

1 New Phytologist Supporting Information
2 Article title: Identification of new marker genes from plant single-cell RNA-seq data using
3 interpretable machine learning methods
4 Authors: Haidong Yan, Jiyoung Lee, Qi Song, Qi Li, John Schiefelbein, Bingyu Zhao, Song Li
5 Article acceptance date: Feb-06-2022

1 **Methods S1.** Additional materials and methods for this study.

2 **Explanation of the machine learning methods used in this study**

3 The machine learning methods for selecting marker genes were based on distinct underlying
4 computational models. In this section, we will briefly summarize and explain the mechanisms of
5 each type of model with non-technical language for readers who are not familiar with the technical
6 details of these machine learning models. We also illustrated these methods using **Figure S1C**.
7 Classification of two cell types was used as our example to explain the general principles of these
8 machine learning methods. More details and exact parameters used in our analyses are provided
9 as a separate section titled “Technical details for the machine learning methods used in this study”
10 in this document.

11
12 **KNN** is a baseline machine learning model to assign cell types (**Figure S1C**, KNN with $K=4$). For
13 each cell, the gene expression similarities between this cell and each individual other cells were
14 calculated based on all genes that were used in the analysis. The K most similar cells were then
15 selected and the cell types of these K nearest neighboring cells were used to assign cell type of the
16 original cell by majority vote. Specifically, for one cell, if majority of the neighboring cells are of
17 type A, then this cell is assigned as type A. Because all genes were used to determine similarity
18 between cell types, we did not use KNN to select marker genes. KNN is fast and simple, which
19 makes it a first choice of machine learning classifier in many cases when computation resource is
20 limited.

21
22 **PCA** is an unsupervised method to group cells (**Figure S1C**, PCA). For all cells, each gene’s
23 expression level was treated as a feature. These cells were projected into a multi-dimensional space
24 where first and second PCA dimensions were typically plotted to demonstrate the grouping of the
25 cells. The dimensions of PCA were based on decreasing variations explained by the data where
26 the first dimension had most variation in the data followed the second dimension. The contribution
27 of each gene to each of the PCA dimension can be extracted by a loading factor and higher absolute
28 loading represented higher impact on the PCA dimension. Therefore, the loadings were used to
29 rank genes and select marker genes.

30
31 **SVM** is based on an optimization method that can find best-separating-hyperplane between
32 different cell clusters (**Figure S1C**, SVM). In SVM, each gene was represented by one dimension.
33 If we had only two genes in all cell types, we could generate a x-y scatter plot with each dot
34 represents one cell and x axis is the expression of gene 1 and y axis is the expression of gene 2.
35 The SVM method, when applied in this hypothetical situation of 2 genes, is essentially a method
36 to find a line that best separate two cell types in this two-dimensional space. In single cell data, we
37 typically have a few thousands of dimensions and the separating line between two cell types
38 becomes a “hyperplane”. The weight of each gene explained how much each gene is contributing
39 to the separating hyperplane and genes with higher weights would be selected as marker genes.
40 This way of classification also partly explains why SVM markers do not have good correlations
41 because these genes are important in determining the boundaries between clusters of cells. The
42 highly correlated genes with specific cell types are more similar to the expression in the centers of
43 each cell clusters.

44

1 **RF** is an ensemble, tree-based methods for classification (**Figure S1C**, RF). The basic component
2 of a RF model is a decision tree. To decide whether a cell is type A or type B, a decision tree would
3 evaluate every single gene to decide a threshold in the expression level for best separating cells
4 into these two classes. For example, in figure S1C, gene G1 was used to make the first decision.
5 After the first step, the cells were split into two clusters based on the single gene expression and
6 associated threshold that best separate two cell types. Each sub cluster was then divided based on
7 the second-best gene (G2). The decision tree could grow to a very large tree with many genes used
8 to make this decision until a stopping criterion was met. RF is a method based on a large number
9 of decision trees and trained on bootstrapped input data. The consensus of all the decision trees
10 was used as a trained model to make predictions. **SHAP** is one of the latest approaches to select
11 features in decision trees that help to improve the interpretability of the tree. The idea of SHAP is,
12 for every cell, the method would calculate how much each gene contributes to the prediction of
13 the cell type. The SHAP value was calculated by summing up the loss of prediction power if the
14 marker gene were excluded from the model and a permutation of all possible combinations of other
15 genes used in the prediction. This is challenging to evaluate explicitly but a computational
16 approximation was applied in our manuscript based on a published Python package (Pedregosa et
17 al., 2011).

18
19 **BNN** is a baseline neural network where three fully connected neural network layers were used
20 (**Figure S1C**, BNN). At each neuron, gene expression from all genes were used as inputs, and an
21 output of this neuron was calculated based on linear regression followed by an activation function
22 (ReLU in our model). At the first layer of the network, 586 neurons were used, thus the gene
23 expression from >20,000 genes were converted into 586 neurons at the first layer. These were
24 converted into a second layer of the neural network by a similar process and then converted into a
25 third layer of the neural network. The final output layer will have two neurons if there were two
26 cell types or multiple neurons based on the number of cell types to be predicted. This neural
27 network model was “trained” using labeled cell types from input data and the weights on each
28 neuron were determined by an algorithm called back propagation. In this training process, the
29 changes in the “goodness of fit” to the labeled data were used to update the weights on the neurons
30 at each layers sequentially in a reversed order.

31
32 **Triplet NN** and **Contrastive NN** are two more complexed neural network architectures as
33 compared to **BNN** (**Figure S1C**, TRINN and CTNN). In BNN, the gene expression from each cell
34 was used to decide whether one cell is in type A or type B. The TRINN and CTNN did not directly
35 assign cell types. In contrast, these two networks were used to learn the similarity between different
36 cells. This is considered as manifold learning in general. The advantage of manifold learning as
37 compared to the BNN was that manifold learning can help to identify rare cell types because the
38 model was trained to learn “distance”, not cell identity. For rare cell types, it is challenging to find
39 enough training data to train classifiers to learn the signatures of different cells. In CTNN, the
40 models were trained by input data where cells from the same cell types would have higher weight
41 if they were predicted to be more similar than cells from different cell types. In TRINN, the models
42 were trained such that each training were evaluated using three cells, two from the same cell type
43 and one from a different cell type. The model parameters were optimized such that cells from same
44 cell types should be classified as similar and simultaneously, cells from different cell types had to
45 be distant from the two cells of the same cell type. In contrast to CTNN, TRINN encourages the
46 model to distinguish different cells while maintaining similarity between similar cells.

1 **Technical details for the machine learning methods used in this study**

2 **KNN.** KNN is a commonly used simple classifier that does not have explicit training process.
3 KNN first computed a distance between the new input vector and every feature vector in the
4 training dataset. Then the top K nearest neighbors were used for new prediction. In the last step,
5 class label of the new input vector was determined by majority vote among the K nearest
6 neighbors. The hyperparameter for KNN is K, the number of top nearest neighbors. K was set to
7 be 50 in our analysis. We also set weights as ‘uniform’ that means all points in each neighborhood
8 were weighted equally and set p to be 2 that means Euclidean distance. All other hyperparameters
9 were set as default.

10 **RF.** RF is an ensemble tree-based machine learning approach. For each decision tree, a subset
11 of training examples was randomly sampled as inputs and a subset of features were randomly
12 sampled to split each tree node. The final class label was determined by majority vote. Number
13 of trees (N) for RF was set as 50 in our analysis.

14 **SVM.** SVM is a machine learning classifier that maximizes the margin between different
15 classes in a high dimensional space transformed by a kernel function. Depending on the kernel
16 function, SVM can be a linear classifier (linear kernel) or a non-linear classifier (e.g., Gaussian
17 kernel). To be able to extract interpretable feature weights, linear kernel was used in our analysis
18 to train SVM classifier.

19 **Baseline NN.** Baseline NN refers to a basic type of neural network that uses densely connected
20 layer as input layer and hidden layers. The output layer has number of neurons equal to number of
21 cell types (ten cell types). Architecture of the base NN is demonstrated in **Figure S23B**. Briefly,
22 input layer has number of neurons equal to number of genes used for classification and three hidden
23 layers were used, of which each has 586, 256, and 100 neurons. The last layer is an output layer to
24 which a softmax is applied to ensure output scores are summed to 1.

25 **Triplet NN.** Triplet NN is the implementation of Siamese neural network with triplet loss
26 function. The use of triplet loss function was discussed in a published study (Alavi et al., 2018).
27 Briefly, Siamese DNN consists of two subnetworks which had identical architecture and weights.
28 The two neural networks connected to the same distance layer which computed a vector of distance
29 between the last two hidden layers in the two subnetworks. The last two hidden layers were lower
30 dimensional embeddings of original feature vectors. Architecture of Siamese NN is demonstrated
31 in **Figure S23**. In this work, number of neurons in input layer was equal to number of genes used

1 for classification (29,929). Numbers of neurons used in three hidden layers were 586, 256, and
 2 100. In training dataset, each scRNA-seq expression profile was an “anchor” that can be paired
 3 with positive example and negative example. Positive examples were those labeled with the same
 4 cell type with anchor and negative examples were those with different cell type. For each anchor,
 5 it would be paired with a positive example and a negative example, which formed a group of
 6 triplets. Then for each group of triplets, anchor-positive and anchor-negative pairs would be
 7 respectively fed into triplet NN. Based on the discussion in (Schroff et al., 2015) and (Alavi et al.,
 8 2018), the loss function of triplet NN can be written as:

$$9 \quad L(D) \max \left\{ 0, \left(\sum_{i=1}^T (D_{a,p}^i)^2 - (D_{a,n}^i)^2 + m \right) \right\}$$

10 Where T is the number of groups of triplets. $D_{a,p}^i$ is the Euclidean distance between anchor and
 11 positive samples and $D_{a,n}^i$ is the Euclidean distance between anchor and negative samples. m is a
 12 hyperparameter that represents the margin between $(D_{a,p}^i)^2$ and $(D_{a,n}^i)^2$.

13 To ensure that triplet NN can be effectively trained, the groups of triplets need to include
 14 anchor-positive pairs with large distances and anchor-negative pairs with small distances. These
 15 are the hard training examples that enforce the model to learn effectively. As discussed in Alavi’s
 16 study (2018), batch hard loss function was used to generate hard training examples. In each
 17 iteration of optimization, M cell types which had K cells in each were sampled to generate a mini-
 18 batch. In this mini-batch, losses of hard training examples were selected and summed up as final
 19 loss value for the mini-batch. A slight modification of batch hard loss function was made in this
 20 study to include more training samples in each mini-batch. Instead of using one pair of hardest
 21 anchor-positive and anchor-negative respectively for each anchor, top k pairs of hardest pairs were
 22 selected for each anchor. The batch hard loss function therefore can be written as:

$$23 \quad L'(D) = \left\{ 0, \sum_{i=1}^M \sum_{j=1}^K [topmax(k, P_j^i) - topmin(k, N_j^i) + m] \right\}$$

24 Where P_j^i is the set of distances between j th cell from i th cell type and all other cells in i th cell
 25 type (anchor-positive pairs) and N_j^i is the set of distances between j th cell from i th cell type and
 26 all other cells not from i th cell type (anchor-negative pairs). $topmax(k, P_j^i)$ selects the top k pairs
 27 with largest distances in P_j^i and sums the selected distances. $topmin(k, N_j^i)$ selects the top k pairs
 28 with smallest distances in N_j^i and sums the selected distances. This gives k pairs of anchor-positive

1 sample pairs and k pairs of anchor-negative sample pairs for each anchor. In our analysis k was
2 set as 10.

3 **Contrastive NN.** Contrastive NN is an implementation of Siamese neural network with
4 contrastive loss function (Alavi et al., 2018). In our work, contrastive NN was constructed using the
5 same neural network architecture as triplet NN (**Figure S15**). The difference here was that
6 contrastive NN uses paired samples which pair the cell assigned with same/different cell types.
7 The idea was to penalize large distances between samples of same cell type and small distances
8 between samples of different cell types. The loss function of Contrastive NN can be written as:

$$9 \quad L(Y, D) = \sum_{i=1}^P (Y^i) \frac{1}{2} (D)^2 + (1 - Y^i) \frac{1}{2} (\{0, m - D\})^2$$

10 Where P represents number of pairs of training samples. $Y^i = 1$ if two samples in the i th pair
11 are assigned with same cell type and $Y^i = 0$ if not. D is the Euclidean distance between the two
12 samples in each pair, computed using the last hidden layers of the two sub-networks. m is a
13 hyperparameter that represents the margin between two samples assigned with different cell types,
14 usually set to 1.

15

16 ***Model evaluation***

17 For the KNN, SVM, RF and baseline NN, the sub-training datasets were used to train the models
18 that can directly predict cell type label. The trained models were then used to predict cell type
19 labels for the independent testing datasets. For triplet NN, and contrastive NN, the sub-training
20 datasets were used to train models that predict neural embeddings of the original feature vectors
21 of in training dataset. For each new input vector from testing dataset, the trained models were first
22 used to predict a neural embedding and this embedding was compared to all neural embeddings of
23 the training dataset. The final cell type label was determined by majority vote of m nearest
24 neighbors. Here we set $m = 50$.

25 To further evaluate the performance of these seven models, we calculated parameters such as
26 numbers of true positive (TP), false positive (FP), true negative (TN), and false negative (FN). The
27 sensitivity (SE), accuracy (AC), specificity (SP), precision (PR), geometric mean (GM) of SE and
28 SP, and Matthews Correlation Coefficient (MCC) were used to evaluate these models. SE, SP,
29 AC, PR, GM, MCC, and F1 were defined as follows:

30 $SE = TP / (TP + FN)$;

1 $SP = TN / (FP + TN);$

2 $PR = TP / (TP + FP);$

3 $GM = \sqrt{SE \times SP};$

4 $AC = (TP + TN) / (TP + TN + FP + FN);$

5 $MCC = \frac{(TP \times TN - FP \times FN)}{\sqrt{(TP+FP)(TP+FN)(TN+FP)(TN+FN)}}$

6 $F1 = 2 \times (PR \times SE / (PR + SE))$

7 Tukey's honestly significant difference test was used as a conservative statistical test to find
8 significant differences in all pairwise comparisons and to control for family wise error rate (Abdi
9 and Williams, 2010). The Mean average precision (MAP) was also used as an evaluation metric
10 for all classification approaches. The MAP works on ranked lists (e.g. a list of nearest neighbor
11 cells in a retrieval database) by calculating the precision at exact-match cutoffs in the list, and then
12 taking the mean of these. We followed the MAP calculation in Alavi's study (Alavi et al., 2018).

13

14 ***Identification of marker genes based on machine learning approaches***

15 For the five-publication datasets, we first selected the top 20 percent of highly variable genes
16 (5,986 genes) using the Seurat package and then identified SHAP and SVMM markers from these
17 5,986 genes. The TreeExplainer in SHAP package was used to calculate feature importance in the
18 RF model (Lundberg et al., 2020). Briefly, this package calculated marginal contribution of each
19 feature of a given observation from all model combinations. For a target feature, each combination
20 contains one model with and another without this feature, and the marginal contribution can be
21 calculated based on the difference yielding between these two models. Due to its local
22 interpretability that each observation can get its own set of SHAP values, we can calculate these
23 values of each gene under each cell type. The higher SHAP value suggests higher contribution of
24 the feature to the classification. The novel marker genes assumed to have higher SHAP values than
25 other genes. The implementation of the SVM model is based on libsvm (Chang and Lin, 2011).
26 The absolute size of the coefficient relative to the other ones gave an indication of how important
27 the feature was for the separation. We assumed the absolute coefficient values represent feature
28 importance. The multiclass support of the SVM was handled according to a one-vs-one scheme.
29 The attributes coefficients had the shape: (number_of_cell_type * (number_of_cell_type -1) / 2,
30 number of features). To identify the feature importance of each gene on cell type, we calculated
31 the average absolute coefficient values from all pairs for a specific cell type for each feature. Each

1 feature had one coefficient of each cell type. The cell type with the highest coefficient was assigned
2 to the feature.

3

4 ***Identification of correlation marker genes***

5 Pearson correlation analysis was conducted between the cell expression of known marker
6 genes and other genes (Benesty et al., 2009). Each of the unknown markers had a correlation score
7 for each cell type. We ranked the ten cell types for each marker based on the correlation score. An
8 unknown marker was assigned to a cell type where this marker achieved the highest correlation
9 score in this cell type.

10

11 ***Cell clustering***

12 The integrated dataset with 25,618 cells and 25,092 genes was used for clustering analysis.
13 The top 30 aligned correlated components were used as input for UMAP dimension reduction and
14 clustering analysis. Clusters were identified using Seurat FindClusters function with default
15 settings. The DoHeatmap and DotPlot function in the Seurat was used to visualize expression
16 patterns of the novel marker genes for the identified clusters.

17 Method for SHAP markers consisting of the top 20 markers in each of ten cell types were used.
18 The same methods were used to select 200 markers in SVM and CORR. To select the top 20
19 BULR and KNOW markers of each cell type, we ranked them based on their expression specificity.
20 To calculate the marker specificity for a specific cell type, we generated a cell vector by labeling
21 cells under this cell type to 1, and all the other cells to 0. The cells under the specific cell type were
22 defined based on the ICI method. The Pearson correlation analysis was used to calculate the
23 correlation rate between the marker expression and the cell vector developed before. The higher
24 absolute correlation rate means higher marker specificity. The top 180 BULR markers in the cell
25 types were selected except for the Protophloem where no marker was detected. In the KNOW
26 markers, the top 161 of them were selected since no Meri_xylem marker and only one Protoxylem
27 marker were found. The 232 ICIM markers were all used in the comparison (**Table S4**).

28

29 ***Assign cell identity***

30 In figure 2H, expression values for each marker genes were normalized across all cells and
31 clusters using AverageExpression function in the Seurat package (Satija et al., 2015). For each of

1 the marker genes, the normalized expressions were ranked from high to low (from 1 to N with N
2 equals the number of clusters, N=17 in our dataset). Finally, the average rankings of all marker
3 genes were used to determine the cluster identity.

4

5 ***Classify cells into different developmental stages using machine learning methods***

6 We used our published single cell data (Ryu et al., 2019), and root hair, non-root hair and
7 lateral root caps cells were extracted. The selected 2,932 cells have been differentiated into nine
8 sub-populations at different development stages (LateralRootCap, Differ_LateralRootCap,
9 Early_Differ_NonHairEpiderm, Differ_NonHairEpiderm, NonHairEpiderm, RootHairEpiderm,
10 Late_Differ_RootHairEpiderm, Mid_Differ_RootHairEpiderm, Early_Differ_RootHairEpiderm)
11 in a previous analysis. This dataset was divided into independent, training, and validation datasets,
12 and five-fold cross validation was conducted described above. The RF and SVM models were
13 trained, and predictions were made on trichoblast and atrichoblast cells (11,904) extracted from
14 the other four datasets (ICI score > 0.5) as well as 1,970 ‘positive’ WER cells labeled by expressed
15 WER-AT gene. Next, we integrated the 2,932 training cells and the 13,874 prediction cells using
16 Seurat (v3.1) multicanonical correlation analysis with top 50 aligned correlated components as
17 input for UMAP dimension reduction.

18

19 ***Compare overlapping ratio with rice markers among six marker types***

20 A recent rice scRNA-seq study (Liu et al., 2021) had listed a number of candidate marker genes.
21 There were three cell types (Cortex, Endodermis, Trichoblast) in that study that exactly matched
22 three of ten cell types used in our study. To compare the ratio of cells where genes are detected in
23 these cell types among the six different marker types, we identified rice orthologs based on
24 information from Phytozome (v12.1) (Goodstein et al., 2012). The previous top 20 markers based
25 on ICI score over 0.9 of each marker type were compared. The marker type with less than three
26 rice orthologs was not considered. To compare the frequency of markers overlapping with the rice
27 set, all the six types of markers corresponding to rice orthologs were used to overlap with these
28 rice candidates. We randomly picked rice genes with the same orthologs number for each marker
29 type. This step was repeated for 100 times to calculate an average overlapping ratio. We performed
30 an exact binomial test (Wagner-Menghin, 2014) by setting this overlapping ratio as hypothesized

1 probability of success, number of overlapping markers as number of successes, number of rice
2 candidate markers as number of trials, and alternative as ‘greater’.

3
4 ***Methods for literature search and additional wet-bench experimental support for newly***
5 ***identified markers.***

6
7 To identify supporting evidence of cell type specificity of new marker genes that are identified in
8 our work, we performed literature search for all SHAP markers (200). We also searched SVM
9 and CORR markers (20 for each type) for trichoblast cell types. Because it is difficult to perform
10 complete literature search automatically, this analysis is done as a demonstration that many newly
11 identified markers have support from published, non-high throughput experiments. We used the
12 following procedure for each marker genes used in this search. We first search the TAIR website
13 using the ATxGxxxxx ID to extract all literature related to this gene ID. We then exclude the
14 publications where more than 20 gene ids were associated with a single publication. This is to
15 avoid finding evidence that were generated by high throughput methods such as RNAseq. For the
16 rest of associated publications for each gene, we check the full text for evidence of promoter
17 GFP/YFP/GUS reporters. The literature search results are provided in supplementary table S10.
18 We found that in 11 cases, the newly identified SHAP markers had published reporter genes. In 2
19 cases, the published reporter genes are not in agreement with the specificity determined by
20 SHAP/scRNA-seq data. For the rest 187 genes, there are no published evidence based on
21 promoter-reporter genes. Here is a list of genes that are supported by published literatures.

22
23 **1. Trichoblast:** AT2G21045 (HAC1). SHAP selected this gene as trichoblast marker. Promoter-
24 GFP shows this gene is expressed in epidermal layers, and in mature trichoblast which is
25 overlapping with WER-GFP expressed mature trichoblast. Fischer et al., Journal of Experimental
26 Botany, 2021. DOI: 10.1093/jxb/eraa465

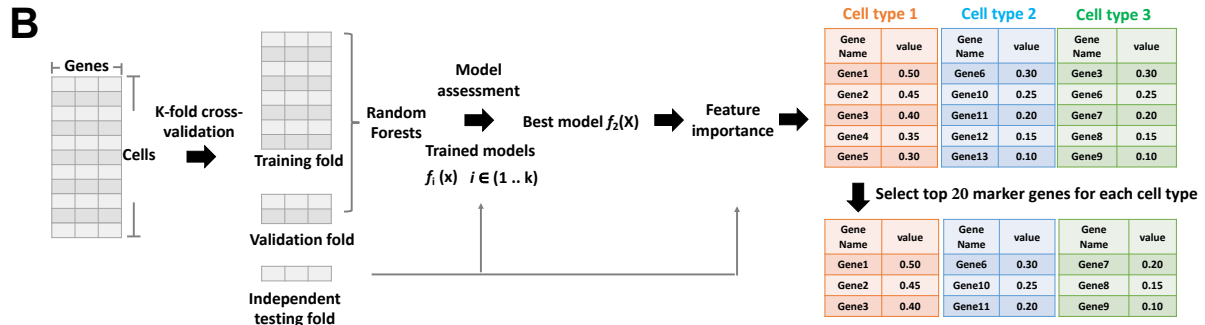
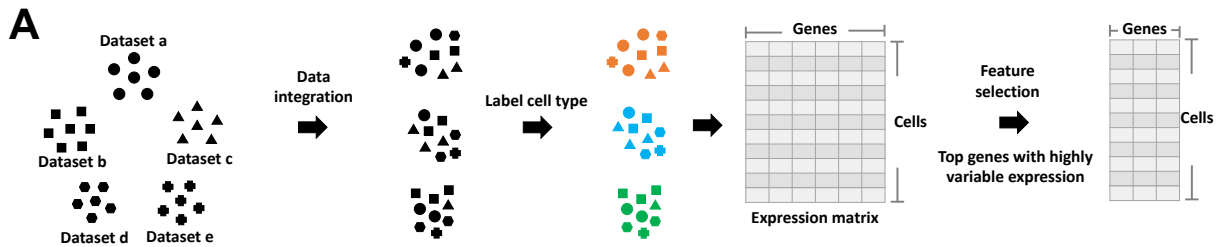
27
28 **2. Atrichoblast:** At2g41800 (DUF642). SHAP and CORR selected this gene as atrichoblast (non-
29 hair epidermal cells) marker. Figure 1 and 4 in this paper shows promoter-GFP fusion of this gene
30 is expressed in epidermal cells. Salazar-Irbe et al., Plant Science, 2016. DOI:
31 10.1016/j.plantsci.2016.10.007

32
33 **3. Endodermis:** AT2G37180 (PIP2;2). SHAP and SVM selected this gene as an endodermis
34 marker. Figure 4 of this publication shows highly specific endodermis expression of promoter-
35 GUS reporter of this gene. Javot et al., The Plant Cell, 2003. DOI: 10.1105/tpc.008888

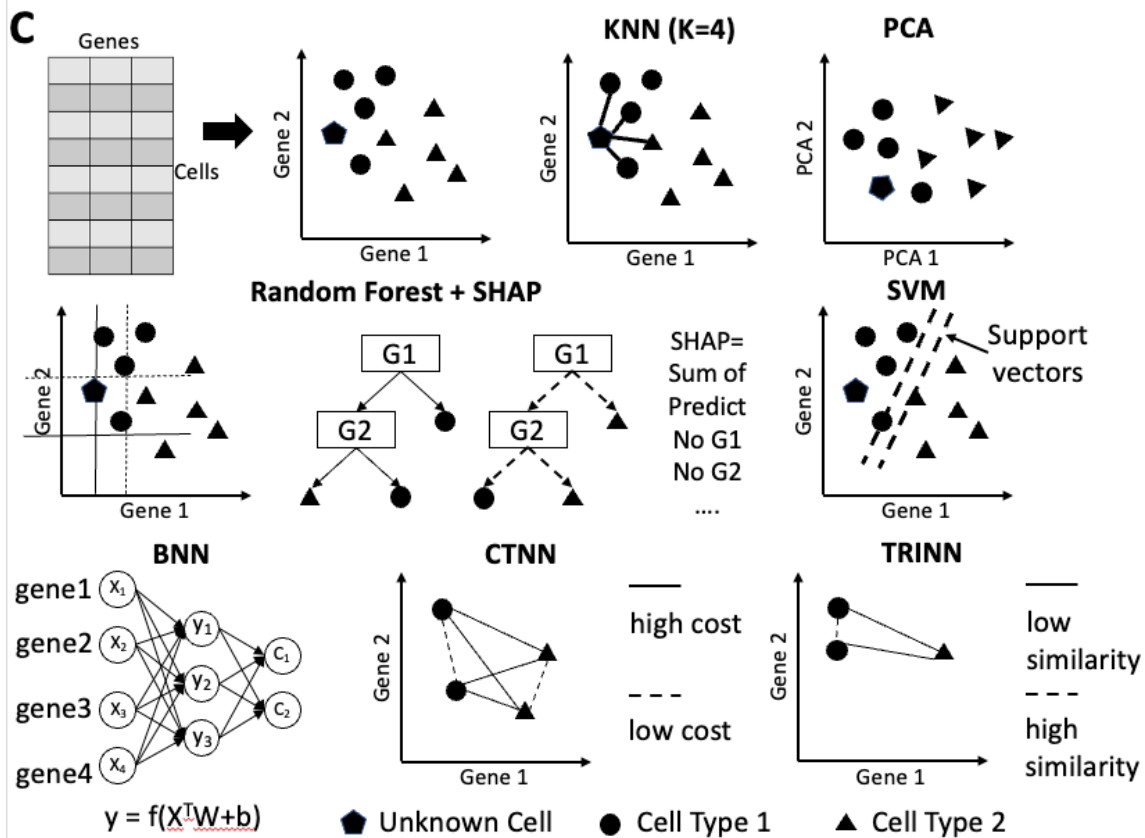
36
37 **4. Endodermis:** AT3G22600 (ATXLP12). SHAP/SVM/CORR selected this gene as an
38 endodermis marker. Figure 5 in this paper shows strong endodermis expression of this marker with

1 some expression in pericycle. Kobayashi et al., Plant and Cell Physiology, 2011. DOI:
2 10.1093/pcp/pcr060
3
4 **5. Protoxylem:** AT4G04460 (PASPA3). SHAP selected this gene as protoxylem marker. Figure
5 5C of this paper shows this gene expressed in differentiating protoxylem. Fendrych et al. Current
6 Biology, 2014. DOI: <https://doi.org/10.1016/j.cub.2014.03.025>.
7
8 **6. Protoxylem:** AT2G40320 (TBL33). SHAP selected this gene as protoxylem marker. Figure
9 1B of this paper shows this gene expressed in protoxylem. Yuan et al. PloS One. 2016. DOI:
10 <https://doi.org/10.1371/journal.pone.0146460>.
11
12 **7. Protophloem:** AT4G29920 (SMXL4). SHAP/CORR selected this gene as protophloem marker.
13 Figure 1 and 2 of this paper shows that this gene promoter is active protophloem differentiation.
14 Wallner et al., Current Biology, 2017. DOI: <https://doi.org/10.1016/j.cub.2017.03.014>.
15
16 **8. Phloem:** AT5G02600 (NAKR1). SHAP/CORR selected this gene as phloem_CC marker.
17 Figure 3F of this paper shows that the promoter-GUS of this gene is expressed in phloem.
18 Shibuta and Abe. Plant Cell Physiology, 2017. DOI: <https://doi.org/10.1093/pcp/pcx133>.
19
20 **9. Phloem:** AT3G21190 (MSR1). SHAP selected this gene as protophloem marker. Figure 4 of
21 this paper shows that the promoter-GUS of this gene is expressed in Phloem. Wang et al. The
22 Plant Journal, 2012. DOI: <https://doi.org/10.1111/tpj.12019>.
23
24 **10. Cortex:** AT2G25810 (TIP4;1) SHAP selected this gene as cortex marker. Figure 1 of this
25 paper shows that the TIP4;1 expressed in both epidermal and cortex. Gattolin et al., BMC Plant
26 Biology, 2009. DOI: 10.1186/1471-2229-9-133
27
28 **11. Cortex:** AT2G45960 (PIP1;2). SHAP selected this gene as cortex marker. Figure 1B of this
29 paper shows that the promoter-GUS of this gene is expressed in cortex, endodermis and stele.
30 Postaire et al., Plant Physiology, 2009. DOI: 10.1104/pp.109.145326
31
32
33

1
2



3



4

5 **Figure S1 Summary of the SPmarker. A. Data processing pipeline.** The different datasets are integrated together.

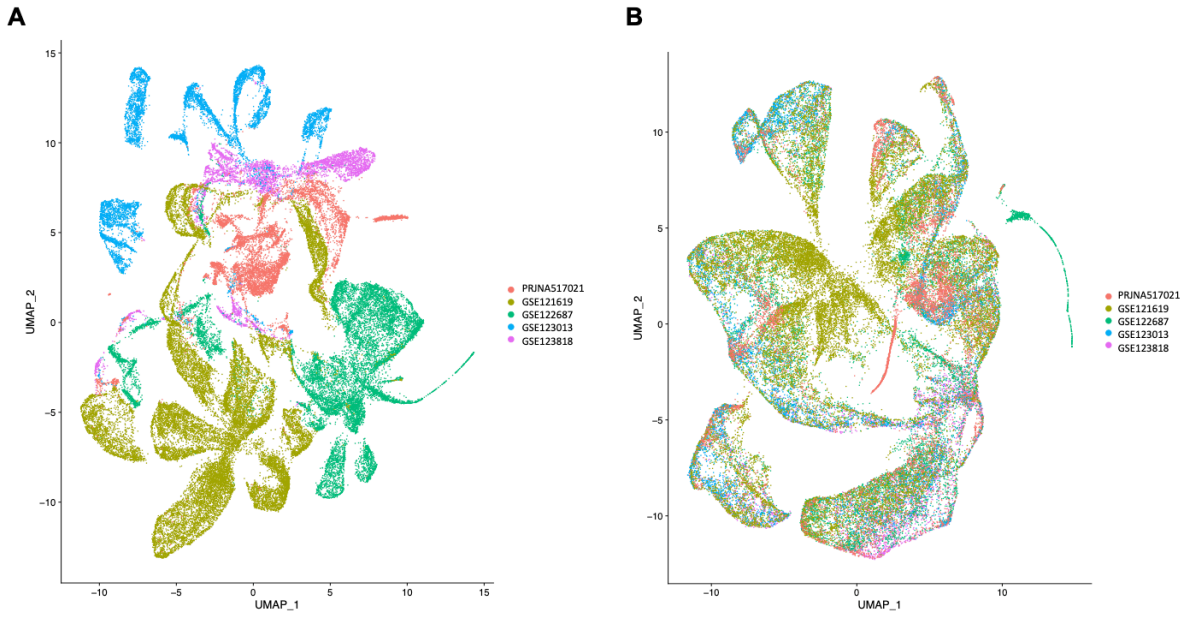
6 After labeling the cells, a gene by cell expression matrix is built. The top genes with highly variable expression are

7 selected to build a new expression matrix. **B. Model training and identification of SHAP marker genes.** The

1 integrated expression matrix was divided into the training dataset (90%) and the independent testing dataset (10%).
2 The independent testing was used to evaluate the prediction performance a $f_1(x)$ model trained with the training dataset.
3 The best model ($f_2(x)$ in this case) was selected to identify the feature importance using the SHAP method. The top
4 SHAP marker genes were selected from each cell type such that each cell type having its own marker genes that are
5 not shared with others. C. Explanations of different machine learning methods used in this study.

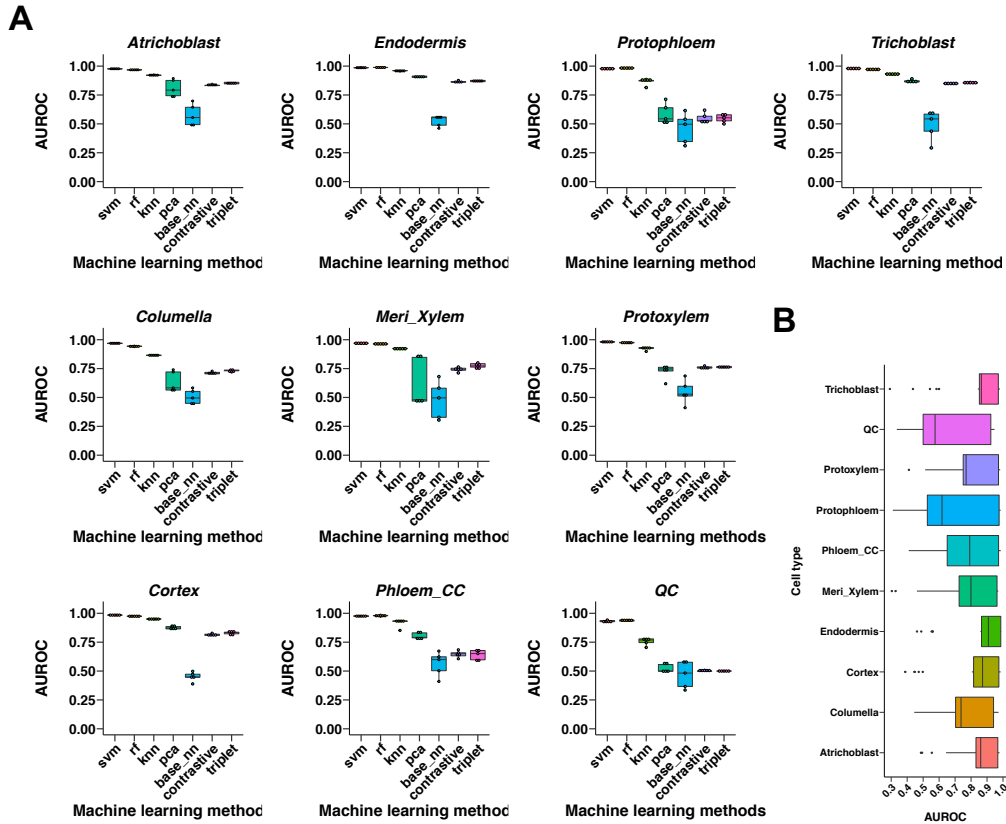
6
7

1

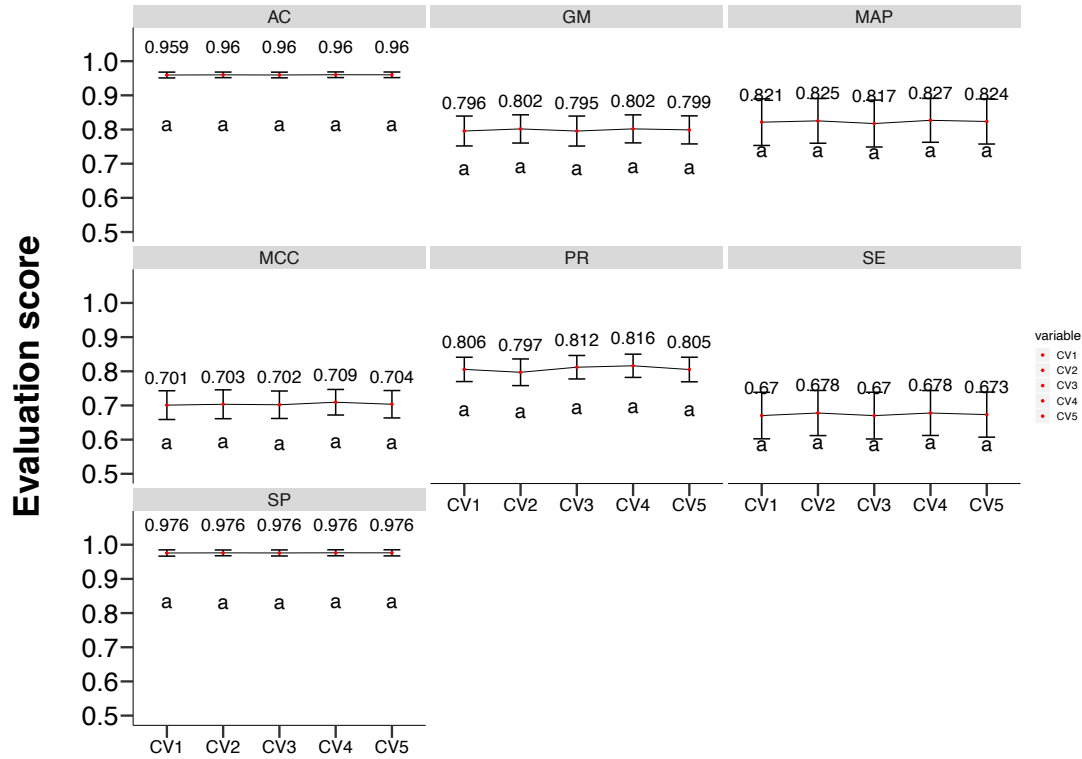


2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17

Figure S2 Integration of five datasets using the canonical correlation function in the Seurat package. A. Data before integration. B. Data after integration.



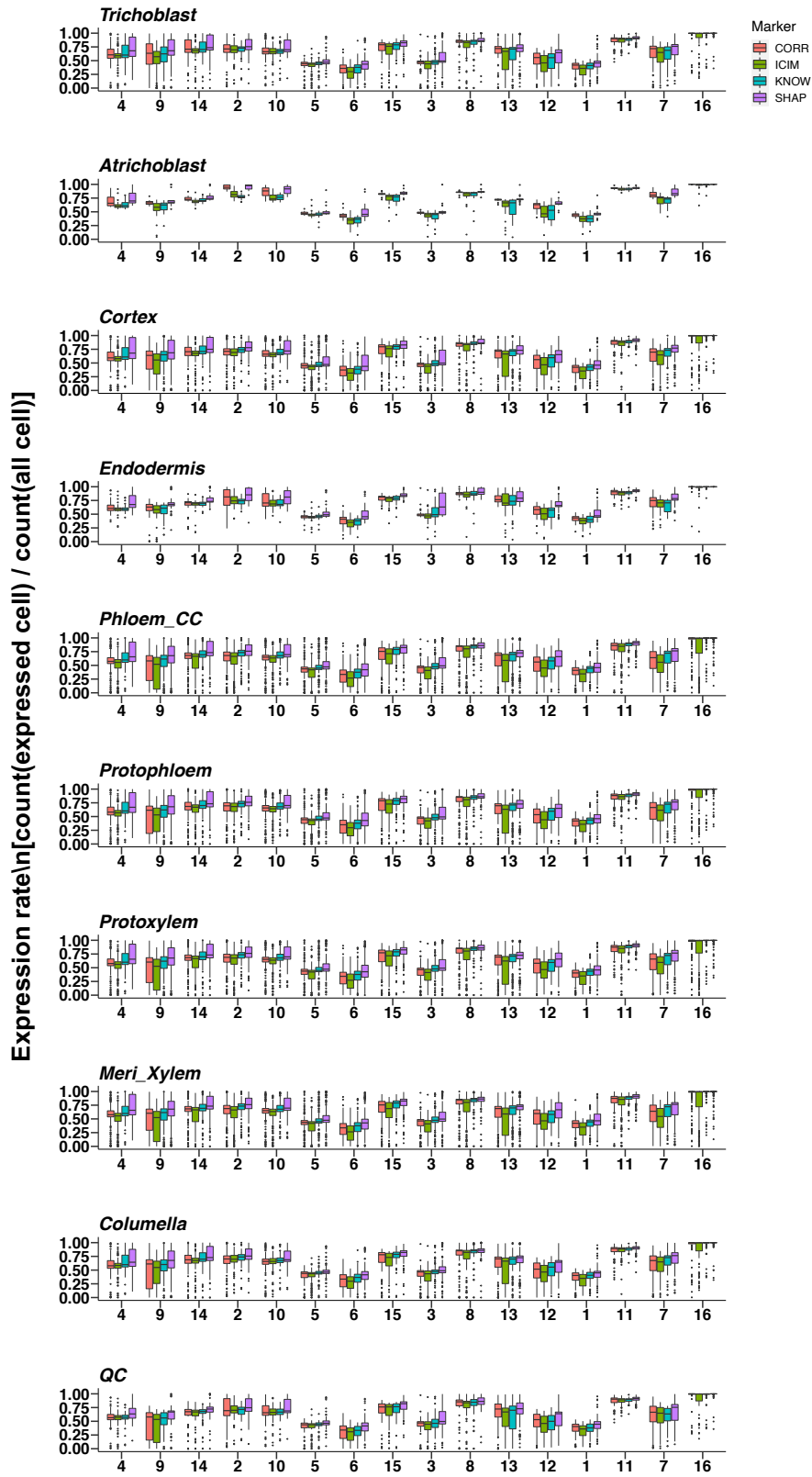
1
 2 **Figure S3 Classification performance (AUROC) of ten root cell types of *Arabidopsis*.** A. comparison of seven
 3 machine learning models on cell type classification. In these boxplots, the mid-horizontal line represents the median
 4 and dots represent data points. B. comparison of classification performance of all the ten cell types. Dots represent
 5 outliers. AUROC means Area Under the Receiver Operating Characteristics. Number of cells used in this figure is the
 6 same as shown in **Figure 1A**.



Software

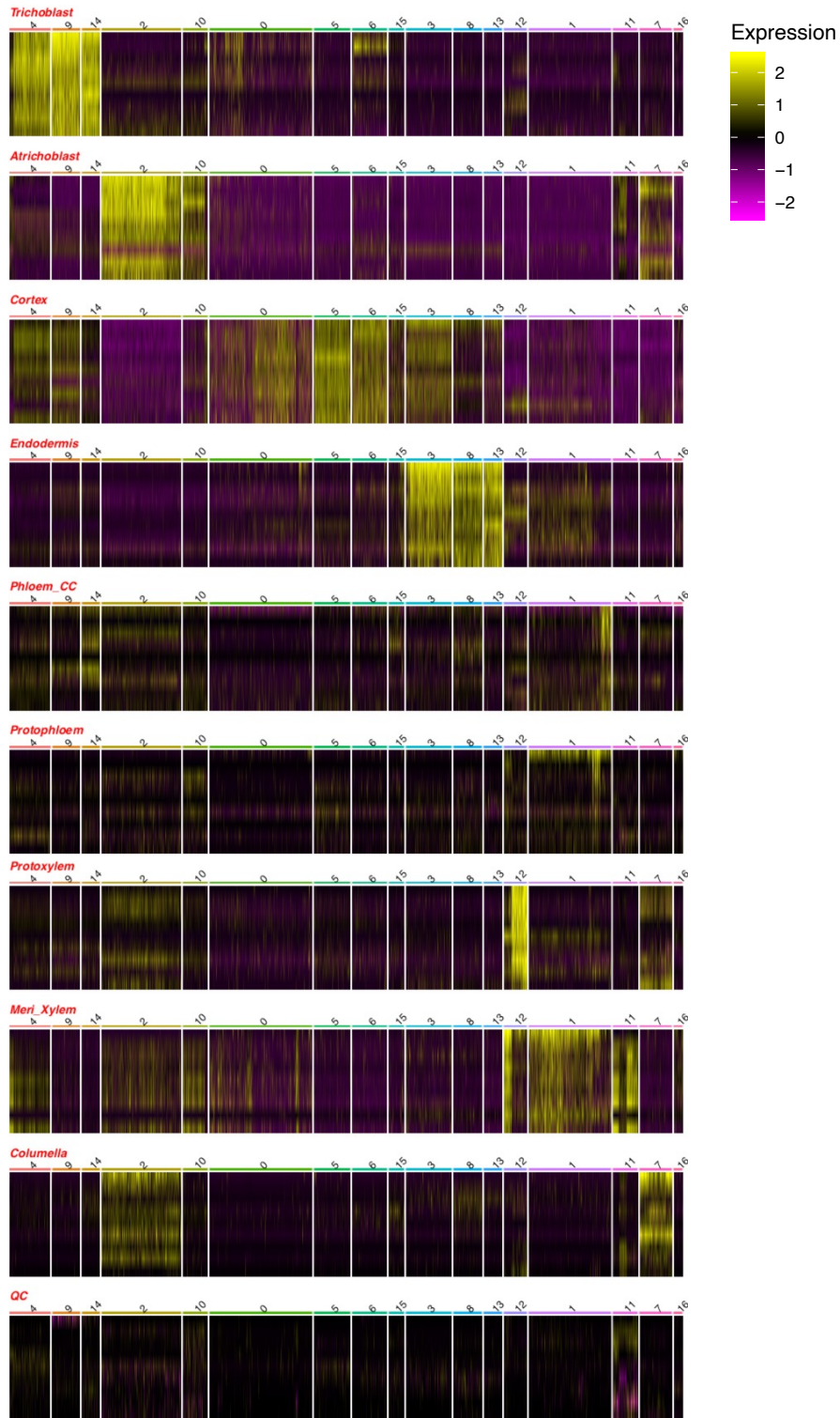
Figure S5 Performance comparisons among five cross validation random forest models. CV1 to CV5 suggest models obtained from the five-fold cross validation. The error bar suggests evaluation score variations of the ten cell types. The evaluation scores include sensitivity (SE), accuracy (AC), specificity (SP), precision (PR), geometric mean (GM), matthews correlation coefficient (MCC), and mean average precision (MAP). All pair wise comparisons are not statistically significant, as represented by the same letter a.

1
2
3
4
5
6
7
8



1 **Figure S6** Comparisons of proportion of expressed cells among the SHAP, CORR, ICIM, and KNOW markers
2 across all the clusters. In these boxplots, the mid-horizontal line represents the median and dots represent data
3 outliers.
4

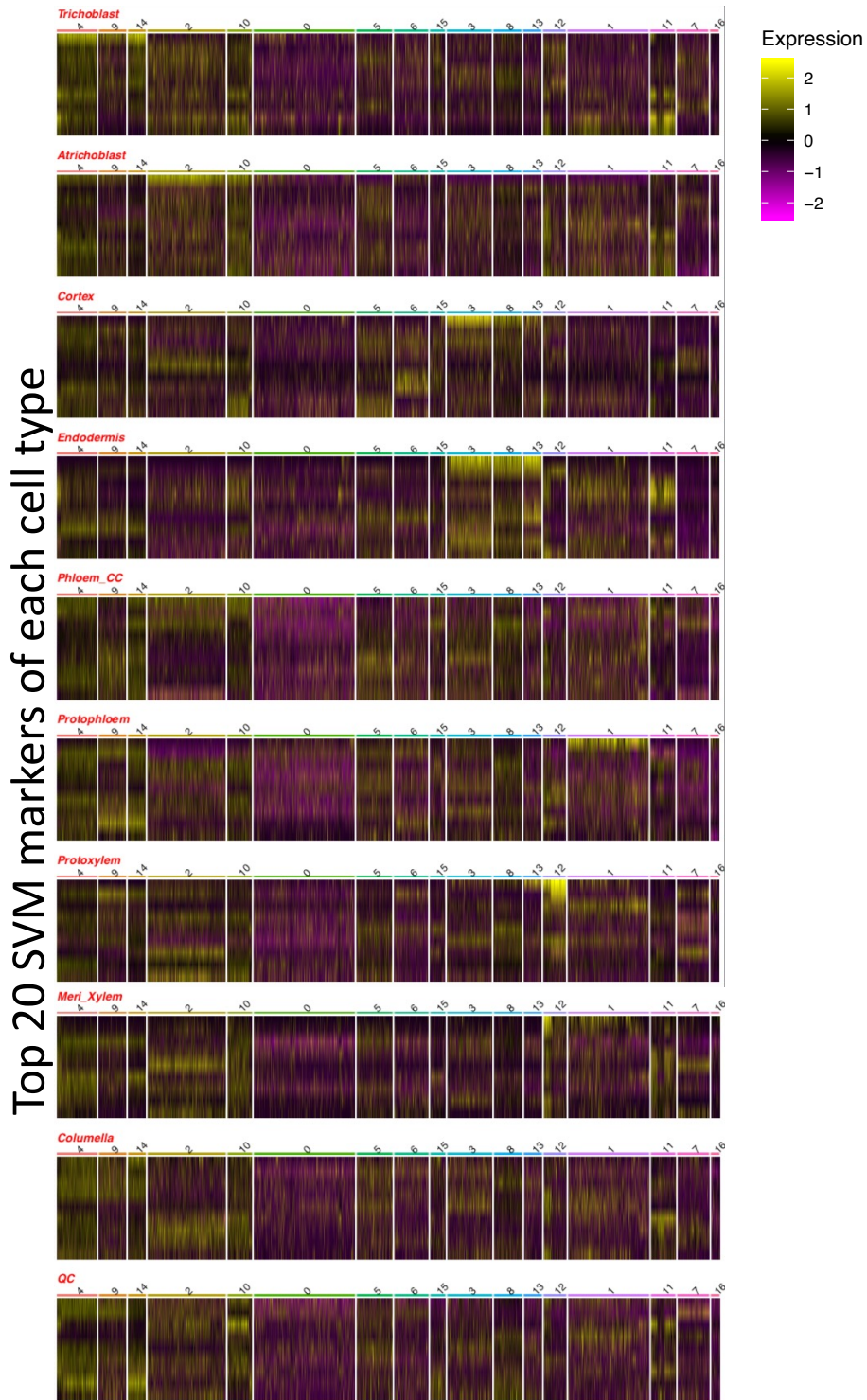
Top 20 SHAP markers of each cell type



1
2
3
4

Figure S7 Heat map of top 20 SHAP markers from each cell type across all the 17 cell clusters.

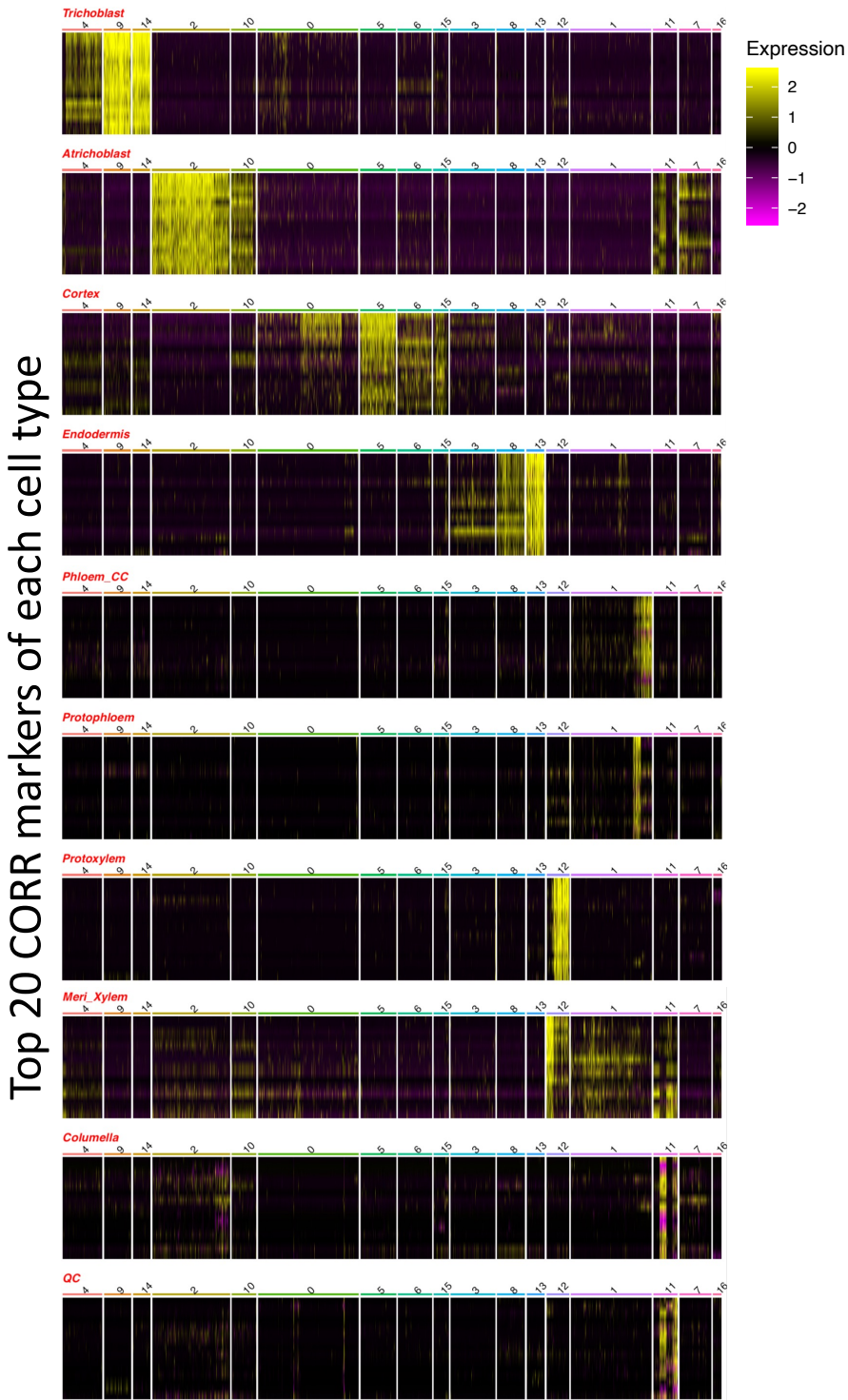
1



2
3
4

Figure S8 Heat map of top 20 SVM markers from each cell type across all the 17 cell clusters.

1



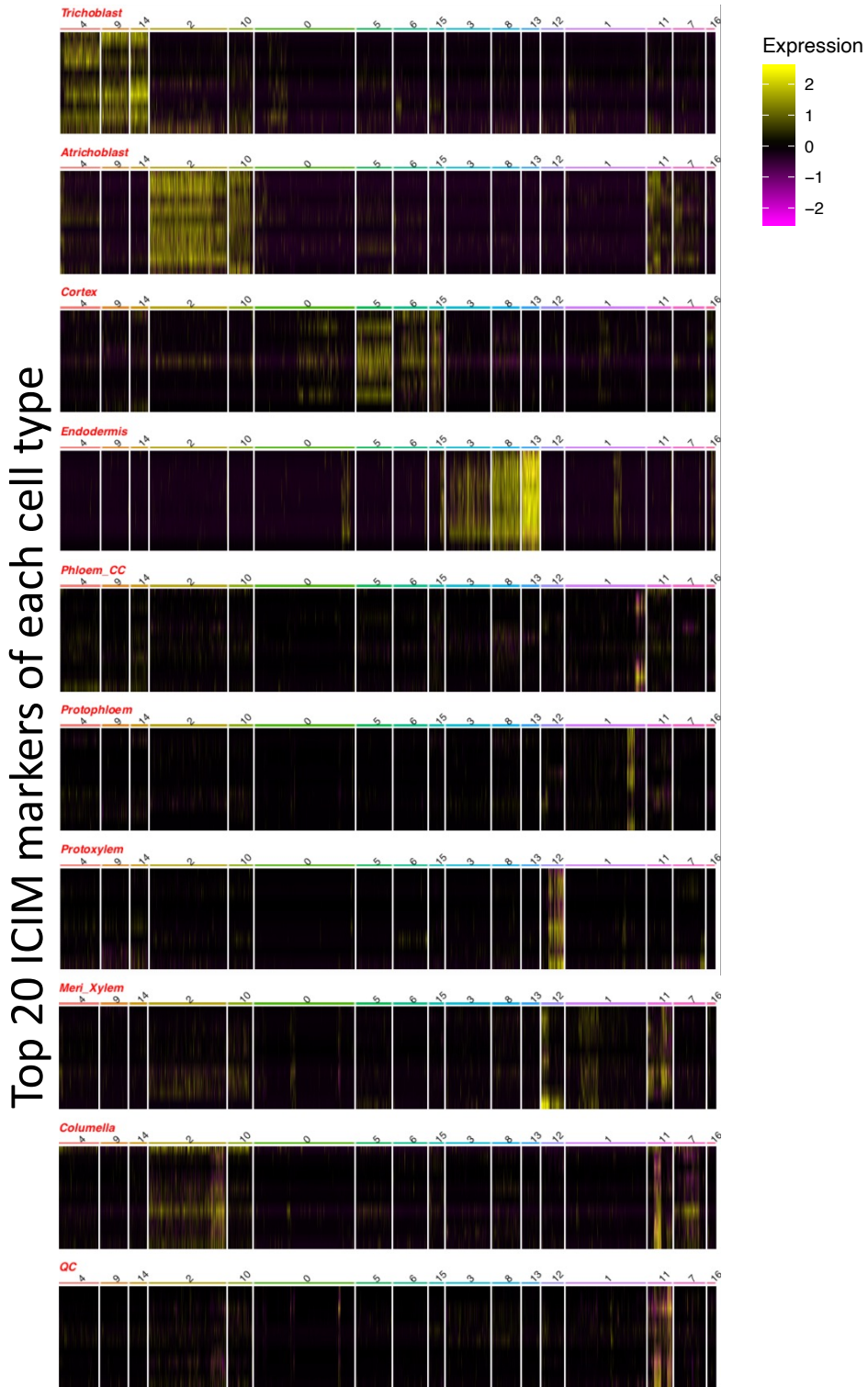
2

3

4

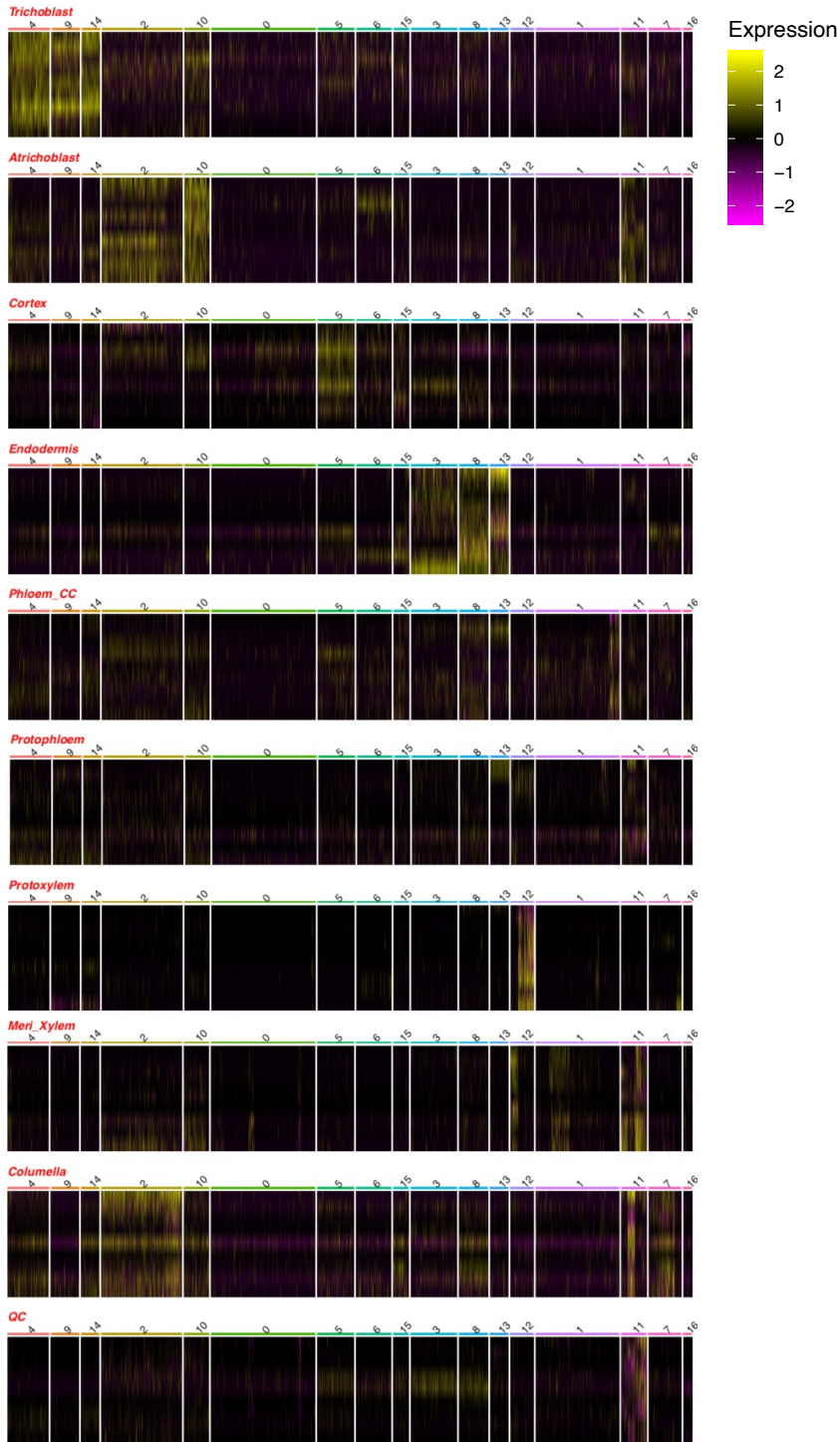
Figure S9 Heat map of top 20 CORR markers from each cell type across all the 17 cell clusters.

5



1
2
3
4
Figure S10 Heat map of all ICIM markers from each cell type across all the 17 cell clusters.

Top 20 KNOW markers of each cell type



1
2

Figure S11 Heat map of top 20 KNOW markers randomly selected from each cell type across all the 17 cell clusters

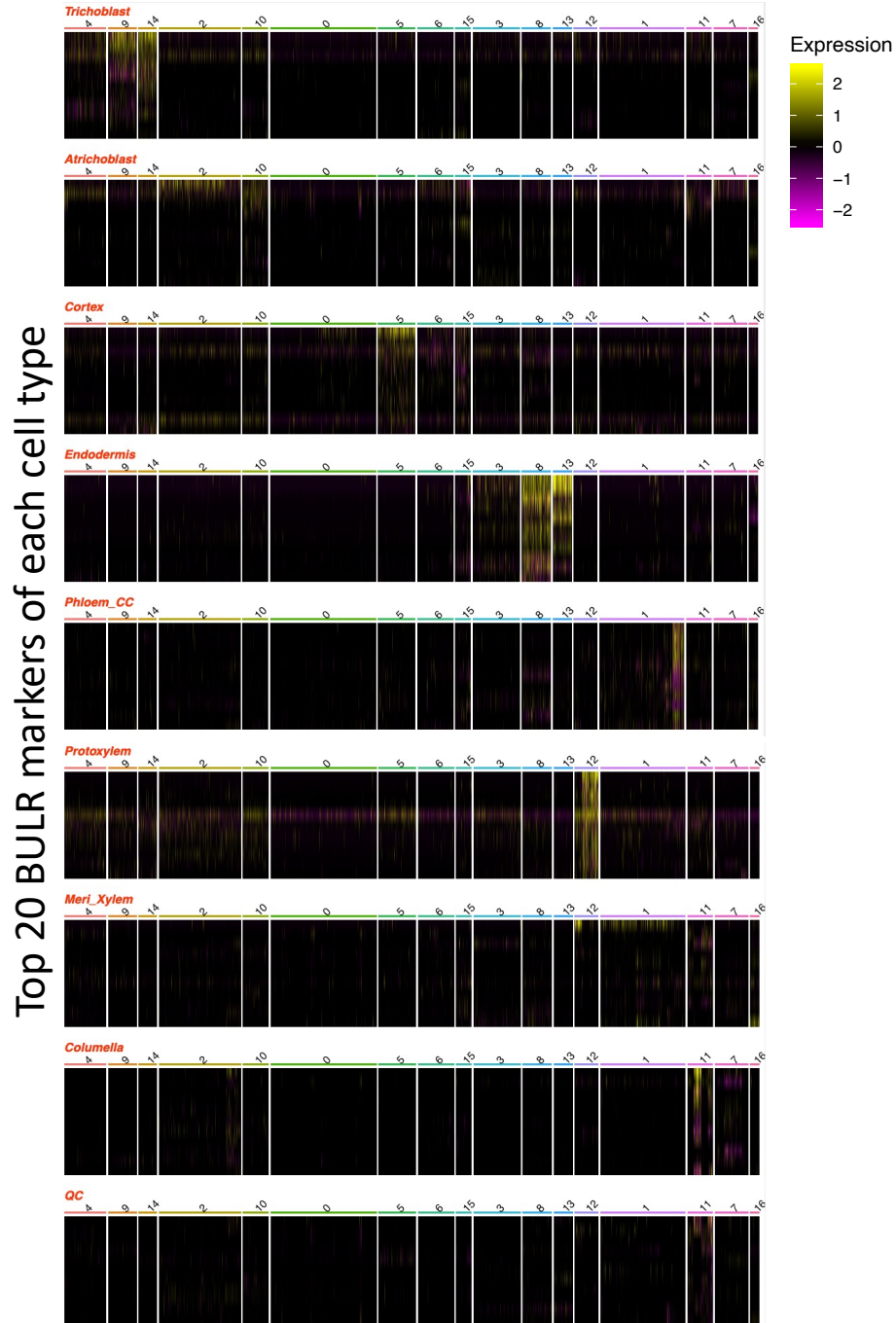


Figure S12 Heat map of top 20 BULR markers from each cell type across all the 17 cell clusters.

1
2
3
4
5

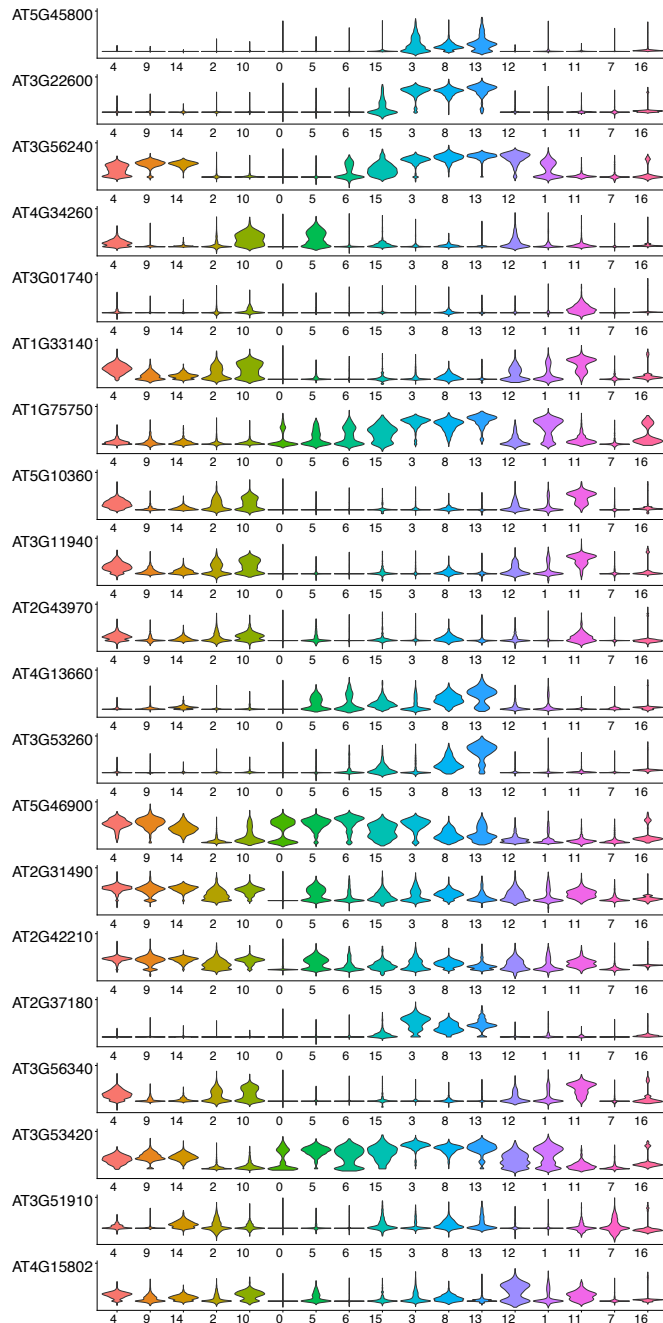
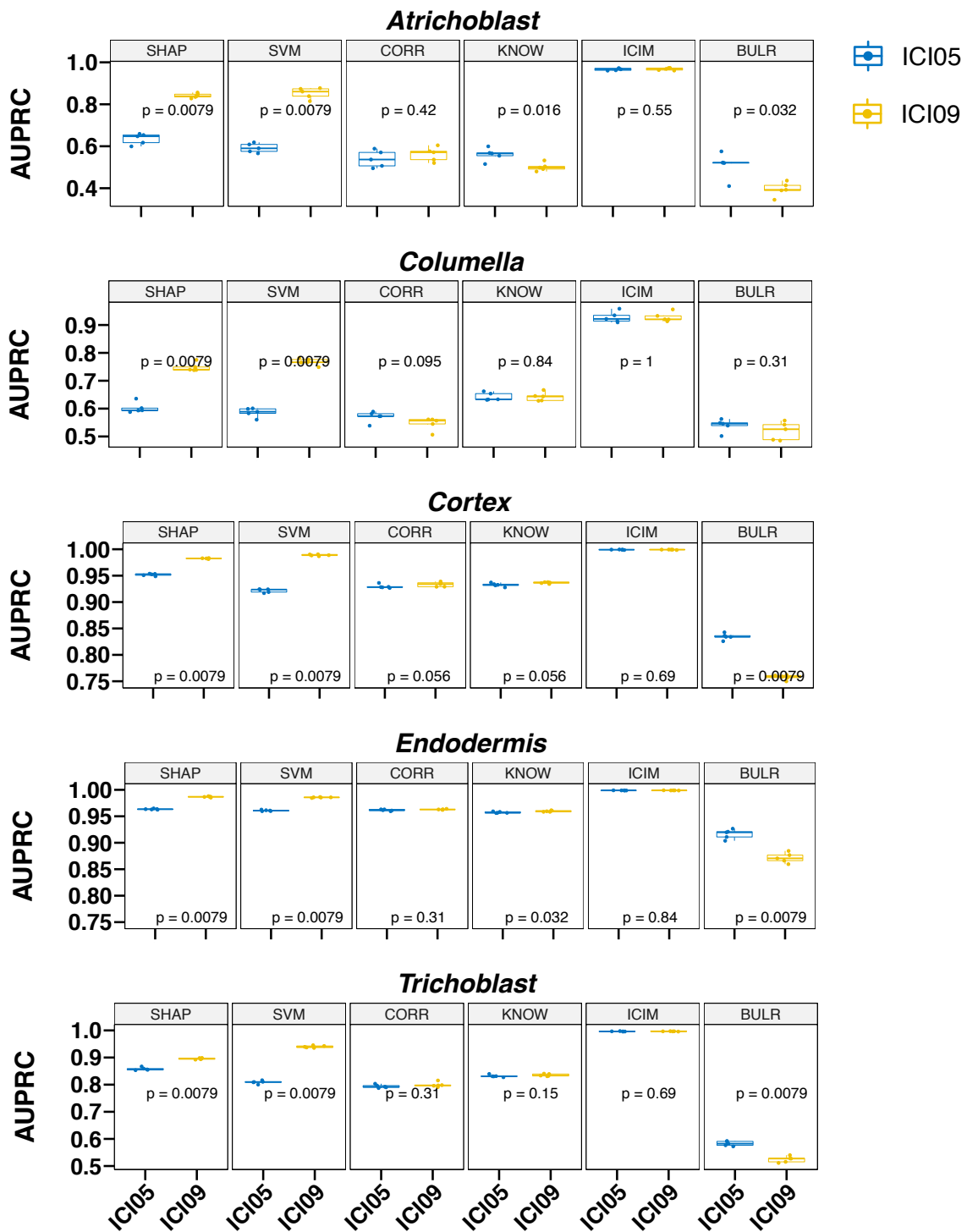


Figure S13 Violin plots of top 20 SVMM markers of endodermis cell type. Violin plots only show the distribution of the data.

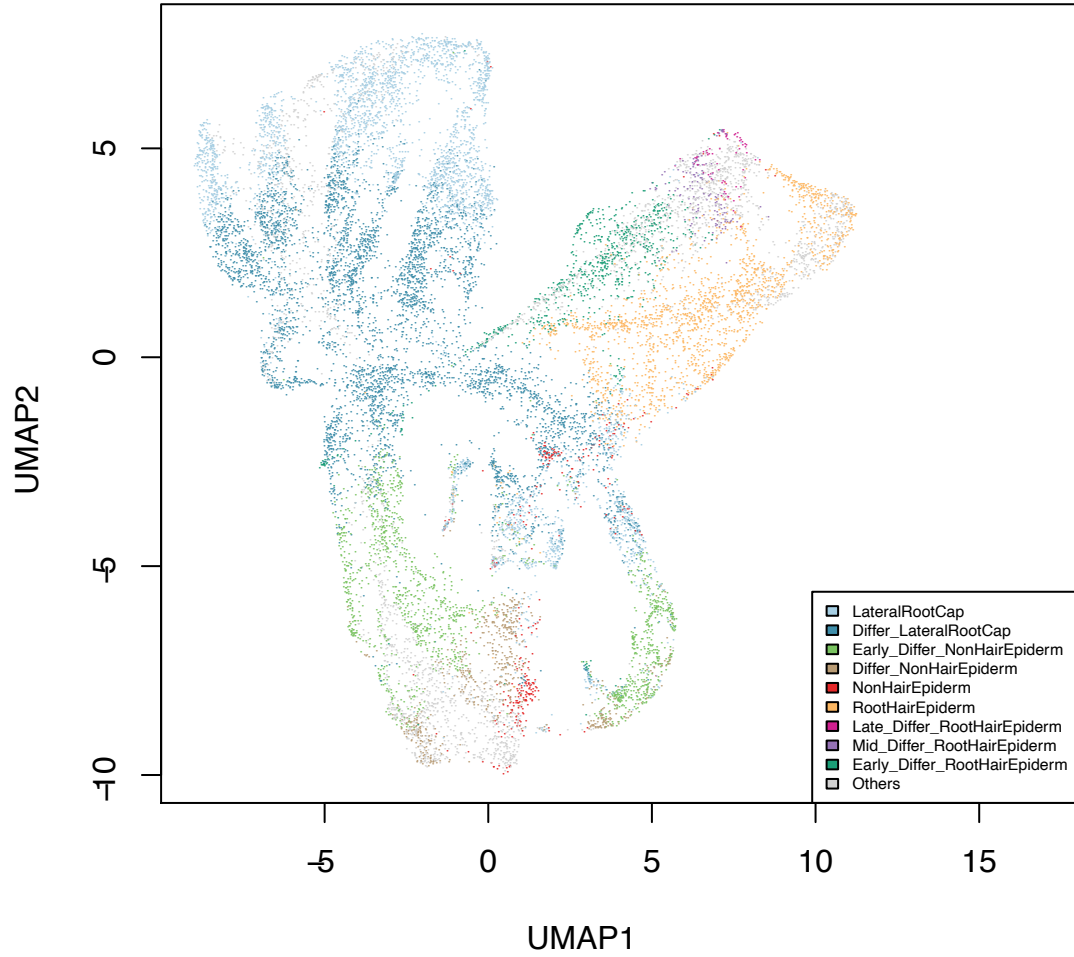
1
2
3
4
5
6
7
8
9
10
11
12

1
2



1
 2 **Figure S14 Comparison of classification performance based on ICI labeling method between 0.5 and 0.9**
 3 **thresholds in five cell types.** p value < 0.05 indicates significant differences between ICI05 and ICI09 groups. In
 4 these boxplots, the mid-horizontal line represents the median and dots represent data points.
 5
 6

1



2
3
4
5
6
7

Figure S15 A UMAP for SVM predicted cell types. Note that a group of cells predicted as non-hair epidermal cells by SHAP random forest model are predicted as lateral root cap cells. ‘Others’ label indicates cells from Ryu’s study (2019).

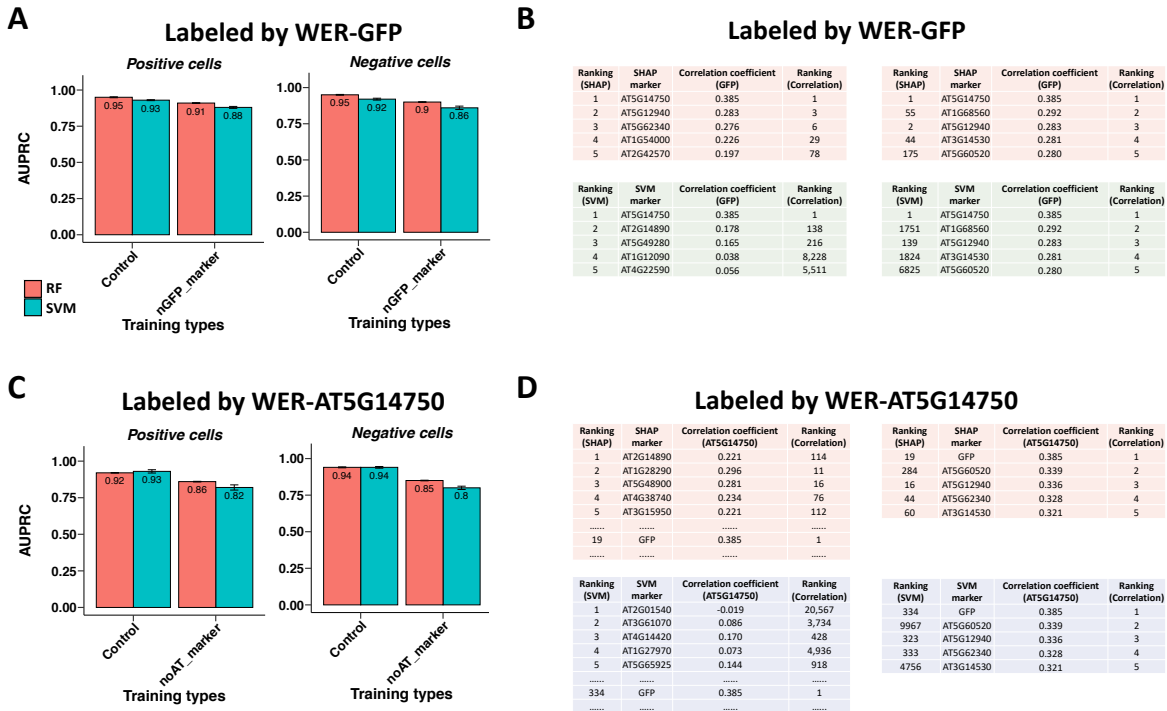
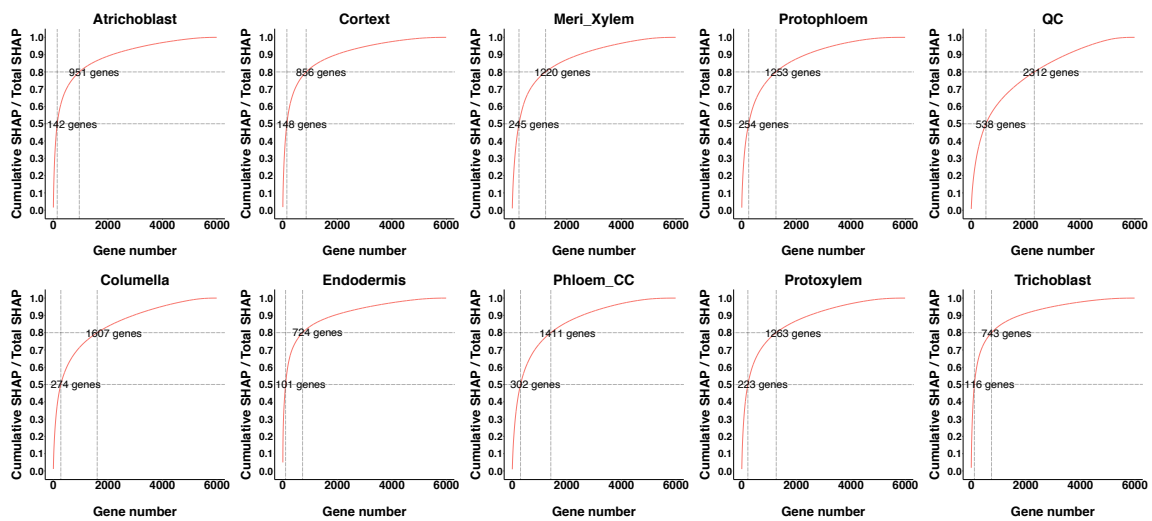


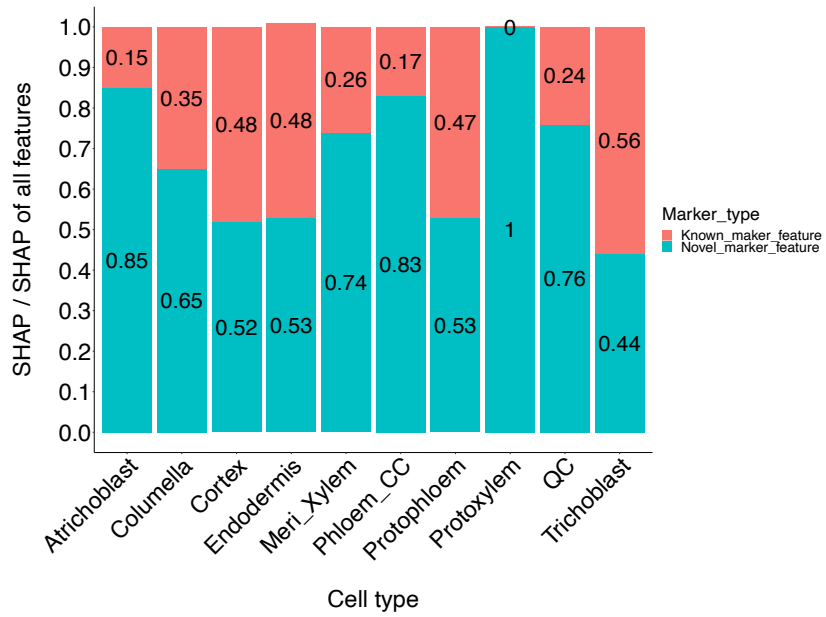
Figure S16. Classification and predicted top markers. A. Comparison of classification performance on GFP-labeled WER cells (positive cells) and none GFP-labeled WER cells (negative cells) between using all genes (control) and genes without GFP marker (nGFP_marker) for both RF and SVM models. B. Ranking of best SHAP and SVM markers to predict WER-GFP positive cells (left two tables). Ranking of genes with top correlation values with the GFP markers (right two tables). C-D show the similar legends as A-B except cells were labeled by AT5G14750. AUPRC means Area Under Precision-Recall Curve. Error bars represent +/- SE.

1
2
3
4
5
6
7
8
9
10
11
12



1
2
3

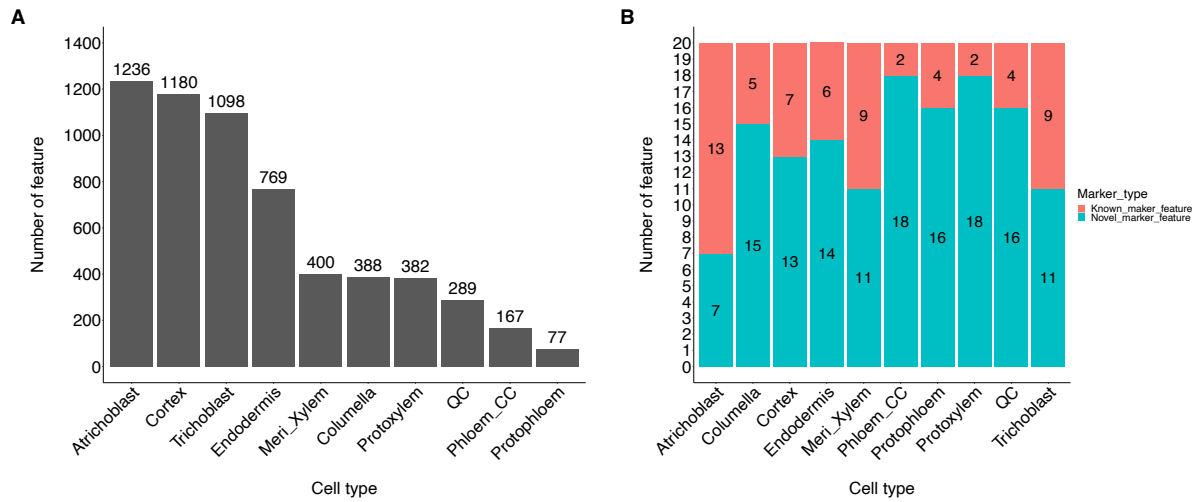
Figure S17. Cumulative SHAP values for all SHAP markers for each of the ten cell types.



1
 2 **Figure S18** Comparison of proportion of cumulative SHAP values from the SHAP to and from the known markers
 3 in the top 20 features with the highest SHAP value in each cell type.

4
 5

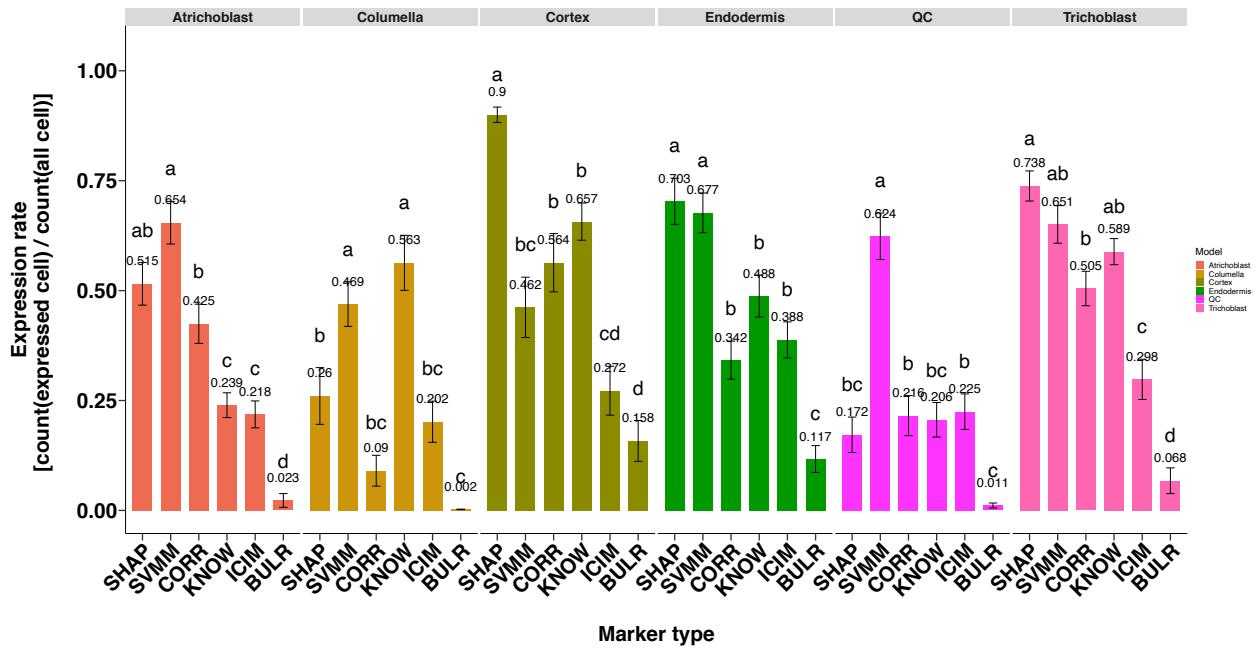
1



2
3
4
5
6
7
8
9
10
11
12

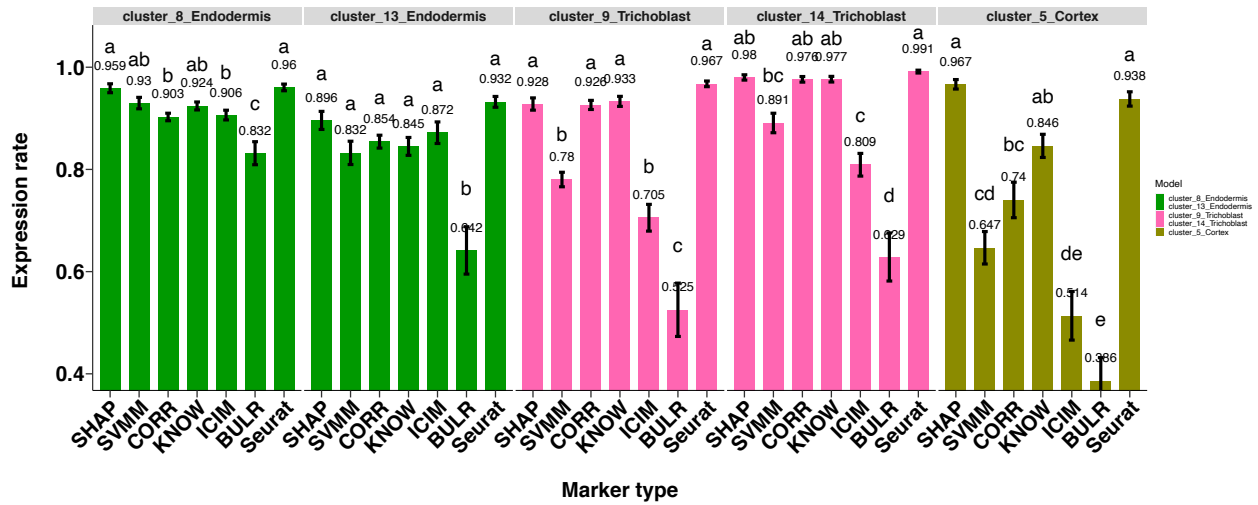
Figure S19 Summary of the SVMM markers. The SVMM genes were determined based on the feature importance estimated from the absolute coefficient values in the SVM model. Each gene has a coefficient for each of the ten cell types, and each gene is assigned to the cell type with the highest coefficient. The number of unique marker genes and novel marker genes determined by SVM are similar to that determined by SHAP. For example, most unique SVMM genes were assigned to atrichoblast (1,236), cortex (1,180) and vascular tissue and QC cells have lower numbers of SVMM markers (A). However, some cell types, such as trichoblast, have approximately twice as many SVMM markers as SHAP markers (1,098 SVMM vs 555 SHAP). In the top 20 genes with the highest coefficient of each cell type, on average, 69.5% of SVMM are novel marker genes (B).

1

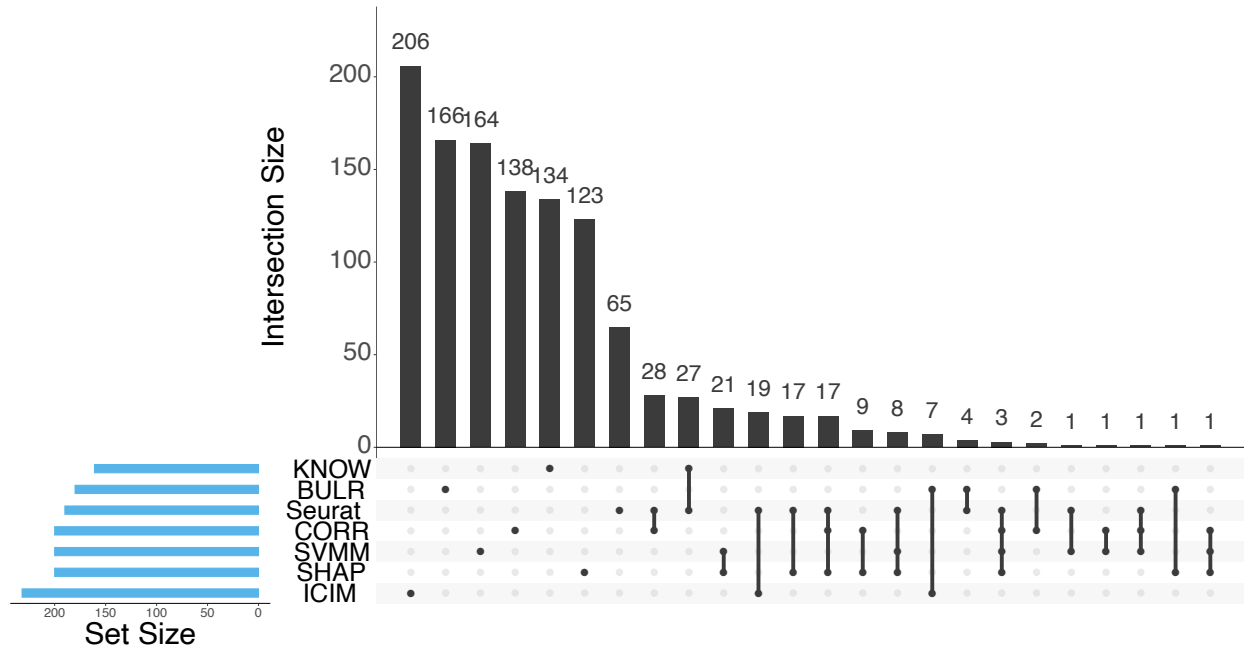


2
3
4
5
6
7

Figure S20 Expression rate of SHAP marker genes in (Shahan et al., 2020). All pair wise comparisons are statistically significant as indicated by different letters (a, b, c, and d). If two bars have the same letter, then they are not significantly different from each other. Error bars represent +/- SE.



1
 2 Figure S21. Expression rate comparison for five clusters between different marker types. Error
 3 bars represent +/- SE. All pair wise comparisons are statistically significant as indicated by
 4 different letters (a, b, c, d, and e).
 5
 6



1
2
3
4
5
6
7
8
9

Figure S22. Number of marker genes uniquely and commonly identified by different methods
 Intersection size means gene count of different marker types. The dots under the bars mean the genes specifically exist in the relative marker type. The line connected between two or more dots under the bars mean genes exist in two or more marker types. If two or more marker types do not have connection, it means these groups do not have shared genes.

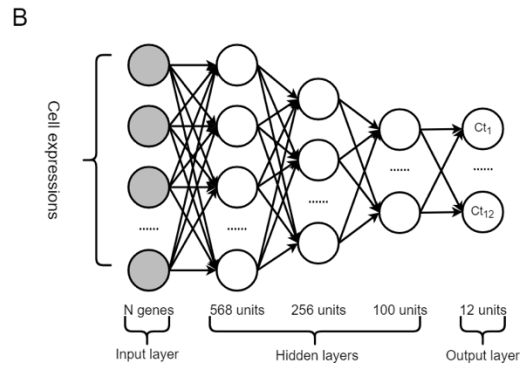
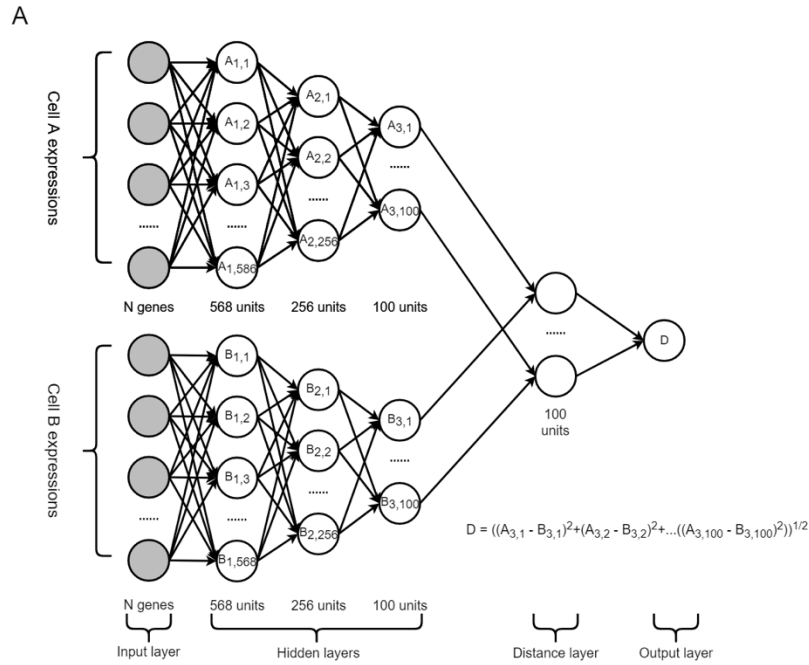


Figure S23 Schematic demonstration of architecture for each type of neural network. A. Architecture of Siamese NN, which was used for both triplet NN and contrastive NN. The distance layer computes a vector of distance between the last two hidden layers A3 and B3. This distance was then used in the objective function of triplet NN and contrastive NN to train cell type classifier. B. Architecture of multi-task NN. Ct represents a cell type (10 cell types were used in total for classification)

1
2
3
4
5
6
7
8
9
10

1 **Reference**

- 2 Abdi, H., and Williams, L.J. (2010). Tukey's honestly significant difference (HSD) test.
3 Encyclopedia of Research Design. Thousand Oaks, CA: Sage, 1-5.
- 4 Alavi, A., Ruffalo, M., Parvangada, A., Huang, Z., and Bar-Joseph, Z. (2018). A web server for
5 comparative analysis of single-cell RNA-seq data. *Nature communications* 9, 1-11.
- 6 Benesty, J., Chen, J., Huang, Y., and Cohen, I. (2009). Pearson correlation coefficient. In *Noise*
7 *reduction in speech processing* (Springer), Editors: Jacob Benesty, Jingdong Chen, Yiteng Huang
8 and Israel Cohen. pp. 1-4.
- 9 Chang, C.-C., and Lin, C.-J. (2011). LIBSVM: A library for support vector machines. *ACM*
10 *transactions on intelligent systems and technology (TIST)* 2, 1-27.
- 11 Goodstein, D.M., Shu, S., Howson, R., Neupane, R., Hayes, R.D., Fazo, J., Mitros, T., Dirks, W.,
12 Hellsten, U., and Putnam, N. (2012). Phytozome: a comparative platform for green plant
13 genomics. *Nucleic acids research* 40, D1178-D1186.
- 14 Liu, Q., Liang, Z., Feng, D., Jiang, S., Wang, Y., Du, Z., Li, R., Hu, G., Zhang, P., and Ma, Y.
15 (2021). Transcriptional landscape of rice roots at the single-cell resolution. *Molecular Plant* 14,
16 384-394.
- 17 Lundberg, S.M., Erion, G., Chen, H., DeGrave, A., Prutkin, J.M., Nair, B., Katz, R.,
18 Himmelfarb, J., Bansal, N., and Lee, S.-I. (2020). From local explanations to global
19 understanding with explainable AI for trees. *Nature machine intelligence* 2, 2522-5839.
- 20 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,
21 Prettenhofer, P., Weiss, R., and Dubourg, V. (2011). Scikit-learn: Machine learning in Python.
22 *the Journal of machine Learning research* 12, 2825-2830.
- 23 Ryu, K.H., Huang, L., Kang, H.M., and Schiefelbein, J. (2019). Single-cell RNA sequencing
24 resolves molecular relationships among individual plant cells. *Plant physiology* 179, 1444-1456.
- 25 Satija, R., Farrell, J.A., Gennert, D., Schier, A.F., and Regev, A. (2015). Spatial reconstruction of
26 single-cell gene expression data. *Nature biotechnology* 33, 495-502.
- 27 Schroff, F., Kalenichenko, D., and Philbin, J. (2015). Facenet: A unified embedding for face
28 recognition and clustering. In *Proceedings of the IEEE conference on computer vision and*
29 *pattern recognition*, pp. 815-823.
- 30 Shahan, R., Hsu, C.-W., Nolan, T.M., Cole, B.J., Taylor, I.W., Vlot, A.H.C., Benfey, P.N., and
31 Ohler, U. (2020). A single cell Arabidopsis root atlas reveals developmental trajectories in wild
32 type and cell identity mutants. *bioRxiv*. DOI: 10.1101/2020.06.29.178863
- 33 Wagner-Menghin, M.M. (2014). Binomial test. *Wiley StatsRef: Statistics Reference Online*.
34