# Expression-based machine learning models for predicting plant tissue identity

Sourabh Palande[1], Jeremy Arsenault[2], Patricia Basurto-Lozada[3], Andrew Bleich[4], Brianna N. I. Brown[4], Sophia F. Buysse[4,5,6], Noelle A. Connors[7], Sikta Das Adhikari[1,8], Kara C. Dobson[5,9], Francisco Xavier Guerra-Castillo[10,11], Maria F. Guerrero-Carrillo[12], Sophia Harlow[7], Héctor Herrera-Orozco[13,14], Asia T. Hightower[4,5], Paulo Izquierdo[15], MacKenzie Jacobs[16,17], Nicholas A. Johnson[5,18], Wendy Leuenberger[5,9], Alessandro Lopez-Hernandez[3,19], Alicia Luckie-Duque[12], Camila Martínez-Avila[20], Eddy J. Mendoza-Galindo[12], David Plancarte[21], Jenny M. Schuster[22,17], Harry Shomer[2], Sidney C. Sitar[15,23,24], Anne K. Steensma[4,17,25], Joanne Elise Thomson[17,22], Damián Villaseñor-Amador[26], Robin Waterman[4,5,6], Brandon M. Webster[4], Madison Whyte[15], Sofía Zorilla-Azcué[27], Beronda L. Montgomery[28], Aman Y. Husbands[29], Arjun Krishnan[30], Sarah Percival[1], Elizabeth Munch[1,31], Robert VanBuren[7,32], Daniel H. Chitwood[1,7,]*, Alejandra Rougon-Cardoso[12,33]*

[1]Department of Computational Mathematics, Science and Engineering, Michigan State University, USA
[2]Department of Computer Science and Engineering, Michigan State University, USA
[3]Laboratorio Internacional de Investigación sobre el Genoma Humano (LIIGH), Universidad Nacional Autónoma de México, México
[4]Department of Plant Biology, Michigan State University, USA
[5]Ecology, Evolution, and Behavior Program, Michigan State University, USA
[6]Kellogg Biological Station, Michigan State University, USA
[7]Department of Horticulture, Michigan State University, USA
[8]Department of Statistics and Probability, Michigan State University, USA
[9]Department of Integrative Biology, Michigan State University, USA
[10]Unidad de Investigación Médica en Inmunología e Infectología, Instituto Mexicano del Seguro Social, México
[11]Programa de Posgrado en Ciencias Biológicas, Facultad de Medicina, Universidad Nacional Autónoma de México, México
[12]Laboratory of Agrigenomic Sciences, Escuela Nacional de Estudios Superiores Unidad León, Universidad Nacional Autónoma de México, México
[13]Posgrado en Ciencias Biológicas, Universidad Nacional Autónoma de México, México
[14]Laboratorio de Ecología Evolutiva y Conservación de Anfibios y Reptiles. Facultad de Estudios Superiores Iztacala, Universidad Nacional Autónoma de México, México
[15]Department of Plant, Soil, and Microbial Sciences, Michigan State University, USA
[16]Department of Biochemistry and Molecular Biology, Michigan State University, USA
[17]Molecular Plant Sciences Program, Michigan State University, USA
[18]Genetics and Genome Sciences, Michigan State University, USA
[19]Computational Population Genetics Group, Universidad Nacional Autónoma de México, México
[20]Colección Nacional de Aves, Posgrado en Ciencias Biológicas, Instituto de Biología, Universidad Nacional Autónoma de México, México
[21]Departamento de Botánica, Posgrado en Ciencias Biológicas, Instituto de Biología, Universidad Nacional Autónoma de México, México
[22]Cell and Molecular Biology, Michigan State University, USA
[23]Plant Breeding, Genetics, and Biotechnology, Michigan State University, USA
[24]Crop and Soil Sciences Program, Michigan State University, USA
[25]MSU-DOE Plant Research Laboratory, Michigan State University, USA
[26]Programa de Posgrado en Ciencias Biológicas, Facultad de Ciencias, Universidad Nacional Autónoma de México, México
[27]Programa de Posgrado en Ciencias Biológicas, Escuela Nacional de Estudios Superiores (ENES), Unidad Morelia, Universidad Nacional Autónoma de México, México
[28]Department of Biology, Grinnell College, USA
[29]Department of Biology, University of Pennsylvania, USA
[30]Department of Biomedical Informatics, Center for Health AI, University of Colorado Anschutz Medical Campus, USA
[31]Department of Mathematics, Michigan State University, USA
[32]Plant Resilience Institute, Michigan State University, USA
[33]Plantecc National Laboratory, ENES-León, México

*Corresponding authors
Dr. Alejandra Rougon-Cardoso, arougon@enes.unam.mx
Dr. Daniel H. Chitwood, dhchitwood@gmail.com

**ABSTRACT**

The selection of *Arabidopsis* as a model organism played a pivotal role in advancing genomic science, firmly establishing the cornerstone of today's plant molecular biology. Competing frameworks to select an agricultural- or ecological-based model species, or to decentralize plant science and study a multitude of diverse species, were selected against in favor of building core knowledge in a species that would facilitate genome-enabled research that could assumedly be transferred to other plants. Here, we examine the ability of models based on *Arabidopsis* gene expression data to predict tissue identity in other flowering plant species. Comparing different machine learning algorithms, models trained and tested on *Arabidopsis* data achieved near perfect precision and recall values using the K-Nearest Neighbor method, whereas when tissue identity is predicted across the flowering plants using models trained on *Arabidopsis* data, precision values range from 0.69 to 0.74 and recall from 0.54 to 0.64, depending on the algorithm used. Below-ground tissue is more predictable than other tissue types, and the ability to predict tissue identity is not correlated with phylogenetic distance from *Arabidopsis*. This suggests that gene expression signatures rather than marker genes are more valuable to create models for tissue and cell type prediction in plants. Our data-driven results highlight that, in hindsight, the assertion that knowledge from *Arabidopsis* is translatable to other plants is not always true. Considering the current landscape of abundant sequencing data and computational resources, it may be prudent to reevaluate the scientific emphasis on *Arabidopsis* and to prioritize the exploration of plant diversity.

**INTRODUCTION**

Historically, plant biology has focused on inferring genetic, molecular, physiological, and ecological mechanisms. Conventionally, through quantifying phenomena and applying statistics, hypotheses are tested and decisions of most likely scenarios are determined. New technologies and computational approaches have caused a shift from hypothesis- to data-driven research (Mazzocchi, 2015). Moreover, plant biology has embraced the inclusion of machine learning methods in addition to traditional statistical approaches (Ij, 2018). Both a deluge of data and new computational methods have allowed for predictive, rather than inferential, methods. Both statistics and machine learning can be used for inference and prediction, but machine learning methods more often classify and predict on class labels rather than inferring statistical parameters of a population. In plant biology, such predictive approaches underlie the frameworks of phenotyping (Coppens et al., 2017), precision agriculture (Zhang et al., 2002), genomic prediction (Crossa et al., 2014), linking transcriptomic profiles to phenotype (Azodi et al., 2020), and protein structure determination (Jumper et al., 2021). Just as inferential statistics has its limitations, the robustness and ability to extrapolate predictive models are also constrained by the empirical context from which the data originates. Although data-driven research is slowly becoming more theoretical and predictive (Hogeweg, 2011), the creation of universal plant models is hindered by their overwhelming diversity. Not only is the phylogenetic diversity among flowering plants immense (The Angiosperm Phylogeny Group et al., 2016), but plants are exceptionally responsive to their environments (Sultan, 2000) and have evolved symbiotic interactions with and defense mechanisms against innumerable microbes (Mitchell et

101 al., 2006). Furthermore, technical variability in data acquisition makes it difficult to exploit the
102 huge amount of expression data archived in databases. The number of ways we sample
103 molecular profiles from plant tissues and the interaction effects that arise between
104 phylogenetically diverse species with environments, stresses, and biotic interactions is
105 countless and prevents extrapolating results between studies.
106
107 Due to the clear advantages of studying a single model species, the early days of the genomics
108 era tended to overlook the importance of prioritizing plant diversity. The candidates considered
109 for the first sequenced genome were either easily transformable (e.g., species within
110 Solanaceae; Knapp et al., 2004) or were already used for genetics (e.g., maize; Strable and
111 Scanlon, 2009), but never was biodiversity considered (Meyerowitz, 2001). Reasons for
112 choosing *Arabidopsis* as the first sequenced plant genome (Arabidopsis Genome Initiative,
113 2000) include ease of transformation (Clough and Bent, 1998), its small genome (Bennett et al.,
114 2003), and life history traits that allow for genetics through crossing, and short generation times
115 (Meyerowitz, 1987). The justification for initially sequencing the genome of a single model
116 species was that such focus would allow unprecedented molecular discoveries that could be
117 translated into other species and improve our understanding of all plants (Bevan and Walsh,
118 2005). The strategy to focus on a single model species was successful, and *Arabidopsis* is the
119 most cited plant in the last 20 years, even surpassing key crops and all other plant species
120 (Marks et al., 2023). Our molecular knowledge in plants was purposefully constructed to focus
121 on *Arabidopsis* over crops and plant genetic diversity. However, such a choice has little
122 relevance in a changing climate with dwindling natural resources and vanishing biodiversity that
123 have become the most pressing concerns of our time. The cultural dynamics that influenced the
124 choice of *Arabidopsis* as the first sequenced genome are reflected in subsequently sequenced
125 plant genomes. Plants intrinsic to Indigenous cultures and territories have been sequenced by
126 colonial powers (Marks et al., 2021; Dyer et al., 2022). While sequencing *Arabidopsis* has
127 certainly expanded our knowledge of molecular processes, due to such an intense focus, our
128 understanding in other species remains limited. This leaves us questioning the extent to which
129 the insights from *Arabidopsis* can be extrapolated to the rest of flowering plants.
130
131 In the 20 years since the release of the *Arabidopsis* genome sequence (Arabidopsis Genome
132 Initiative, 2000), the number of sequenced plant genomes has dramatically risen (Michael and
133 Jackson, 2013; Li and Harkess, 2018; Marks et al., 2021) leading to a greater understanding of
134 the evolutionary origin and genetic mechanisms underlying numerous traits across the green
135 lineage. Next-generation sequencing, for example, has enabled unprecedented surveys of
136 genome-scale features across species, tissue types, environments, and interactions between
137 plants with abiotic and biotic factors. There are currently over 300,000 public gene expression
138 datasets spanning thousands of diverse plant species (Lim et al., 2022). Cross-species
139 comparisons of gene expression across plants have usually been limited by the number of
140 species analyzed (Proost and Mutwil, 2018) or their sampling breadth. Most studies have
141 generated datasets from scratch (Julca et al., 2021) instead of leveraging public repositories.
142 Databases and datasets curating and making vast amounts of gene expression profiles and
143 their associated metadata have been created. For example, an *Arabidopis* RNA-seq database
144 (ARS) compiles 20,068 publicly available *Arabidopsis* RNA-Seq libraries (Zhang et al., 2020),

145 and the Plant Public RNA☐seq Database has ~45,000 maize, rice, wheat, soybean and cotton
146 samples (Yu et al. 2022). Previously we had curated a dataset of 2,671 publicly available gene
147 expression profiles from 54 flowering plant species across 7 developmental tissue types and
148 nine stresses (Palande et al., 2023). More than 20 years after the release of the *Arabidopsis*
149 genome, not only have we accumulated enough data across plants to ask unprecedented
150 questions but new computational tools are available that permit comparative approaches to
151 analyze such massive amounts of data.
152
153 Here, building upon large, curated databases of *Arabidopsis* (Zhang et al., 2020) and flowering
154 plant gene expression profiles (Palande et al., 2023), we examine how predictive *Arabidopsis* is
155 as a model species relative to the rest of the flowering plants and to what degree we can
156 extrapolate our knowledge from model organisms to diverse plant species. Dimension reduction
157 through principal component analysis (PCA) reveals that biotic stress response and tissue type
158 are primary, orthogonal sources of structure in gene expression data from *Arabidopsis*, and
159 while angiosperm data projected onto this space retains some structure, the regions occupied
160 between tissue types become less distinct. We next compare the performance of different
161 machine learning models. The k-nearest neighbor (KNN) method yields precision and recall
162 values of up to 0.99 with models trained and tested on *Arabidopsis* data. Model performance
163 drops significantly, with higher precision than recall values, when data from across flowering
164 plants is tested using models trained on *Arabidopsis* data. Below-ground tissue is more
165 separated from and predictable than other tissue types, and phylogenetic distance from
166 *Arabidopsis* does not appear to influence prediction rates. We end with a discussion of the
167 implications of our results for the current structure of the plant science community,
168 acknowledging that the past focus on *Arabidopsis* as a model organism based on decisions
169 decades ago was effective at that time; however, we now advocate for a shift in approach due
170 to changing circumstances, particularly in light of the pressing issue of biodiversity loss. We
171 argue for a more decentralized and inclusive research framework that better encompasses the
172 diversity of plants and the human cultures that represent them, adapting to current
173 environmental and scientific challenges.
174
175 **MATERIALS AND METHODS**
176
177 *Datasets*
178
179 We used two curated databases in this analysis. The first contained 28,165 *Arabidopsis* gene
180 expression profiles across 37,334 genes (Zhang et al. 2020). The second contained 2,671
181 flowering plant expression profiles across 6,327 orthogroups (Palande et al. 2023). Metadata
182 labels for each sample from both of the databases was assigned one of four tissue type labels
183 (above-ground, below-ground, whole plant, or other). The categories are purposefully
184 encompassing and chosen to facilitate accurate assignment across the broad categories of
185 experimental data we analyzed, focusing on above-ground and below-ground tissue identity as
186 one of the simplest cases to test tissue predictability. After removing samples with missing
187 metadata and samples with low unique mapped rate (<75%), the *Arabidopsis* database was left
188 with 19,415 samples. A conserved *Arabidopsis* database was also constructed by keeping only

189   the genes mapped to the orthogroups from the flowering plant database. The conserved
190   *Arabidopsis* database contained the same number of samples, but with much smaller
191   expression profiles across only the 6,327 orthogroups shared with the angiosperm dataset.
192
193   *Classification models*
194
195   Classification is a common machine learning task where, given data points belonging to two or
196   more classes, the goal is to *learn* a function that best differentiates between points from different
197   classes. Then, given a new data point, the function can be used to decide which class the point
198   belongs to. The classifier function can be learned in many different ways, leading to various
199   types of machine learning models. For each classifier model in this study, we employed the
200   following modeling methods:
201
202   *Linear support vector classifier (SVC)*: In linear classification, each point is viewed as a vector in
203   *k*-dimensional space (Cortes and Vapnik, 1995). The goal is to find *(k-1)*-dimensional
204   hyperplanes that separate the points belonging to different classes. There are many possible
205   choices for hyperplanes that can classify the points. A reasonable choice is to find the ones that
206   maximize the separation between points from different classes. These are known as maximum-
207   margin hyperplanes. Geometrically, the max-margin hyperplanes are defined by the points that
208   lie closest to them; therefore, such points are called support vectors.
209
210   *Multi-layer perceptron (MLP)*: The SVC model assumes that the classes are linearly separable,
211   which may not be true. MLPs are a class of artificial neural networks (Haykin, 1998) with three
212   or more layers of "perceptrons" with non-linear activation. An MLP consists of an input and an
213   output layer, with one or more hidden layers of neurons. We experimented with one and two
214   hidden-layer MLPs and used rectified linear unit (ReLU) activation in all cases. In ReLU, a
215   neuron's activation is the weighted sum of its inputs, if the sum is non-negative, and zero
216   otherwise. Even with this simple nonlinear activation function, MLPs are able to outperform the
217   linear SVC.
218
219   *Random forest (RF)*: Random forests (Ho, 1995) perform classification by constructing an
220   ensemble of decision trees. Each decision tree outputs a class label for the given sample and
221   the output of the RF is the class label predicted by the majority of the trees. In a decision tree,
222   each internal node is labeled by an input feature and the leaf nodes are labeled by the class
223   labels. Starting from the root node, the input set is recursively partitioned into children nodes
224   using the input feature associated with the node. The recursion ends when all data points in the
225   node belong to the same class, or some pre-specified termination criteria, such as maximum
226   depth of the tree, are met. Which feature to split the data on at each level is determined using
227   information criteria such as gini impurity or entropy that measure how consistent the subsets are
228   with respect to the class labels after the split.
229
230   *Histogram-based gradient boosting (HGB)*: Gradient boosting (Mason et al., 1999) is another
231   class of methods that uses a large ensemble of decision trees. In histogram-based boosting, the
232   real-valued input features are first discretized into a few (typically 256) bins using histograms.

5

233    This allows the training algorithm to run much more efficiently and construct a much larger
234    ensemble of decision trees to support the classification.

235

236    *K-nearest neighbor (KNN) classifier*: In KNN classifiers (Cover and Hart, 1967) class labels are
237    assigned based on a majority vote of the K nearest training points. The distance metric and the
238    number of neighbors are specified by the user. In our experiments, correlation distance between
239    the expression profiles was used to train the KNN classifier.

240

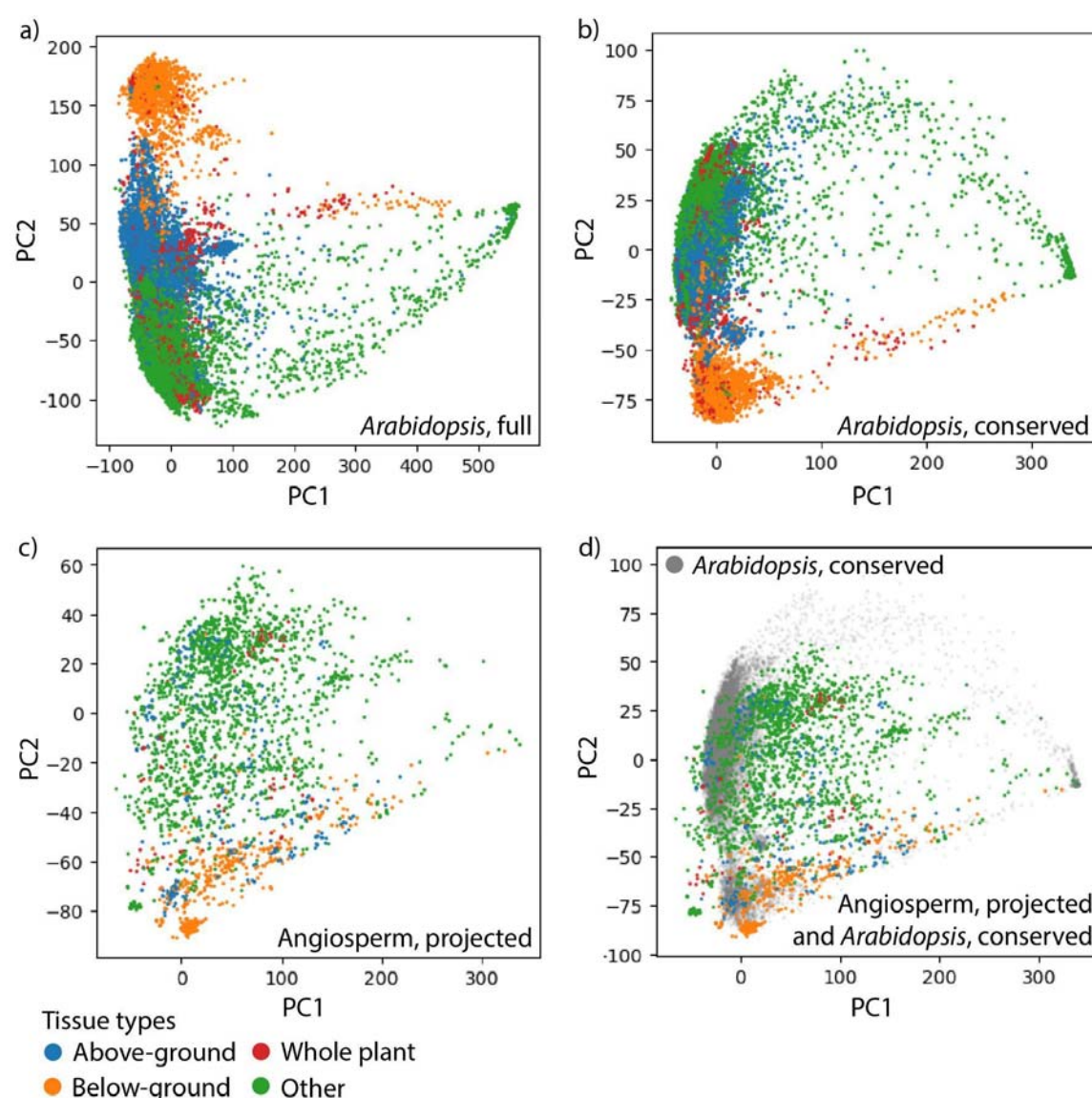241    *Experimental design*

242

243    To establish the utility of gene expression profiles in predicting tissue type, we trained the
244    supervised machine learning models to classify the *Arabidopsis* data by tissue types (**Table 1**).
245    The database was split into training and test sets (70%-30% split). To ensure comparability, all
246    five models were trained and tested on the same training and test sets. Next, we wanted to
247    examine how predictive *Arabidopsis* is to the rest of the flowering plants (**Table 2**). To test this,
248    we used a set of conserved *Arabidopsis* transcripts with orthogroups across angiosperms, split
249    into training and test sets (70%-30% split) as before. The same five machine learning models
250    were trained on the conserved *Arabidopsis* training set. The performance of these models was
251    first tested on the conserved gene *Arabidopsis* test set to make sure that the models were still
252    able to predict the tissue types with a significantly smaller number of features. We then used the
253    same models to classify the angiosperm data to test how well they extrapolate to species other
254    than *Arabidopsis*. Each machine learning model employed in our experiments requires
255    additional hyperparameters that need to be tuned to optimize model performance. We used the
256    Bayesian optimization procedure implemented in the hyperopt package in Python (Bergstra et
257    al., 2013). To gain insights into the functional annotation and enrichment of our gene list, we
258    performed a Gene Ontology term analysis using the DAVID Functional Annotation Clustering
259    tool (version 2021) from the web interface http://david.ncifcrf.gov (Huang et al., 2009). We
260    filtered the 200 genes with the most positive and negative PC1 loading values. The annotation
261    was performed using TAIR IDs and selecting Gene Ontology terms from levels 3 and 4 of
262    Molecular Function and Biological Process categories. All data and code to reproduce the
263    results in this manuscript are available at https://github.com/PlantsAndPython/arabidopsis-gene-
264    expression.

265

266    **RESULTS**

267

268    *Dimension reduction and alignment between* Arabidopsis *and angiosperm gene expression*
269    *datasets*

270

271    A principal component analysis (PCA) performed on the full dataset of 19,415 *Arabidopsis*
272    RNAseq samples shows a clear separation by tissue type (**Fig. 1a**). For simplicity, we
273    categorized samples into bins of above-ground, below-ground, whole plant, and other. The
274    above-ground, below-ground, and other tissue types are well-separated from each other, but the
275    below-ground tissue has the least overlap with other tissues. The whole plant tissue type,
276    composed of different combinations of the other tissues, is not well separated, as we would

277    expect. The separation of tissues occurs along a gradient defined by PC2, demonstrating that
278    tissue type is not the primary source of variance in the data. Rather, a small proportion of
279    samples are strewn across PC1 in an additive, orthogonal manner, preserving the separation of
280    tissue types defined by PC2. To investigate the underlying cause responsible for the primary
281    source of variation in the data, we performed GO enrichment on genes with the most extreme
282    PC1 loading values that are most responsible for defining PC1. In the full *Arabidopsis* dataset
283    (**Fig. 1a**), high PC1 values, which include a small number of samples that contribute to a
284    disproportionate amount of variance in the data, are defined by high expression of genes
285    associated with response to biotic stress and oxidative damage GO terms (**Table S1**). Low PC1



**Figure 1: Principal Component Analysis (PCA) of gene expression profiles.** PCAs with gene expression profiles colored by above-ground (blue), below-ground (orange), whole plant (red), and other (green) tissue types for **a)** the full *Arabidopsis* dataset, **b)** the conserved *Arabidopsis* data set, **c)** the angiosperm dataset projected onto the conserved *Arabidopsis* PCA from b), and **d)** the same as c), but with conserved *Arabidopsis* gene expression profiles in the background (transparent gray).
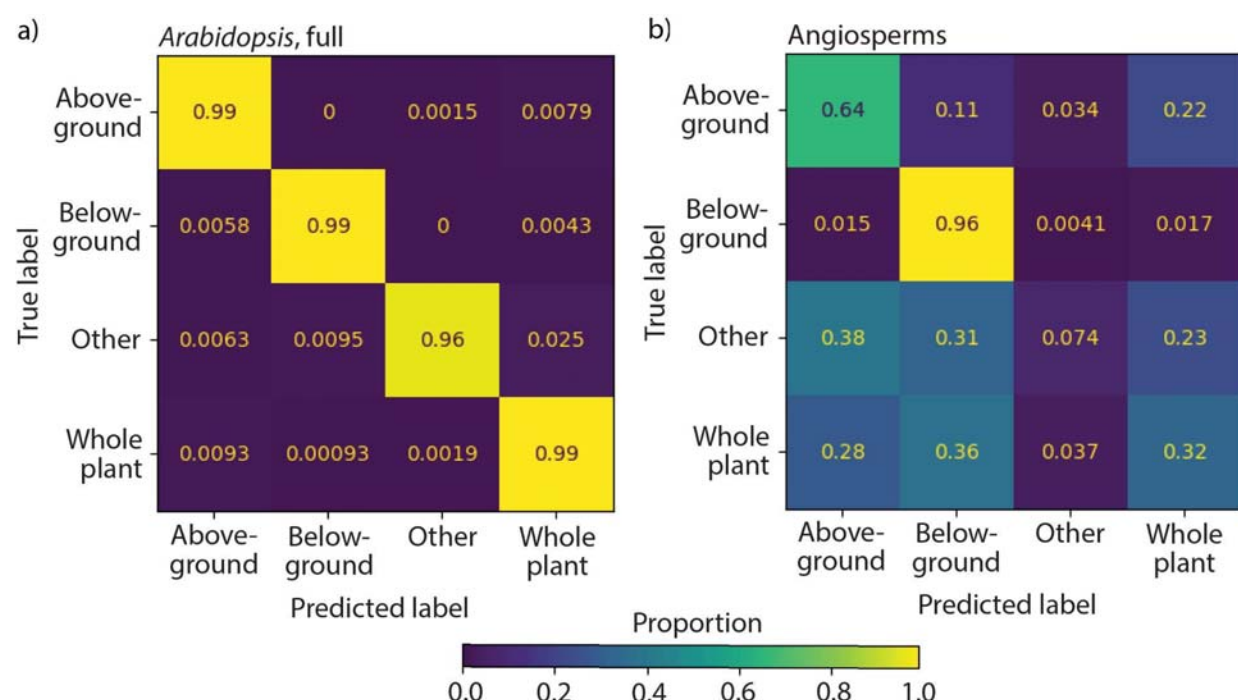
7

286   values, which include a majority of samples across tissues and which we assume arise from
287   plants grown under regular conditions associated with the less stress, are defined by high
288   expression of genes with GO terms associated with biosynthesis, biogenesis, and cell growth.
289   Remarkably, in the full *Arabidopsis* dataset, negative PC1 loading values are enriched for
290   glucosinolate biosynthetic and other metabolic processes (FDR <0.05) .
291
292   From these large-scale datasets, we developed a predictive model to test if tissue type could be
293   inferred from expression patterns alone and if this *Arabidopsis*-trained model could be
294   transferred to other flowering plants. We previously created a set of 6,328 low copy orthogroups
295   that are deeply conserved across flowering plants (Palande et al., 2023) and used a set of 6,327
296   *Arabidopsis* genes corresponding to these orthogroups for all downstream analyses. A PCA
297   performed on this subset of 6,327 conserved flowering plant genes shows mostly the same
298   structure as the analysis with all *Arabidopsis* genes included (**Fig. 1b**). However, while the
299   below-ground tissue type remains distinct from the rest of the data, the above-ground tissue
300   type overlaps more with whole plant and other tissue types. Note that the sign of principal
301   components is arbitrary, which explains the "flip" of PC2 values relative to the full set of
302   *Arabidopsis* genes. An analysis of the enriched GO terms for PC1 loading values from the
303   conserved gene PCA reveals that high PC1 values are associated with biotic responses, but
304   also with anther- and pollen-related GO terms (**Table S1**). Low PC1 values are associated
305   overwhelmingly with photosynthesis. Because the two datasets have corresponding orthogroup
306   features, we are able to project the angiosperm dataset onto the PCA defined by the conserved
307   gene *Arabidopsis* dataset (**Fig. 1c-d**). While the overall structure defining the distributions of
308   tissue types is maintained in the projected angiosperm data, there is substantial overlap
309   between above-ground and below-ground tissue types. We conclude that indeed there is
310   conservation of tissue-specific expression between *Arabidopsis* and the rest of the flowering
311   plants, but that as expected, the alignment of the underlying structures of gene expression
312   patterns defining tissue type identity are not identical.
313
314   *Predictive modeling of plant tissue from gene expression*
315
316   We used supervised learning classifiers to test if gene expression profiles could predict tissue
317   type in *Arabidopsis* and if these *Arabidopsis* trained models could be applied more broadly to
318   flowering plants. We first split the *Arabidopsis* data into testing and training sets with samples
319   split into four classes of above-ground, below-ground, whole-plant, or other as described above.
320   Models trained on *Arabidopsis* expression data and used to predict tissue type in *Arabidopsis*,
321   whether the full or conserved gene datasets, achieved high precision and recall scores. The
322   highest f1-scores (the harmonic mean of precision and recall) for the full and conserved
323   datasets were achieved using a K-Nearest Neighbors algorithm (KNN) (0.99 and 0.99,
324   respectively; **Tables 1 and 2**) and the lowest using Linear Support Vector Classification (SVC)
325   (0.78 and 0.75). Histogram-Based Gradient Boosting (HGB) also achieved high f1-scores (0.98
326   and 0.97) while the results for Random Forest (RF) (0.83 and 0.86) and Multilayer Perceptron
327   (MLP) (0.83 and 0.82) were intermediate. When used to predict *Arabidopsis* data, the precision
328   and recall values for each model were similar to each other, indicating similar positive prediction
329   value (precision, true positives divided by true positives and false positives) and sensitivity

330    (recall, true positives divided by true positives and false negatives). The relative prediction rates
331    of different tissue types to each other were equivalent for the full *Arabidopsis* dataset (**Fig. 2a**).



**Figure 2: Confusion matrices using the KNN-classifier.** Confusion matrices showing true label identity (vertical axis) and the proportion of samples assigned to predicted label identities (horizontal axis) for **a)** the full *Arabidopsis* dataset and **b)** the angiosperm dataset. Proportion indicated by viridis color scale.

332
333    The projection of gene expression patterns from across flowering plants onto a PCA using a
334    conserved set of genes from *Arabidopsis* shows considerable variability (**Fig. 1c-d**). Using
335    models trained on *Arabidopsis* data and tested on flowering plants, prediction rates are more
336    similar to each other using different algorithms than *Arabidopsis* alone but perform much worse,
337    and with higher precision than recall rates (**Table 2**). For KNN, HGB, RF, MLP, and SVC
338    methods, precision values were 0.73, 0.74, 0.75, 0.73, and 0.70, respectively, whereas the rates
339    of recall were 0.64, 0.57, 0.57, 0.55, and 0.58. Although these rates are moderately high, they
340    must be interpreted in the context of using only four tissue type labels. The relatively higher
341    precision rates compared to recall indicate that when a sample is retrieved, there is a higher
342    rate of the models calling a true positive (positive prediction value) compared to the fraction of
343    relevant samples retrieved (sensitivity). The prediction rates across tissue types were not evenly
344    distributed (**Fig. 2b**). Below-ground tissue was accurately classified, at a rate of 0.96, while
345    above-ground tissue was only correctly predicted at a rate of 0.64. Other and whole plant tissue
346    types were classified poorly (0.074 and 0.32, respectively), and almost no samples were
347    predicted as other tissue type, including other samples themselves. Although the prediction
348    accuracy varies considerably across plant families (**Fig. 3**), from around 0.4 to 0.8, we could not
349    identify any phylogenetic signal or find any support that prediction of tissue identity is inversely
350    correlated with distance of a plant family from *Arabidopsis* in the Brassicaceae.
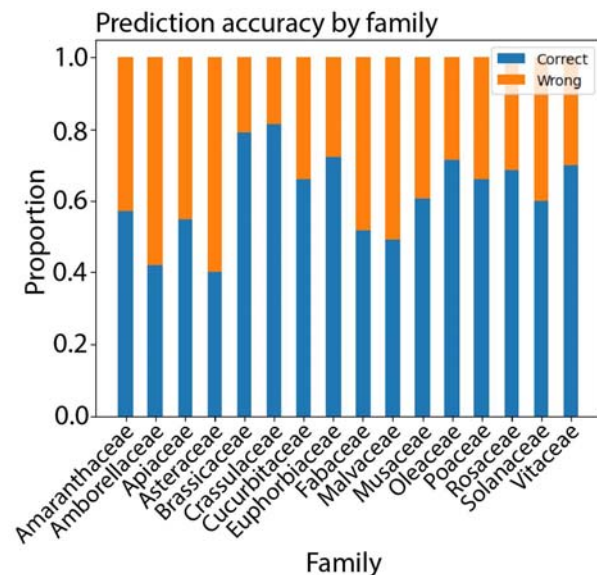
9

351

## DISCUSSION

353

*Arabidopsis-only models are highly accurate*

355

Although we focus on tissue identity in this study, we note that the strongest source of variance (PC1) in publicly available *Arabidopsis* gene expression profiles is a signature associated with biotic defense (**Table S1**) and that it acts in an additive, orthogonal manner with respect to tissue type which is the next strongest source of variance (PC2). Not only are higher prediction rates expected for the *Arabidopsis*-only models because the same dataset is being used for training and testing, but because of the data structure itself that separates the main factors we are testing—above and below ground tissues—as visualized in a PCA (**Fig. 1a-b**). From this perspective, it is perhaps not surprising that KNN is the best performing algorithm, based on the overall



**Figure 3: Prediction accuracy by plant family.** Using KNN-classifier on the angiosperm dataset, the proportion of samples correctly (blue) and wrongly (orange) predicted is shown as a stacked bar plot.

distance-based proximity of gene expression profiles for each label to each other (**Table 1**). The other methods, based on decision trees or neural networks, by focusing on individual gene expression values as parameters, fail to account for overall distance. The focus on individual gene expression values instead of the overall signature or profile is reminiscent of the molecular biology concept of "biomarkers" to indicate the tissue or stress from which a sample arises. The outperformance of KNN over other algorithms we tested may suggest that gene expression signatures (rather than focusing on individual gene expression values) are more valuable to create models for tissue and cell type prediction.

Arabidopsis *gene expression as a model for other flowering plants may not be the most suitable approach*

384

Lower prediction rates are expected when testing a model on different data than its training set (**Table 2**). However, the lower precision and recall scores when a model trained on *Arabidopsis* is tested on gene expression samples across the flowering plants undermines the foundational argument for using model species: that data from *Arabidopsis* would be predictive for plants in general. This is not to say that there is not substantial conservation of tissue-specific gene expression patterns. Our own work (Palande et al., 2023) and that of others (Julca et al., 2021) strongly supports conserved tissue-specific gene expression patterns across flowering plants, as is true of animals as well (Fukushima and Pollock, 2020). Rather, the ability to leverage and predict tissue identity from conserved gene expression profiles is diminished when building a model from a single, arbitrary species.

395

396    Details of the performance of our model hint at underlying biological considerations when using

397    model species data. Not all tissue types are equally predictable, and the prediction of below-

398    ground tissue outperforms other tissue types (**Fig. 2**). We hypothesized that the ability to predict

399    tissue identity from *Arabidopsis* may be inversely correlated with phylogenetic distance of a

400    sample from Brassicaceae, but we found no evidence to support this idea (**Fig. 3**). Additionally,

401    the precision values for predicting tissue type of flowering plant data from *Arabidopsis* are much

402    higher than recall values (**Table 2**). This may indicate that models are relatively better at calling

403    samples with conserved tissue-specificity with *Arabidopsis* (a true positive) over those without (a

404    false negative). These results may also be a product of our classification scheme. For example,

405    above and whole plant tissues are often more similar to each other than below ground tissue

406    because they are missing roots, and might more easily be misclassified with each other. The

407    other category is composed of diverse tissues which may not have clear predictive features.

408    These factors should be considered when evaluating the classification results (**Fig. 2**).

409

410    Our results potentially arise not only from genes with evolutionary differences in tissue-specific

411    expression compared to *Arabidopsis*, but ones that may indeed have conserved expression but

412    differ in the ways we have culturally constructed our developmental descriptions of plant

413    species. Such a circumstance might arise when the cell type-specific expression of a gene is

414    truly conserved, but that evolved differences in functional morphology between species lead us

415    to apply different tissue descriptors (for example, between an herbaceous annual and a woody

416    perennial, or a CAM succulent compared to a weedy C3 plant). The misalignment of tissue

417    labels extends to more quantitative descriptors and to the molecular level, including Gene

418    Ontology (GO) and Kyoto Encyclopedia of Genes and Genomes (KEGG) terms that ultimately

419    become biased to plants with sequenced genomes (Provart et al., 2016). For example, in our

420    analysis of genes corresponding to the most positive and most negative PC1 loading values,

421    there was a noticeable enrichment of genes associated with the glucosinolate biosynthetic and

422    metabolic pathways in *Arabidopsis* samples (**Table S1**). However, this enrichment was absent

423    in broader angiosperm samples, as these compounds are found almost exclusively in

424    Brassicaceae. Glucosinolates are a diverse group of secondary metabolites that play a critical

425    role in plant defense against herbivores and pathogens. Beyond their defensive role, they seem

426    to be involved in growth, development, microbiota interactions, and phosphate nutrition

427    (Kopriva, 2021). Focusing on a single organism, or small group of model species to predict

428    attributes of all plants is flawed from both biological (arising from evolutionary novelty) as well as

429    philosophical (due to semantic, ontological, and cultural differences in how we socially construct

430    plants) perspectives.

431

432    *Moving forward and embracing plant and cultural diversity*

433

434    *Arabidopsis* was selected as a model species unilaterally, over raised objections, decades ago

435    arising from mostly genetic and molecular biology considerations (Meyerowitz, 1987; Clough

436    and Bent, 1998; Arabidopsis Genome Initiative, 2000; Bennett et al., 2003; Bevan and Walsh,

437    2005). Arguments in favor of plant diversity or selecting agricultural or ecological models were

438    ignored. These past decisions have led to continued focus on *Arabidopsis* and there is

439 continuing advocacy for a plant model species and to fund *Arabidopsis* research at the expense
440 of plant diversity to this current day (Provart et al., 2016; Parry et al., 2020). Since then, data
441 science and computational approaches have begun to grow. Retrospectively, after which
442 decades of sequencing data across flowering plants has allowed us to objectively ask if the
443 focus on a single, arbitrary plant allows us to predict the biology of other flowering plants better
444 than if we had studied all plants equally from the start, the answer is no (**Table 2**). Using a data
445 science approach and building machine learning models on *Arabidopsis* gene expression data
446 to predict the tissue identity of gene expression samples from across flowering plants as we
447 have done here, does not preclude the consideration of other, more important qualitative
448 arguments against the model species concept that continue to limit the potential of the plant
449 science community. Beyond just *Arabidopsis*, there is still a focus on agriculturally important
450 species at the expense of all plants (Marks et al., 2023). More insidiously, the social construct of
451 plants and their diversity arises from colonialism, evidenced not only by the gaze of the Global
452 North and the plants we have chosen to research and document and how we do so, but in ways
453 that can be quantified related to the specific discussion of *Arabidopsis* here, specifically which
454 plant genomes have been sequenced and by whom (Marks et al., 2021), usually through
455 extinguishing and stealing the cultural knowledge of Indigenous people (Dwer et al., 2022).
456
457 Useful discoveries and insights have arisen from *Arabidopsis* (Arabidopsis Genome Initiative,
458 2000). Rather than advocating for continued focus and funding for a single model species
459 (Provart et al., 2016; Parry et al., 2020), it is long past due that we address the historical
460 inequities that have led to our current construction of the plant sciences and that we avoid a
461 biased focus and embrace the biological and cultural diversity of the plant world.
462

475

476 **REFERENCES**
477

478 Angiosperm Phylogeny Group, Chase, M.W., Christenhusz, M.J., Fay, M.F., Byng, J.W., Judd,
479 W.S., Soltis, D.E., Mabberley, D.J., Sennikov, A.N., Soltis, P.S. and Stevens, P.F., 2016. An
480 update of the Angiosperm Phylogeny Group classification for the orders and families of
481 flowering plants: APG IV. *Botanical Journal of the Linnean Society*, *181*(1), pp.1-20.
482

483  Arabidopsis Genome Initiative, 2000. Analysis of the genome sequence of the flowering plant
484  Arabidopsis thaliana. *Nature*, *408*(6814), pp.796-815.

485

486  Azodi, C.B., Pardo, J., VanBuren, R., de Los Campos, G. and Shiu, S.H., 2020. Transcriptome-
487  based prediction of complex traits in maize. *The Plant Cell, 32*(1), pp.139-151.

488

489  Bennett, M.D., Leitch, I.J., Price, H.J. and Johnston, J.S., 2003. Comparisons with
490  Caenorhabditis (approximately 100 Mb) and Drosophila (approximately 175 Mb) using flow
491  cytometry show genome size in Arabidopsis to be approximately 157 Mb and thus
492  approximately 25% larger than the Arabidopsis genome initiative estimate of approximately 125
493  Mb. *Annals of Botany*, *91*(5), pp.547-557.

494

495  Bergstra, J., Yamins, D. and Cox, D.D., 2013. Hyperopt: A python library for optimizing the
496  hyperparameters of machine learning algorithms. In *Proceedings of the 12th Python in Science
497  Conference* (Vol. 13, p. 20).

498

499  Bevan, M. and Walsh, S., 2005. The Arabidopsis genome: a foundation for plant research.
500  *Genome Research*, *15*(12), pp.1632-1642.

501

502  Clough, S.J. and Bent, A.F., 1998. Floral dip: a simplified method for Agrobacterium-mediated
503  transformation of Arabidopsis thaliana. *The Plant Journal, 16*(6), pp.735-743.

504

505  Coppens, F., Wuyts, N., Inzé, D. and Dhondt, S., 2017. Unlocking the potential of plant
506  phenotyping data through integration and data-driven approaches. *Current Opinion in Systems
507  Biology*, *4*, pp.58-63.

508

509  Cortes, C. and Vapnik, V., 1995. Support-vector networks. *Machine Learning*, *20*, pp.273-297.

510

511  Cover, T. and Hart, P., 1967. Nearest neighbor pattern classification. *IEEE transactions on
512  Information Theory*, *13*(1), pp.21-27.

513

514  Crossa, J., Perez, P., Hickey, J., Burgueno, J., Ornella, L., Cerón-Rojas, J., Zhang, X.,
515  Dreisigacker, S., Babu, R., Li, Y. and Bonnett, D., 2014. Genomic prediction in CIMMYT maize
516  and wheat breeding programs. *Heredity*, *112*(1), pp.48-60.

517

518  Dwyer, W., Ibe, C.N. and Rhee, S.Y., 2022. Renaming Indigenous crops and addressing
519  colonial bias in scientific language. *Trends in Plant Science*.

520

521  Fukushima, K. and Pollock, D.D., 2020. Amalgamated cross-species transcriptomes reveal
522  organ-specific propensity in gene expression evolution. *Nature Communications*, *11*(1), p.4459.

523

524  Haykin, S., 1998. *Neural networks: a comprehensive foundation*. Prentice Hall PTR.

525

526    Ho, T.K., 1995, August. Random decision forests. In *Proceedings of 3rd international*
527    *conference on document analysis and recognition* (Vol. 1, pp. 278-282). IEEE.

529    Hogeweg, P., 2011. The roots of bioinformatics in theoretical biology. *PLoS Computational*
530    *Biology*, *7*(3), p.e1002021.

532    Huang, D.W., Sherman, B.T. and Lempicki, R.A., 2009. Systematic and integrative analysis of
533    large gene lists using DAVID bioinformatics resources. *Nature Protocols*, *4*(1), pp.44-57.

535    Ij, H., 2018. Statistics versus machine learning. *Nat Methods*, *15*(4), p.233.

537    Julca, I., Ferrari, C., Flores-Tornero, M., Proost, S., Lindner, A.C., Hackenberg, D.,
538    Steinbachová, L., Michaelidis, C., Gomes Pereira, S., Misra, C.S. and Kawashima, T., 2021.
539    Comparative transcriptomic analysis reveals conserved programmes underpinning
540    organogenesis and reproduction in land plants. *Nature Plants*, *7*(8), pp.1143-1159.

542    Jumper, J., Evans, R., Pritzel, A., Green, T., Figurnov, M., Ronneberger, O., Tunyasuvunakool,
543    K., Bates, R., Žídek, A., Potapenko, A. and Bridgland, A., 2021. Highly accurate protein
544    structure prediction with AlphaFold. *Nature*, *596*(7873), pp.583-589.

546    Kopriva, S. (2021). Glucosinolates revisited—A follow-up of ABR volume 80: Glucosinolates. In
547    Advances in Botanical Research (Vol. 100, pp. 249-274). Academic Press.

549    Knapp S, Bohs L, Nee M, Spooner DM., 2004. Solanaceae—a model for linking genomics with
550    biodiversity. Comparative and Functional Genomics. Apr;5(3):285-91.

552    Li, F.W. and Harkess, A., 2018. A guide to sequence your favorite plant genomes. *Applications*
553    *in Plant Sciences*, *6*(3), p.e1030.

555    Lim, P.K., Zheng, X., Goh, J.C. and Mutwil, M., 2022. Exploiting plant transcriptomic databases:
556    resources, tools, and approaches. *Plant Communications*, p.100323.

558    Marks, R.A., Hotaling, S., Frandsen, P.B. and VanBuren, R., 2021. Representation and
559    participation across 20 years of plant genome sequencing. *Nature Plants*, *7*(12), pp.1571-1578.

561    Marks, R.A., Amézquita, E.J., Percival, S., Rougon-Cardoso, A., Chibici-Revneanu, C., Tebele,
562    S.M., Farrant, J.M., Chitwood, D.H., VanBuren, R., 2023. A critical analysis of plant science
563    literature reveals ongoing inequities. *Proc Natl Acad Sci USA*

565    Mason, L., Baxter, J., Bartlett, P. and Frean, M., 1999. Boosting algorithms as gradient descent.
566    *Advances in Neural Information Processing Systems*, *12*.

568    Mazzocchi, F., 2015. Could Big Data be the end of theory in science? A few remarks on the
569    epistemology of data-driven science. *EMBO Reports*, *16*(10), pp.1250-1255.

570

571     Meyerowitz, E.M., 1987. Arabidopsis thaliana. *Annual Review of Genetics*, *21*(1), pp.93-111.

572

573     Meyerowitz, E.M., 2001. Prehistory and history of Arabidopsis research. *Plant Physiology*,
574     *125*(1), pp.15-19.

575

576     Michael, T.P. and Jackson, S., 2013. The first 50 plant genomes. *The Plant Genome*, *6*(2).

577

578     Mitchell, C.E., Agrawal, A.A., Bever, J.D., Gilbert, G.S., Hufbauer, R.A., Klironomos, J.N.,
579     Maron, J.L., Morris, W.F., Parker, I.M., Power, A.G. and Seabloom, E.W., 2006. Biotic
580     interactions and plant invasions. *Ecology Letters*, *9*(6), pp.726-740.

581

582     Palande, S., Kaste, J.A., Roberts, M.D., Aba, K.S., Claucherty, C., Dacon, J., Doko, R.,
583     Jayakody, T.B., Jeffery, H.R., Kelly, N. and Manousidaki, A., 2023. The topological shape of
584     gene expression across the evolution of flowering plants. *PLOS Biology*.

585

586     Parry, G., Provart, N.J., Brady, S.M., Uzilday, B., Multinational Arabidopsis Steering Committee,
587     Adams, K., Araújo, W., Aubourg, S., Baginsky, S., Bakker, E. and Bärenfaller, K., 2020. Current
588     status of the multinational Arabidopsis community. *Plant Direct*, *4*(7), p.e00248.

589

590     Proost, S. and Mutwil, M., 2018. CoNekT: an open-source framework for comparative genomic
591     and transcriptomic network analyses. *Nucleic Acids Research*, *46*(W1), pp.W133-W140.

592

593     Provart, N.J., Alonso, J., Assmann, S.M., Bergmann, D., Brady, S.M., Brkljacic, J., Browse, J.,
594     Chapple, C., Colot, V., Cutler, S. and Dangl, J., 2016. 50 years of Arabidopsis research:
595     highlights and future directions. *New Phytologist*, *209*(3), pp.921-944.

596

597     Sultan, S.E., 2000. Phenotypic plasticity for plant development, function and life history. *Trends
598     in Plant Science*, *5*(12), pp.537-542.

599

600     Strable J, Scanlon MJ., 2009. Maize (Zea mays): a model organism for basic and applied
601     research in plant biology. *Cold Spring Harb Protoc.* Oct 1;10(2009):pdb-emo132.

602

603     Yu, Y., Zhang, H., Long, Y., Shu, Y. and Zhai, J., 2022. Plant public RNA☐seq database: a
604     comprehensive online database for expression analysis of~ 45 000 plant public RNA☐seq
605     libraries. *Plant Biotechnology Journal*, *20*(5), p.806.

606

607     Zhang, N., Wang, M. and Wang, N., 2002. Precision agriculture—a worldwide overview.
608     *Computers and Electronics in Agriculture*, *36*(2-3), pp.113-132.

609

610     Zhang, H., Zhang, F., Yu, Y., Feng, L.I., Jia, J., Liu, B.O., Li, B., Guo, H. and Zhai, J., 2020. A

611     comprehensive online database for exploring~ 20,000 public Arabidopsis RNA-seq libraries.

612     *Molecular Plant*, *13*(9), pp.1231-1233.

**Tables**

**Table 1: Classification performance of models trained on the full *Arabidopsis* dataset.**

| Model | Precision | Recall | f1-score |
|---|---|---|---|
| SVC | 0.765131 | 0.80103 | 0.777531 |
| MLP | 0.843599 | 0.844979 | 0.832854 |
| RF | 0.845664 | 0.826609 | 0.833746 |
| HGB | 0.976665 | 0.976481 | 0.976319 |
| KNN | 0.98921 | 0.989185 | 0.989193 |

**Table 2: Classification performance of models trained on the conserved *Arabidopsis* dataset and tested on conserved *Arabidopsis* or Angiosperm datasets.**

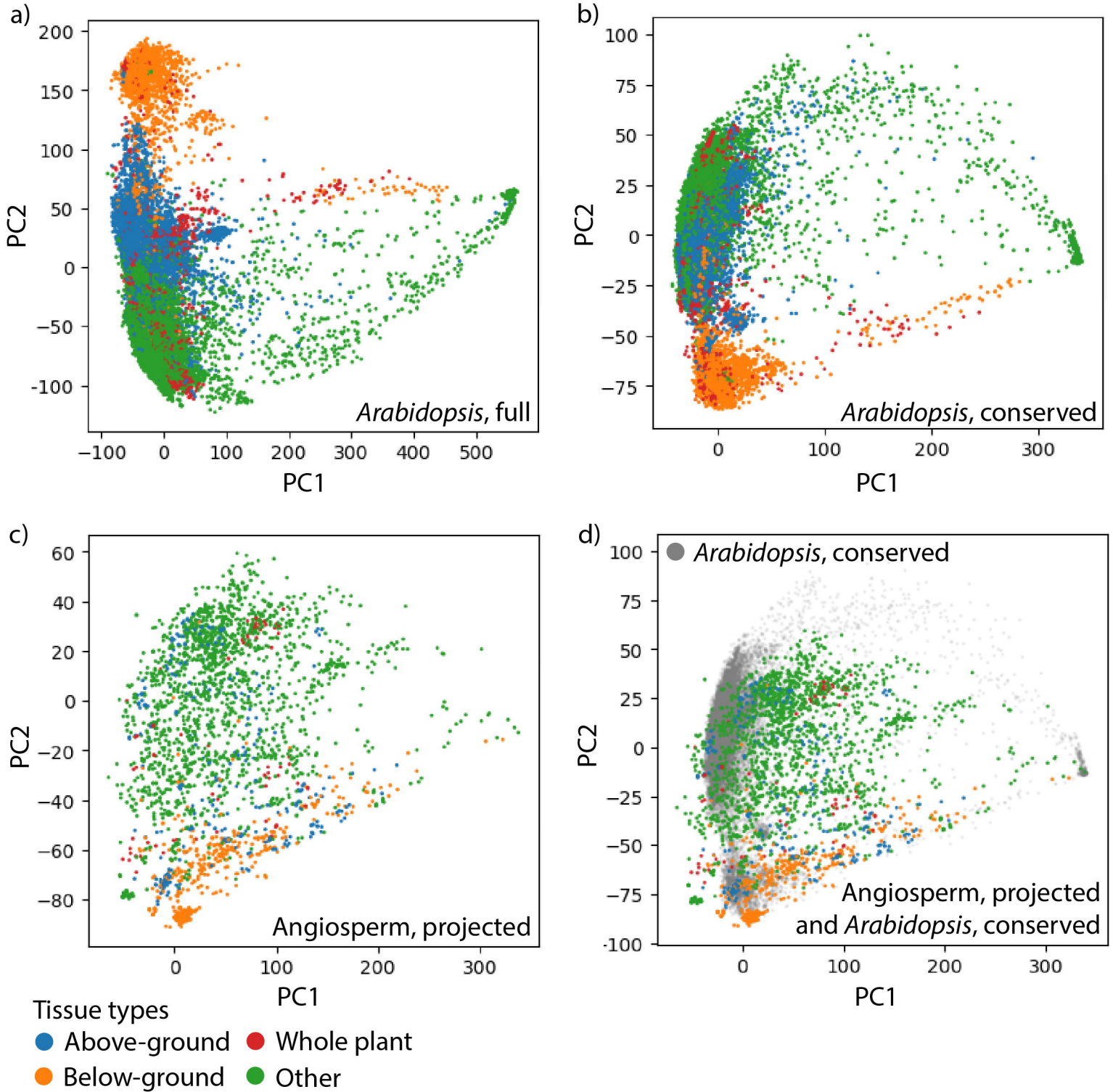| Model | Test Set | Precision | Recall | f1-score |
|---|---|---|---|---|
| SVC | Arabidopsis | 0.740855 | 0.778026 | 0.754276 |
|  | Angiosperm | 0.695691 | 0.576189 | 0.591683 |
| MLP | Arabidopsis | 0.822682 | 0.828155 | 0.824351 |
|  | Angiosperm | 0.734603 | 0.547361 | 0.611767 |
| RF | Arabidopsis | 0.862941 | 0.864721 | 0.861927 |
|  | Angiosperm | 0.747272 | 0.569075 | 0.622122 |
| HGB | Arabidopsis | 0.971034 | 0.970987 | 0.970574 |
|  | Angiosperm | 0.741902 | 0.567952 | 0.640741 |
| KNN | Arabidopsis | 0.987804 | 0.987811 | 0.987803 |
|  | Angiosperm | 0.733478 | 0.643205 | 0.663313 |

**Figure Legends**

**Figure 1: Principal Component Analysis (PCA) of gene expression profiles.** PCAs with gene expression profiles colored by above-ground (blue), below-ground (orange), whole plant
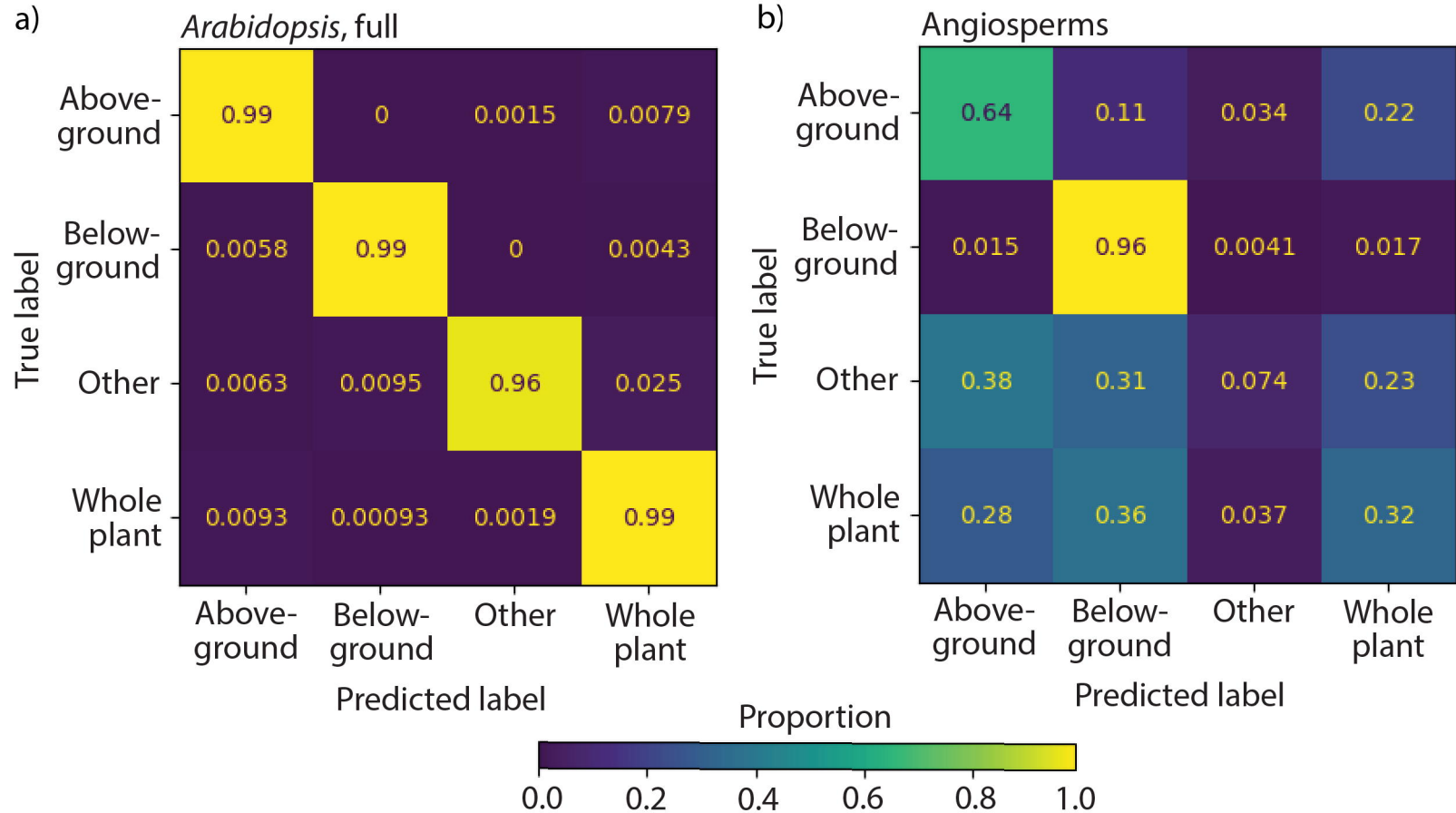
16

626    (red), and other (green) tissue types for **a)** the full *Arabidopsis* dataset, **b)** the conserved

627    *Arabidopsis* data set, **c)** the angiosperm dataset projected onto the conserved *Arabidopsis* PCA

628    from b), and **d)** the same as c), but with conserved *Arabidopsis* gene expression profiles in the

629    background (transparent gray).

630

631    **Figure 2: Confusion matrices using the KNN-classifier.** Confusion matrices showing true

632    label identity (vertical axis) and the proportion of samples assigned to predicted label identities

633    (horizontal axis) for **a)** the full *Arabidopsis* dataset and **b)** the angiosperm dataset. Proportion

634    indicated by viridis color scale.

635

636    **Figure 3: Prediction accuracy by plant family.** Using KNN-classifier on the angiosperm

637    dataset, the proportion of samples correctly (blue) and wrongly (orange) predicted from

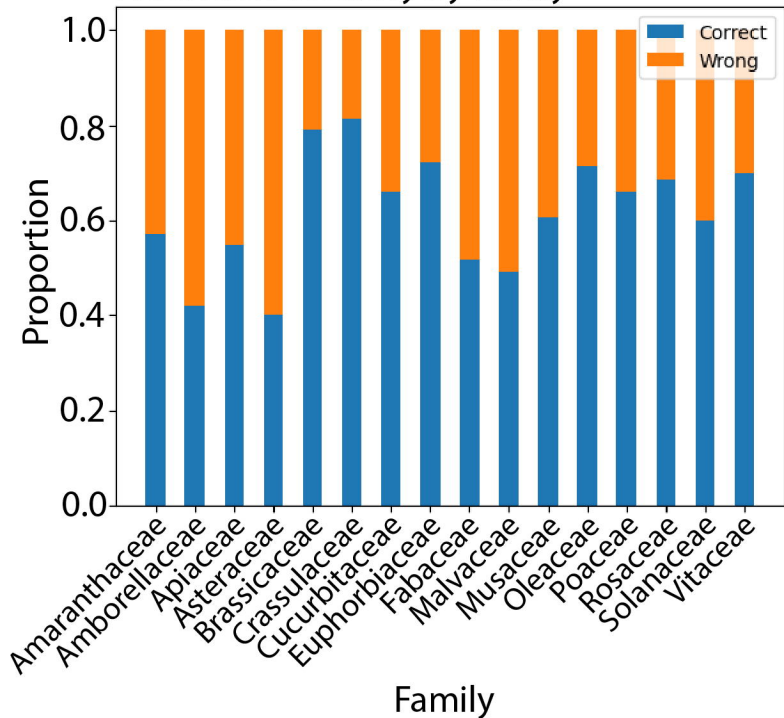638    *Arabidopsis* data is shown as a stacked bar plot.

**Figure 1: Principal Component Analysis (PCA) of gene expression profiles.** PCAs with gene expression profiles colored by above-ground (blue), below-ground (orange), whole plant (red), and other (green) tissue types for **a)** the full *Arabidopsis* dataset, **b)** the conserved *Arabidopsis* data set, **c)** the angiosperm dataset projected onto the conserved *Arabidopsis* PCA from b), and **d)** the same as c), but with conserved *Arabidopsis* gene expression profiles in the background (transparent gray).

**Figure 2: Confusion matrices using the KNN-classifier.** Confusion matrices showing true label identity (vertical axis) and the proportion of samples assigned to predicted label identities (horizontal axis) for **a)** the full *Arabidopsis* dataset and **b)** the angiosperm dataset. Proportion indicated by viridis color scale.

**Figure 3: Prediction accuracy by plant family.** Using KNN-classifier on the angiosperm dataset, the proportion of samples correctly (blue) and wrongly (orange) predicted is shown as a stacked bar plot.