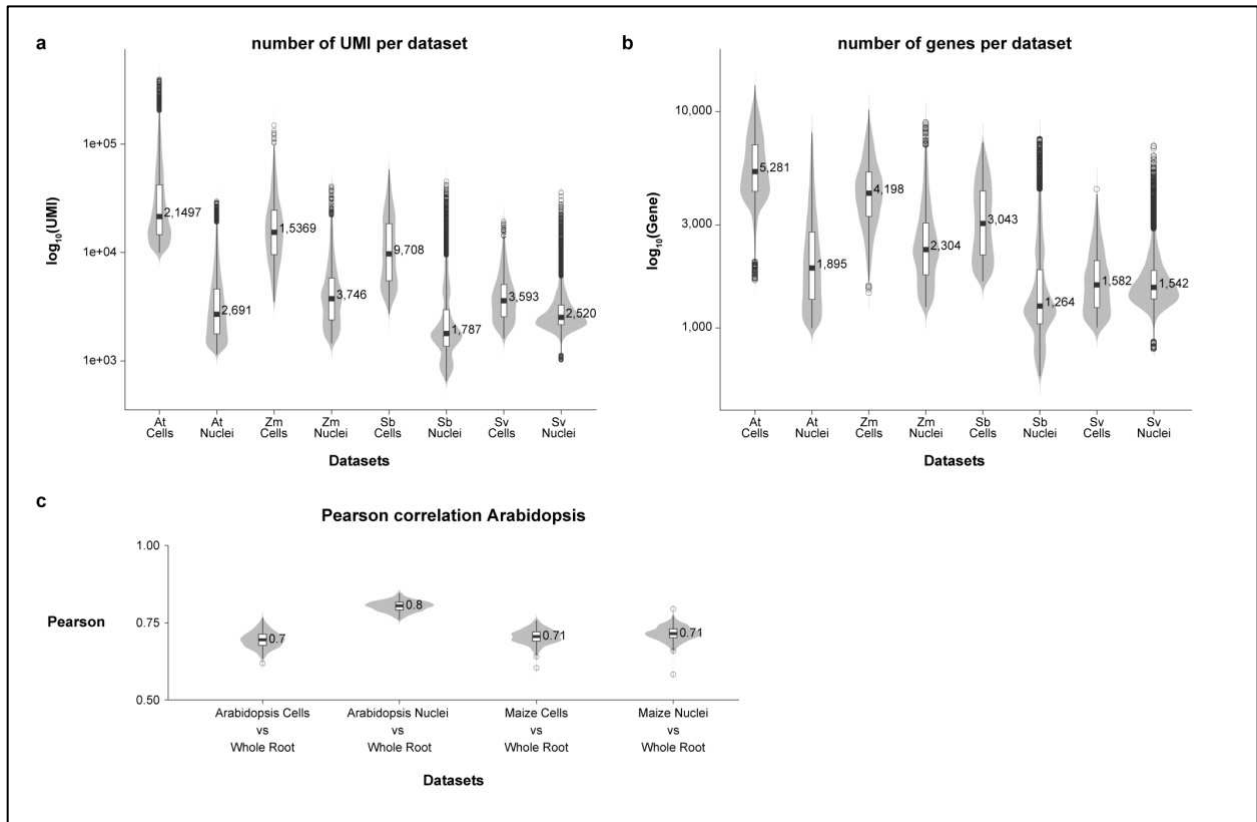
- 379 22. Farmer, A., Thibivilliers, S., Ryu, K. H., Schiefelbein, J. & Libault, M. Single-nucleus RNA
380 and ATAC sequencing reveals the impact of chromatin accessibility on gene expression
381 in Arabidopsis roots at the single-cell level. *Mol. Plant* **14**, 372–383 (2021).
- 382 23. Long, Y. *et al.* FlsnRNA-seq: protoplasting-free full-length single-nucleus RNA profiling in
383 plants. *Genome Biol.* **22**, 66 (2021).
- 384 24. Marand, A. P., Chen, Z., Gallavotti, A. & Schmitz, R. J. A cis-regulatory atlas in maize at
385 single-cell resolution. *Cell* **184**, 3041–3055.e21 (2021).
- 386 25. Ortiz-Ramírez, C. *et al.* Ground tissue circuitry regulates organ complexity in maize and
387 Setaria. *Science (80-.)*. **374**, 1247–1252 (2021).
- 388 26. Allen, A. M. *et al.* A single-cell transcriptomic atlas of the adult Drosophila ventral nerve
389 cord. *Elife* **9**, (2020).
- 390 27. Ding, J. *et al.* Systematic comparison of single-cell and single-nucleus RNA-sequencing
391 methods. *Nat. Biotechnol.* **38**, 737–746 (2020).
- 392 28. Narsai, R. *et al.* Genome-Wide Analysis of mRNA Decay Rates and Their Determinants
393 in Arabidopsis thaliana. *Plant Cell Online* **19**, 3418–3436 (2007).
- 394 29. Sorenson, R. S., Deshotel, M. J., Johnson, K., Adler, F. R. & Sieburth, L. E. Arabidopsis
395 mRNA decay landscape arises from specialized RNA decay substrates, decapping-
396 mediated feedback, and redundancy. *Proc. Natl. Acad. Sci. U. S. A.* **115**, E1485–E1494
397 (2018).
- 398 30. Hafemeister, C. & Satija, R. Normalization and variance stabilization of single-cell RNA-
399 seq data using regularized negative binomial regression. *bioRxiv* 1–15 (2019).
400 doi:10.1101/576827
- 401 31. Muyle, A., Marais, G. A. B., Bačovský, V., Hobza, R. & Lenormand, T. Dosage
402 compensation evolution in plants: theories, controversies and mechanisms. *Philos. Trans.
403 R. Soc. B Biol. Sci.* **377**, (2022).
- 404 32. Coate, J. E., Farmer, A. D., Schiefelbein, J. W. & Doyle, J. J. Expression Partitioning of
405 Duplicate Genes at Single Cell Resolution in Arabidopsis Roots. *Front. Genet.* **11**, (2020).
- 406 33. Hughes, T. E., Langdale, J. A. & Kelly, S. The impact of widespread regulatory
407 neofunctionalization on homeolog gene evolution following whole-genome duplication in
408 maize. *Genome Res.* **24**, 1348–1355 (2014).
- 409 34. Yanai, I. *et al.* Genome-wide midrange transcription profiles reveal expression level
410 relationships in human tissue specification. *Bioinformatics* **21**, 650–659 (2005).
- 411 35. Birchler, J. A. & Veitia, R. A. The Gene Balance Hypothesis: From Classical Genetics to
412 Modern Genomics. *Plant Cell* **19**, 395–402 (2007).
- 413 36. Ohno, S. *Evolution by Gene Duplication*. *Teratology* **1**, (Springer Berlin Heidelberg,
414 1970).

- 415 37. Xu, X. *et al.* Single-cell RNA sequencing of developing maize ears facilitates functional
416 analysis and trait candidate gene discovery. *Dev. Cell* **56**, 557-568.e6 (2021).
- 417 38. Rastogi, S. & Liberles, D. A. Subfunctionalization of duplicated genes as a transition state
418 to neofunctionalization. *BMC Evol. Biol.* **5**, 28 (2005).
- 419 39. Lu, Z. *et al.* The prevalence, evolution and chromatin signatures of plant regulatory
420 elements. *Nat. Plants* **5**, 1250–1259 (2019).
- 421 40. Crow, M., Paul, A., Ballouz, S., Huang, Z. J. & Gillis, J. Characterizing the replicability of
422 cell types defined by single cell RNA-sequencing data using MetaNeighbor. *Nat.*
423 *Commun.* **9**, 884 (2018).
- 424 41. Lee, J., Shah, M., Ballouz, S., Crow, M. & Gillis, J. CoCoCoNet: Conserved and
425 comparative co-expression across a diverse set of species. *Nucleic Acids Res.* **48**,
426 W566–W571 (2021).
- 427 42. Galloway, A. F., Knox, P. & Krause, K. Sticky mucilages and exudates of plants: putative
428 microenvironmental design elements with biotechnological value. *New Phytol.* **225**, 1461–
429 1469 (2020).
- 430 43. Roy, S. S., Mitra, B., Sharma, S., Das, T. K. & Babu, C. R. Detection of root mucilage
431 using an anti-fucose antibody. *Ann. Bot.* **89**, 293–299 (2002).
- 432 44. Werker, E. & Kiselev, M. Mucilage on the root surface and root Hairs of Sorghum:
433 Heterogeneity in structure, manner of production and site of accumulation. *Ann. Bot.* **42**,
434 809–816 (1978).
- 435 45. Voiniciuc, C., Guenl, M., Schmidt, M. H.-W. & Usadel, B. Highly Branched Xylan Made by
436 IRX14 and MUC121 Links Mucilage to Arabidopsis Seeds. *Plant Physiol.* **169**,
437 pp.01441.2015 (2015).
- 438 46. Efroni, I., Ip, P.-L., Nawy, T., Mello, A. & Birnbaum, K. D. Quantification of cell identity
439 from single-cell gene expression profiles. *Genome Biol.* **16**, 9 (2015).
- 440 47. Stuart, T. *et al.* Comprehensive Integration of Single-Cell Data. *Cell* **177**, 1888-1902 e21
441 (2019).
- 442 48. Hernández Coronado, M. *et al.* Repel or Repair: Plant Glutamate Receptor-Like
443 Channels Mediate a Defense vs. Regeneration Tradeoff. *SSRN Electron. J.* (2021).
444 doi:10.2139/ssrn.3818443
- 445 49. Jackson, D., Veit, B. & Hake, S. Expression of maize KNOTTED1 related homeobox
446 genes in the shoot apical meristem predicts patterns of morphogenesis in the vegetative
447 shoot. *Development* **120**, 405–413 (1994).

448

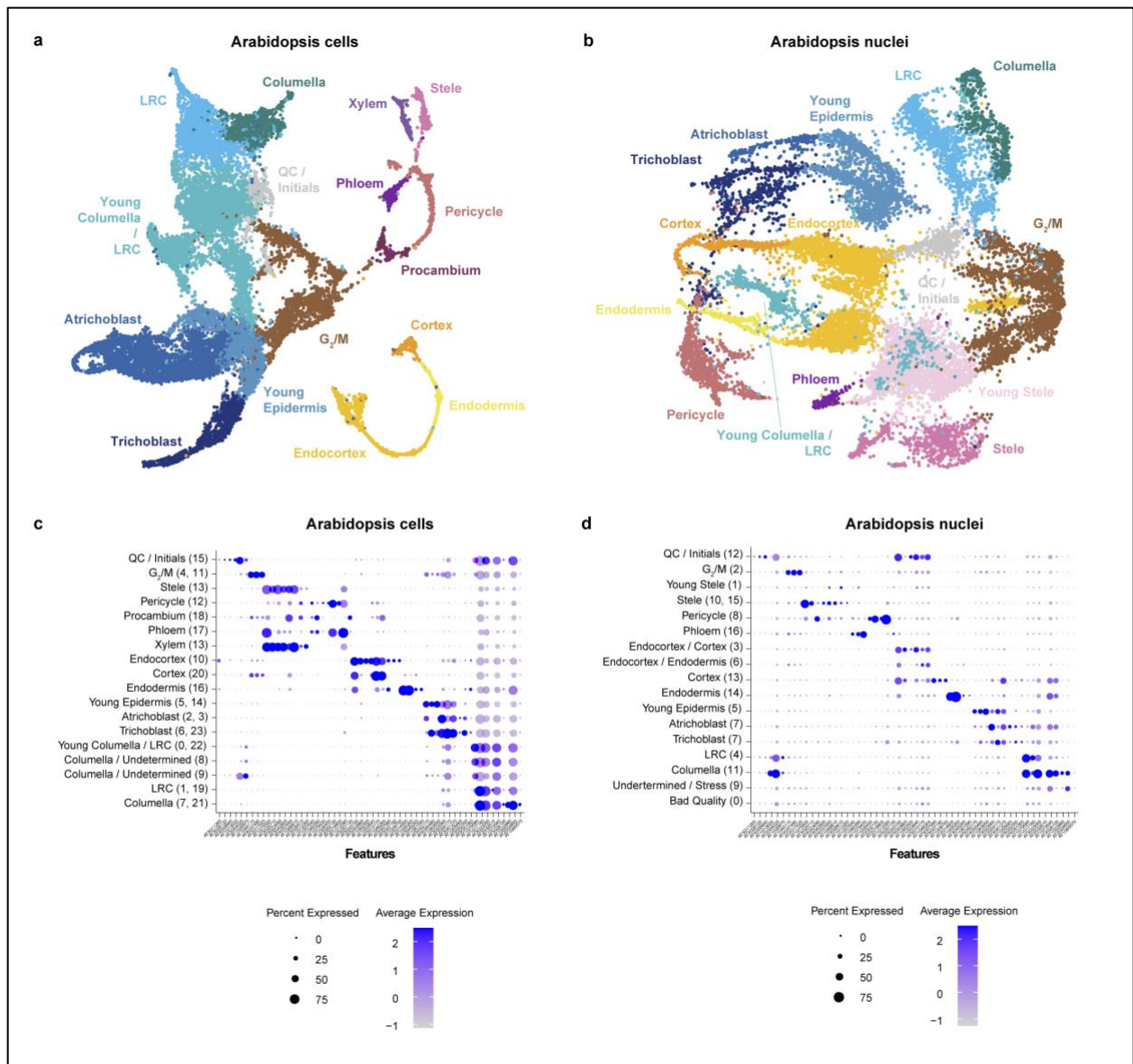
449

450 **Extended Data**



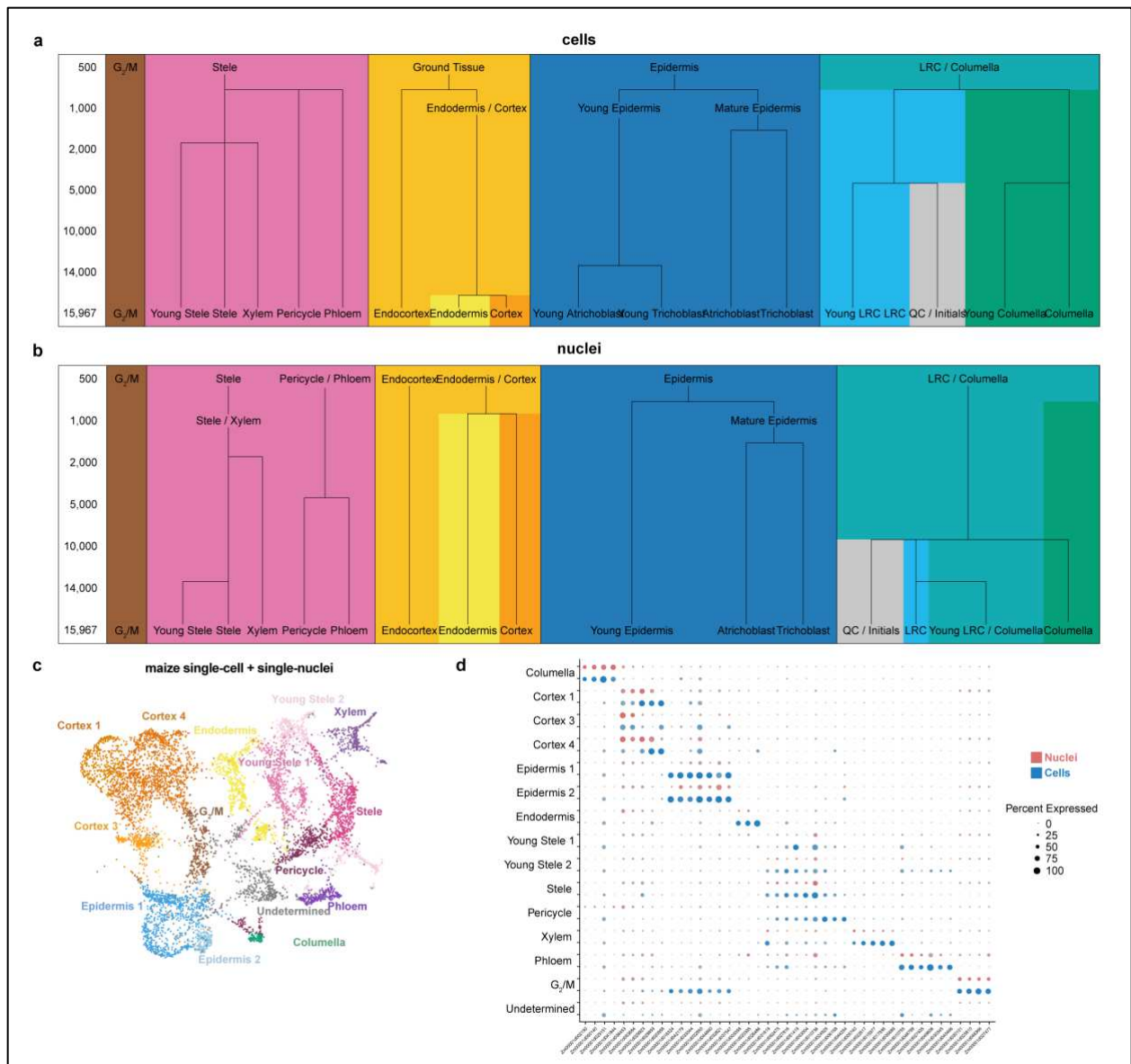
451

452 **Extended Data Fig. 1:** **a** Violin plot distribution of the number of UMI detected among cells vs.
 453 nuclei in Arabidopsis, maize, sorghum, and *Setaria*. **b** Violin plot distribution of the number of
 454 genes detected among cells vs. nuclei in the same species as in (a). Black bar is median, box
 455 represents 25th and 75th percentile, and vertical line is the 5th and 95th percentile. **c** Pearson
 456 correlation distributions of gene expression from single-cell or single-nuclei compared to whole-
 457 root RNAseq in Arabidopsis and maize, performed on replicate runs of single-cell or single-
 458 nucleus profiles.
 459



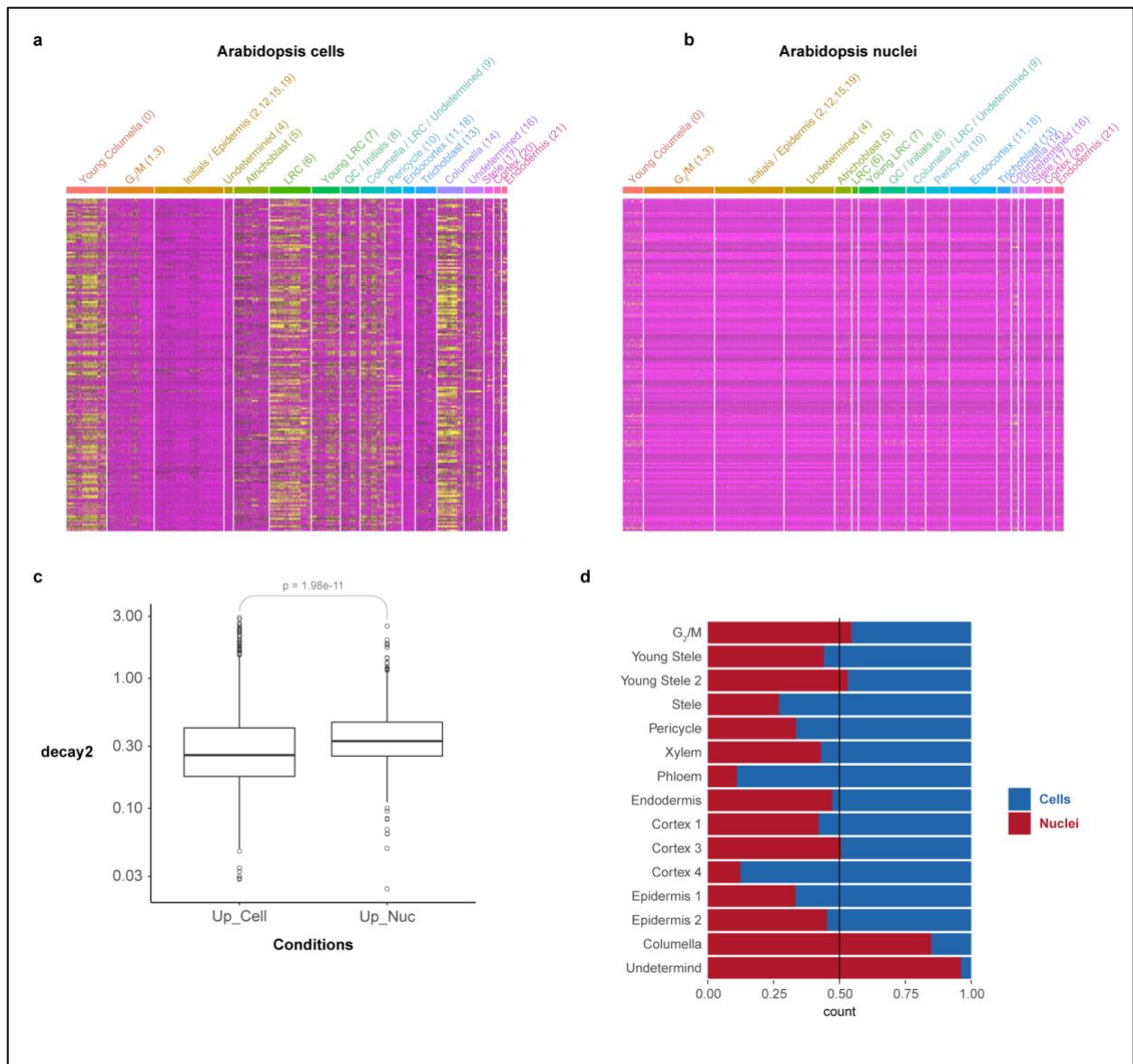
460
461
462
463
464
465
466
467
468

Extended Data Fig. 2: a, b UMAP clustering in Arabidopsis single-cell (**a**) and single-nucleus (**b**) datasets clustered independently, showing clusters with the same assigned cell identities. **c, d** Dot plots showing cluster-wise average expression levels and percent of cells detected for a set of Arabidopsis cell type markers in cells (**c**) vs. nuclei (**d**). The plot shows that markers for a given cell type identified in one profile type (cells or nuclei) show largely the same enriched expression in the second profile type.



469
 470
 471
 472
 473
 474
 475
 476
 477
 478
 479
 480

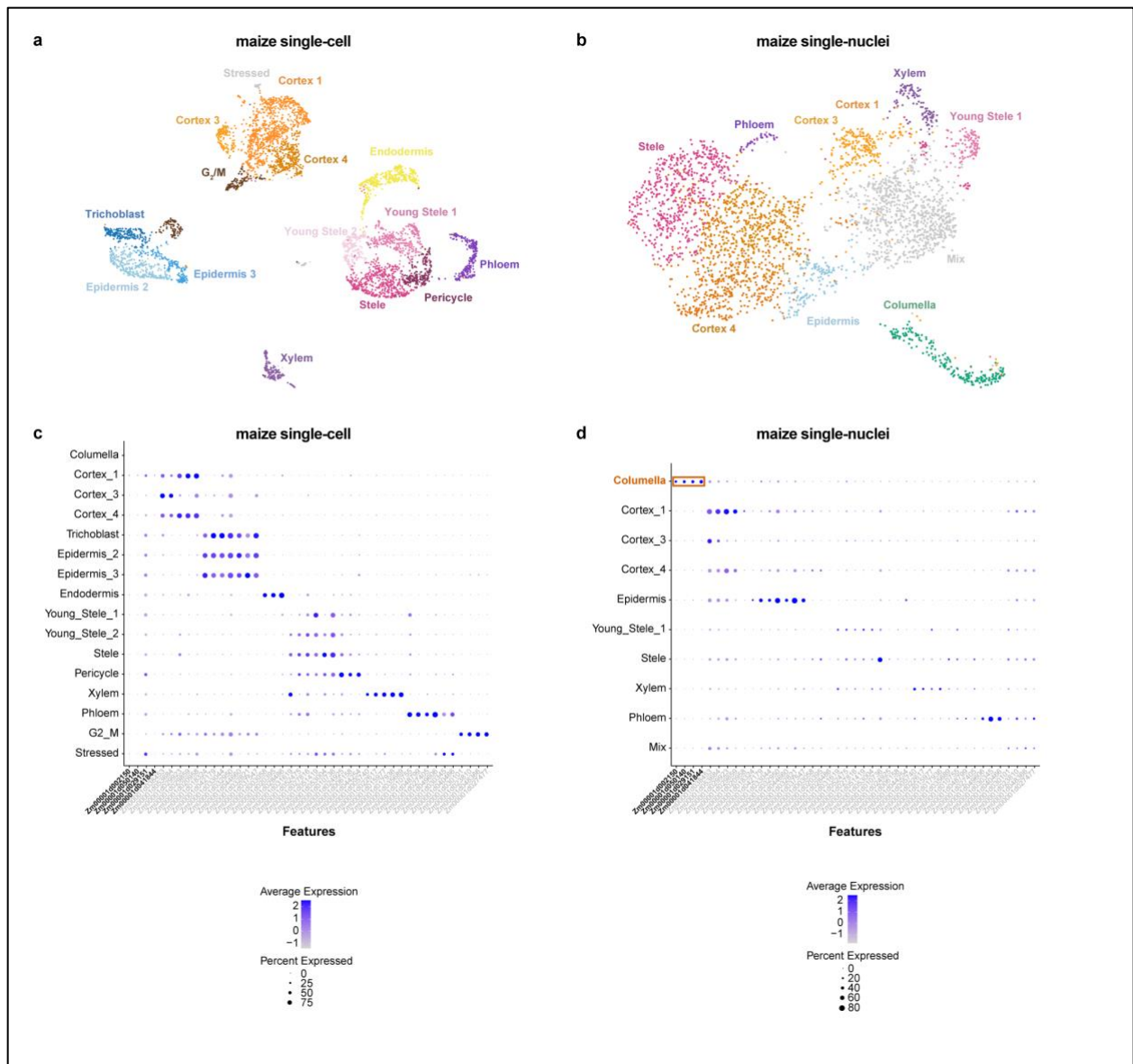
Extended Data Fig. 3: **a** Arabidopsis down-sampling analysis shows the number of cells needed to resolve clusters into different cell types within tissues. A branch signifies the number of cells needed for Seurat to distinguish a new cluster from a cloud of points when analyzing lower numbers of cells. **b** A similar analysis using the single-nucleus dataset, showing more nuclei are needed to resolve clusters compared to cell profiles in (a). Tracking the branches of graphs in (a) vs. (b) leads to a rule-of-thumb that two-fold more nuclei than cells are needed to identify clusters. **c** UMAPs of the combined maize single-cell and single-nucleus datasets with clusters colored by cell type. **d** Dot plot of maize marker genes in cells (blue) or in nuclei (red), showing concordance of marker expression in the two datasets.



481
482

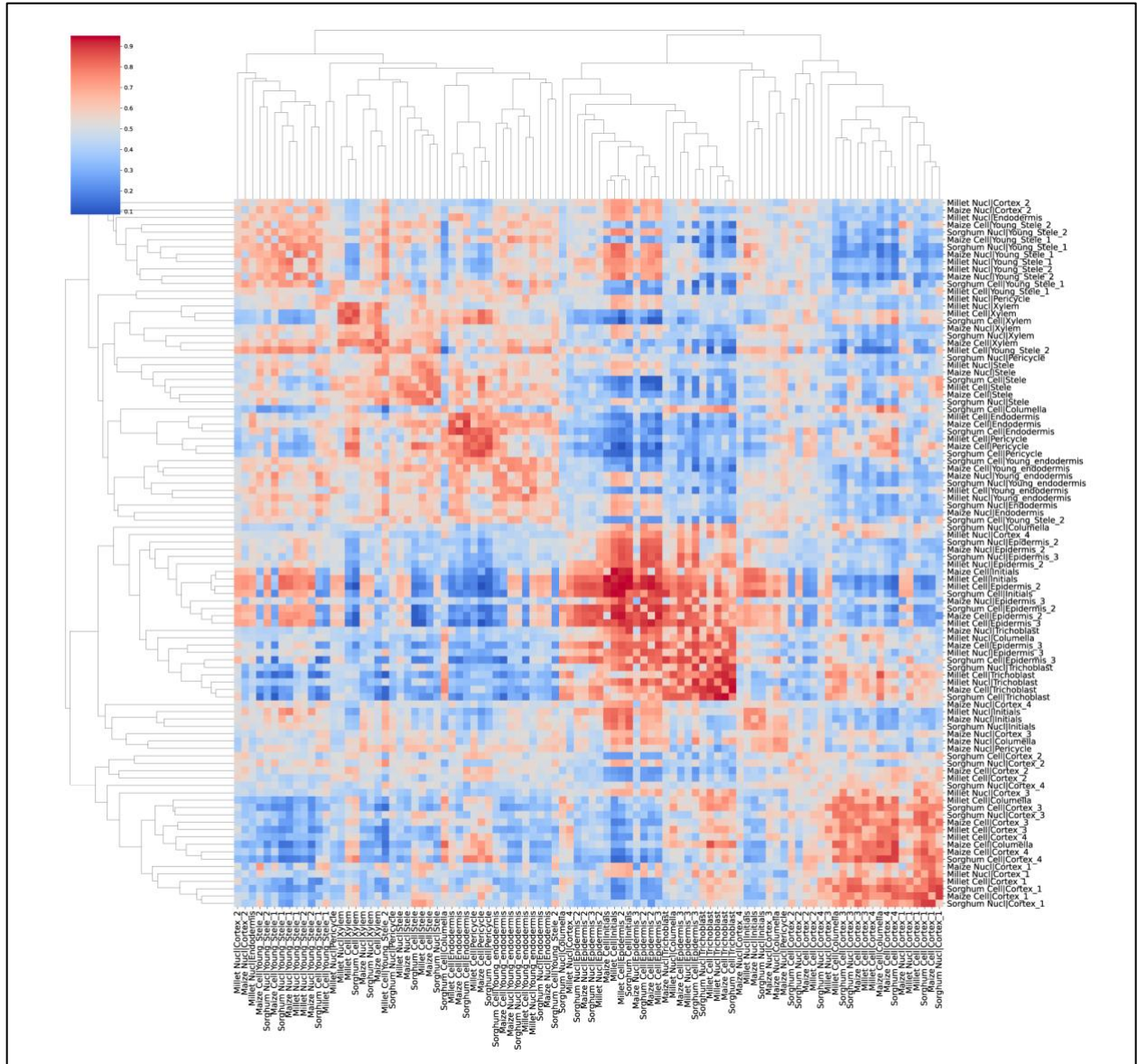
483 **Extended Data Fig. 4: a, b** Heatmaps of Arabidopsis genes known to be induced by protoplast
 484 generation (Birnbaum et al., 2003) showing their expression in cells (a) vs. nuclei (b). The analysis
 485 shows that stress-induced genes also have higher expression in cells vs. nuclei, with particular
 486 induction in specific cell types. **c** Distribution of expression levels of genes annotated for mRNA
 487 decay in cells or in nuclei. The same genes are analyzed in both cells and nuclei. A significant
 488 increase in mRNA decay gene expression level was detected in nuclei (Wilcoxon rank sum test,
 489 p -value = $1.98e-11$). **d** Proportion of cells (blue) vs nuclei (red) present in each cell type cluster
 490 for maize.

491



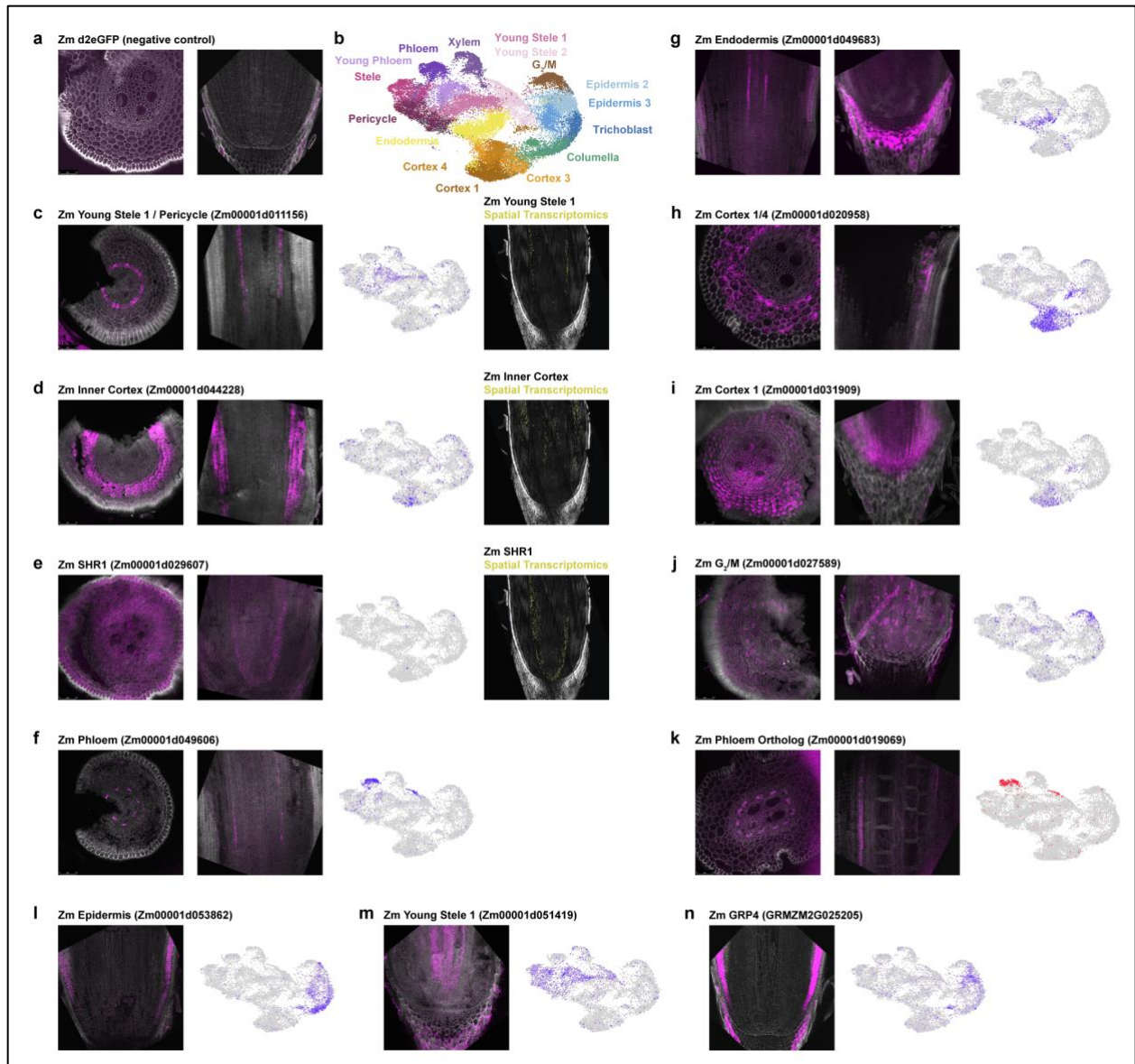
492

493 **Extended Data Fig. 5: a, b** UMAPs of maize single-cell and single-nucleus datasets clustered
 494 independently. Only the single-nucleus dataset displays a cluster annotated as columella, which
 495 is absent in the single-cell dataset. **c, d** Dot plot of maize marker genes for each cell type cluster,
 496 showing expression in cells (**c**) and in nuclei (**d**) datasets independently. Markers for columella
 497 outlined in the red box are only present in the nuclei dataset.
 498



499
 500
 501
 502
 503
 504
 505

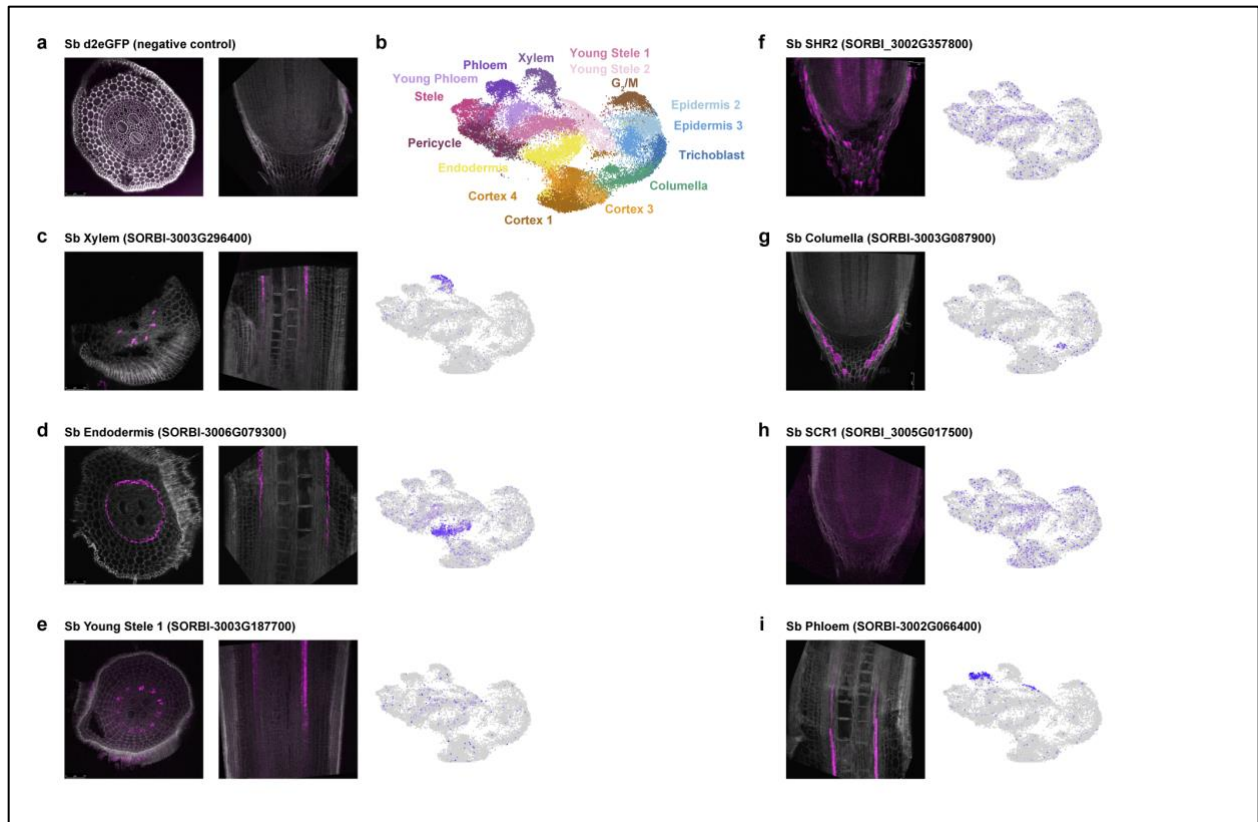
Extended Data Fig. 6: AUROC test as performed in MetaNeighbor comparing every cell type in maize, sorghum, and *Setaria* using single-cell and single-nucleus datasets separately, showing that cells and nuclei largely group by cell type and not by either species or profiling method (cells vs. nuclei).



506

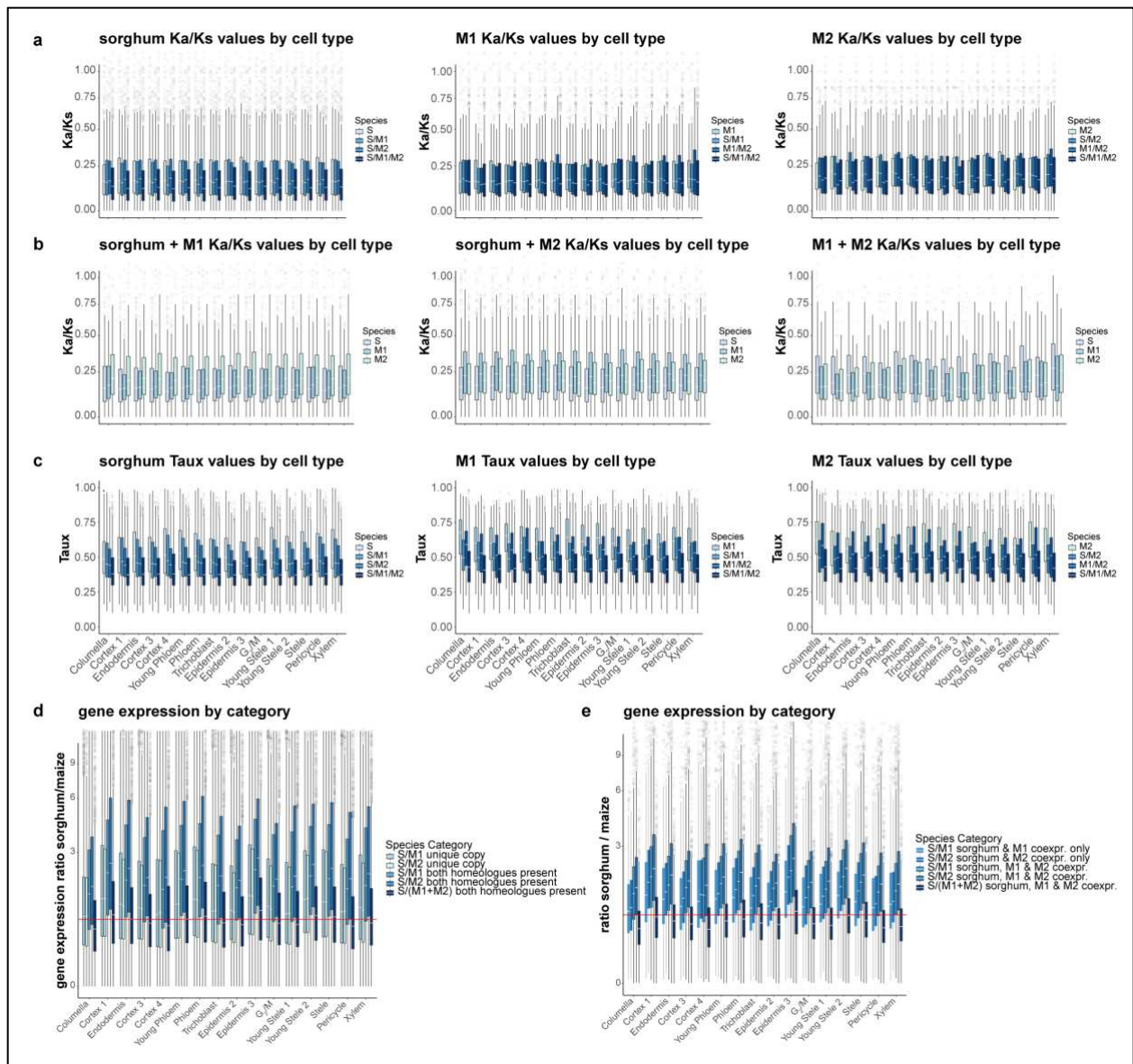
507 **Extended Data Fig. 7:** a-n *in situ* hybridization using Hairpin Chain Reaction (HCR) probes
 508 labeling various transcripts in maize. UMAPs showing each transcript's cluster localization are
 509 shown next to each probe's fluorescent image. Additionally, spatial transcriptomics imaging data
 510 is shown for a three probes in maize (c-e), further validating the cluster annotations. The
 511 minimum/maximum values for each fluorescence channel (grey: autofluorescence, magenta:
 512 HCR probes) have been adjusted to show the localization more clearly in the merged image.

513



514

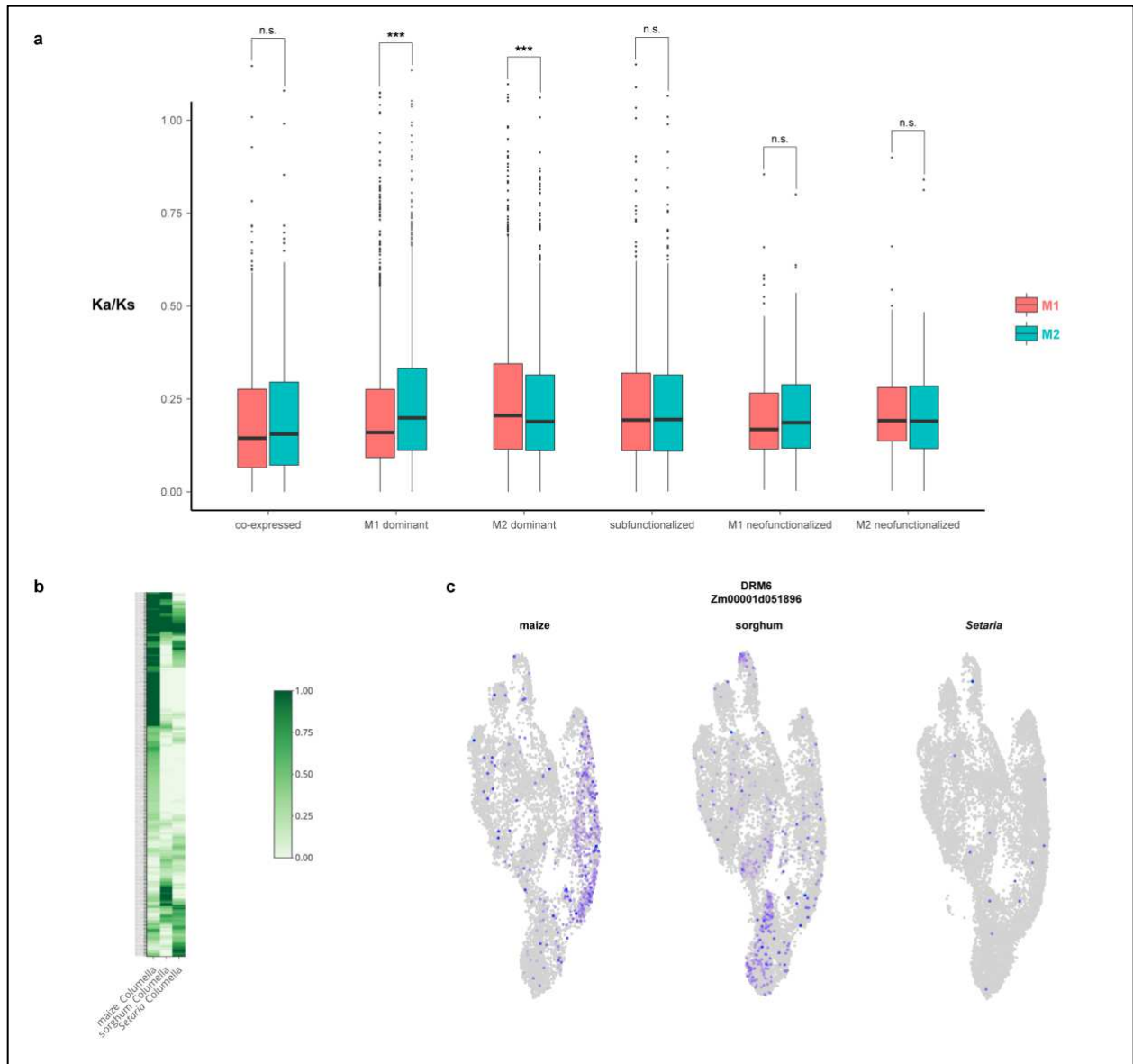
515 **Extended Data Fig. 8:** a-i *in situ* hybridization using Hairpin Chain Reaction (HCR) probes
 516 labeling various transcripts in sorghum. UMAPs showing each transcript's cluster localization are
 517 shown next to each probe's fluorescent image. The minimum/maximum values for each
 518 fluorescence channel (grey: autofluorescence, magenta: HCR probes) have been adjusted to
 519 show the localization more clearly in the merged image.
 520



521
522

523 **Extended Data Fig. 9:** a Distribution of Ka/Ks values (relative to *Setaria*) of genes or homeologs
 524 among the different genomes or subgenomes across cell types (labeled on bottom graph in (c)),
 525 with expression pattern by cell type shown in subcategories. Letters represent co-expression
 526 categories of sorghum genes and maize homeologs in a given cell type (i.e. S = only sorghum
 527 gene expressed in cell type, SM1 = sorghum and M1 genes expressed in cell type, etc.). b
 528 Distribution of Ka/Ks values in cases of specific co-expression patterns (listed in title of graph) for
 529 each cell type (labeled in (c)). For example, in the first graph on the left, when sorghum and M1
 530 are co-expressed in a given cell type, the M2 homeolog always has a slightly higher Ka/Ks value,
 531 showing they are under less stringent purifying selection. c Distribution of Taux (cell specificity)
 532 measure of genes in different cell types according to comparative expression patterns. For
 533 example, in the first graph on the left, sorghum genes have a higher Taux (i.e. more specificity)
 534 when they are not expressed in the same cell type as M1 or M2. d Analysis similar to Fig. 2d
 535 showing expression ratios between sorghum and M1 or M2 homeologs, here broken down by cell
 536 type and genomic status of M1 or M2 homeolog (retained or lost). The trends show the strongest
 537 dosage compensation at the cell-type level when both copies are retained and expressed. e

538 Analysis similar to Fig. 2d on only the retained homeologs showing ratios of expression between
539 sorghum and M1 and/or M2 homeologs broken down by cell type depending on their co-
540 expression with sorghum. The trends show dosage compensation by cell type only when the two
541 homeologs are expressed in that cell type.



542

543 **Extended Data Fig. 10: a** Ka/Ks distribution of co-expressed vs dominant vs neofunctionalized
 544 genes. Co-expressed genes display higher purifying selection. However, only the non-dominant
 545 ortholog shows an increase in Ka/Ks distribution compared to the dominant ortholog, while other
 546 categories don't display significant differences in purifying selection. Pairwise Wilcoxon Test ***
 547 $p < 0.001$. **b** Heatmap represents the 443 most divergent genes across species in the CoCoCoNet
 548 analysis, showing expression differences in the columella of the three grass species. **c** Example
 549 of the gene DMR6 switching its expression between columella in maize vs epidermis / cortex in
 550 sorghum.

551

552

553 **Methods**

554 **Plant growth conditions**

555 Seeds of *Arabidopsis thaliana* var Col0, *Zea maize* B73 and *Sorghum bicolor* Btx623, *Setaria*
556 *viridis* (Accession PI 669942, U.S. National Plant Germplasm System)

557 *Arabidopsis* seeds were imbibed for 48h at 4°C before being surface sterilized and placed on a
558 nylon mesh (110µm) and agar plates-1/2 × Murashige and Skoog salts (Sigma M5524), 0.5%
559 sucrose, 0.8% Agar (Sigma A1296). Plant were transferred vertically in growth chambers set to
560 23°C and a 16 h light/8 h dark cycle (400 µmol m⁻² s⁻¹). Root tips were collected 7 days after
561 transfer, cut with a feather scalpel at 150µm from the tip and directly transferred to either the
562 protoplast solution at room temperature or the nuclei lysis buffer at 4°C.

563
564 Maize, *Sorghum* seeds were sterilized using bleach (1.5% active chloride) and 0.001% tween
565 20 for 20mins then 4% chloramine T for 20mins. *Setaria* seed germination was induced by
566 incubation in 4% liquid smoke (Colgin, Authentic Natural Hickory) at 29°C for 24h. Then *setaria*
567 seeds were sterilized using bleach (1.5% active chloride) and 0.001% tween 20 for 20mins. All
568 seeds were placed between two layers of brown paper (Anchor Paper&Cie., 38# regular), rolled
569 and covered with aluminum foil to prevent roots from exposure to direct light. Rolls were placed
570 in a bucket of tap water at 28/24°C at 16 h light/8 h dark cycle (250 µmol m⁻² s⁻¹) for 7 days
571 (15days for *Setaria*) before harvesting the root tips. Primary and seminal root tips were cut using
572 a fine scalpel at 0.5 cm from the tip for Maize and *Sorghum*, 0.2cm from the tip for *Setaria* and
573 transferred either to the pre-incubation solution for single cell or to the nuclei lysis buffer.

574

575 **Protoplasting**

576 Protoplasts were generated from primary and seminal roots as described previously⁴⁶. For Maize,
577 *Sorghum* and *Setaria*, roots were cut above the meristem as describe above and placed in
578 pretreatment solution containing L-cysteine for 40 min (3% sorbitol, 2.5mM L-cysteine, 20mM
579 MES, pH 5.8 with Tris) to improve enzyme efficiency and cell wall digestion. Cell walls were
580 digested for 90 min in an enzyme solution optimized for monocot roots (Mannitol 8%, 400mM,
581 MES 20mM, KCl 20mM, CaCl₂ 40mM, pH 5.8 with Tris, BSA 100ug/ml, 2% cellulase “Onozuka”
582 RS, 1.2% cellulase “Onozuka” R10, 0.4% macerozyme R-10 (all three Yakult Pharmaceutical
583 Industry CO.), 0.36% pectolyase Y-23 (MP Biomedicals)). Protoplast were then filtered through a
584 40 µm cell strainer and transferred to microcentrifuge tubes for centrifugation.

585 For *Arabidopsis*, roots were cut above the meristem as described above and placed in an enzyme
586 solution optimized for *Arabidopsis* (Mannitol 8%, 400mM, MES 20mM, KCl 20mM, CaCl₂ 40mM,
587 pH 5.8 with Tris, BSA 100ug/ml, 1.2% cellulase “Onozuka” R10, 0.4% macerozyme R-10 (all three
588 Yakult Pharmaceutical Industry CO.)). Protoplast were then filtered through a 20 µm cell strainer
589 and transferred to microcentrifuge tubes for centrifugation.

590

591 Protoplast were centrifuge for 3 min at 500 x g and the pellets were washed and resuspended in
592 washing solution twice (Mannitol 8%, MES 20mM, KCl 20mM, CaCl₂ 10mM, pH 5.8 with Tris,
593 BSA 100ug/ml) and used immediately for single-cell RNAseq.

594 An aliquot of protoplasts was stained with trypan blue (0.2% final) and check on hemacytometer
595 under microscope to determine viability and cell concentration before loading into the 10x
596 chromium.

597

598 **Nuclei extraction**

599 For all species, root tips are directly transferred in prechilled lysis buffer (0.3M sucrose, 15mM
600 Tris HCl pH8, 60mM KCl, 15mM NaCl, 2mM EDTA, Spermine 0.5mM, Spermidine 0.5mM, 15mM
601 MES, 0.1% Triton, 5mM DTT*, 1mM PMSF* , 1% Plant Protease Inhibitors* 1ml(Sigma P9599),
602 BSA 0.4%*, RNase inhibitor 0.2u/ul*, (* add last minute)). Roots are chopped on ice with scalpel

603 blades for 5-10mins and transferred into a pre-chilled dounce homogenizer (Kimble, 885302).
604 Pestle is moved for 10 back and forth, sample keep on ice for 10mins then additional 10 back and
605 forth. Then root extracts are filtered at 20 μ m into a centrifuge tube and centrifuge for 10mins at
606 500g (Maize, Sorghum, Setaria) or 1000g (Arabidopsis). Pellet is washed once with washing
607 buffer (0.3M sucrose, 15mM Tris HCl pH8 , 60mM KCl, 15mM NaCl, Spermine 0.5mM,
608 Spermidine 0.5mM, 15mM MES, 5mM DTT*,1mM PMSF*, 1% Plant Protease Inhibitors*
609 1ml(Sigma P9599), BSA 0.4%*, RNase inhibitor 0.2u/ul*, (* add last minute)). Finally, nuclei are
610 resuspended into final buffer (0.3M sucrose, 15mM Tris HCl pH8 , 60mM KCl, 15mM NaCl,
611 Spermine 0.5mM, Spermidine 0.5mM, 15mM MES, 5mM DTT*, 1% Plant Protease Inhibitors*
612 1ml(Sigma P9599), BSA 0.4%*, RNase inhibitor 0.2u/ul*,(* add last minute)) and filtered using a
613 10 μ m filter. A nuclei aliquot is stained with DAPI for quality check and nuclei counting under
614 microscope and used immediately for single-nuclei RNAseq.

615

616 **Single cell RNA-seq**

617 16,000 cells or nuclei were loaded in a Single Cell B Chip (10x Genomics) per replicate. Single-
618 cell libraries were then prepared using the Chromium Single Cell 3' library kit, following
619 manufacturer instructions. Libraries were sequenced with an Illumina NextSeq 550 platform using
620 a 1x150 high-output (2 libraries per chip) or Novaseq 6000 chip SP V2.5, 4 libraries per chip. Raw
621 scRNAseq data was analyzed by Cell Ranger 5.0.1 (10x Genomics) to generate gene-cell
622 matrices. Gene reads were aligned to Arabidopsis TAIR10.38, Maize B73 v4, Sorghum bicolor v3
623 and Setaria viridis v2 reference genome.

624

625 **UMAP and ICI analysis**

626 Replicates (see supplementary Table 1) were integrated and cells mapped using the Seurat
627 package v3.0⁴⁷ as follows: first, genes with counts in fewer than three cells were excluded from
628 the analysis and their counts were removed. Second, low quality cells were removed using
629 threshold variable depending on the libraries quality (see supplementary Table 1). Clustering of
630 cells or nuclei separately were done by log-normalized raw counts and the 2000 most variable
631 genes were identified for each replicate using the "vst" method in Seurat. Next, we used the
632 *FindIntegrationAnchors* function to identify anchors between the three datasets, using 20
633 dimensions. A new profile with an integrated expression matrix containing cells from all replicates
634 was produced with the *IntegrateData* function. For dimensionality reduction, the integrated
635 expression matrix was scaled (linear transformed) using the *ScaleData* function, and Principal
636 Component analysis (PCA) performed. The top 30 principal components were selected. Cells or
637 nuclei were clustered using a K-nearest neighbor (KNN) graph, which is based on the Euclidean
638 distance in PCA space. The *FindNeighbors* and *FindClusters* function with a resolution of 0.5.
639 was applied. Next, non-linear dimensional reduction was performed using the UMAP algorithm
640 with the top 30 PCs. Co clustering of cells and nuclei was performed using the SCT approach.
641 First raw reads were normalized using the *SCTransform* function, then *SelectIntegrationFeatures*
642 was used to identify anchors between the three datasets, using 3000 features. For multiple
643 species clustering, all orthologous genes names were replace by their corresponding maize ID in
644 sorghum and setaria raw features.tsv.gz files. Anchors are combined using PrepSCTIntegration
645 and selected using *FindIntegrationAnchors*. For clustering of maize, sorghum and setaria
646 together, maize was selected as a reference dataset in the *FindIntegrationAnchors* function, for
647 single species integration all datasets are considered equally for the integration. Finally, a
648 Principal Component analysis (PCA) is performed using the first 100 principal components and a
649 non-linear dimensional reduction was performed using the UMAP algorithm with the top 100 PCs.

650

651 **GO enrichments.**

652 All GO enrichment were performed using shinyGO V0.61 (<http://bioinformatics.sdstate.edu/go/>)
653 with an FDA of 0.05.

654
655
656
657
658
659
660
661

Gene expression analysis across species.

Whole-root transcriptomes were obtained from Ortiz-Ramírez *et al.*, 2021²⁵ for maize and Hernández Coronado *et al.*, 2021⁴⁸ for arabidopsis.

Gene expression was normalized for each species using the *NormalizedData* function from Seurat. Then the average expression per cluster was calculated using *AverageExpression* from Seurat. Ka and Ks values were provided upon request by J.C. Schnable from the lab publication Zhang *et al.*, 2017

$$\tau = \frac{\sum_{i=1}^N (1 - x_i)}{N - 1},$$

662 Tau t was calculated as describe in (Yanai *et al.*, 2.005) where N is the total
663 number of cell type and xi is the expression profile component normalized by the maximal
664 component value.

665
666
667

Cell type prediction across species and technologies.

668 To determine how well the cell clusters characterized the shared identities of cells in their own
669 clusters and the overlaps with the identities of all other cells, we utilized the MetaNeighbor
670 package in Python (<https://github.com/gillislab/pyMN>) (Fischer *et al.* 2021; Crow *et al.* 2018).
671 MetaNeighbor measures the replicability of cell-types by learning a model in one dataset (or
672 subset) and testing for its ability to reconstruct cell-type clusters in the other dataset. First, we
673 labeled all cells and nuclei by the technology used to sequence the transcriptome, by the cluster
674 identity, and by the plant species to which they belonged. Then, we used the
675 `PyMN.variable_genes` function from MetaNeighbor to subset the gene list to variable genes. This
676 generates a list of genes that are variable across the technology and species. Next, we employed
677 the `PyMN.MetaNeighborUS` function to measure how well the transcriptional profiles of cells from
678 clusters in one division of the dataset (e.g., technology) predict the identities of cell clusters in the
679 other fraction of the data. This generates pairwise AUROCs for each combination of clusters. To
680 generate the heatmaps, the `PyMN.plotMetaNeighborUS` was used with a Brown Blue-green color
681 map. This plots the pairwise AUROCs generated previously.

682
683

Co-expression Conservation between maize subgenomes and sorghum.

684 To generate co-expression conservation scores between the two maize sub genomes and
685 sorghum, we use our existing aggregated co-expression networks (Lee *et al.* 2020). In brief,
686 these networks are built by taking all publicly available data and calculating average correlations
687 between genes pairs within experiments, standardizing within experiments, and then averaging
688 to construct robust meta-analytic networks. We filtered these networks to a previously generated
689 list of gene triplet pairs for maize sub genomes and sorghum. Next, for each gene, we compare
690 the top co-expression partners across species to determine the degree of functional conservation,
691 as described in more detail in previous work (Crow *et al.* 2022). We calculate this by taking the
692 ranks of a gene's co-expression strength to all other genes in one species and using it to predict
693 that gene's top 10 co-expressed partners in the second species. This is then done again in the
694 reverse direction, and the two scores are averaged (calculated as an AUROC).

695 with the lowest co-expression scores ($0.34 < FC.Score$) and highest cell specificity ($\square > 0.8$) in
696 the root cap (Supplementary Table 7; Extended Data Fig. 10b). Similar to trends in other cell
697 types, these highly cell-type specific genes

Dominance vs. partition score:

699 To calculate the Dominance vs. partition score, for each ortholog triplet (S, M1, M2) we calculated
700 the number of cells in which M1 or M2 was dominant or co-expressed together in the same cells
701 where the sorghum ortholog was expressed.

702 $Score = (number\ of\ cells\ in\ which\ M1\ is\ dominant * number\ of\ cells\ in\ which\ M2\ is\ dominant) -$
703 $(number\ number\ of\ cell\ of\ the\ dominant\ ortholog - number\ of\ cell\ of\ the\ non\ dominant\ ortholog)$

704 If the score is negative, the score is normalized by

$$705 \text{ NormScore} = \frac{Score}{\# \text{ of cell in which M1 and M2 are expressed}}$$

706 If the score is positive, the score is normalized by dividing it by:

$$707 \text{ NormScore} = \frac{Score}{(\# \text{ of cell in which M1 and M2 are expressed} * 0.5)^2}$$

708 **Statistical analysis:**

709 Each species marker genes were identified using *FindAllmarkers* functions from Seurat, log.FC=
710 0.25, pt.1 > 0.750 pt.2 < 0.250. Differential gene expression was done using the *Findmarkers*
711 function from Seurat with default parameter function. For Fig2 e, Extended Data Fig. 4 c, 10 a,
712 statistical analysis was performed on R using a pairwise Wilcoxon test with p.adjust method "BH"
713 as data is not normally distributed.

714
715 Correlation analysis on Extended Data Fig 1 c was performed using Pearson correlation function
716 on R between whole-root data coming from and single cell or single nuclei. Briefly averaged gene
717 expression was calculated for each gene while combining every cell types using the
718 *AverageExpression* function from Seurat.

719
720 For Fig 4 a, to generate p-values for evaluating the significance of the differences between each
721 pair of AUROCs generated by MetaNeighbor, we utilized the Hanley McNeil test, which produces
722 a Z-score for the difference (Hanley and McNeil 1983). As each MetaNeighbor AUROC is the
723 averaged AUROC from two reciprocal tests between a pair of cell clusters, we chose the smaller
724 of the two clusters as the number of true positives (NTP) to generate the most conservative p-
725 value. The number of true negatives was the total number of cells, less the number of true
726 positives. Following the calculation of Z-scores for each pairwise combination of AUROCs, we
727 utilized the *scipy.stats.norm.sf* function in Python to convert the Z-scores into p-values for a two
728 tailed test.

729
730 Hanley, J. A., and B. J. McNeil. 1983. "A Method of Comparing the Areas under Receiver
731 Operating Characteristic Curves Derived from the Same Cases." *Radiology* 148 (3): 839–43.
732 <https://doi.org/10.1148/radiology.148.3.6878708>.

733 734 **"Half mount" *in situ* hybridization:**

735 Probes (Hairpin Chain Reaction (HCR) RNA-FISH) and reagents (including the Probe
736 Hybridization Buffer, Probe Wash Buffer and Amplification buffer) are ordered from Molecular
737 Instruments (<https://www.molecularinstruments.com/shop>)(**Supplementary Table 9**).

738
739 For fixation, germination paper containing 7 day old maize or sorghum roots are unrolled and
740 small volume of fixative FAA (4% formaldehyde, 5% glacial acetic acid, 50% ethanol in RNase

741 free water) is pipetted onto each root. Then longitudinal sectioning of root tips is performed using
742 a 15° microscalpel. Roots are cut up to ~3cm from the tip, then immediately fixed by transferring
743 to FAA in 5ml screw caps and put under vacuum several times until they no longer float. Roots
744 are then agitated at RT for at least 1 hour in a tube revolver. (All washes in the protocol are
745 performed in a tube revolver or stated otherwise.)

746 Samples are dehydrated in a series of washes at RT: 70% ethanol for 15 min, 90% ethanol for 15
747 min, 100% ethanol 2x for 15 min each, 100% methanol 2x for 15 min each. Samples can then be
748 stored at -20°C for several weeks. Samples are washed 2x for 15 min in 100% ethanol at RT
749 before being permeabilized for 30 min in 50% Histo-Clear II / 50% EtOH at RT. Then they are
750 incubated 2x for 30 minutes in a solution of 100% Histo-Clear II at RT. Each time, vacuum is
751 applied for the first 10 minutes.

752
753 Samples are rehydrated through a series of washes: 50% Histo-Clear II / 50% EtOH for 15 min,
754 100% EtOH for 15 min, 50% EtOH / 50% DPBS-T (0.1% Tween20, 1x DPBS) for 15 min (roots
755 will float up then settle after a few minutes), 100% DPBS-T 2x for 15 min (roots will float up again).
756 Samples are incubated with Proteinase K (0.1 M Tris-HCl (pH 8), 0.05 M EDTA (pH 8), Proteinase
757 K 80 µg ml⁻¹ final) at RT under vacuum for 5 min then digested with Proteinase K for 25 min in a
758 37°C water bath with manual agitation every 5-10 minutes (roots should turn a little yellow after
759 this step). Samples are washed 2x for 15 min in DPBS-T at RT then incubated with Fixative II (4%
760 formaldehyde in DPBS-T) under gentle vacuum for 10 min then in a tube revolver for 30 mins at
761 RT. They are then washed 2x for 15 min each in DPBS-T at RT. Roots are aliquoted into 2 mL
762 Eppendorf tubes and incubated in 500 µL of HCR Probe Hybridization Buffer, vacuum is applied
763 for 10 mins then roots are incubated for 1 hour at 37°C in a thermomixer with agitation (1000
764 rpm).

765
766 Samples can then be stored in Probe Hybridization Buffer at -20°C up to several weeks.
767 Probe buffers are made by adding 0.8 pmol of each probe set (e.g. 2 µL of the 1 µM stock) to 500
768 µL of HCR Probe Hybridization Buffer at 37°C. Pre-hybridization solution is removed and replaced
769 with probe solution. Samples are hybridized by incubating overnight (~20h) at 37°C in a
770 thermomixer with agitation (1000 rpm). The following day, excess probes are removed by washing
771 4x for 15 min each with 1 mL of HCR Probe Wash Buffer at 37°C in a thermomixer with agitation.
772 Samples are washed 2x for 5 min each with 1 mL of 5x SSC-T (25% 20x SSC, 0.1% Tween20)
773 at RT in a thermomixer with agitation. SSC-T is replaced with 500 µL of amplification buffer, gentle
774 vacuum is applied in a fume hood for 10 minutes and then samples are pre-amplified by incubating
775 in a tube rotator at RT for 50 min. While samples pre-amplify, 6 pmol of hairpin h1 and 6 pmol of
776 hairpin h2 (i.e. 5 µL of the 3 µM stocks) are prepared, each in its own separate tube. Hairpins are
777 snap-cooled by heating at 95°C for 90 seconds then kept in a dark drawer at RT for 30 min.
778 Amplification solution is prepared by combining snap-cooled h1 and h2 hairpins in 250 µL of HCR
779 Amplification Buffer at RT. Pre-amplification solution is removed and replaced with
780 amplification buffer containing hairpin solution overnight (~20h) in the dark at RT in a thermomixer
781 with agitation (1000 rpm). Excess hairpins are removed by washing with 1 mL of 5x SSC-T at RT
782 in a thermomixer with agitation, 2x for 5 min each, then 2x for 30 min each, 1x for 5 min. Samples
783 are transferred onto a glass slide (in 5x SSC-T) and cut using a 30° microscalpel and arranged
784 so that the cut face of the roots is facing upwards. They are then covered with coverslip and
785 imaged on confocal microscope.

786 **Spatial transcriptomics:**

787 Tissue fixation and embedding was performed as described in⁴⁹.

788

789 **Sample slide preparation:** Formaldehyde-fixed paraffin-embedded tissue sections (10 μm) were
790 placed within capture areas on Resolve Bioscience slides and incubated on a hot plate for 10 min
791 at 60 $^{\circ}\text{C}$ to attach the samples to the slides. Slides were treated to allow deparaffinization,
792 permeabilization, acetylation, and refixation. After complete dehydration of the samples, a few
793 drops of SlowFade-Gold Antifade reagent (Invitrogen) were added to the sections and covered
794 with a thin glass coverslip to prevent damage during shipment to Resolve BioSciences (Germany).
795

796 **Sample pre-treatment and priming: In preparation for hybridization,** the coverslip is removed
797 and the mounting reagent is washed twice in 1x PBS for 30 min 4 $^{\circ}\text{C}$, followed by one min washes
798 in 50% Ethanol and 70% Ethanol at room temperature. Samples were primed, after the aspiration
799 of ethanol, by the addition of buffer BST1 for optimal hybridization of probes during the Molecular
800 CartographyTM procedure, which uses a combination of probes and single-molecule fluorescence
801 in-situ hybridization to identify 100 separate transcripts. Tissues were hybridized overnight at a
802 constant temperature with all probes specific to the target genes. Samples were washed the next
803 day to remove excess probes and fluorescently labeled in a two-step procedure. Regions of
804 interest were imaged as described below and fluorescent signals were removed after imaging via
805 a decolorization procedure. Color development, imaging, and decolorization were repeated over
806 several cycles to develop a unique combinatorial code for every target gene that was derived from
807 raw images as described below.
808

809 **Probe design:** The probes for 100 genes were designed based on full-length protein-coding
810 transcript sequences (Supplementary Table 9). Probe design is based the manufacturer's
811 proprietary algorithm, with probes available from the Resolve. After screening to generate probe
812 candidates and discard ambiguous ones, the probes were mapped to the background
813 transcriptome using *ThermonucleotideBLAST*, and probes with stable off-target hits were
814 discarded.
815

816 **Imaging:** Samples were imaged on a Zeiss Celldiscoverer 7, using the 50x Plan Achromat
817 water immersion objective with an NA of 1.2 and the 0.5x magnification changer, resulting in a 25x
818 final magnification. Standard CD7 LED excitation light source, filters, and dichroic mirrors were
819 used together with customized emission filters optimized for detecting specific signals. Excitation
820 time per image was fixed at 1000 ms for each channel, 20 ms for DAPI, and 1 ms for Calcofluor
821 White. A z-stack was taken at each region with a distance per z-slice according to the Nyquist-
822 Shannon sampling theorem. A custom CD7 CMOS camera (Zeiss AxioCam Mono 712, 3.45 μm
823 pixel size) was used. The imaging for the cell-wall specific stain, Calcofluor White, was done at
824 the end of all primary imaging. Before the preprocessing of the images, all images were corrected
825 for background fluorescence. Based on the raw data image, the 20% darkest local pixel values
826 and positions were determined and copied to a new empty image (background image) having the
827 same size as the image to be corrected. The remaining 80% of pixels of the background image
828 were generated based upon the surrounding existing pixel values using a distance-weighted
829 average value. Finally, the background-corrected image (bc-image) was created by subtracting
830 the background image values from the raw data image values.
831

832 **Extraction of features:** In the first step, a target value for the allowed number of maxima was
833 calculated based on the area of the slice in μm^2 multiplied by an empirically optimized factor
834 (0.5x). The resulting target value was used to adapt the threshold for the algorithm iteratively
835 searching local 2D-maxima. The threshold leading to the closest number of maxima equal to or
836 smaller than the target value was used for further steps and the respective maxima were stored
837 in a reiterative process for every image slice independently. Maxima that did not have a
838 neighboring maximum in an adjacent slice (termed as z-group) within a radius of one pixel were
839 excluded. For the resulting list of maxima, the absolute brightness (Babs), the local background

840 (Bback), and the average brightness of the pixels surrounding the local maximum (Bperi) were
841 measured and stored. The resulting maxima list was further filtered in an iterative loop by adjusting
842 the allowed thresholds for (Babs-Bback) and (Bperi-Bback) to reach a feature target value based
843 on the total volume of the 3D image. Only maxima still in a z-group with a size of at least 2 passed
844 this stringent filter step. Each z-group was counted as one hit and the members of the z-groups
845 with the highest absolute brightness were used as features to resemble 3D point clouds.
846

847 **Determination of transformation matrices, pixel evaluation, and decoding:** To align the raw
848 data images from different imaging rounds, these images had to be corrected for the 6 degrees
849 of freedom in 3D-space. The extracted feature point clouds were used to find the transformation
850 matrices to align the raw data images. Based on the transformation matrices, the corresponding
851 images were processed by a rigid transformation using trilinear interpolation. The aligned images
852 were used to create a profile for each pixel, which were then filtered for a variance from zero
853 normalized by the total brightness of all pixels in the profile. Matched pixel profiles with the highest
854 score were assigned as an ID to the pixel to further group the neighboring pixel with the same ID.
855 The local 3D-maxima of the groups were determined as potential final transcript locations, which
856 were additionally evaluated by the number of maxima in the raw data images where a maximum
857 was expected. The finalized maxima were decoded by the fit to the corresponding code to be
858 written to the results file and considered to resemble transcripts of the corresponding gene. The
859 ratio of signals matching to codes used in the experiment and signals matching to codes not used
860 in the experiment were used as estimation for specificity (false positives). Final image analysis
861 was performed in ImageJ using the Polylux tool plugin from Resolve BioSciences to examine
862 specific Molecular Cartography signals.
863

864 All R scripts related to models and statistical analyses are available upon request.
865

866 All raw RNA-seq data will be deposited in GEO upon publication.
867
868

869

Supplementary Files

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

- [SupplementaryTable1Summaryruns.xlsx](#)
- [SupplementaryTable2ArabidopsisclusteringrepartitionofNucleiorCellaloneandmergeddataset.csv](#)
- [SupplementaryTable3DifferentiallyexpressedgenesbetweenAthNucleiandCells.xlsx](#)
- [SupplementaryTable4CellTypemarkersMaizeSorghumSetaria.xlsx](#)
- [SupplementaryTable5GOEnrichmentSSM1SM2M1M2SM1M2.csv](#)
- [SupplementaryTable6OrthologsDominance.xlsx](#)
- [SupplementaryTable7CococoNetAnalysisallfunctionalscorescalculated.csv](#)
- [SupplementaryTable8ProberefInsituhybridizationandmolecularcartography.xlsx](#)