

# Maize Yield Prediction Accuracy Increased By Inclusion of Genetics, Environment, and Management Interactions With Deep Learning

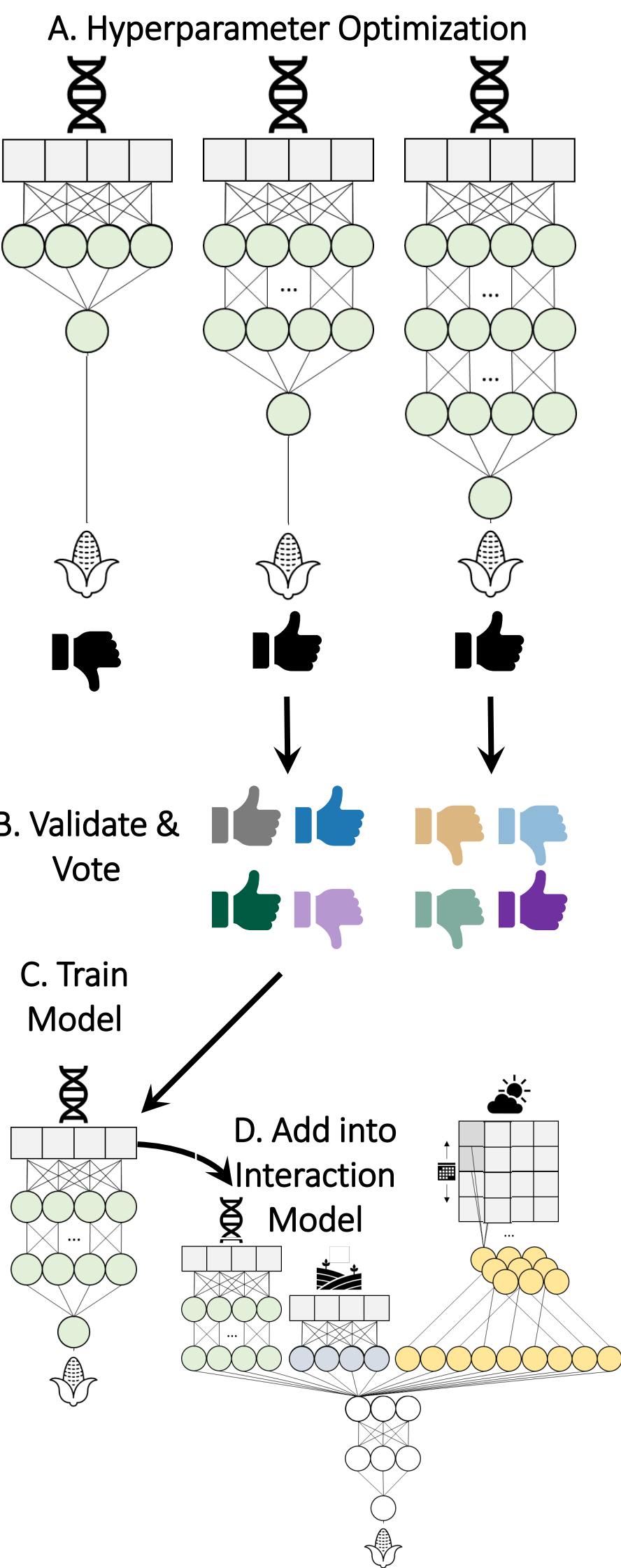
## Main Questions:

1. Is phenotypic prediction via deep learning more accurate than other models?
2. Does optimization strategy influence performance?
3. Does inclusion of GxExM effects change the importance of input variables?

## Abstract

Phenotypic prediction for polygenic traits and those with interaction effects between genetic, environmental, and management factors remains challenging. We developed a deep learning model with interactions between these factors to better predict maize (*Zea mays*) yield across diverse environments in the continental United States. This model and the modeling approach used is of potential use for genomic selection, management optimization, and forecasting. We find the optimization strategy used influences model performance with the best performance coming from consecutively optimizing submodels, each processing a single data group or interactions between data groups, rather than optimizing all aspects of the model's structure simultaneously. Furthermore, we find deep learning can but is not guaranteed to outperform other model types when all data groups (genomic, soil, and daily weather and management data) are provided. When restricted to a single data group, the model is less accurate. Lastly, we observe that interactions between data groups substantially influence the importance of different variables, reducing the influence of weather events at the end of the season, approximately 200 days following planting.

## Methods



## Simultaneous Optimization (SO)

Steps 1-3 above are repeated with interactions between data types allowed at the onset. Using same hyperparameter search space.

## Measuring Model Performance

Root mean squared error for the test set is used to assess model performance ( $RMSE = \sqrt{\frac{\sum \hat{y}_i - y_i}{n}}$ ). A deep neural network's performance can be sensitive to it's randomly initialized parameter values. To account for this each final architecture was initialized with 10 replicates. To reduce the time required to fit deep learning models, variables were centered and scaled by the training set mean and standard deviation. In the case of yield (bushels/acre) the transformation is:  $y = \frac{y_{original} - 147.397}{48.169}$ .

## Benchmarking Models: Linear and Machine Learning Models

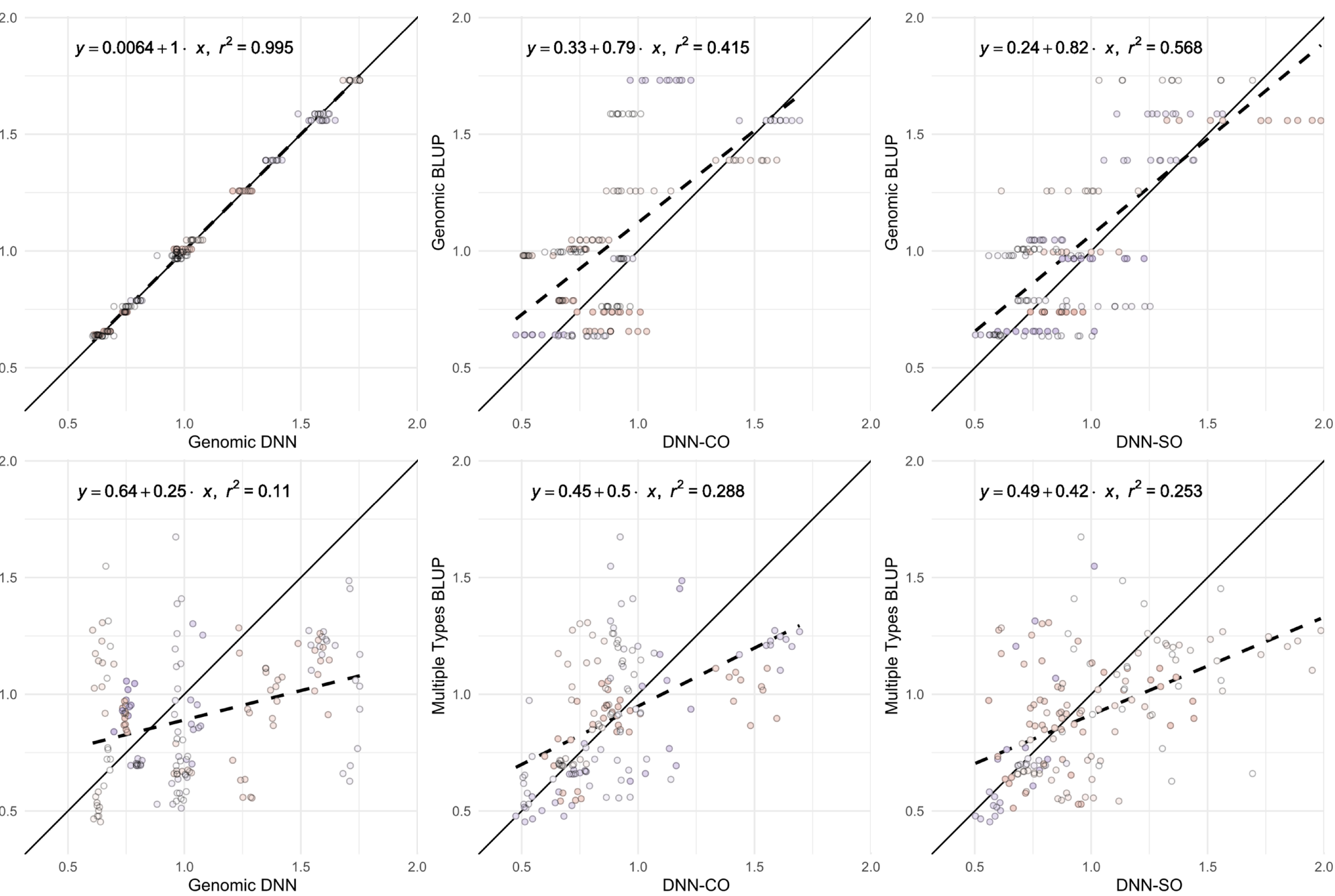
To assess model performance, linear models with fixed were trained on the same groups of data. These are not amenable to combining variables with one value per season (genetics, soil) with those variable within a season (weather and management). The latter variables were clustered into categories before use. For each variable, a time series k-means with dynamic time warping was fit for values of k from 2-40. k was selected as the value one less than the first k where the silhouette score decreased.

For each neural network's input data (exclusively genomic, soil, weather and management, or all three) Best linear unbiased predictor (BLUP) and machine learning models were created to contextualize model performance. BLUP were modeled on those in Washburn et al. 2021. In the model using multiple data types, genomic by soil and genomic by weather and management interaction effects were included. 4 machine learning models were considered: K Nearest Neighbors (KNN), Radius Neighbors Regression (RNR), Random Forest (RF), Support Vector Machine with a linear kernel (SVR). These models' hyperparameters which were optimized before use.

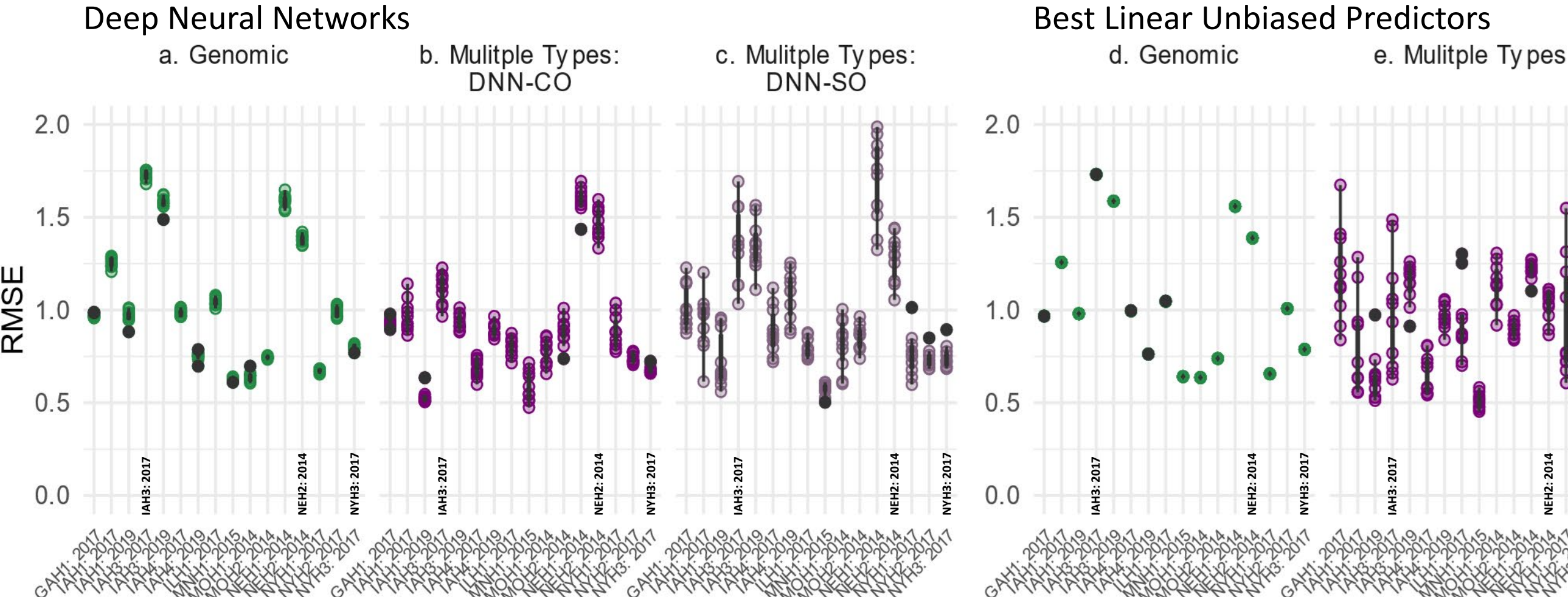
## Feature Importance: Saliency

The input variables that are most influential in the predicted yield for a given observation can be determined by calculating each variable's saliency. Saliency is calculated based off the derivative of each input variable for a given observation. Individual measures of saliency are averaged to produce the saliency maps shown here.

## Site-by-Year Group Average RMSE Between DNN and BLUP Models



## Differing Model Performance Across Site-by-Year Groups



Performance in different site-by-year groups differs with respect to model type (DNN vs BLUP). Lower correlation between model errors (e.g. Multiple Types BLUP vs DNN-CO) suggests potential benefit from ensembling model estimates. Point color indicates correlation within site-by-year group (lavender > 0).

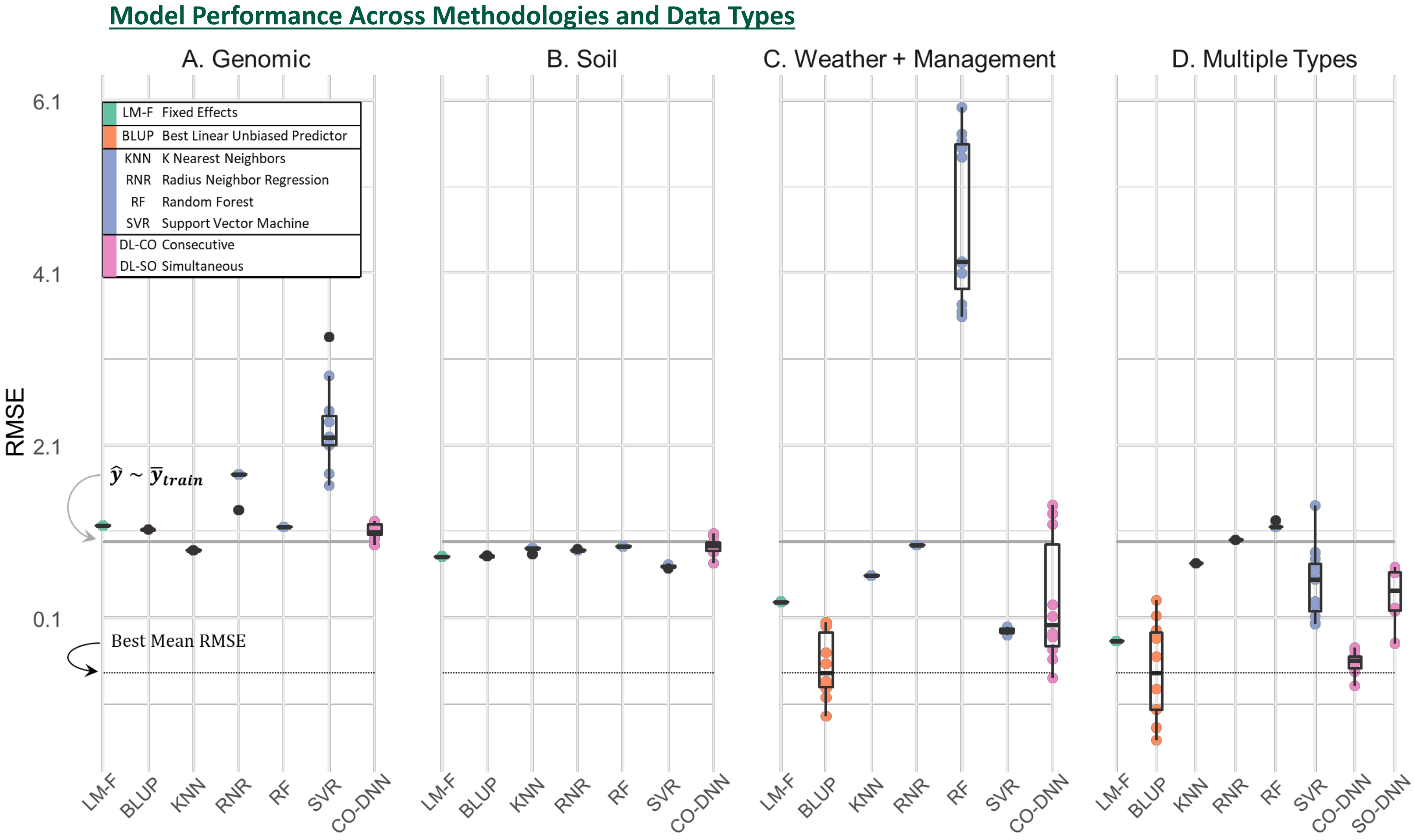
Daniel R Kick<sup>1,2\*</sup>, Jacob D. Washburn<sup>1,2</sup>  
USDA-ARS, Plant Genetics Research Unit, Columbia, MO<sup>1</sup>,  
Division of Plant Sciences, University of Missouri, Columbia, MO, USA<sup>2</sup>.



## Main Findings:

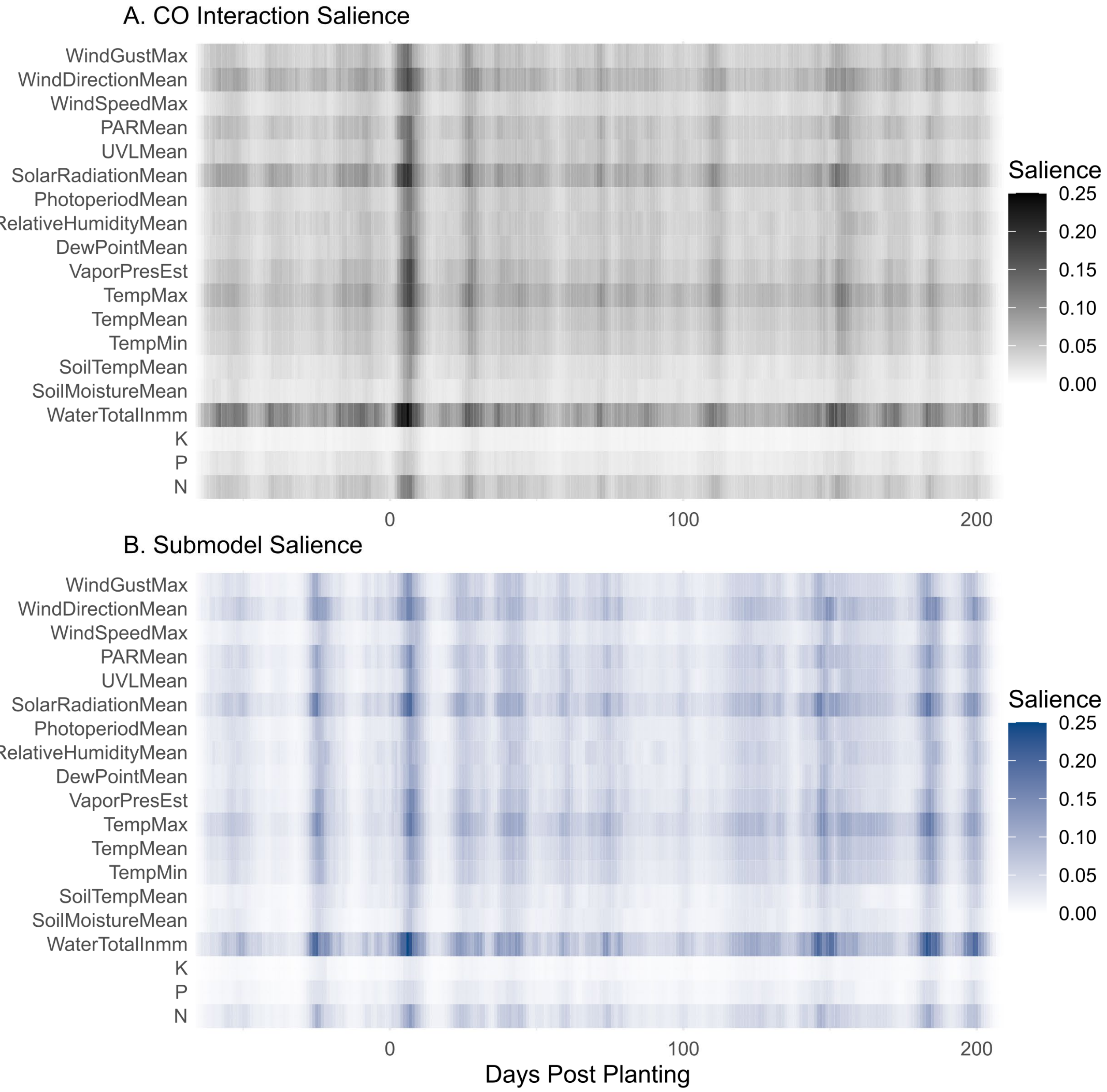
1. BLUP outperforms Deep Learning, but with greater variability in performance.
2. Optimizing sequentially improves performance.
3. GxExM interactions increase accuracy & alters variables' importance.

## Main Results



The root mean squared error (RMSE) of the testing set is shown for each data group (panels **A. – D.**) and class of model (linear models: green, BLUPs: orange, machine learning: blue, deep learning: pink). Lower values indicate better performance. Deep learning models are divided by whether they were part of the consecutive optimization strategy (DNN-CO) or the simultaneous optimization strategy (DNN-SO). LM-F use all only main effects except in D where interactions between PC1-8 and all weather and management and soil variables are included. BLUPs in D contain genome by soil and genome by weather and management interactions.

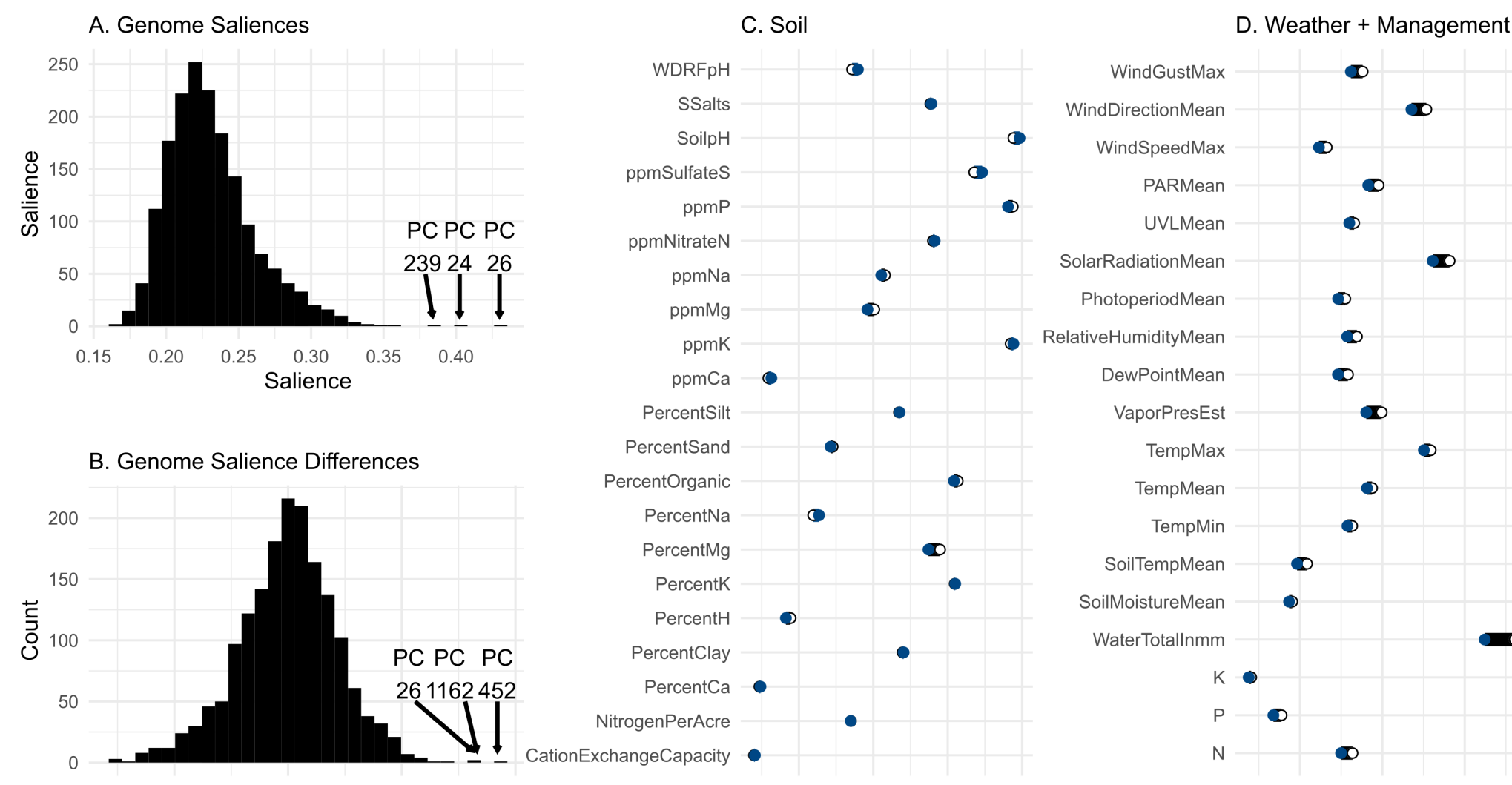
## Influence of Interaction Effects on Daily Feature Saliency



**A.** Average saliency for each day. Interaction model values shown. **B.** The same values for the weather submodel. Saliency peaks shortly after planting in both. The submodel contains salient dates prior to planting and near the end of the date range. The interaction-containing model appears to place greater importance for certain features, e.g., irrigation & rainfall, represented as “WaterTotalInmm”. The difference between the two saliency maps indicates additional times of sensitivity in the submodel (approximately -25, +180, +195) that the interaction model is relatively insensitive to.

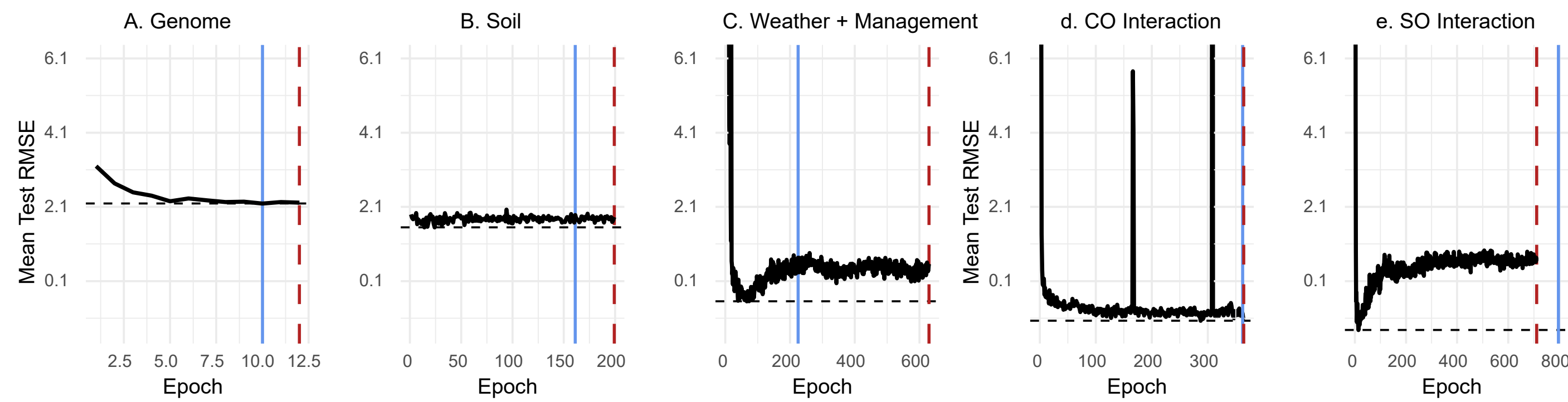
## Additional Results

## Influence of Interaction Effects on Feature Saliency in Aggregate



**A.** A histogram of the saliencies for Genomic PCs for the interaction containing DNN-CO model does not indicate a clear preference for PCs that explain more variance. Note that the two most salient PCs were PC 26 and PC 24 which account for 0.350% and 0.392% of the total variance respectively. **B.** The difference between saliencies of the full model and genomic submodel are shown. Those with the highest difference, explain little of the total variance. **C.** Saliency of soil variables is shown for the interaction model (black open circles) and the submodel (blue filled circles). The difference is shown in black if saliency is higher in the interaction model and blue if not. **D.** Saliency values for weather and management variables for the interaction model and submodel are shown as average saliency over the season. This indicates an overall similarity in saliency across features, with the most notable difference being “WaterTotalInmm”. Since that these values do not indicate the full effect of the weather and management variables in influencing predictions as these are daily values, not the total for the 210 day time window, resulting in smaller values in D than A - C.

## Optimization Strategy Results in Different architectures; degree of overfitting



Mean test set RMSE across 10 replicates (**A. – E.**). The horizontal dashed line is the minimum error. The vertical lines indicate the epochs selected by the chosen heuristic - minimizing total validation error (red dashed line). The mean plus standard deviation of validation error (solid blue line). Both strategies resulted in apparent overfitting in the Weather and Management submodel (C.) and the SO model (E.).

## Acknowledgements

This project was enabled with the mentorship of Dr. Jacob D. Washburn<sup>1,2</sup>. **USDA Project: 5070-21000-041-000-D** Provided funding for this study while the Genomes to Fields Intuitive (genomes2fields.org) and the DAYMET database (daymet.ornl.gov) provided the data for this study.

For their contributions in support of this study we would like to thank :  
Jason G. Wallace<sup>3</sup>, James C. Schnable<sup>4</sup>, Judith M. Kolkman<sup>5</sup>, Baris Alaca<sup>6, 7</sup>, Timothy M. Beissinger<sup>6, 7</sup>, Jode W. Edwards<sup>8</sup>, David Ertl<sup>9</sup>, Sherry Flint-Garcia<sup>1</sup>, Joseph L. Gage<sup>10</sup>, Candice N. Hirsch<sup>11</sup>, Joseph E. Knoll<sup>12</sup>, Natalia de Leon<sup>13</sup>, Dayane C. Lima<sup>14</sup>, Danilo Moreta<sup>5</sup>, Maninder P. Singh<sup>15</sup>, Teclemariam Weldekidan<sup>16</sup>

<sup>1</sup> United States Department of Agriculture, Agricultural Research Service, Columbia, MO 65211, USA  
<sup>2</sup> Division of Plant Sciences, University of Missouri, Columbia, MO 65211, USA  
<sup>3</sup> Department of Crop & Soil Science, University of Georgia, Athens GA 30602  
<sup>4</sup> Center for Plant Science Innovation & Department of Agronomy and Horticulture, University of Nebraska-Lincoln, Lincoln, NE 68583 USA  
<sup>5</sup> School of Integrative Plant Science, Cornell University, Ithaca NY 14853  
<sup>6</sup> University of Goettingen, Division of Plant Breeding Methodology, Department of Crop Science Breeding Research  
<sup>7</sup> United States Department of Agriculture, Agricultural Research Service, Ames, IA 50011, USA  
<sup>8</sup> Iowa Corn Promotion Board, Johnston, IA 50131 USA  
<sup>9</sup> Department of Crop and Soil Science, North Carolina State University, Raleigh, NC 27695, USA  
<sup>10</sup> Department of Agronomy and Plant Genetics, University of Minnesota, St. Paul, MN 55108  
<sup>11</sup> USDA-ARS Crop Genetics and Breeding Research Unit, Tifton, GA 31793 USA  
<sup>12</sup> Department of Agronomy, University of Wisconsin, Madison  
<sup>13</sup> Department of Agronomy - Plant Breeding & Plant Genetics University of Wisconsin - Madison Plant Breeding and Plant Genetics Program Madison, WI 53706  
<sup>14</sup> Plant, Soil and Microbial Sciences Dept., Michigan State University, East Lansing, MI 48824  
<sup>15</sup> University of Delaware Newark, DE 19716 USA