

1 **Gene expression analysis provides insight into the physiology of the important staple food**
2 **crop cassava**

3
4 Mark C. Wilson¹, Andrew M. Mutka¹, Aaron W. Hummel^{2,3}, Jeffrey Berry¹, Raj Deepika
5 Chauhan¹, Anupama Vijayaraghavan¹, Nigel J. Taylor¹, Daniel F. Voytas², Daniel H. Chitwood¹
6 and Rebecca S. Bart^{1,4}

7
8 ¹Donald Danforth Plant Science Center, 975 North Warson Road. St. Louis, MO, 63132

9 ²Department of Genetics, Cell Biology & Development and Center for Genome Engineering,
10 University of Minnesota, Minneapolis, Minnesota 55455, USA

11 ³Current address: KWS SAAT SE, Gateway Research Center, St. Louis, MO, USA

12 ⁴Correspondence: 314-587-1696, rbart@danforthcenter.org

13

14 Total Word Count: 3753

15 Summary: 134

16 Introduction: 211

17 Materials and Methods: 1641

18 Results: 1606

19 Discussion: 170

20 Acknowledgements: 103

21 Number of Figures: 5 (all color)

22 Number of Supporting Files: 11

23 Data is deposited in GEO repository: Accession number GSE82279

24

25 **Summary**

- 26 • Cassava (*Manihot esculenta*) feeds approximately 800 million people worldwide. Although
27 this crop displays high productivity under drought and poor soil conditions, it is susceptible
28 to disease, postharvest deterioration and the roots contain low nutritional content.

- 29 • Here, we provide molecular identities for eleven cassava tissue types through RNA-
30 sequencing and develop an open access, web-based interface for further interrogation of the
31 data.
- 32 • Through this dataset, we report novel insight into the physiology of cassava and identify
33 promoters able to drive specified tissue expression profiles. Specifically, we focus on
34 identification of the transcriptional signatures that define the massive, underground storage
35 roots used as a food source and the favored target tissue for transgene integration and
36 genome editing, friable embryogenic callus (FEC).
- 37 • The information gained from this study is of value for both conventional and
38 biotechnological improvement programs.
- 39
- 40 • Key words: biotechnology, cassava, food security, friable embryogenic callus, gene
41 expression, organized embryogenic structures, RNAsequencing

42

43 Introduction

44 Cassava (*Manihot esculenta*) is the food security crop that feeds approximately 800
45 million people worldwide (Liu *et al.*, 2011; Howeler *et al.*, 2013). Although this crop displays
46 high productivity under drought and poor soil conditions, it is susceptible to disease, postharvest
47 deterioration and the roots contain low nutritional content (Gegios *et al.*, 2010; Stephenson *et al.*,
48 2010; Patil *et al.*, 2015). Cassava improvement programs are focused on addressing these
49 constraints but are hindered by the crop's high heterozygosity, difficulty in synchronizing
50 flowering, low seed production and a poor understanding of the physiology of this plant
51 (Ceballos *et al.*, 2004). Among the major food crops, cassava is unique in its ability to develop
52 massive, underground storage roots. Despite the importance of these structures, their basic
53 physiology remains largely unknown, especially the molecular genetic basis of storage root
54 development. Similarly, in cassava, the favored target tissue for transgene integration and
55 genome editing is a friable embryogenic callus (FEC) (Taylor *et al.*, 2001; Bull *et al.*, 2009;
56 Taylor *et al.*, 2012; Zainuddin *et al.*, 2012; Nyaboga *et al.*, 2013). Little is known concerning

57 gene expression in this tissue, or its relatedness to the somatic organized embryogenic structures
58 (OES) from which it originates (Gresshoff & Doy, 1974; Taylor *et al.*, 2012; Chauhan *et al.*,
59 2015). Here, we provide molecular identities for eleven cassava tissue types through RNA-
60 sequencing and develop an open access, web-based interface for further interrogation of the data.
61 Through this dataset, we report novel insight into the physiology of cassava and identify
62 promoters able to drive specified tissue expression profiles.

63

64 **Materials and Methods**

65 *Plant material and tissues sampled*

66 Plant tissues were sampled from 3-month-old TME 204 cassava plants, grown in a
67 greenhouse at the Donald Danforth Plant Science Center (St. Louis, MO). Plants were
68 established from in vitro micropropagated plants (Taylor *et al.*, 2012) and grown in a 12 h
69 light:12 h dark photoperiod (250-500 $\mu\text{mol s}^{-1} \text{m}^{-2}$ irradiance Daytime day, temperatures ranged
70 from 28°-32°C with 70% relative humidity, and night time, temperatures ranged from 25°-27°C
71 with 70% relative humidity. The following tissues were sampled from these plants at 2pm: leaf
72 blade, leaf midvein, petiole, stem, lateral buds, shoot apical meristem (SAM), storage roots,
73 fibrous roots, and root apical meristem (RAM). For non-meristem tissues, approximating 100 mg
74 of tissue was collected in 3 separate biological replicates. For the SAM and RAM, 6 meristems
75 were dissected and pooled for each of 3 biological replicates. All samples were frozen in liquid
76 nitrogen after collection. Samples of TME 204 OES and FEC were generated as described
77 previously (Chauhan *et al.*, 2015). The OES induced on DKW/Juglans basal salts
78 (*PhytoTechnology Laboratories*, Kansas, USA) containing Murashige and Skoog (MS) vitamins
79 and supplemented with 2% w/v sucrose, 50 μM picloram was sampled after four weeks of
80 culture. The OES was separated from the non-embryogenic tissues and collected in 2ml sampling
81 tubes. FEC tissues were sampled after three weeks of culture on Gresshoff and Doy basal
82 medium supplemented with 2% w/v sucrose, 50 μM picloram, 500 μM tyrosine and 50 mg/l
83 moxalactam. Approximately 200 to 250 mg of tissue was collected and the tubes containing the
84 tissues were immediately placed on dry ice.

85

86 *Preparation of RNA-seq libraries and Illumina sequencing*

87 For non-meristem tissues, total RNA was isolated with the Spectrum Plant Total RNA
88 Kit (Sigma). For SAM and RAM tissues, total RNA was isolated with the Arcturus PicoPure
89 RNA Isolation Kit. RNA quality was assessed on an Agilent Bioanalyzer. For library preparation
90 with tissues other than SAM and RAM, 5 µg of RNA was used as input. For SAM and RAM, six
91 samples were pooled to obtain a total of 500-600 ng each. The NEBNext Poly(A) mRNA
92 Magnetic Isolation Module (New England BioLabs) was used to isolate mRNA, which was then
93 used for library prep using NEBNext mRNA Library Prep Master Mix Set for Illumina (New
94 England BioLabs) with 13 cycles of PCR amplification. Standard library prep protocol was
95 followed for all samples, except for the SAM and RAM in which 1 µL of fragmentation enzyme
96 was used instead of 2 µL, and 0.5 µL of random primer instead of 1 µL. Library quality was
97 assessed with the Agilent Bioanalyzer. In total, 32 RNA-seq libraries were made from 11
98 different tissue types with 3 biological replicates each, except for storage root which only had 2
99 biological replicates. All libraries were multiplexed into one lane of Illumina HiSeq2500.

100

101 *Read mapping and gene expression analysis*

102 Illumina RNA-seq reads from each replicate were cleaned using Trimmomatic version
103 0.32 (Bolger *et al.*, 2014). Using TopHat2 version 2.1.0 (Trapnell *et al.*, 2009), these cleaned
104 reads were then mapped to the version 6.1 draft assembly of *Manihot esculenta* AM560-2
105 provided on Phytozome10.3 (<http://phytozome.jgi.doe.gov/pz/portal.html>). The read mapping
106 output was linked to candidate gene models for each sample using Cufflinks version 2.2.1
107 (Trapnell *et al.*, 2010). The gene models from all samples of the experiment were merged into
108 one gene model file using Cuffmerge version 2.2.1. Using the output from Cuffmerge and the
109 read mapping files from each replicate, a differential expression analysis between tissue types
110 was performed using Cuffdiff version 2.2.1. Quality checks were performed on the Cuffdiff
111 output using cummeRbund version 2.6.1 in R (R Core Team, 2015). The output of Cuffdiff was
112 processed in Python with the pandas, numpy, and seaborn packages to visualize the expression
113 data (McKinney, 2010).

114

115 *Multivariate statistics*

116 Analysis of transcript expression profiles began with those transcripts with 1) FPKM
117 values above a threshold of 1 FPKM and 2) those transcripts significantly differentially
118 expressed in at least one pairwise comparison of all tissues types against each other. For
119 Principal Component Analysis (PCA) of tissue replicates (**Fig. 3A-B**), the `prcomp()` function in
120 R was used with scaled FPKM values across transcripts as input. For the PCA of transcript
121 profiles (**Fig. 3C-D**), the `prcomp()` function was again used with scaled mean FPKM values
122 across tissues as input. Self-Organizing Maps (SOMs) were performed using the Kohonen
123 package in R. Scaled mean transcript expression values across tissues were assigned to four
124 nodes in a 2x2 hexagonal topology over 100 training iterations. To focus on those transcripts
125 with expression profiles closest to over-represented patterns of variance in the dataset, only those
126 transcripts with distances from their respective nodes less than the median for the overall dataset
127 were subsequently used and projected back onto the transcript PCA space. Data visualization for
128 the above was carried out with the `ggplot2` package in R, using `geom_point()` and `geom_line()`
129 functions, among others. The color space for the above was determined using palettes from
130 colorbrewer2.org.

131

132 *Gene Ontology (GO) enrichment analysis*

133 GO enrichment analysis was completed using the Python `goatools` package
134 (<https://github.com/tanghaibao/goatools>). `goatools` was run with the `--fdr` flag to calculate the
135 False Discovery Rate (FDR) error corrected p-value, and the `--no_propagate_counts` flag to
136 prevent nodes at the root of the GO tree from being included in the analysis. GO terms for each
137 gene were used from the annotation file provided on Phytozome. GO enrichment output was then
138 filtered to include only enriched GO terms with a FDR error corrected p-value < 0.001. For SOM
139 node GO enrichment, each SOM node identified above was processed separately. The genes
140 identified as part of the SOM node were used as the study group, and all genes expressed greater
141 than 1 FPKM in at least one tissue with significant differential expression in at least one pairwise
142 comparison were used as the population or background group. For pairwise tissue comparison
143 GO enrichment, genes identified as significantly upregulated with a $|\log_2(\text{fold_change})| > 2$ in
144 one tissue were treated as one study group to look at each tissue separately. This resulted in two

145 GO enrichment analyses for each pairwise comparison. Genes with at least 1 FPKM in either
146 tissue were used as the background dataset.

147

148 *Identification of genes with strong, constitutive, and tissue-specific expression patterns*

149 Custom Python code was used in a Jupyter notebook using the Pandas, NumPy, Seaborn,
150 and SciPy packages to organize, process, and display the data (**SuppFile3**). Genes with strong
151 expression across all tissue types were identified using expression values from the gene_exp.diff
152 file produced by Cuffdiff. The genes were first checked for functional annotations, then
153 shortened to a list of genes with a minimum expression of 300 FPKM in each tissue sampled.
154 Specifically and constitutively expressed genes were identified using expression values from
155 each replicate in the genes.read_group_tracking file produced by Cuffdiff. Genes used were
156 annotated in the AM560-2 v6.1 assembly on Phytozome10.3. For specifically expressed genes,
157 this list was then subset by selecting genes with expression greater than 10 FPKM in the tissue(s)
158 specifically expressing the gene, and no more than 1 FPKM in all other tissues. For a more
159 relaxed analysis, genes were required to be expressed greater than 8 FPKM in the tissue(s)
160 specifically expressing the gene, and no more than 4 FPKM in all other tissues. Constitutively
161 expressed genes were identified using the replicate expression data. This list was filtered to
162 include only genes with greater than 40 FPKM in all replicates, and then the coefficient of
163 variation was calculated across all replicates for each gene.

164

165 *Data availability*

166 A graphical user interface was created using R Shiny (v 0.13.2) to explore the tissue-
167 specific data and discover trends therein. This application uses data from RNA-seq differential
168 expression analysis completed with the Tuxedo Suite pipeline (v 2.2.1), functional gene
169 annotations from the Joint Genome Institute's Phytozome, and analysis from principle
170 components (prcomp in R "stats" package v 3.2.3) and self-organizing maps (som in R
171 "kohonen" package v 2.0.19). The application has two main features: 1) gene discovery based on
172 gene expression patterns across tissues, and 2) creation of a tissue-specific heatmap of known or
173 newly discovered genes for visualizing expression patterns. Detailed instructions are included in
174 the application. The application can be found at: shiny.danforthcenter.org/cassava_atlas/.

175 Additional R packages used in this application include: png (v 0.1-7), grid (v 3.2.3), ggplot2 (v
176 2.1.0), shinyBS (v 0.61), shinydashboard (v 0.5.1), DT (v 0.1), stringr (v 1.0.0), mailR (v 0.4.1),
177 and shinyjs (v 0.5.2).

178

179 *In planta expression assays*

180 Promoter fragments, listed in SuppFile4, were cloned from cassava variety TME419 into
181 a pCAMBIA vector upstream of GUS. Constructs were transformed into *Agrobacterium*
182 *tumefaciens* strain LBA4404. Strains carrying the reporter constructs were re-suspended in IM
183 media (10 mM MES, pH5.6; 10 mM MgCl₂; 150 μM Acetosyringone), incubated at room
184 temperature for 3 hours and then infiltrated into *Nicotiana benthamiana* leaves at an OD600 =
185 0.1. 48 hours post inoculation, leaves were detached and placed in a petri dish. GUS staining
186 solution (0.1 M NaPO₄ pH7; 10 mM EDTA; 0.1% Triton X-100; 1 mM K₃Fe(CN)₆; 2 mM X-
187 Gluc) was pipetted on to the detached leaf and a glass tube rolled across the leaf surface to
188 lightly crush the tissue. Leaves were incubated overnight at 37C⁰. Prior to imaging, leaves were
189 cleared of chlorophyll through several washes in 95% EtOH. To quantify GUS staining, multiple
190 image processing steps were implemented using ImageJ to obtain the pixel statistics that are
191 reported. The original RGB image was converted to HSL colorspace using the “Color
192 Transform” plugin and the lightness channel was extracted. The image look-up table was
193 changed to “thermal” and a manually defined circular ROI was created whose size and shape
194 remain constant when gathering the mean and standard deviation of the pixel intensities for each
195 of the strains. Using the same ROI, the image was cropped for each of the strains to display the
196 exact regions sampled.

197

198 *Cassava transformation*

199 Reporter constructs were introduced to cassava FEC cells by LBA4404 following our
200 published methods (Chauhan *et al.*, 2015).

201

202 **Results**

203 To shed light on the development and physiology of cassava plants from a gene
204 expression perspective, eleven tissue/organ types from cassava cultivar TME 204 were sampled

205 for transcriptome profiling (**Fig. 1**). Tissue type relatedness was assessed based on Jensen-
206 Shannon (JS) distances (**Fig. 2**) and principal component analysis (PCA) (**Fig. 3**). Biological
207 replicates clustered closely together confirming the high quality of the dataset (**Fig. S1A, Fig.**
208 **3A-B**). Both analyses divided the 11 tissues into three major groups: aerial (leaf, midvein,
209 petiole, stem, lateral bud, and shoot apical meristem (SAM)), subterranean (storage root, fibrous
210 root and root apical meristem (RAM)), and embryogenic (OES and FEC). Leaf and midvein,
211 petiole and stem, lateral bud and SAM, and OES and FEC samples cluster together within the
212 dendrogram (**Fig. 2B**), and occupy similar positions across the first four principal components
213 (PCs), which collectively explain 67.3% of transcript expression level variance (**Fig. 3A-B**).
214 These groupings are expected, representing leaf blade, vascular, shoot meristem, and callus-
215 associated tissues. The root tissues show more complicated relationships. Figure 2 indicates
216 storage roots as distant from fibrous roots and RAM (**Fig. 2B**). Similarly, whereas the RAM,
217 storage root and fibrous root samples cluster closely together when projected onto PCs 1 and 2
218 (**Fig. 3A**), these tissues occupy more disparate positions when evaluated by PCs 3 and 4 (**Fig.**
219 **3B**). This indicates that while root samples share common gene expression patterns, tissue
220 specific signatures differentiate storage roots from the other subterranean tissues.

221 Two tissue comparisons within the dataset: OES vs FEC and storage vs fibrous roots, are
222 particularly intriguing, given how little is known about the features distinguishing their
223 physiology. The results of a PCA on the expression profiles of individual transcripts was
224 considered. A self-organizing map (SOM) was used to identify four main clusters of transcripts
225 with similar expression profiles across tissues, which was then projected back onto the PCA
226 transcript space (**Fig. 3C-D**). To determine the identities of transcripts with shared expression
227 profiles, we performed Gene Ontology (GO) enrichment analysis for each SOM node (**Fig. 3F**).
228 In addition, we directly examined the genes most highly differentially expressed between each
229 comparison (**Fig. 2, Fig. 4**).

230 First, we used the above approach to examine gene expression patterns for well
231 characterized tissues and comparisons. Node 4 transcripts (teal) are highly expressed in the
232 photosynthetic tissues of the leaf and mid vein (**Fig. 3E**). Similarly, comparison of leaf and
233 fibrous roots revealed ~4900 genes differentially expressed greater than 4-fold
234 ($|\log_2(\text{fold_change})| > 2$) between tissues (**Fig. 2C, SuppFile1a**). A similar number were up-

235 regulated in each tissue and consistent with the GO term analysis presented in Figure 3F, the
236 most highly up-regulated genes in leaf tissue pertained to photosynthesis activity while genes
237 induced in fibrous root were related to lignin, ion binding, and transcription. We highlight that
238 these analyses are complementary. The former takes an unbiased approach to identify variability
239 within the dataset. The latter, directly looks at genes with maximum expression differences.

240 OES and FEC tissue are closely related with the latter generated from the former by a
241 simple switch in the basal medium (Taylor *et al.*, 2012). FEC tissues are highly disorganized and
242 ultra-juvenile in nature, consisting of proliferating, sub-millimeter sized pre-embryo units from
243 which somatic embryos will regenerate on removal of auxin. Efficacy of FEC production from
244 the OES is genotype dependent and can be challenging in some farmer-preferred varieties,
245 though this recalcitrance is poorly understood (Liu *et al.*, 2011). Node 3 transcripts (burnt
246 orange) are highly expressed in both callus tissues, but especially the FEC (**Fig. 3E**). Node 3
247 transcripts are associated with GO terms related to epigenomic reprogramming (DNA
248 methylation and histone modification). Over two thousand genes were identified as differentially
249 expressed between OES and FEC tissues (**Fig. 2D, SuppFile1b**). Genes up-regulated in OES
250 tissue are associated with GO tags for heme, iron and tetrapyrrole binding and oxidoreductase
251 activity. In contrast, genes upregulated in FEC tissue are associated with sulfur and sulfate
252 transport (**SuppFile2**). Overall, our analyses emphasize the striking similarity between the two
253 tissue types.

254 What distinguishes storage roots from other subterranean structures is ambiguous. A
255 recent anatomical examination of these structures revealed that roots develop from the cut base
256 of the stem cutting (basal) and from buried nodes (nodal), but that only the nodal roots will
257 develop to form storage roots (Chaweewan & Taylor, 2015). Once initiated the storage roots
258 develop by massive cell proliferation from the cambium to generate the central core that consists
259 largely of xylem parenchyma in which starch is synthesized and stored. Node 1 transcripts
260 (magenta) are highly expressed in the RAM and somewhat in the fibrous root while Node 2
261 transcripts (lavender) are highly expressed in the storage root suggesting that storage roots
262 exhibit distinct gene expression patterns relative to RAM and fibrous roots. Node 1 transcripts,
263 highly expressed in the RAM and the fibrous root, are enriched for GOs related to translation,
264 proteolysis, and intracellular transport that might be expected for a tissue undergoing growth.

265 Node 2 transcripts highly expressed in the storage root, are associated with zinc ion and
266 phosphatidylinositol binding GO terms. In contrast to differentially expressed gene comparisons
267 for leaf versus fibrous roots, and OES versus FEC, comparison of fibrous and storage roots
268 revealed a significant shift towards gene induction in the former (**Fig. 2E, SuppFile1c**). Taken
269 together, these data and analyses demonstrate that OES and FEC are highly similar tissue types
270 and suggest that their difference may come mostly from the media on which they are cultured. In
271 contrast, fibrous and storage roots appear as inherently distinct on a transcriptional level.

272 Promoters capable of driving gene expression in one or more defined tissue/organ types
273 is essential for the successful application of biotechnology to improve crop plants. Currently, a
274 limited set of promoters are available to achieve desired expression patterns for cassava *in*
275 *planta*. For example, the root-specific patatin promoter from *Solanum tuberosum* has been used
276 to overexpress transgenes that enhance iron and zinc levels in cassava storage roots (Gaitan-Solis
277 *et al.*, 2015; Narayanan *et al.*, 2015). De Souza *et. al.* has characterized the Pt2L4 gene
278 (Manes.09G108300) and confirmed preferential expression in cassava storage roots but also in
279 stems (de Souza *et al.*, 2006; de Souza *et al.*, 2009). This previously published expression pattern
280 is consistent with the current dataset. To identify cassava promoters capable of tissue specific
281 expression, we queried the dataset for genes expressed in a single tissue type, henceforth referred
282 to as uniquely expressed genes. To identify uniquely expressed genes, FPKM values of 1 and 10
283 were chosen to represent ‘below the limit of detection’ and ‘expressed’, respectively. These
284 cutoffs were determined by investigating read mapping coverage for individual genes within our
285 datasets. An FPKM value of less than one generally correlated with less than 1x coverage across
286 a coding sequence. Genes expressed at greater than 10 FPKM had read mapping across the entire
287 coding sequence. In addition, we choose an expression value of ≥ 300 FPKM as the cutoff for
288 highly expressed genes which encompasses approximately the top 2% of expression values
289 across our dataset. Below the limit of detection, expressed, and highly expressed cutoffs within
290 the context of the entire dataset are displayed in Figure 4b. Uniquely expressed genes were
291 identified as those expressed at greater than 10 FPKM in one tissue, and less than 1 FPKM in all
292 other tissues (**Fig. 4a**). Applying the cutoff criteria, unique gene expression was observed for
293 FEC, fibrous root, RAM and SAM, but not for the other seven tissues. Using less stringent cutoff
294 FPKM values (OFF<4; ON>8), we were able to identify uniquely expressed genes for all

295 additional tissues (**Fig. 4a**). In addition, we considered expression that would be constrained to
296 the major groupings from the dendrogram in Figure 2. Storage root was excluded from the
297 subterranean group because of its distinct gene expression patterns (**Fig. 2B, Fig. 3**).

298 In addition to identification of uniquely expressed genes, the data was queried to identify
299 candidate promoters for driving strong gene expression within all surveyed tissue types
300 (constitutive). We identified genes that showed expression values of ≥ 300 FPKM across our
301 entire dataset. This analysis resulted in a list of 31 genes (**Fig. 5a**). In order to test the *in silico*
302 analysis, promoters from five of the 31 putative constitutively expressed genes were cloned and
303 functionally validated by fusing to the *uidA* (GUS) reporter gene. These constructs were
304 expressed transiently in *Nicotiana benthamiana* leaves and stably transformed into cassava FEC
305 cells. All five tested promoters were confirmed to drive GUS expression in cassava FEC cells
306 while one promoter fusion, Manes.G035300, failed to drive expression in *N. benthamiana* for
307 unknown reasons (**Fig. 5b, Fig. S2**).

308 A small collection of ‘housekeeping genes’ are routinely used for internal controls in
309 quantitative reverse transcription polymerase chain reaction (qRT-PCR) experiments. Three
310 cassava genes have previously been used for this purpose, *GTPb*, *PP2A*, and *UBQ10* (Moreno *et*
311 *al.*, 2011). However, data from the present study show that all three genes display significant
312 variance between tissue types (**Fig. S3**). The datasets described here were queried to identify
313 candidate genes displaying medium level expression with low variance across the tissue types.
314 We identified genes with expression greater than 40 FPKM in all replicates with the lowest
315 coefficient of variation in order to normalize for magnitude of expression. Figure S3 shows the
316 top 10 candidates from our analysis in comparison to the three genes previously used.

317 To facilitate future analyses, a web application has been developed wherein users can
318 specify a desired gene expression pattern across all tissue types and receive a list of candidate
319 genes. This application also allows users to visualize a heatmap of expression values for any
320 gene of interest across each tissue type. The queried gene is displayed in the PCA and overlaid
321 SOM nodes. This application can be accessed here: shiny.danforthcenter.org/cassava_atlas/

322

323 **Discussion**

324 To assist cassava improvement efforts, various genomic, transcriptomic and epigenomic
325 resources have previously been described (Prochnik *et al.*, 2012; Wang *et al.*, 2014; Wang *et al.*,
326 2015). Our study provides a unique resource: we characterize the cassava transcriptome across a
327 wide range of tissue types. Comparison of gene expression patterns revealed a dramatic
328 similarity between OES and FEC tissue. Storage roots were found to be significantly different
329 from the other root tissues, and closer examination of the data suggest that the majority of this
330 difference comes from a lack of gene expression, consistent with the role of this organ as a sink.
331 Our study provides new insight into cassava physiology, and the data will serve as a valuable
332 resource for cassava researchers. In addition, we identify both genes that are constitutively
333 expressed as well as those that are highly tissue specific. The promoters of these genes may be
334 useful for diverse biotechnological applications, including those that seek to alter cassava
335 metabolism and improve the value of cassava as a source of food for a large fraction of the
336 world's population.

337

338 **Acknowledgements**

339 This research was supported by the Bill and Melinda Gates Foundation. Sequencing was
340 performed at the Genome Technology Access Center in the Department of Genetics at
341 Washington University School of Medicine. The Center is partially supported by NCI Cancer
342 Center Support Grant #P30 CA91842 to the Siteman Cancer Center and by ICTS/CTSA Grant#
343 UL1 TR000448 from the National Center for Research Resources (NCRR), a component of the
344 National Institutes of Health (NIH), and NIH Roadmap for Medical Research. This publication is
345 solely the responsibility of the authors and does not necessarily represent the official view of
346 NCRR or NIH.

347

348 **Author Contributions:**

349 M.C.W. analyzed data and co-wrote the manuscript. A.M.M. designed experiments and isolated
350 tissues for RNAseq analysis and co-wrote the manuscript. A.W.H. designed experiments and
351 made constructs. J.B created the shiny application. R.D.C. created transgenic cassava plants.
352 A.V. generated RNAseq libraries. N.J.T and D.F.V. supervised the study and edited the

353 manuscript. D.H.C. performed statistical analyses and co-wrote the manuscript. R.S.B. designed
354 experiments, supervised the study, and co-wrote the paper.

355

356 **Bolger AM, Lohse M, Usadel B. 2014.** Trimmomatic: a flexible trimmer for Illumina sequence
357 data. *Bioinformatics* **30**(15): 2114-2120.

358 **Bull SE, Owiti JA, Niklaus M, Beeching JR, Gruissem W, Vanderschuren H. 2009.**

359 *Agrobacterium*-mediated transformation of friable embryogenic calli and regeneration of
360 transgenic cassava. *Nat Protoc* **4**(12): 1845-1854.

361 **Ceballos H, Iglesias CA, Perez JC, Dixon AGO. 2004.** Cassava breeding: opportunities and
362 challenges. *Plant Mol Biol* **56**: 503-516.

363 **Chauhan RD, Beyene G, Kalyaeva M, Fauquet CM, Taylor N. 2015.** Improvements in
364 *Agrobacterium*-mediated transformation of cassava (*Manihot esculenta* Crantz) for large-
365 scale production of transgenic plants. *Plant Cell, Tissue and Organ Culture (PCTOC)*
366 **121**(3): 591-603.

367 **Chaweewan Y, Taylor N. 2015.** Anatomical Assessment of Root Formation and Tuberization in
368 Cassava (*Manihot esculenta* Crantz). *Tropical Plant Biology* **8**(1-2): 1-8.

369 **de Souza CR, Aragao FJ, Moreira EC, Costa CN, Nascimento SB, Carvalho LJ. 2009.**

370 Isolation and characterization of the promoter sequence of a cassava gene coding for
371 Pt2L4, a glutamic acid-rich protein differentially expressed in storage roots. *Genet Mol*
372 *Res* **8**(1): 334-344.

373 **de Souza CR, Carvalho LJ, de Almeida ER, Gander ES. 2006.** A cDNA sequence coding for
374 a glutamic acid-rich protein is differentially expressed in cassava storage roots. *Protein*
375 *Pept Lett* **13**(7): 653-657.

376 **Gaitan-Solis E, Taylor NJ, Siritunga D, Stevens W, Schachtman DP. 2015.** Overexpression
377 of the transporters AtZIP1 and AtMTP1 in cassava changes zinc accumulation and
378 partitioning. *Front Plant Sci* **6**: 492.

379 **Gegios A, Amthor R, Maziya-Dixon B, Egesi C, Mallowa S, Nungo R, Gichuki S, Mbanaso**

380 **A, Manary MJ. 2010.** Children consuming cassava as a staple food are at risk for
381 inadequate zinc, iron, and vitamin A intake. *Plant Foods Hum Nutr* **65**(1): 64-70.

- 382 **Gresshoff PM, Doy CH. 1974.** Derivation of a haploid cell line from *Vitis vinifera* and the
383 importance of the stage of meiotic development of the anthers for haploid culture of this
384 and other genera. *Z Pflanzenphysiol* **73**: 132-141.
- 385 **Howeler R, Litaladio N, Thomas G. 2013.** *Save and Grow: Cassava - A guide to sustainable*
386 *production intensification*. Rome: Food and Agriculture Organization of the United
387 States of America.
- 388 **Liu J, Zheng Q, Ma Q, Gadidasu KK, Zhang P. 2011.** Cassava genetic transformation and its
389 application in breeding. *J Integr Plant Biol* **53**(7): 552-569.
- 390 **McKinney W. 2010.** Data structures for statistical computing in python. *Proceedings of the 9th*
391 *Python in Science Conference*: 51-56.
- 392 **Moreno I, Gruissem W, Vanderschuren H. 2011.** Reference genes for reliable potyvirus
393 quantitation in cassava and analysis of Cassava brown streak virus load in host varieties.
394 *J Virol Methods* **177**(1): 49-54.
- 395 **Narayanan N, Beyene G, Chauhan RD, Gaitan-Solis E, Grusak MA, Taylor N, Anderson P.**
396 **2015.** Overexpression of Arabidopsis VIT1 increases accumulation of iron in cassava
397 roots and stems. *Plant Sci* **240**: 170-181.
- 398 **Nyaboga E, Njiru J, Nguu E, Gruissem W, Vanderschuren H, Tripathi L. 2013.** Unlocking
399 the potential of tropical root crop biotechnology in east Africa by establishing a genetic
400 transformation platform for local farmer-preferred cassava cultivars. *Front Plant Sci* **4**:
401 526.
- 402 **Patil BL, Legg JP, Kanju E, Fauquet CM. 2015.** Cassava brown streak disease: a threat to
403 food security in Africa. *J Gen Virol* **96**(Pt 5): 956-968.
- 404 **Prochnik S, Marri PR, Desany B, Rabinowicz PD, Kodira C, Mohiuddin M, Rodriguez F,**
405 **Fauquet C, Tohme J, Harkins T, et al. 2012.** The Cassava Genome: Current Progress,
406 Future Directions. *Tropical Plant Biology* **5**(1): 88-94.
- 407 **R Core Team 2015.** R: A language and environment for statistical computing.: R Foundation for
408 Statistical Computing, Vienna, Austria URL <http://www.R-project.org/>.
- 409 **Stephenson K, Amthor R, Mallowa S, Nungo R, Maziya-Dixon B, Gichuki S, Mbanaso A,**
410 **Manary M. 2010.** Consuming cassava as a staple food places children 2-5 years old at

- 411 risk for inadequate protein intake, an observational study in Kenya and Nigeria. *Nutr J* 9:
412 9.
- 413 **Taylor NJ, Gaitán-Solís E, Moll T, Trauterman B, Jones T, Pranjal A, Trembley C,**
414 **Abernathy V, Corbin D, Fauquet CM. 2012.** A high-throughput platform for the
415 production and analysis of transgenic cassava (*Manihot esculenta*) plants. *Trop Plant Biol*
416 **5(1): 127-139.**
- 417 **Taylor NJ, Masona MV, Carcamo R, Ho T, Schopke C, Fauquet CM. 2001.** Production of
418 embryogenic tissues and regeneration of transgenic plants in cassava (*Manihot esculenta*
419 Crantz). *Euphytica* **120: 25-34.**
- 420 **Trapnell C, Pachter L, Salzberg SL. 2009.** TopHat: discovering splice junctions with RNA-
421 Seq. *Bioinformatics* **25(9): 1105-1111.**
- 422 **Trapnell C, Williams BA, Pertea G, Mortazavi A, Kwan G, van Baren MJ, Salzberg SL,**
423 **Wold BJ, Pachter L. 2010.** Transcript assembly and quantification by RNA-Seq reveals
424 unannotated transcripts and isoform switching during cell differentiation. *Nat Biotechnol*
425 **28(5): 511-515.**
- 426 **Wang H, Beyene G, Zhai J, Feng S, Fahlgren N, Taylor NJ, Bart R, Carrington JC,**
427 **Jacobsen SE, Ausin I. 2015.** CG gene body DNA methylation changes and evolution of
428 duplicated genes in cassava. *Proceedings of the National Academy of Sciences of the*
429 *United States of America* **112(44): 13729-13734.**
- 430 **Wang W, Feng B, Xiao J, Xia Z, Zhou X, Li P, Zhang W, Wang Y, Møller BL, Zhang P, et**
431 **al. 2014.** Cassava genome from a wild ancestor to cultivated varieties. *Nature*
432 *Communications* **5: 5110.**
- 433 **Zainuddin IM, Schlegel K, Gruissem W, Vanderschuren H. 2012.** Robust transformation
434 procedure for the production of transgenic farmer-preferred cassava landraces. *Plant*
435 *methods* **8(1): 24.**

436

437 **Figure Legends**

438 **Figure 1. Cartoon and pictures of cassava tissues sampled for gene expression atlas.** Eleven
439 tissue types were dissected by hand and frozen in liquid nitrogen prior to processing for RNA
440 sequencing library preparation.

441
442 **Figure 2. Comparison of global gene expression patterns across 11 cassava tissue types. (A)**
443 Heatmap showing every pairwise comparison for the 11 tissue types sampled, as produced by
444 CummeRbund's csDistHeat() method. Lighter colors correspond to more closely related tissue
445 types. The numbers in each cell represent the Jensen-Shannon distance between those two tissues
446 using the mean expression values of the biological replicates. **(B)** Output of CummeRbund's
447 csDendro() method. This dendrogram is created using the Jensen-Shannon distances calculated
448 between the consensus expression values of genes for each tissue type. A low squared coefficient
449 of variation for biological replicates was observed indicating the high quality of this dataset (Fig.
450 S1A). **(C-E)** Volcano plots showing the differential expression of genes from leaf to fibrous
451 roots **(C)**, OES to FEC **(D)** and storage root to fibrous root **(E)** using FDR corrected p-value as
452 the y axis. Number of genes significantly up-regulated in each tissue type, for each comparison,
453 are listed. Red vertical lines: $\pm \log_2(\text{fold_change}) = 2$, red horizontal line: $\log \text{score} = 1.3$. The
454 green points indicate significantly differentially expressed genes based on these cutoffs.

455
456 **Figure 3: Transcript expression profiles across tissues. (A,B)** Principal component analysis
457 (PCA) performed on replicates of tissue samples, using transcript expression levels. **(C,D)** PCA
458 performed on transcript profiles, across tissue samples. Colors correspond to self-organizing map
459 (SOM) nodes used to find transcripts with similar expression profiles, which cluster together in
460 the PCA space. **(E)** Scaled transcript expression profiles of SOM nodes across tissue types. **(F)**
461 Heatmap of genes showing gene expression pattern corresponding to the nodes in Fig. 3C,D.
462 Gene Ontology terms associated with these genes are listed on the right.

463
464 **Figure 4. Identification of genes specifically expressed in a single or subset of tissue types.**
465 **(A)** Identification of genes with specified expression patterns. (top; middle) Heatmap of the most
466 highly, uniquely expressed genes in each tissue. Requirements for below and above the limits of
467 detectable expression are listed on the left (OFF and ON, respectively). (bottom) Genes
468 expressed highly in a subset of tissues are reported. No genes specifically expressed across all
469 subterranean tissues (storage root, fibrous root and root apical meristem (RAM)) were identified
470 so storage root was excluded from that group. **(B)** Cumulative distribution plot of FPKM values

471 of functionally annotated genes with expression in at least one tissue type. The vertical lines
472 represent 3 cutoffs used in the analysis. < 1 FPKM (maroon line) = below the limit of detection;
473 < 4 FPKM (yellow line) = below the limit of detection (relaxed); > 8 FPKM (red line) = detected
474 expression (relaxed); >10 FPKM (green line) = detected expression; > 300 FPKM (blue line) =
475 highly expressed genes.

476

477 **Figure 5. Identification of highly expressed genes in all tissue types.** (A) Heatmap displaying
478 expression for 31 annotated genes expressed above 300 FPKM in all tissues. Values greater than
479 7000 are condensed within the heatmap scale. (B) Ability of promoters to drive gene expression
480 was assessed transiently in *Nicotiana benthamiana*. In addition, reporter gene constructs were
481 stably transformed into cassava FEC cells. BF: bright field image of GUS staining intensity;
482 HSL: images were converted to HSL color space; Mean intensity across infiltrated spot shown
483 +/- standard deviation. Genes tested are highlighted in blue in A. Full image is shown in Figure
484 S2.

485

486 **Figure S1. Assessing variation among biological replicates.** (A) The squared coefficient of
487 variation of replicates in each of the 11 tissue types is shown here to be low and reasonably
488 uniform across all tissue types. This is indicative of replicates being closely related, limiting the
489 possibility for error in sampling. (B) Distribution of all FPKM values greater than 1 in
490 functionally annotated genes in each tissue type plotted against a \log_2 scale on the y-axis. This
491 demonstrates the similar expression of each tissue type across annotated genes.

492

493 **Figure S2. Original images that accompany Figure 5b.**

494

495 **Figure S3. Identification of constitutively expressed genes and assessment of expression**
496 **variation across tissue type.** Expression profile across eleven tissue types was investigated for
497 three housekeeping genes: GTPb, PP2A and UBQ10. The dataset was queried for genes that
498 showed medium level expression (greater than 40 FPKM) and low variability (low coefficient of
499 variation) across all tissues. Top ten genes are displayed.

500

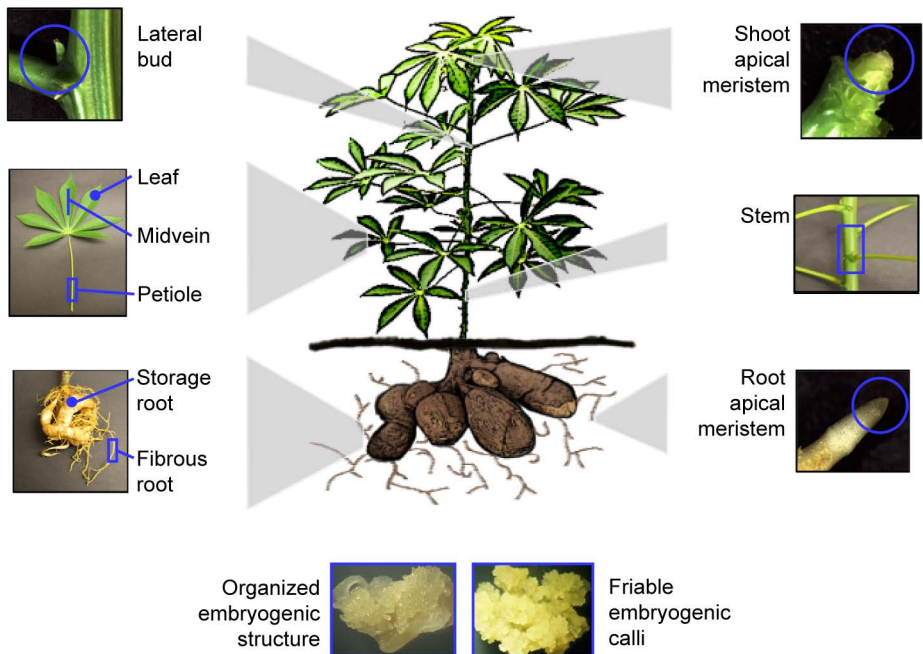


Figure 1. Cartoon and pictures of cassava tissues sampled for gene expression atlas. Eleven tissue types were dissected by hand and frozen in liquid nitrogen prior to processing for RNA sequencing library preparation.

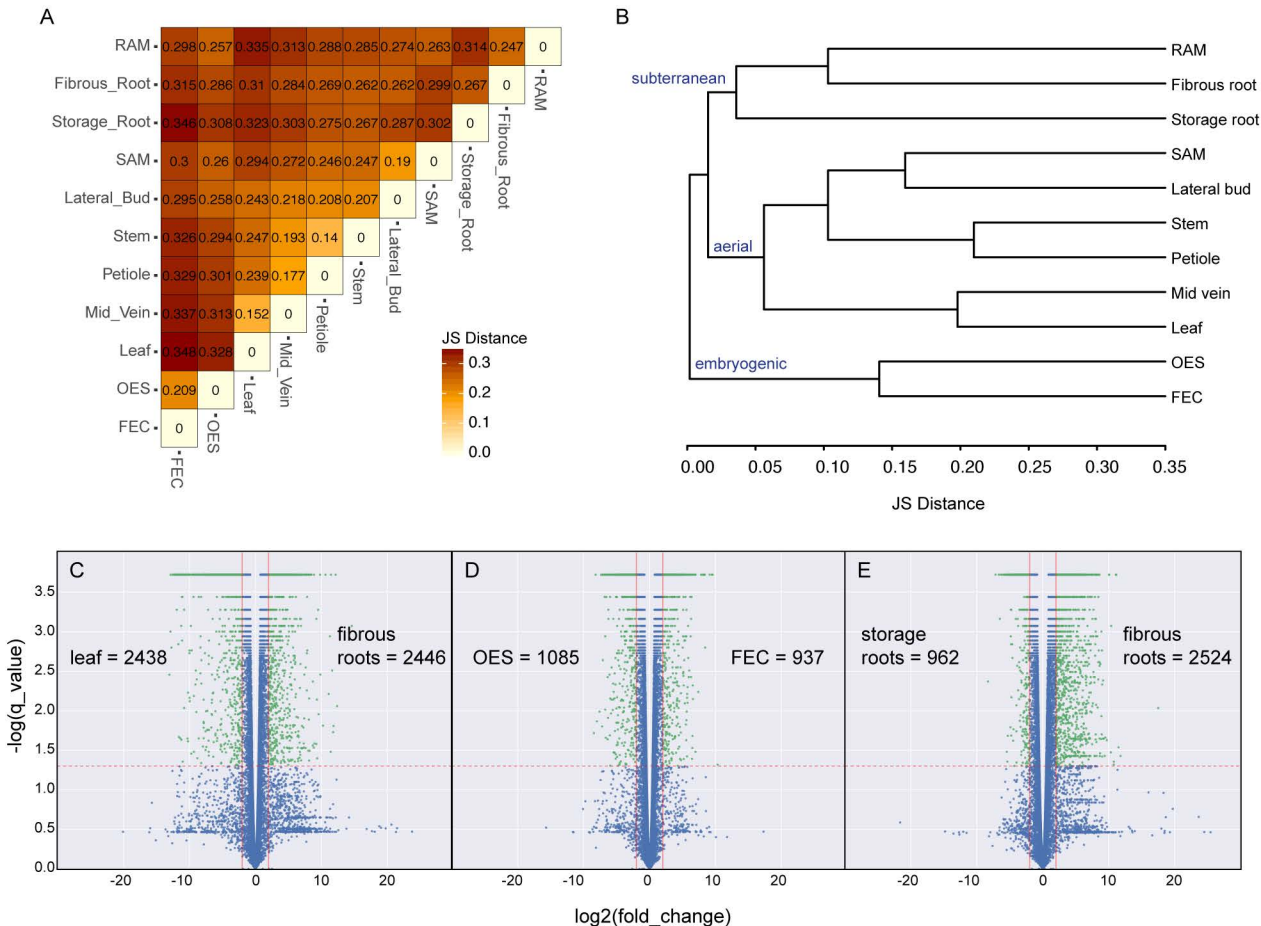


Figure 2. Comparison of global gene expression patterns across 11 cassava tissue types. (A) Heatmap showing every pairwise comparison for the 11 tissue types sampled, as produced by CummeRbund's `csDistHeat()` method. Lighter colors correspond to more closely related tissue types. The numbers in each cell represent the Jensen-Shannon distance between those two tissues using the mean expression values of the biological replicates. (B) Output of CummeRbund's `csDendro()` method. This dendrogram is created using the Jensen-Shannon distances calculated between the consensus expression values of genes for each tissue type. A low squared coefficient of variation for biological replicates was observed indicating the high quality of this dataset (Fig. S1A). (C-E) Volcano plots showing the differential expression of genes from leaf to fibrous roots (C), OES to FEC (D) and storage root to fibrous root (E) using FDR corrected p-value as the y axis. Number of genes significantly up-regulated in each tissue type, for each comparison, are listed. Red vertical lines: $\pm \log_2(\text{fold_change}) = 2$, red horizontal line: $\log \text{score} = 1.3$. The green points indicate significantly differentially expressed genes based on these cutoffs.

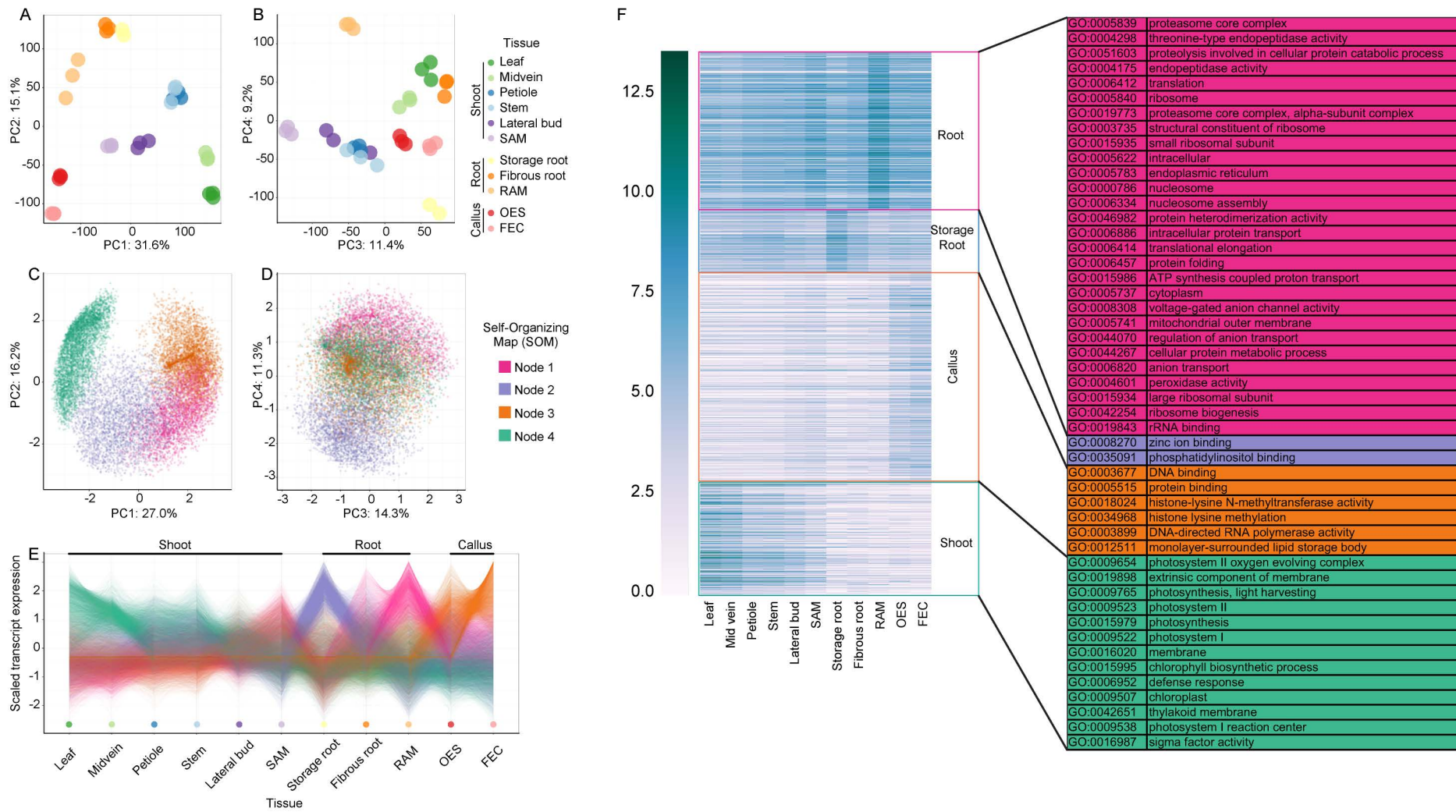


Figure 3: Transcript expression profiles across tissues. (A,B) Principal component analysis (PCA) performed on replicates of tissue samples, using transcript expression levels. (C,D) PCA performed on transcript profiles, across tissue samples. Colors correspond to self-organizing map (SOM) nodes used to find transcripts with similar expression profiles, which cluster together in the PCA space. (E) Scaled transcript expression profiles of SOM nodes across tissue types. (F) Heatmap of genes showing gene expression pattern corresponding to the nodes in Fig. 3C,D. Gene Ontology terms associated with these genes are listed on the right.

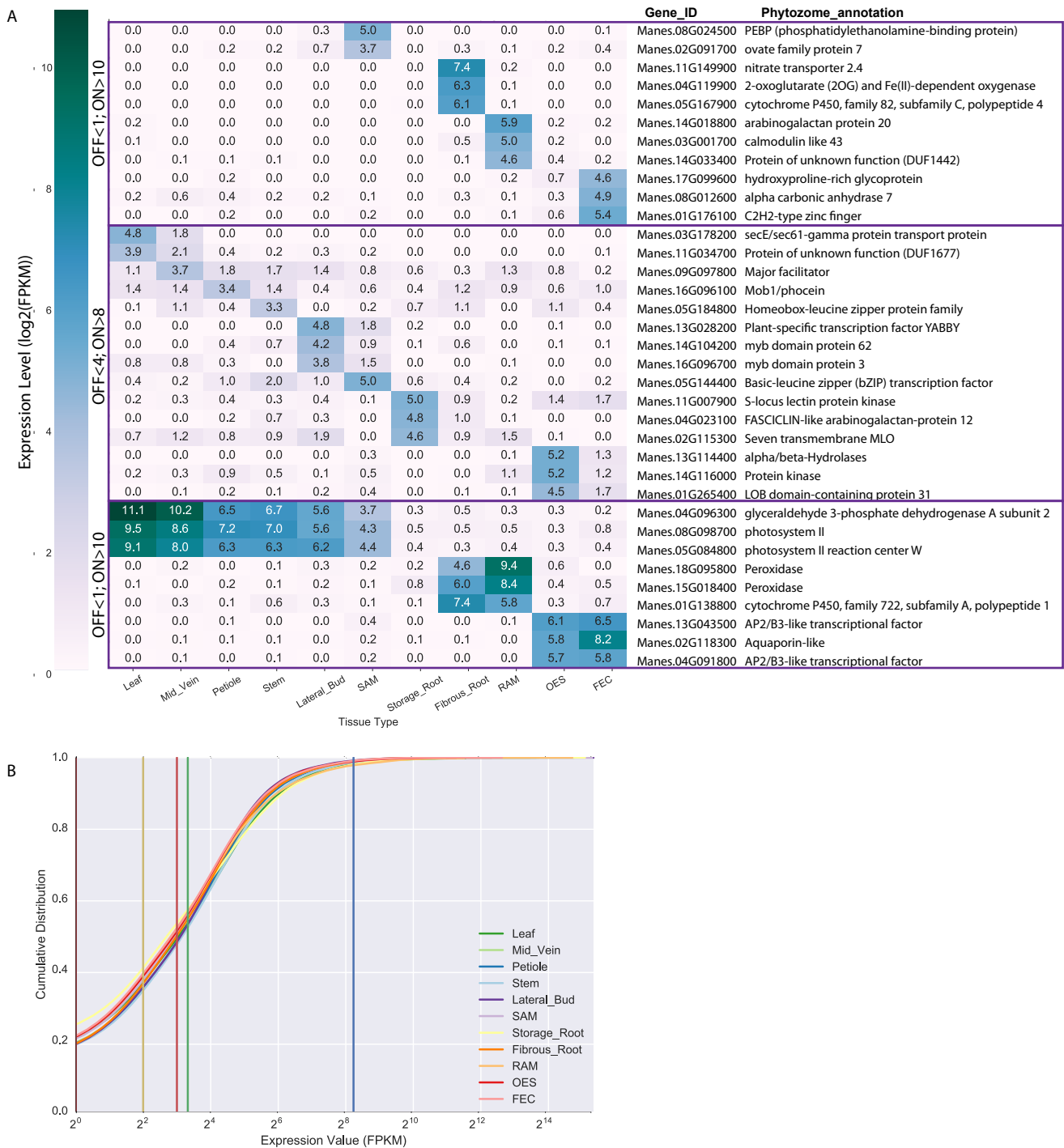


Figure 4. Identification of genes specifically expressed in a single or subset of tissue types. (A) Identification of genes with specified expression patterns. (top; middle) Heatmap of the most highly, uniquely expressed genes in each tissue. Requirements for below and above the limits of detectable expression are listed on the left (OFF and ON, respectively). (bottom) Genes expressed highly in a subset of tissues are reported. No genes specifically expressed across all subterranean tissues (storage root, fibrous root and root apical meristem (RAM)) were identified so storage root was excluded from that group. (B) Cumulative distribution plot of FPKM values of functionally annotated genes with expression in at least one tissue type. The vertical lines represent 3 cutoffs used in the analysis. < 1 FPKM (maroon line) = below the limit of detection; < 4 FPKM (yellow line) = below the limit of detection (relaxed); > 8 FPKM (red line) = detected expression (relaxed); >10 FPKM (green line) = detected expression; > 300 FPKM (blue line) = highly expressed genes.

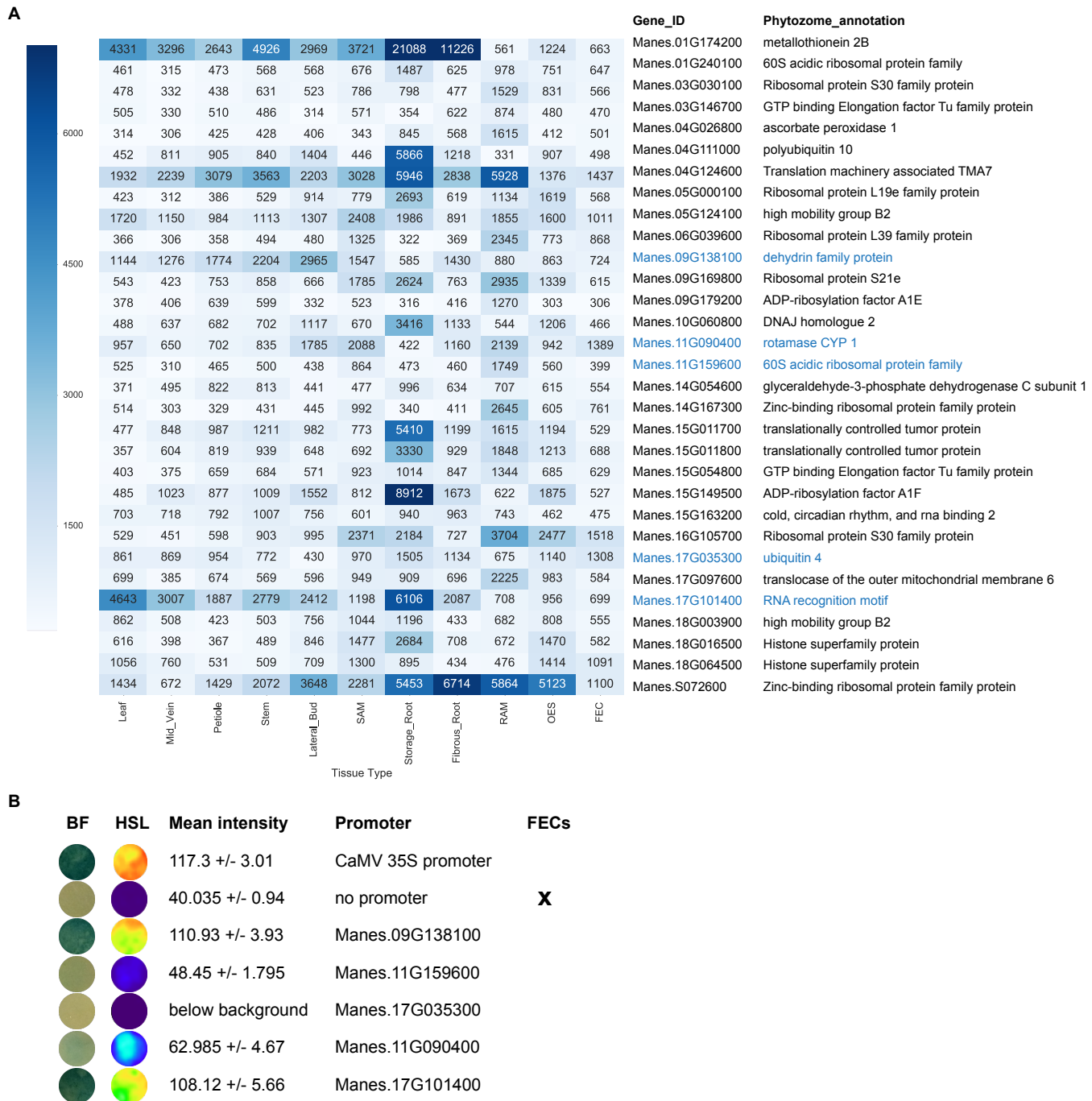


Figure 5. Identification of highly expressed genes in all tissue types. (A) Heatmap displaying expression for 31 annotated genes expressed above 300 FPKM in all tissues. Values greater than 7000 are condensed within the heatmap scale. (B) Ability of promoters to drive gene expression was assessed transiently in *Nicotiana benthamiana*. In addition, reporter gene constructs were stably transformed into cassava FEC cells. BF: bright field image of GUS staining intensity; HSL: images were converted to HSL color space; Mean intensity across infiltrated spot shown +/- standard deviation. Genes tested are highlighted in blue in A. Full image is shown in Figure S2.

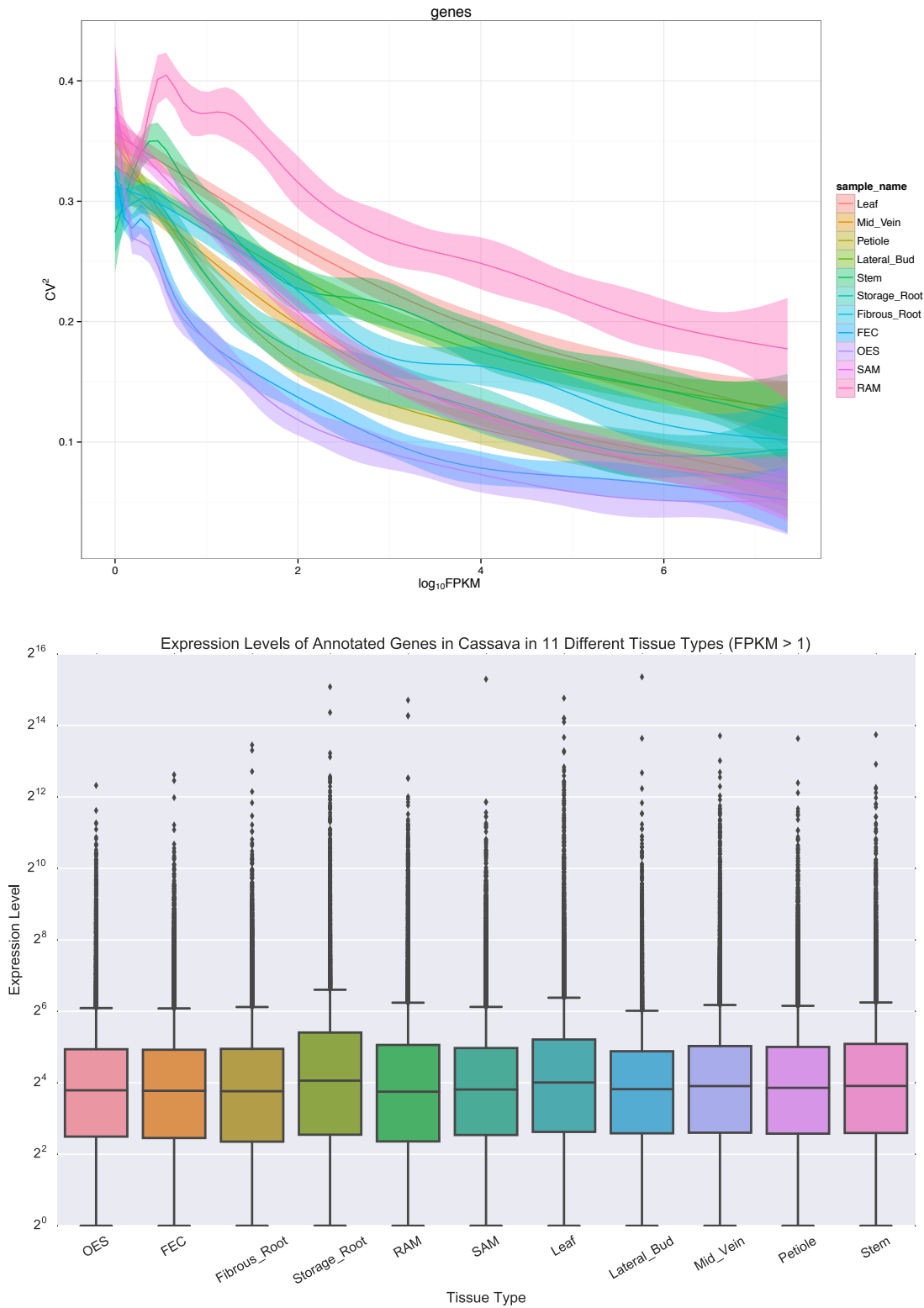


Figure S1. Assessing variation among biological replicates. (A) The squared coefficient of variation of replicates in each of the 11 tissue types is shown here to be low and reasonably uniform across all tissue types. This is indicative of replicates being closely related, limiting the possibility for error in sampling. (B) Distribution of all FPKM values greater than 1 in functionally annotated genes in each tissue type plotted against a \log_2 scale on the y-axis. This demonstrates the similar expression of each tissue type across annotated genes.

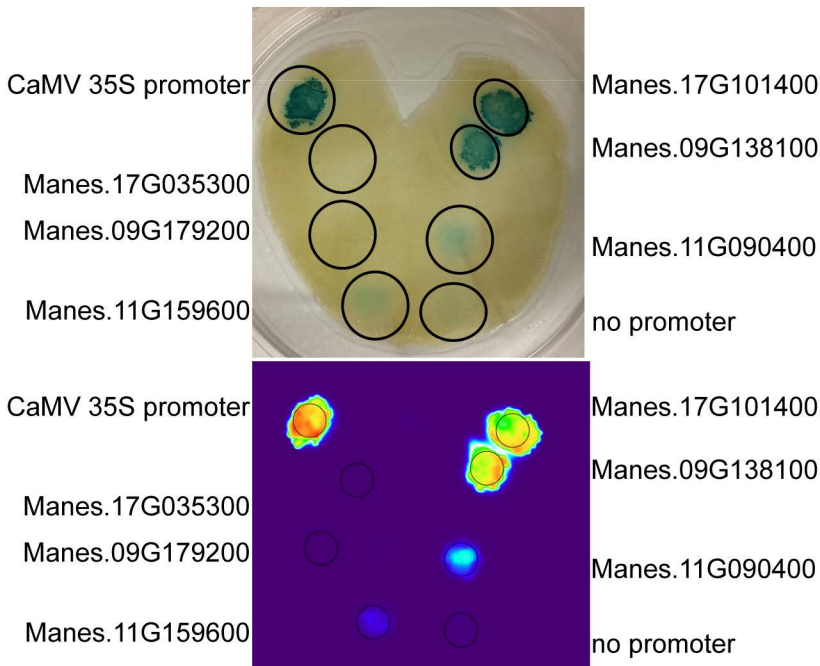


Figure S2. Original images that accompany Figure 5b.

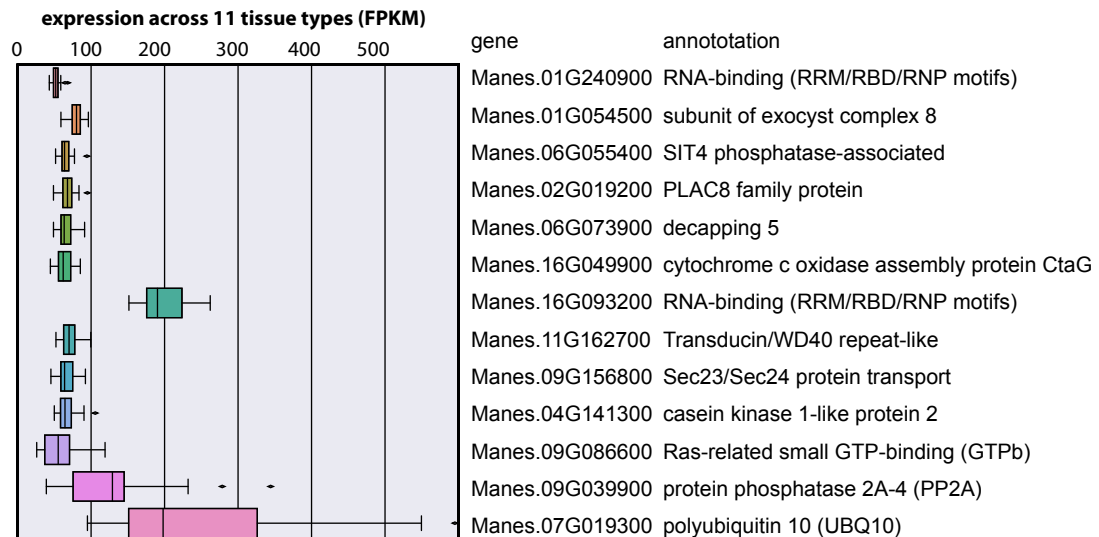


Figure S3. Identification of constitutively expressed genes and assessment of expression variation across tissue type. Expression profile across eleven tissue types was investigated for three housekeeping genes: GTPb, PP2A and UBQ10. The dataset was queried for genes that showed medium level expression (greater than 40 FPKM) and low variability (low coefficient of variation) across all tissues. Top ten genes are displayed.