# Supplementary Information for

## Predicting transcriptional responses to cold stress across plant species

**Xiaoxi Meng, Zhikai Liang, Xiuru Dai, Yang Zhang, Samira Mahboub, Daniel W. Ngu, Rebecca L. Roston, and James C. Schnable**

**James C. Schnable.**
**E-mail: schnableunl.edu**

**This PDF file includes:**

Figs. S1 to S12
Legend for Dataset S1
SI References

**Other supplementary materials for this manuscript include the following:**
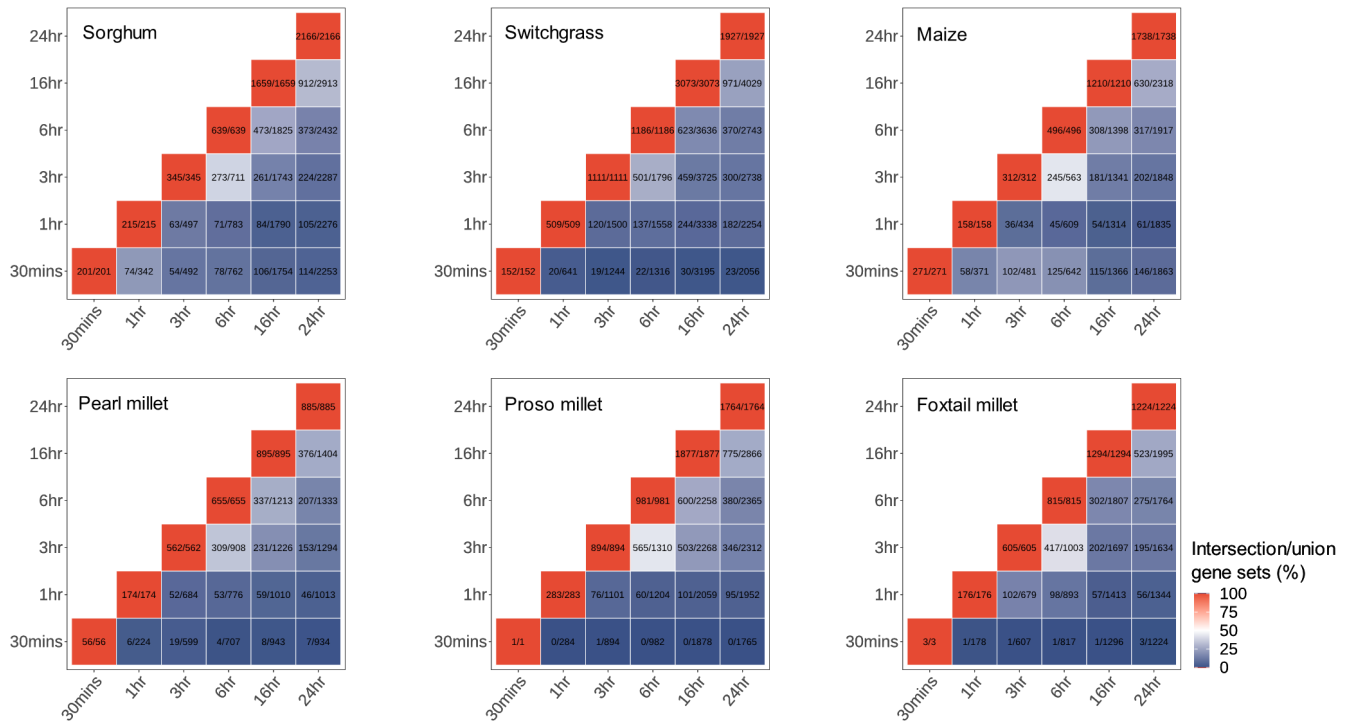
Dataset S1

**Fig. S1. Overlap of cold-responsive genes among time points in each species.** Abundance of intersecting differentially expressed genes as a percent of the union of differentially expressed genes between pairs of time points in each grass species analyzed in this study. In each cell the numerator indicates the intersection of the sets of differentially expressed genes identified at the two time points and the indicates the union of the same two sets of genes.
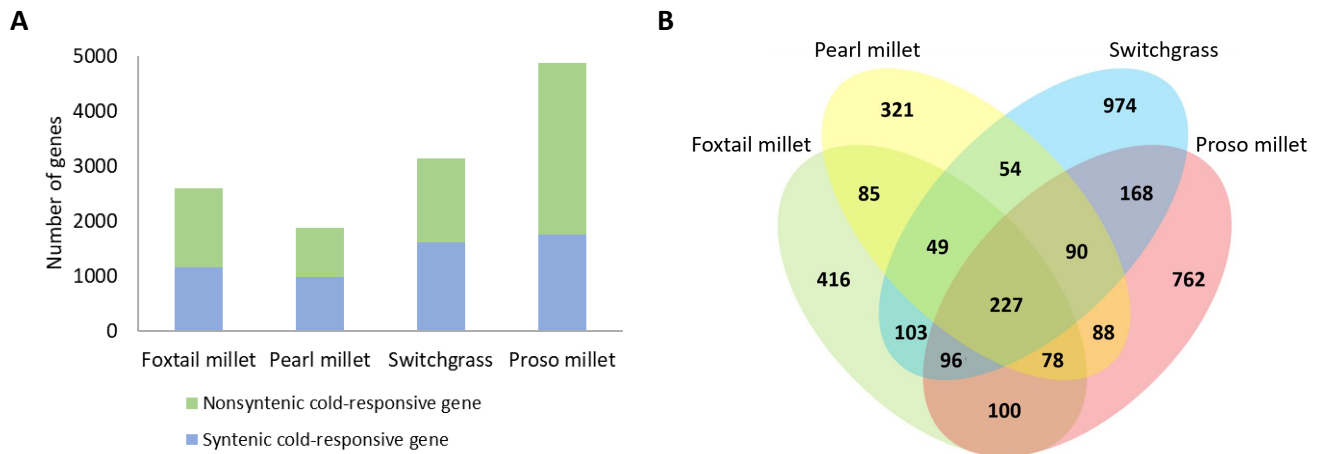
**A**



**B**



**Fig. S2. Conserved cold-responsive genes across foxtail millet, pearl millet, switchgrass, and proso millet.** A. Proportions of syntenic orthologous genes among cold-responsive genes; B. Overlapping cold-responsive syntenic orthologs among the four species.
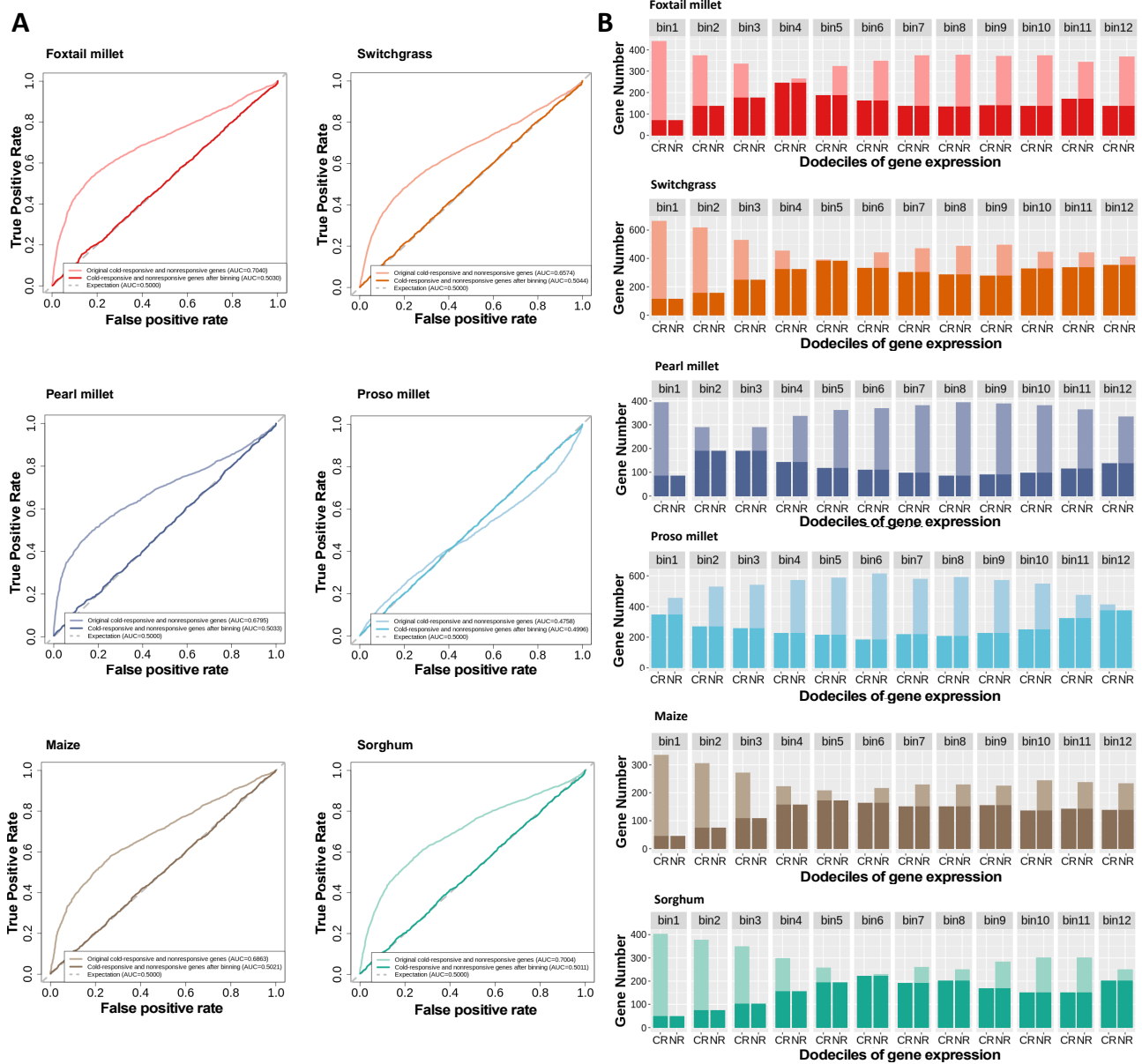
**Fig. S3. Baseline expression controls.** A. Accuracy of genes being scored as cold-responsive genes solely based on average FPKM values before and after baseline expression control; B. Distribution of average FPKM values of cold-responsive genes (CR) and nonresponsive genes (NR), and training sets resampled from genes in dodeciles with balanced gene expression levels (darker color).
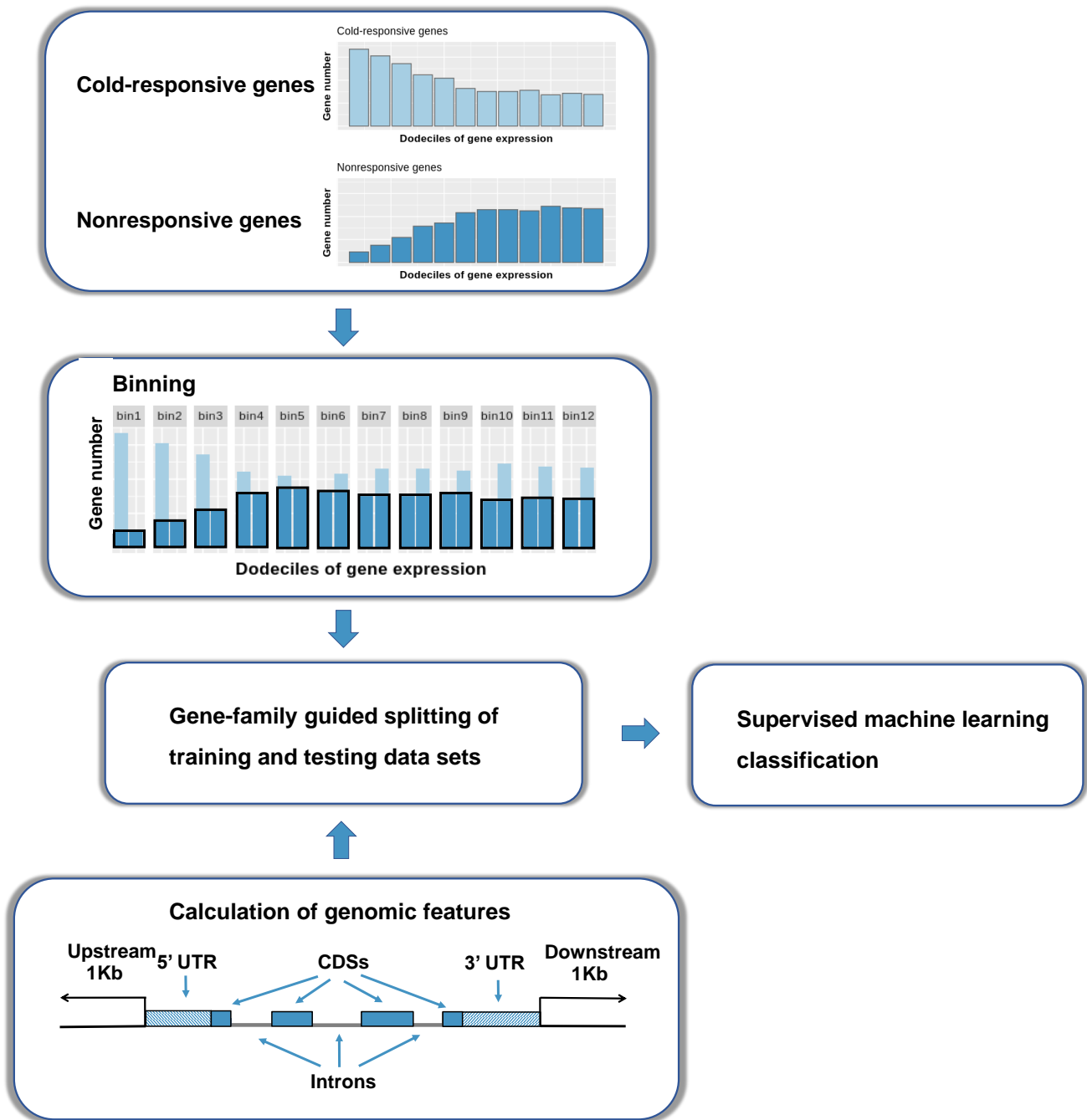
**Xiaoxi Meng, Zhikai Liang, Xiuru Dai, Yang Zhang, Samira Mahboub, Daniel W. Ngu, Rebecca L. Roston, and James C. Schnable**

**Fig. S4. Workflow of the supervised machine classification model for predicting cold-responsive genes.** For within species predictions, gene-family guided subsampling and splitting consisted of first subsampling each gene family present in the species and then dividing into training/validation and testing data. For cross-species predictions, gene-family guided subsampling and splitting consisted of first dividing gene families into training/validation and testing data and then subsampling one gene per family per species.
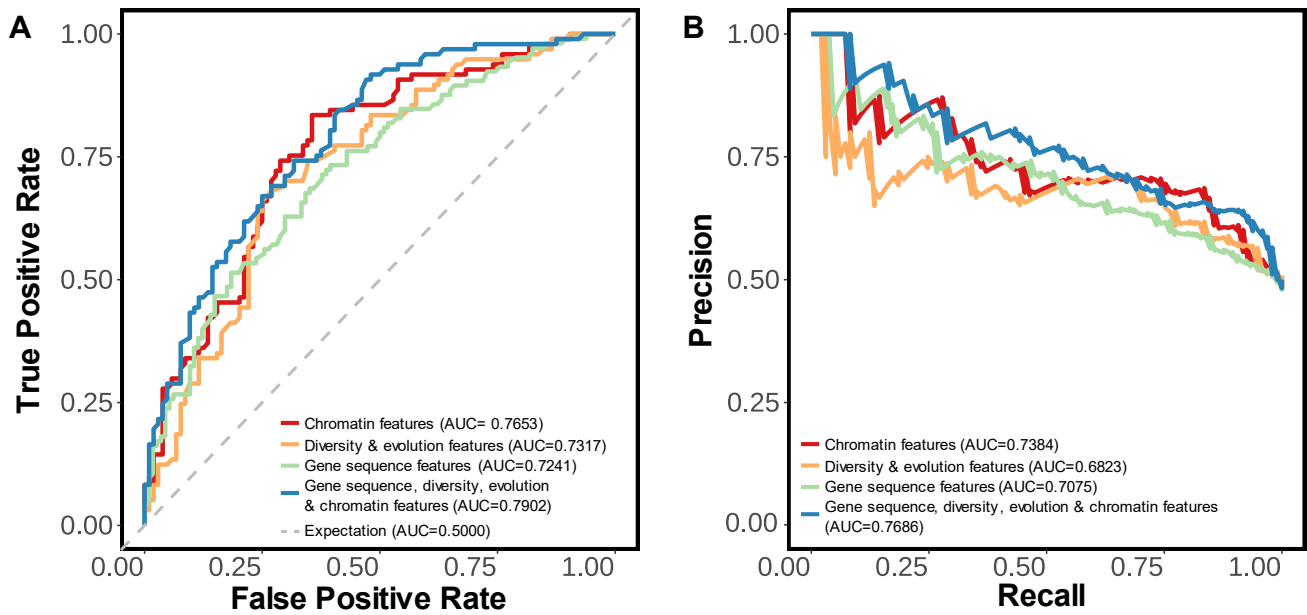
**Fig. S5. Cold-responsive gene predictions in maize using different subsets of features** A. Receiver operating characteristic (ROC) curves shows the classification on holdout test data; B. Precision-recall (PR) curves shows the classification on holdout test data.
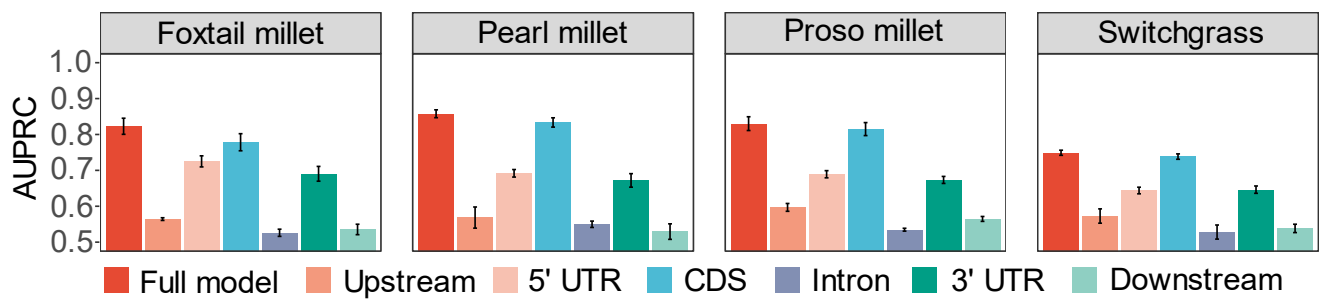
**Fig. S6. AUPRCs of supervised machine learning models for Paniceae grass species based on gene sequence features.** Bar plot showing AUPRCs achieved by the full gene sequence models and single feature group models for foxtail millet, pearl millet, switchgrass and proso millet. Standard error (se) was calculated from five independent predictions.
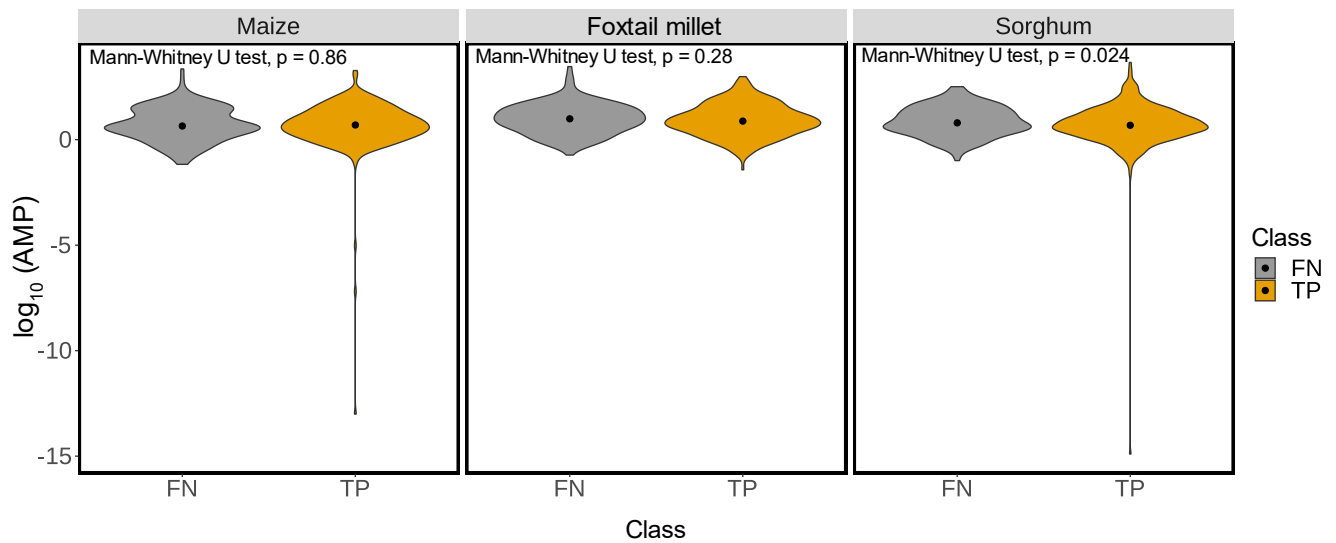
**Fig. S7. Model performance evaluations on predicting cold-responsive genes and circadian genes.** The trained machine learning model in each species (Maize, Sorghum and Foxtail millet) using gene sequence features was applied on predicting holdout test data. From prediction results, genes in holdout test data were split into true positive (TP) and false negative (FN) sets. Amplitudes of unique genes were considered together in one group and the black dot indicate the median value of amplitudes in each group. Mann-whitney U test was applied on comparing raw amplitudes between groups in each species. AMP represents amplitude.
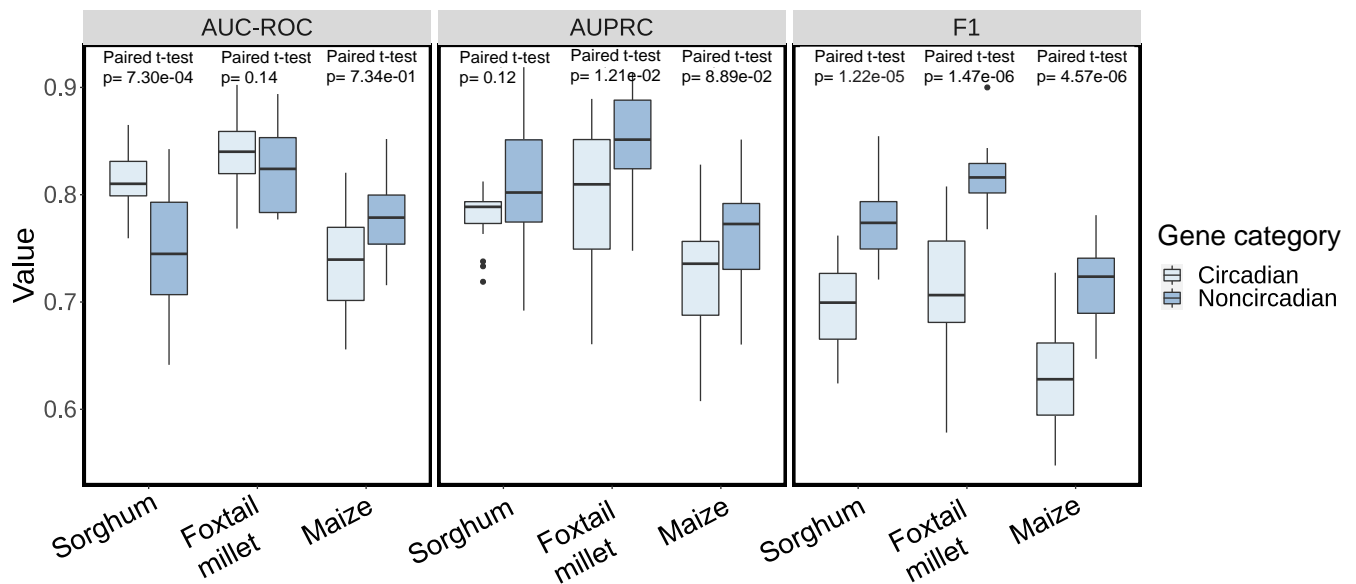
**Fig. S8. Model performance evaluations on predicting circadian genes.** The trained machine learning model in each species (Maize, Sorghum and Foxtail millet) using gene sequence features was applied on predicting diurnal cycling genes. AUC-ROC, AUPRC and F1 values were calculated on 10% holdout test data based on 20 prediction models per species. Paired t-test was used to evaluate statistical significance between samples.
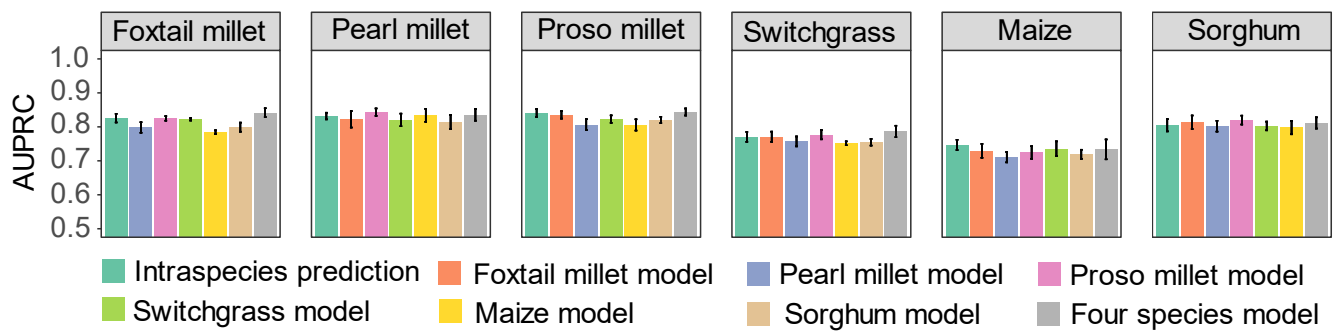
**Fig. S9. AUPRC of models trained for cross-species prediction.** Areas under Precision-Recall Curves (AUPRC) show the classification on holdout test data in machine learning models constructed in different species. Standard error (se) was calculated from five independent predictions. All predictions shown here, including intraspecies predictions, were made using the cross-species prediction framework for partitioning hold out test data (see methods).
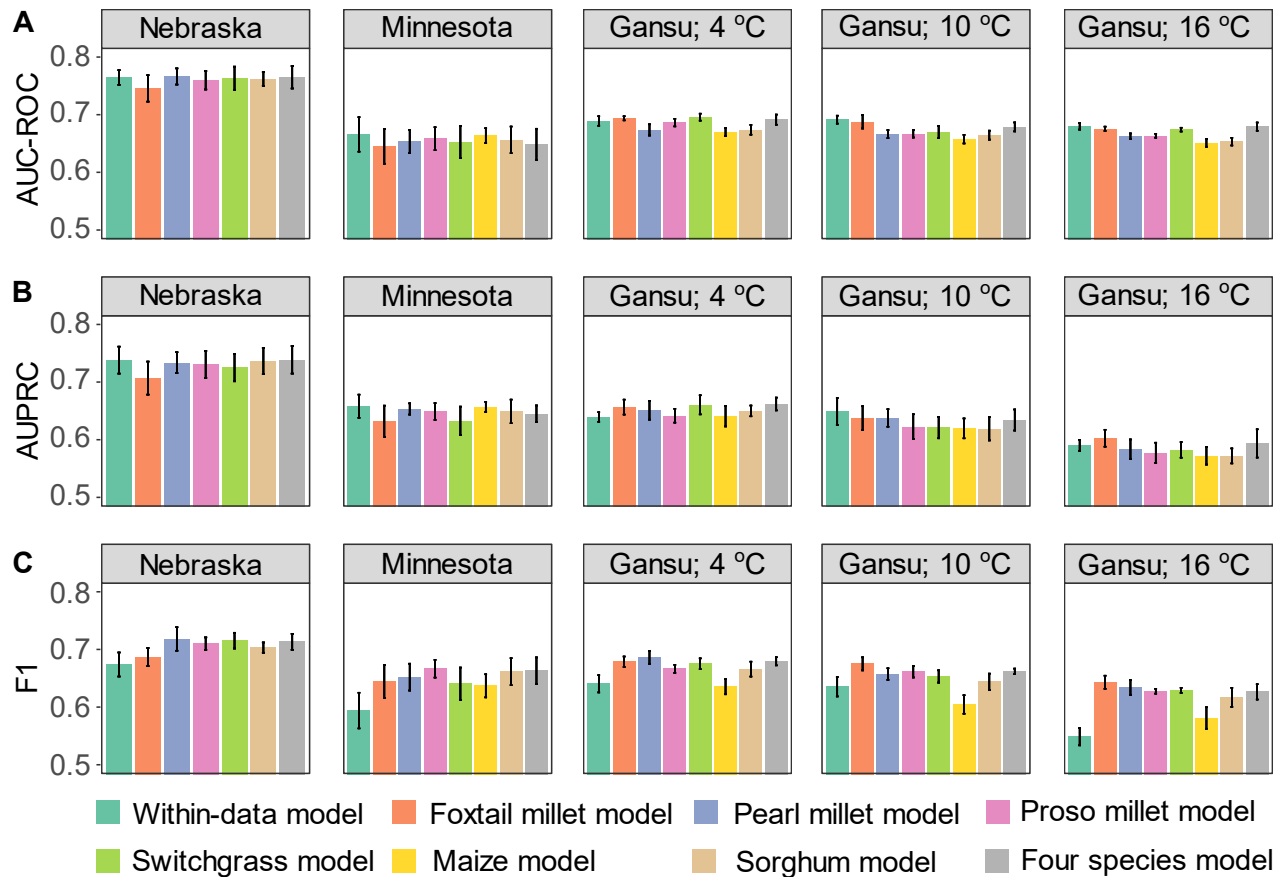
**Fig. S10. Comparisons of model performance in predicting cold-responsive genes in maize identified using RNA seq data collected by different research groups.** A. Areas under the receiver operating characteristic curves (AUC-ROCs) for predicting cold-responsive genes identified in different maize data sets. Standard error (se) was calculated from five independent predictions. Minnesota indicates data generated by Liang *et al.*, 2020 (1). Gansu (China) indicates data generated by Li *et al.*, 2020 (2) with different cold temperatures as labeled. Nebraska and "Maize model" indicate the same data employed in this study; B. Areas under Precision-Recall Curves (AUPRC) show performance of predictions on cold-responsive genes; C. F1 scores show performance of predictions on cold-responsive genes.
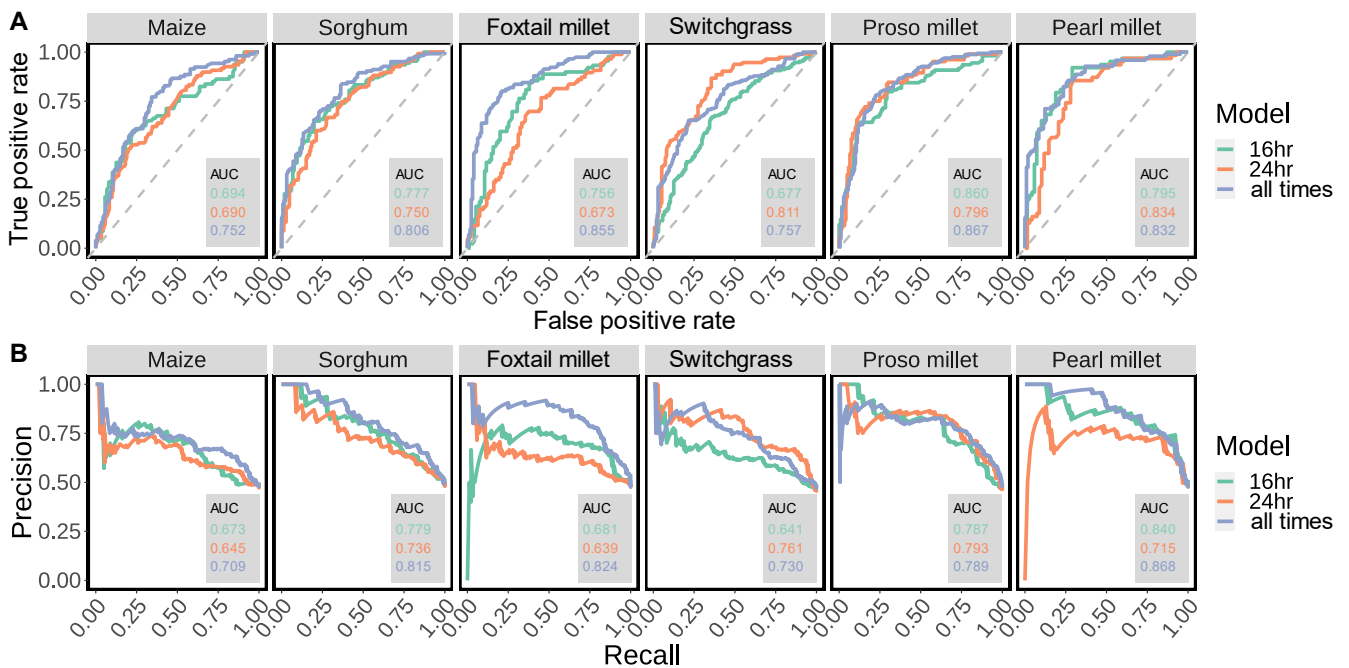
Xiaoxi Meng, Zhikai Liang, Xiuru Dai, Yang Zhang, Samira Mahboub, Daniel W. Ngu, Rebecca L. Roston, and James C. Schnable

**Fig. S11. Comparisons of model performance on predicting cold-repsonsive genes defined from single and multiple time points in six grass species.** A. Receiver operating characteristic (ROC) curves show classifications on holdout test data. Data collected from 16 hr and 24 hr represented single time point data. All times are the same as we defined cold-responsive genes in the study. Values of AUC were indicated within each species with same colors for corresponding models; B. Precision-recall (PR) curves show the classification on holdout test data.
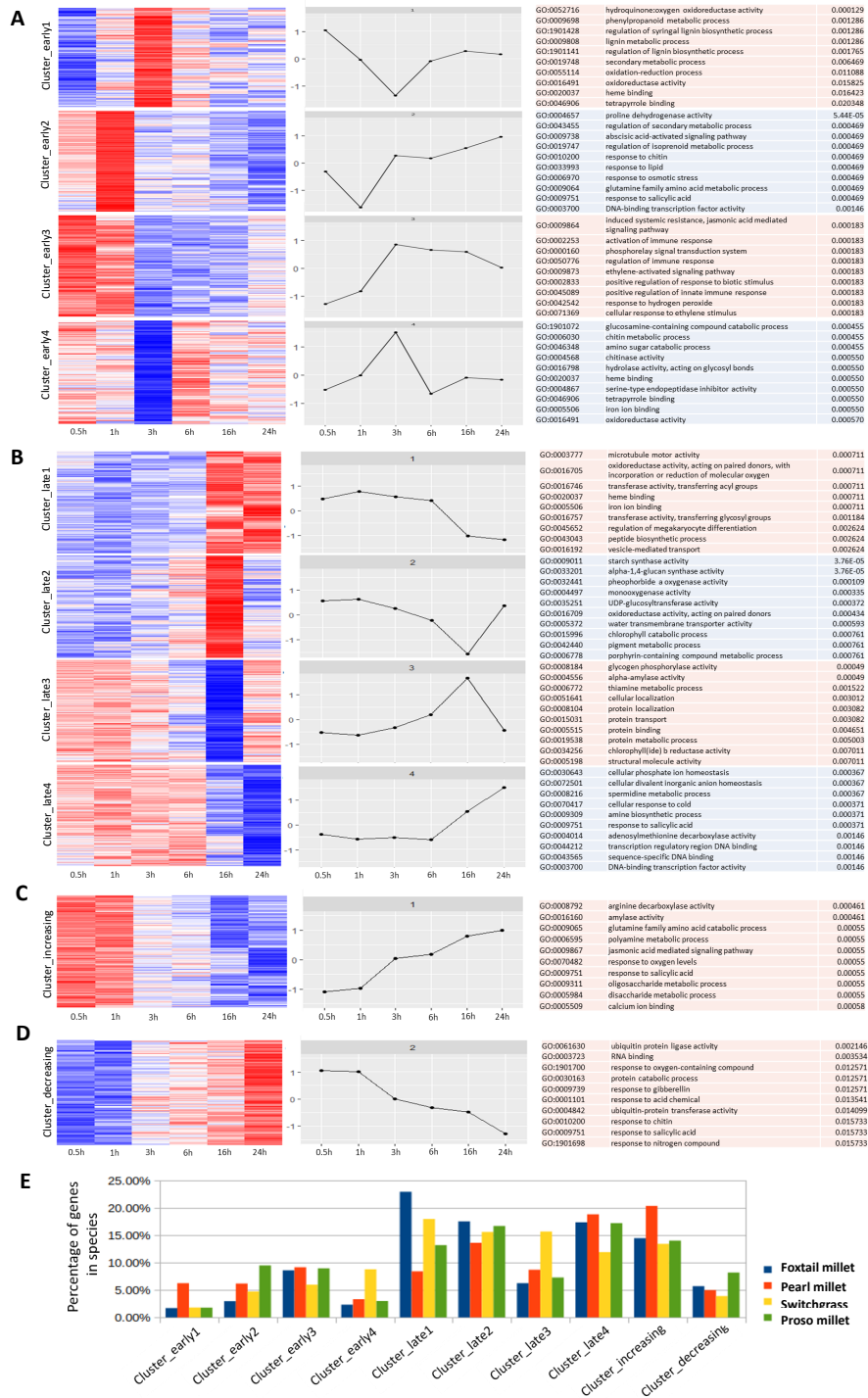
Xiaoxi Meng, Zhikai Liang, Xiuru Dai, Yang Zhang, Samira Mahboub, Daniel W. Ngu, Rebecca L. Roston, and James C. Schnable

**Fig. S12. Gene expression clusters analysis** A-D. Cold-responsive genes from foxtail millet, pearl millet, switchgrass, and proso millet were analyzed using k-means clustering. This process identified eight major groups, as shown in heat map and graphical format, based on patterns of gene expression at different time points. (A) Clusters containing genes of early transcriptional responses to cold (30 min to 3 h); (B) clusters with late responded genes to cold (responded after 6h); (C and D) genes with continuously increasing or decreasing transcriptional levels within 24hrs. Enriched GO terms within clusters were shown in the last column. E. Percentage of genes of each of the four species distributed in clusters.

Xiaoxi Meng, Zhikai Liang, Xiuru Dai, Yang Zhang, Samira Mahboub, Daniel W. Ngu, Rebecca L. Roston, and James C
Schnable

**SI Dataset S1 ()**

**Dataset S1, Tab1** Number of cold-responsive genes identified in the four grass species

**Dataset S1, Tab2** Performance metrics on predicting maize cold responsive genes by models with or without considering evolutionary relatedness and baseline expression.

**Dataset S1, Tab3** Gene sequence features, chromatin features, and diversity/evolutionary features used in supervised machine learning classification.

**Dataset S1, Tab4** The top 20 random forest feature importance presented by Mean Decrease Accuracy for each intraspecies prediction of transcriptional responses to cold stress.

**Dataset S1, Tab5** Model performance on predicting cold responsive gene sets identified across species.

**Dataset S1, Tab6** Model performance on predicting maize cold responsive gene sets identified by different experiments.

**Dataset S1, Tab7** Syntenic gene list among *S. bicolor*, *S. italica*, *P. glaucum*, *P. miliaceum*, and *P. virgatum*.

## References

1. Z Liang, et al., Genetic and epigenetic contributions to variation in transposable element expression responses to abiotic stress in maize. *bioRxiv* (2020).
2. Y Li, et al., Transcriptomic analysis revealed the common and divergent responses of maize seedling leaves to cold and heat stresses. *Genes* **11**, 881 (2020).