# Title: Yield Prediction Through Integration of Genetic, Environment, and Management Data Through Deep Learning

Authors:

Daniel R. Kick[1,2] (0000-0002-9002-1862),

Jason G. Wallace[3] (0000-0002-8937-6543),

James C. Schnable[4] (0000-0001-6739-5527),

Judith M. Kolkman[5] (0000-0001-7388-7245),

Barış Alaca[6,7] (0000-0001-7542-6514),

Timothy M. Beissinger[6,7] (0000-0002-2882-4074),

David Ertl[8] (0000-0003-1374-2970),

Sherry Flint-Garcia[1] (0000-0003-4156-5318),

Joseph L. Gage[9] (0000-0001-5946-4414),

Candice N. Hirsch[10] (0000-0002-8833-3023),

Joseph E. Knoll[11] (0000-0002-7477-4039),

Natalia de Leon[12] (0000-0001-7867-9058),

Dayane C. Lima[13] (0000-0002-5327-2160),

Danilo Moreta[5] (0000-0001-6866-9200),

Maninder P. Singh[14] (0000-0001-6166-2452),

Teclemariam Weldekidan[15] (0000-0003-4427-0099),

Jacob D. Washburn[1,2] (0000-0003-0185-7105)

1 United States Department of Agriculture, Agricultural Research Service, Columbia, MO 65211, USA

2 Division of Plant Sciences, University of Missouri, Columbia, MO 65211, USA

3 Department of Crop & Soil Science, University of Georgia, Miller Plant Sciences Building, 120 Carlton Street, Athens GA 30602

4 Center for Plant Science Innovation & Department of Agronomy and Horticulture, University of Nebraska-Lincoln, Lincoln, NE 68588 USA

5 School of Integrative Plant Science, 236 Tower Rd., Cornell University, Ithaca NY 14853

6 University of Goettingen, Division of Plant Breeding Methodology, Department of Crop Science

7 University of Goettingen, Center for Integrated Breeding Research

8 Iowa Corn Promotion Board, 5505 NW 88th Street, Johnston, IA 50131 515-225-9242

9 Department of Crop and Soil Sciences, North Carolina State University, Raleigh, NC 27695, USA

10 Department of Agronomy and Plant Genetics, University of Minnesota, 411 Borlaug Hall, 1991 Upper Buford Circle, St. Paul, MN 55108

11 USDA-ARS Crop Genetics and Breeding Research Unit, 115 Coastal Way, Tifton, GA 31793 USA

12 Department of Agronomy, University of Wisconsin, Madison

13 Department of Agronomy - Plant Breeding & Plant Genetics University of Wisconsin – Madison Plant Breeding and Plant Genetics Program 1575 Linden Drive Madison, WI 53706

14 Plant, Soil and Microbial Sciences Dept., Michigan State University, Plant and Soil Sciences Building 1066 Bogue Street, Room A286 East Lansing, MI 48824

15 University of Delaware Newark, DE 19716 USA

Correspondence: Jacob D. Washburn jacob.washburn@usda.gov

48

51

# Abstract

Accurate prediction of the phenotypic outcomes produced by different combinations of genotypes, environments, and management interventions remains a key goal in biology with direct applications to agriculture, research, and conservation. The past decades have seen an expansion of new methods applied towards this goal. Here we predict maize yield using deep neural networks, compare the efficacy of two model development methods, and contextualize model performance using linear models, which are the conventional method for this task, and machine learning models We examine the usefulness of incorporating interactions between disparate data types. We find a deep learning model with interactions has the best average performance. Optimizing submodules for each datatype improved model performance relative to optimizing the whole model for all data types at once. Examining the effect of interactions in the best performing model revealed that including interactions altered the model's sensitivity to weather and management features, including a reduction of the importance scores for timepoints expected to have limited physiological basis for influencing yield – those at the extreme end of the season, nearly 200 days post planting. Based on these results, deep learning provides a promising avenue for phenotypic prediction of complex traits in complex environments and a potential mechanism to better understand the influence of environmental and genetic factors.

# Introduction

Prediction of an organism's phenotype is a key challenge for biology, especially when integrating the effects of genetics, environmental factors, and human intervention. For many traits, prediction is complicated by interactions between these factors. For example, within a large multi-site, multi-genotype maize (*Zea mays*) study, more variation in grain yield is

3

76      explained by interactions between genetic and environmental factors than by genetic main

77      effects (Rogers *et al.* 2021). Including interaction effects between environmental and genomic

78      data can improve predictive accuracy in novel environments or for new cultivars (Li *et al.* 2021;

79      Jarquin *et al.* 2021).

80              Within agriculture, diverse methods have been applied to the task of predicting

81      phenotype ranging from classical statistics (Jarquin *et al.* 2021; Rogers *et al.* 2021; Rogers and

82      Holland 2021), machine learning (Westhues *et al.* 2021), physiological crop growth models

83      (Technow *et al.* 2015), to combinations of these and other methods (Messina *et al.* 2018;

84      Shahhosseini *et al.* 2021). Each model contains limitations such as lacking the capacity to

85      model complex non-linear responses (linear models) or interactions between factors,

86      interpretability within a biological framework (machine learning models), or dependence on

87      costly, low throughput data for calibration (crop growth models). Often simplifying assumptions

88      are introduced into the model (e.g. linearity), into the data (dimensionality reduction, feature

89      engineering), or into the experimental design (e.g. considering exclusively genetic,

90      environmental, or managerial effects to the exclusion of all others). While this approach creates

91      more manageable statistical models and enables a sufficiently powered study to be achieved

92      with fewer resources, it limits the capacity of a model to generalize to new genotypes,

93      environments, or management schemes. Furthermore, which factors are treated as "nuisance"

94      variables varies between communities within agriculture: geneticists often restrict management

95      regimes, while agronomists usually consider only a few cultivars. These approaches make it

96      difficult to investigate the interactions between genetic, environmental, and management

97      factors.

98              To predict an organism's phenotype across genotypes, environments, and management

99      strategies simultaneously requires a dataset containing many combinations of these features.

100     Collecting such a dataset requires a large multi-site, multi-condition, experiment featuring

101  diverse genetic backgrounds. The Genomes to Fields Initiative (McFarland *et al.* 2020) seeks to

102  accomplish this aim. To date it has collected measurements of grain yield and other phenotypic

103  traits (plant height, days to silking, stalk lodging, and kernel row number, etc.) from about

104  180,000 plots planted at more than 160 environments. Environments are characterized using a

105  WatchDog 2700 Weather station (Spectrum Technologies, Inc.) which collects continuous

106  weather data thought the season and collaborator submitted soil samples. Across the initiative,

107  over 2,500 maize hybrids have been tested, with Genotyping by Sequencing performed on

108  inbred parental lines used. Beyond the data collection, a means of effectively incorporating

109  diverse data types (genomics, management, soil measurements, weather, etc.) is needed,

110  particularly one that avoids simplifying assumptions where possible.

111      One method with the potential to accomplish this is that of deep neural networks (DNNs)

112  which have the capacity to approximate any function, provided they are sufficiently complex and

113  have sufficient examples to learn from. This capability is present regardless of whether they are

114  composed of dense fully connected (Hornik *et al.* 1989) or convolutional layers (Zhou 2020).

115  Additionally, DNNs "learn" directly from the data provided which enables reduced feature

116  engineering and dimensionality reduction. The methodology is also flexible with respect to data

117  type, allowing combination of variables that are static over a growing season (e.g. genotype)

118  and those that are dynamic (e.g. temperature) in a single model (Washburn *et al.* 2021). While

119  neural networks have been applied to the problem of predicting yield since at least 2001 (J. Liu

120  *et al.* 2001) this field is rapidly developing, with advances in theory, software, and hardware

121  enabling deeper and more accurate networks. Several recent studies have applied these

122  methods with a relatively large dataset either with (Washburn *et al.* 2021) or without (Khaki *et al.*

123  2020) a genetic component into the model, with little feature engineering performed. Both relied

124  instead on DNN's capacity to learn useful data transformations from the data directly.

125

126       Despite their promise, DNNs are not a panacea for prediction. DNNs are prone to

127    overfitting to training data resulting in poor performance. Even when performing well, the

128    complexity of these models can obscure what aspects of the data the model is using. Advances

129    in deep learning have produced methods which reduce these limitations. For example, the use

130    of convolutional layers minimizes the potential of overfitting because they perform well with

131    fewer parameters relative to fully connected layers. Where fully connected layers are used,

132    overfitting can be reduced by randomly removing neurons from a layer with a certain "dropout"

133    percentage. While the inner workings of DNNs remain far less interpretable than simpler models

134    (e.g., Genomic BLUP or physiological models), methods have been developed to aid in

135    interpretation through identifying the importance of different features in the data which can be

136    applied. These methods include salience (Simonyan *et al.* 2014), guided backpropagation

137    (Khaki *et al.* 2020), and permutation based metrics (Shahhosseini *et al.* 2021) among others

138    (Samek *et al.* 2017). Here we use salience to illuminate the operation of the DNNs generated in

139    this study.

140       Here we, leverage DNNs' capacity to determine feature importance from the data which

141    permits us to remain agnostic as to which features, or combinations of features are most

142    relevant. Furthermore, since DNNs are robust to lower--quality data and benefit from an

143    abundance of data, we employ a strategy of minimal feature transformation and curation and

144    maximal inclusion of observations. Using a minimally transformed dataset we begin the search

145    space considered in (Washburn *et al.* 2021), expand the space under consideration, and detail

146    a sequence of reproducible steps and objective heuristics which produced the models under

147    consideration. DNNs require an abundance of data for training. We begin by detailing a

148    workflow that incorporates a wider number of years from the Genomes to Fields Initiative in than

149    previous studies (Rogers *et al.* 2021; Washburn *et al.* 2021; Rogers and Holland 2021), while

150    also limiting the effect of errant and absent measurements. Improving on past studies, we

151   propose a new approach to model optimization whereby the model is broken into sub-modules

152   for each data type and interactions between them, then each submodule is consecutively

153   optimized, using a bayesian optimization procedure to find a suitable structure based on the

154   data itself. As far as we are aware, previous studies using deep learning for phenotypic

155   prediction have instead employed simultaneous optimization of all model components

156   (Washburn *et al.* 2021) or informal inductive tinkering. We compared models developed through

157   consecutive and simultaneous optimization and tested them against a variety of classic machine

158   learning and statistical methods to determine which performed best. To fairly assess model

159   performance we detail a strategy of constructing testing, training, and validation sets stratified

160   by season and location that is broadly useful to assessing model performance, while avoiding

161   overfitting the model to any location

# Materials and Methods

162

## Data Preparation

163

164       We used data from the Genomes to Fields (G2F) initiative for years 2014-2019

165   (McFarland *et al.* 2020), focusing on the sites within the continental United States. Each year's

166   data are publicly available (https://www.genomes2fields.org/resources/), including weather and

167   soil data for field sites, genomic data, management schedules (e.g., application of fertilizer,

168   herbicides, irrigation) and yield (in addition to other phenotypic variables). We augmented this

169   through additional genomic and weather data. Weather data retrieved from Daymet (Thornton *et*

170   *al.* 2020) was used in quality control as discussed below and to infer data for locations which

171   lacked a functional weather station for some or all of the season. Daymet data was retrieved

172   through wget (Techtonik 2015).

173  These data are provided with some variability in format. Custom scripts were used to

174 aggregate and standardize terminology across years. Rather than itemizing each operation, we

175 restrict ourselves to those which are likely to be of interest to those working with similar data

176 sets. The scripts used are available through Bitbucket

177 (https://bitbucket.org/washjake/maizemodel and https://bitbucket.org/daniel_kick/maizemodel/ ).

178 Scripts were written in Python (Van Rossum and Drake 2009 p. 3) and rely on scientific and

179 common general libraries (Seabold and Perktold 2010; Pedregosa *et al.* 2011; *fuzzywuzzy*

180 2017; Virtanen *et al.* 2020; team 2020; Harris *et al.* 2020; Da Costa-Luis *et al.* 2022) along with

181 plotting libraries for exploratory visualizations (Hunter 2007 p. 200; Inc 2015; Waskom 2021;

182 Kibirige *et al.* 2021). We used Anaconda ("Anaconda Software Distribution" 2021) to manage

183 the virtual environment.

184  We reduced the dimensionality of the genomic data with principal components analysis

185 (PCA) before use. Provided genomes were loaded into TASSEL version 5.2.74 (Bradbury *et al.*

186 2007) and filtered. Filter parameters were intended to be relatively unrestrictive, while still

187 reducing the data enough for PCA to complete with the memory available. We arrived at these

188 filtering parameters through inductive tinkering (i.e., iteratively increasing filtering until the

189 dataset was sufficiently small for the process to not exhaust available memory). First, we

190 restricted the loci to those having a heterozygosity of at least 0.001. Next, we converted the

191 data to a numerical genotype, and imputed missing values are imputed with the mean. The

192 quality of observations varied through the dataset. We discarded samples if they had missing

193 values for 90% or more of loci. Once the data was reduced, the genomes were PCA

194 transformed. We find that 31% of the variance is explainable by the first 8 principal components

195 (PCs), 50% is explainable by the first 50 PCs, and >99% of the variance is explainable by 1725

196 PCs.

197        Each hybrid's coordinates in PC space were estimated as the average between its

198    parent's coordinates. This was done rather than creating simulated hybrids due to hardware and

199    software constraints. If simulated hybrids are generated in TASSEL the number of observations

200    input into PCA is substantially increased (i.e., all observed combinations of genomes are

201    transformed, not all observed genomes). Completion of PCA would thus requires stricter filtering

202    for the analysis to be completed. By estimating hybrid values after PCA transformation we retain

203    a greater number of loci per observation influencing the PCs.

204        Environmental data required preprocessing as well. The soil dataset contains many

205    missing values, having an average completion rate of 47% across all site-by-year combinations.

206    For each variable in the soil dataset, missing values were first linearly interpolated across years

207    with respect to location. Locations with no observations for any years were imputed using k-

208    nearest neighbors based on the nearest 5 neighbors for remaining missing values. Within the

209    reported weather data, we observed evidence of equipment malfunction and imputed or

210    adjusted values using linear models.

211        First, we removed outliers and extreme observations by limiting the data to those days

212    where the difference between the minimum daily value estimated by Daymet and the measured

213    value were in the central 60% of the distribution. This removes values from the G2F reported

214    weather data (based on more inexpensive and error prone equipment) that are inconsistent with

215    the Daymet data. Next, linear models were generated using the 4 Daymet estimated variables

216    with the strongest correlations with the target variable across the dataset. All combinations of

217    additive models containing 0-4 of these predictors plus an optional site-by-year term were fit and

218    ranked using Akaike information criterion corrected for small sample size (AICc). The best

219    performing model for each metric was used to impute missing weather data. In the case of

220    temperature measurements, which exhibited the most apparent evidence of instrumental error,

221    observed values were also replaced if the difference between the expected and observed value

222    was greater than 1.5x the interquartile range of those differences.

223        The representation of management data was refined. Fertilizer applications were

224    decomposed into the quantity of nitrogen, phosphorus, and potassium applied. Where fertilizer

225    applications were lacking an application date, we estimated the time difference relative to the

226    planting date with K-NN imputation (k = 5) to cluster based on application quantity. To define the

227    time window under consideration, we used the earliest within-season fertilizer application and

228    the day of the latest harvest to bound the weather and management data. This resulted in

229    selecting 75 days prior to planting, 1 planting day, and the 204 subsequent days (210 total

230    days).

231        Weather and management time-series data were clustered to reduce their

232    dimensionality for use in machine learning and linear models. For each variable we used time

233    series k-means with dynamic time warping implemented through the tslearn library (Tavenard *et*

234    *al.* 2020). K was optimized by calculating the silhouette score for k between 2 and 40 then

235    selecting the k one less than the lowest k in which the silhouette score decreased. Where

236    needed clusters were represented categorically through one hot encoding.

## Defining Training, Validation, and Test sets

238        We generated train/test splits randomly, with the constraint that any location-year

239    combination could appear in only the testing or training set. Nearby experimental sites were

240    grouped for the purpose of generating training and testing sets. The method of generating splits

241    was: (1) one site group was selected at random and added to the testing set. (2) Each site-

242    group-by-year combination was down sampled so that it had no more observations than were

243    found in the smallest site-group-by-year combination in the testing set. This prevented

244    overrepresentation of any one group in the training data. (3) This process was repeated until the

245    testing set accounted for 10% to 15% of the total observations *and* at least 40,000 observations

246    remained across the training and test set. If these conditions were not met the process was

247    repeated with a different seed value.

248        During deep neural network hyperparameter selection, we limited overfitting to a single

249    validation set by using Monte Carlo cross-validation, stratified by site-group-by-year. Folds were

250    created while controlling for year by location groups by drawing the same number of site-group-

251    by-year groups as were present in the test set (3 for the test/train split used here) but without

252    further down sampling. The selected groups constituted the validation set. This was repeated 10

253    times and the membership of the folds preserved throughout training or hyperparameter search

254    for a single model. To enable reproducibility, random number generator seeds were used in

255    train/validate split searching and Monte Carlo cross-validation fold generation.

256        Prior to hyperparameter selection and training the input data was centered and scaled

257    based on the mean (~147.397 bushels per acre) and standard deviation (~48.169 bushels per

258    acre) of the yield in the training data, i.e., $y = \frac{y_{Original} - 147.397}{48.169}$. Transforming the data prior to

259    determining cross validation folds potentially introduces an information leak within the

260    hyperparameter selection process (i.e., the data used in evaluation, here across cross validation

261    folds, by influencing the mean and standard deviation used in centering and scaling the data).

262    However, this does not create an issue for final evaluation of the model's performance because

263    the test set data were not used in these calculations.

# Model Preparation

264

**Overview**

265

266        We sought to model genotype by environment by management interaction effects (GEM

267    effects) in maize yield and to determine utility of doing so. To this end we optimized DNNs to

268    predict yield with a single data modality (i.e., only genomic data, soil characteristics, or time

269    series data each by itself). We use one dimensional convolutional layers to capture the time

270    dependent features of weather data, which have previously been used in yield prediction for this

271    task (Khaki *et al.* 2020; Washburn *et al.* 2021). We used dense, fully connected layers for the

272    other submodules of the DNN.

273        We pursued two strategies for tuning and training GEM models: Consecutive Optimization

274    (CO) and Simultaneous Optimization (SO). CO tunes the hyperparameters of networks predicting

275    yield from a single data modality (genomic data, soil data, or weather and management time

276    series data). Next, the prediction neurons are discarded and the output of the penultimate layer

277    of each single modality network enters a set of layers to permit interactions between data

278    modalities. Hyperparameters for the interaction layers are then tuned. The SO strategy by

279    contrast allows for all hyperparameters to be selected concurrently, both those which affect the

280    processing of a single data modality and those influencing interactions between modalities.

281    **Hyperparameter Search and Training**

282        We selected model architecture through a hyperparameter search using the

283    `BayesianOptimization` tuner provided within the `keras-tuner` package (O'Malley *et al.* 2019).

284    Models were written in Keras (Chollet and others 2015) with Tensorflow as a backend (Martín

285    Abadi *et al.* 2015) and run in a Singularity container (Kurtzer *et al.* 2017; SingularityCE

286    Developers 2021). The subnetworks processing exclusively genomic and exclusively soil data,

287    along with the interaction subnetwork, are constructed exclusively of dense (i.e., fully

288    connected) layers, each subject to batch normalization and dropout. The weather/management

289    processing subnetwork is composed of two one-dimensional convolution layers followed by

290    batch normalization and a pooling layer. The output of this subnetwork is flattened before

291    entering the interaction subnetwork. Hyperparameter ranges explored for each network are

292    listed in Table 1. To avoid overfitting to the validation data we used a custom subclassed

293    version of the tuner to randomly select one of the previously defined validation folds. This is

294    done rather than using mean loss over all folds to avoid increasing the computational cost 10-

295    fold while still preventing overfitting to a single validation set.

296    For all DNNs, a maximum of 40 hyperparameter sets were explored. In cases where no

297    convolution layers were being varied (CO models with only genomic, or only soil data, and the

298    interaction layers trained for the same strategy) hyperparameters were trained for a maximum of

299    1000 epochs with an early stopping patience of 7 epochs or more. For cases in which

300    convolution layers were varied (Sequential Optimization model with only weather and

301    management data, Concurrent Optimization model) hyperparameters were trained for a

302    maximum of 500 epochs with an early stopping patience of 5 epochs. This difference is due to

303    practical rather than theoretical reasons as the convolutional networks required notably more

304    time per epoch to fit. Regardless of network type, if the hyperparameters optimization had not

305    concluded by 290 hours after the script began, the process was terminated and the

306    hyperparameter sets completed by that point were considered.

307    To ensure the selected hyperparameter set performs well across validation sets, the top

308    4 hyperparameter sets for each model were trained for 1000 epochs and evaluated on all 10

309    defined testing/validation set splits. Next, the validation losses over the duration of training were

310    used to calculate the mean and standard deviation for each epoch. Then the training duration

311    was split into 10 bins and the average of the sum of validation loss mean and standard deviation

312    was calculated, i.e. $loss_{bin} = \frac{\sum_{i=1}^{n} \overline{l_i} + s_i}{n}$ where $i$ is epoch relative to the beginning of the bin, $\overline{l_i}$ is

313    the mean validation loss across cross validation folds at the i[th] epoch and $s_i$ is the standard

314    deviation of the same. The hyperparameter set with the lowest value for the most bins was used

315    going forward.

316    For the best hyperparameter set, we selected a training duration from the validation

317    losses. For each fold we calculated a rolling mean of validation loss with a window size of 20

318    epochs. Next, for each epoch we calculated the sum of the mean and standard deviation of the

13

319    rolling mean and the total rolling validation loss. Then we found the epochs which minimized

320    these two values (subtracting 10 from the epoch number to account for the window size). The

321    disagreement between the epochs which minimized these values ranged from 2 epochs in the

322    case of the CO Genomic model and CO interaction model up to 404 epochs for the CO weather

323    and management model. We decided to use total rolling validation loss to decide on the epoch

324    number for each model. This metric resulted in more training epochs for all models except the

325    SO model. Incorporating a more sophisticated method for selecting training duration is a

326    possible improvement for future studies. With the selected hyperparameters and training

327    duration we fit each model 10 times to account for random initialization and saved each

328    replicate and its training history.

## Benchmarking Models

## Overview

331        To contextualize the performance of the generated deep neural networks we use the

332    same training data to fit linear and classic machine learning models. These models often require

333    fewer resources and time to train than deep neural networks. For linear models we consider a

334    small collection of models that varied with respect to the independent variables present, whether

335    interactions are included, and whether effects are fixed or random. For supervised machine

336    learning models, we selected and optimized four methods: k-nearest neighbor (KNN), radius

337    neighbor regression (RNR), random forest (RF), and support vector regression with a linear

338    kernel (SVR).

## Linear Models

340        To aid in evaluating the efficacy of the models produced, we constructed linear models

341    varying in the scope of included data and model complexity. The simplest model was an

342     intercept model, i.e., every predicted yield equals the mean yield in the training set. We

343     considered three models using the genomic data alone: fixed effects for PCs 1-8 (31% variance

344     explained), fixed effects for PCs 1-50 (50% variance explained), and random effects for the first

345     8 PCs. For soil data we considered two models, one with all factors as fixed effects and one

346     with all factors as random effects. From weather and management data we produced three

347     models, using all factors as fixed effects, using clusterings of the top five most salient features

348     (i.e., Water total, solar radiation mean, maximum temperature, mean wind direction, vapor

349     pressure) identified in the weather and management data (averaging over time points) as fixed

350     effects, and as random effects. Most salient features were taken from the deep neural network

351     with the lowest average test set RMSE (CO Interaction). All weather and management data

352     were represented as categorical clusters as described in "Data Preparation".

353         We evaluated five models using a combination of data sources. In three fixed effect

354     models we incorporated: (1) PCs 1-8 and the five most salient weather features' clusterings, (2)

355     the same plus all soil features, and (3): the effects in "1" plus interactions between each PC and

356     weather factor cluster. The two random-effects models we fit using PCs 1-8 and the selected

357     weather clusters excluding or allowing interactions. We fit Fixed effect models with the linear

358     model function in R (R Core Team 2021) and random effect models with lme4 (Bates *et al.*

359     2015). This analysis was aided by common data wrangling and convenience libraries (Wickham

360     *et al.* 2019; Bache and Wickham 2020; Müller 2020; Izrailev 2021) and feather file read/write

361     capabilities through arrow (Richardson *et al.* 2021).

362 ## Classical Machine Learning Models

363         Additional machine learning models were implemented through scikit-learn (Pedregosa

364     *et al.* 2011; Buitinck *et al.* 2013) and hyperparameters for each were optimized through the

365     hyperopt library (Bergstra *et al.* 2013) run within a Docker container. In a workflow similar to that

366     of the deep neural network models, we generated models for each data modality indpendently,

367    and with all data available. Time series data was represented as clusters as described in "Data

368    Preparation". For each model we allowed the following hyperparameters to vary as described:

369    (1) K Nearest Neighbors (KNN): neighbors = 1-250, weights = 'uniform' or 'distance'; (2) Radius

370    Neighbors Regressor (RNR): radius = 0.01-2000, weights = 'uniform' or 'distance'; (3) Random

371    Forest (RF), maximum depth = 2-200, Minimum samples per leaf = 0-0.5; and (4) Support

372    vector machine with a linear kernel (SVR): Loss = 'epsilon_insensitive' or

373    'squared_epsilon_insensitive', C =1-5 (log uniformly distributed).

374        Cross validation folds matching those as described previously and average loss across

375    all folds was measured. We tested a minimum of 115 combinations for each model and selected

376    the best performing hyperparameters for each input dataset, reported in Table 5. Following

377    selection, we trained each model and produced predictions on the testing and training data. This

378    was repeated 10 times to account for randomness in model fitting.

## 379    Model evaluation

380        For every model described above we calculate predicted yields for the test set and

381    calculate root mean squared error ($\text{RMSE} = \sqrt{\frac{\sum_{i}^{n} Prediction_i - Observation_i}{n}}$ ), normalized RMSE

382    percent ($\text{nRMSEnRMSE} = 100\% * \frac{RMSE}{\left(\frac{\sum_i^n Observation_i}{n}\right)}$ ) and $R^2$ using SciPy (Virtanen *et al.* 2020).

383    Unless stated otherwise in the text RMSE and nRMSE will refer to the average value across

384    replicates. Two observations were not predictable using the fit radius neighbors regressor and

385    were predicted as the training set mean. For the best performing DNN we calculated and

386    visualized the salience of features for each data modality. To examine the influence of allowing

387    interactions we contrast these saliences with the saliences of SO single modality DNNs.

388    Saliences were calculated by Tf-keras-vis (Kubota 2021). Visualizations were created with the

389    use of rjson (Couture-Beil 2018), patchwork (Pedersen 2020), and ggplot2 (Wickham *et al.*

390    2019).


# Data Availability Statement

392         Data for maize phenotypes, genotypes, field site soil properties and on location weather

393    recordings from the Genomes to Fields Initiative data (McFarland *et al.* 2020) is publicly

394    available through the CyVerse Discovery Environment. We used data from 2014 to 2019 which

395    correspond to the following DOIs: 2014 – 2017 (0.25739/frmv-wj25), 2018 (10.25739/anqq-

396    sg86), and 2019 (10.25739/t651-yy97). Additional genomic data was provided by Natalia de

397    Leon, Dayane Lima, and Cinta Romay worked with Joseph Gage through personal

398    communication. These data will be available through CyVerse. Following public release, a

399    version of these data containing all genomic data used in this study will be available through

400    Zenodo (10.5281/zenodo.6916775). At time of writing, this repository contains the eigenvectors

401    resulting from the principal components analysis are provided to enable transformation of

402    provided genomes. Additional weather measurements were retrieved from Daymet (Thornton *et*

403    *al.* 2020). Custom python scripts for downloading, aggregating and processing these data are

404    available on bitbucket (https://bitbucket.org/washjake/maizemodel and

405    https://bitbucket.org/daniel_kick/maizemodel/ ) in the notebooks directory (files with the prefix

406    0.0 to 0.5).

407

# Results

## Deep Neural Networks can--but do not necessarily--outperform

## competing model types

When all data sources are incorporated, the CO DNN achieves the best (i.e., lowest) average RMSE (when not otherwise specified, values refer to the average across replicates), followed by a fixed effect model with interaction effects (nRMSE 14.6% vs 14.9%, RMSE 0.948 vs 0.959) (Figure 2, Table 6). Despite this, with 2/10 replicates of the CO DNN model underperform this fixed effect model. Variability in DNN performance across replicates can be caused by the random initialization of the weights at the start of training.

Following CO DNN and a linear model with interaction effects, an exclusively additive linear model incorporating all data modalities ranked third (nRMSE 15.0%, RMSE 0.980) and one excluding soil data ranked fourth (nRMSE 15.1%, RMSE 0.981). Random effects models with and without soil data (nRMSE 15.2%, 15.3%, RMSE 0.991, 0.994) and SO DNN (nRMSE 15.7%, RMSE 1.024) followed. Of the machine learning models only support vector regression with a linear kernel (SVR) and K Nearest Neighbor (KNN) outperformed a simple intercept model. We find similar results for $R^2$ (Supplemental Figure 2, Table 6).

A DNN is not the best performing model when data is restricted to a single modality. When restricted to genomic data, KNN, followed by linear models with fixed or random effects are the only ones which outperform the intercept model (nRMSE 16.5%, 16.6%, 16.7%, 16.7%, RMSE 1.078, 1.084, 1.085, and 1.088, KNN, linear fixed effects, linear random effects, intercept model). SVR performed particularly poorly on this data (nRMSE 18.7%, RMSE 1.212) -- nRMSE 2% or RMSE of 0.131 above the intercept model. Incorporating only soil data SVR performed best (nRMSE 16.3%, RMSE 1.059) with all other models being within nRMSE 0.261% or 0.017 RMSE of the intercept model. Most models performed better when instead trained on weather/

18

432    management data with the exception of the random forest (RF) which had an nRMSE 5.729%,

433    RMSE of 0.373 *above* the intercept model. SVR (nRMSE 15.1%, RMSE 0.985) and a fixed

434    effects model (nRMSE 15.2%, RMSE 0.993) performed remarkably well. CO DNN is capable of

435    outperforming these methods, but not uniformly. 2/10 replicates underperformed the intercept

436    model resulting in a nRMSE 15.6%, RMSE of 1.018 while the median values were nRMSE

437    15.2% and RMSE is 0.992.

# Consecutive Optimization resulted in a larger, more accurate final network.

440        Two hyperparameter selection strategies were employed, Consecutive Optimization

441    (CO) and Simultaneous Optimization (SO), have the same range of possible networks

442    (hyperparameter ranges are listed in Table 1), the same data driving network selection and both

443    use bayesian optimization. Despite this, the strategy applied resulted in notably different final

444    architectures. A visual summary of the relative differences between network hyperparameters is

445    shown in Figure 1, with the hyperparameter values listed in Tables 2 and 3. Supplementary

446    Figure 1 provides a visual overview of the network architecture. We consider the effect of CO vs

447    SO on each of the four subnetworks (processing exclusively genomic, soil, or

448    weather/management factors or interactions between data modalities), listed in decreasing

449    order of approximate similarity.

450        The output of this subnetwork is flattened before entering the interaction subnetwork

451    genomic subnetworks resulting from CO and SO are both two layers, but the CO model widens

452    somewhat (layer 1 = 83 units, 16% dropout, layer 2 = 133 units 23% dropout) while the SO

453    model begins over twice as wide and constricts more (layer 1 = 196 units, 15% dropout, layer 2

454    = 47 units 6% dropout). The interaction subnetworks contained a similar number of layers (CO:

455    5 vs SO: 6), but while CO resulted in layers with similar widths before constricting at the last

19

456    layer (units = 152, 207, 206, 188, 44, dropout percentages = 19%, 29%, 0.5%, 20%, 24%), SO

457    resulted in layers with very few units initially which are later expanded (units = 10, 25, 126, 204,

458    45, 134, dropout percentages = 10%, 15%, 2%, 16%, 24%, 19%). The soil subnetwork resulting

459    from CO is notably deeper than the one from SO (7 and 2 dense layers respectively) but also

460    narrows more by the last processing layer (2 vs 27 units). Finally, in the weather and

461    management subnetwork CO resulted in a notably deeper network (6 pairs vs 2 pairs of

462    convolution layers) but used a similar number of filters in the final convolution layer pairs (CO

463    294 vs SO 303).

464        The performance of these networks differs as well. The CO network was better at

465    predicting yield in the testing set. It achieved a lower mean RMSE (CO: 0.948 vs SO: 1.024)

466    and was more consistently accurate across replicates (standard deviation CO: 0.013 vs SO:

467    0.035). Similar results were seen in the normalized errors (nRMSE CO: 14.6% SO: 15.7%,

468    standard deviation CO: 0.197%, SO: 0.531%). Similarly, average $R^2$ was higher in the CO

469    network (CO: 0.171 vs SO: 0.032) and more consistent across replicates as well (standard

470    deviation CO: 0.022 vs SO: 0.065).

471        Model performance differences are due, in part, to the heuristic used to select the

472    number of training epochs and different tendencies for these models to overfit. The heuristic

473    used to select the number of training epochs (sum of the rolling validation loss) and alternate

474    heuristic considered (mean plus standard deviation of the rolling validation loss) resulted in

475    networks with comparable performance, having on average 0.001 less RMSE. With the

476    exception of the SO DNN, this also resulted in longer training durations. These ranged from an

477    additional 2 epochs in the cases of the CO genomic and interaction models and as many as 404

478    epochs in the case of the CO weather/management, as shown in Table 4.

479        These training durations were often considerably longer than the optimal values as seen

480    in Figure 1B. Furthermore, the length of overtraining appears loosely proportional to the present

20

481 minimum average RMSE each model achieved. The SO and CO weather models had the

482 largest differences between optimal and used epoch numbers – differences of 697 and 563

483 epochs respectively and achieved 121% and 110% of the minimum possible RMSE. The CO

484 soil model trained an excess 185 epochs but only had RMSE at 102% minimum. The two

485 training durations closest to the optimum were the CO genomic model (2 epochs over) and the

486 SO model (77 epochs over). These models performed at just 100.2% and 101% minimum.

487 The SO model overfits faster and to a greater extent than the full CO model, which does

488 not show evidence of substantial overfitting (Figure 1B d, e). The SO model achieves a loss

489 lower than the CO model, and the accuracy worsens rapidly with further training. The different

490 network sizes (CO containing more layers) may account for this difference. Improved heuristics

491 for training duration could represents an opportunity for future refinements, which these results

492 suggest could both increase goodness of fit and reduce the computational resources needed to

493 train these models.

494

495 # Model performance generally improves through incorporating

496 # multimodal data and interactions

497 Incorporating multiple data sources and allowing interactions between data types

498 generally appears to improve accuracy. Within the tested DNNs, allowing interactions increased

499 performance relative to single modality models. The potential exception to this is the SO DNN

500 (nRMSE 15.7%, RMSEs 1.024) and the CO weather/management model (nRMSE 15.6%,

501 RMSE 1.018). Despite this, the former's distribution had lower dispersion with a standard

502 deviation of RMSEs 0.035 relative to 0.074. Within the linear models tested, allowing

503 interactions increased accuracy in the fixed effect model by 0.3% nRMSE or 0.023 RMSE but

504 *decreased accuracy* in random effects models by 0.04% nRMSE or 0.003 RMSE.

21

505    In purely additive linear models, incorporating additional data modalities decreases error.

506    The largest difference in fixed effect models is for genomic data (improvement of 1.916%

507    nRMSE, 0.125 RMSE), followed by soil data (improvement of 1.722% nRMSE, 0.112 RMSE),

508    and weather/management data (improvement of 0.531% nRMSE, 0.035 RMSE). The same

509    trend is seen in models with random effects models, albeit with less variation in improvements

510    (improvements of 1.436% nRMSE, 0.094 RMSE genomic, 1.349% nRMSE, 0.089 RMSE soil,

511    1.132% nRMSE, 0.074 RMSE weather/management).

512    This pattern does not hold for the machine learning methods tested. For KNN, the model

513    trained on exclusively weather data performed best (0.218% nRMSE, 0.014 RMSE better than

514    using all data) although using all data sources did improve accuracy relative to only genomic or

515    soil data. SVR follows a similar pattern but is more exaggerated with using exclusively weather

516    data resulting in an improvement of 0.862% nRMSE or 0.059 RMSE relative to all data whereas

517    all data represented an improvement relative to genomic and soil data. Random forests did not

518    follow this trend– Genomic and soil models performed better than all data by 0.026% and

519    0.375% or 0.002 and 0.024 RMSE respectively, whereas the weather and management model

520    performed 5.439% nRMSE or 0.354 RMSE worse. Finally, radius neighbor regression (RNR)

521    performance was worst using only genomic data (1.112% nRMSE, 0.072 RMSE worse than all

522    data) but using only soil or only weather data improves model accuracy by 0.185% and 0.090%

523    nRMSE or 0.012 and 0.006 RMSE respectively.

## 524    Which factors are most important to the CO DNN?

525    Among the genomic data we observe no major trend in salience with respect to PC

526    (Supplementary Figure 3 A.). The two most salient PCs are PC 26 (0.423) and PC 24 (0.402)

527    which account for 0.350% and 0.392% of the total genomic variance respectively. Given that

528    these saliences are relative to principal components, using salience to implicate specific genes

22

529    or gene loci is infeasible. Among the soil factors we find that the five with the highest average

530    salience were soil pH (0.488), phosphorus ppm (0.487), potassium ppm (0.485), sulfate ppm

531    (0.436), and percent organic matter (0.413) (Supplementary Figure 3 C).

532    Within the weather and management data, considering the average salience across the

533    season (Supplementary Figure 3 D) five factors achieved an average salience greater than

534    0.140 – Total water (0.245), average solar radiation (0.198), maximum temperature (0.175),

535    average wind direction (0.174), and estimated vapor pressure (0.173). The majority of factors

536    had an average salience between 0.140 and 0.10 with six falling below this threshold – average

537    soil temperature (0.095), maximum wind speed (0.084), average soil moisture (0.076),

538    phosphorus applied (0.052) and potassium applied (0.033). Additionally, we find specific time

539    points which appear to be salient broadly with the most salient region of time is within the first

540    few days of planting, indeed 8 of the 10 days with the highest average salience are days 2-9

541    following planting.

## 542  How is factor importance altered by inclusion of interactions?

543    The full CO model, in addition to performing best (albeit by a small margin), presents an

544    opportunity to directly compare the influence of interactions between data modalities on the

545    salience of factors because the single modality subnetworks are identical except for the

546    prediction layer. The salience of genomic factors differs notably between the two networks

547    (Supplementary Figure 3 B). Salience of PCs differs by as much as 0.432 (PC 24), with the

548    difference in salience of the first 8 PCs (31% variance explained) ranging from 0.200 (PC1) to

549    0.309 (PC7). We find comparatively small differences in the salience of soil factors being

550    between -0.011 and 0.0156 (Supplementary Figure 3 C).

551    In general, the salience map of the weather and management data features fewer

552    broadly salient timepoints when interactions are included (Figure 3 A) than when they are not

23

553 (Figure 3 B). The weather and management CO model contains a broadly salient time point

554 around 25 days before planting and 6 days after planting. The SO model also appears to have

555 peaks of salience around 150, 183, and 199 days after planting. When interactions are included

556 the majority of the salient time points become less so with the exception of the peak 6 days after

557 planting as highlighted through subtraction of the two salience maps (Figure 3 C).

558

# Discussion

## Assumptions, Potential Sources of Error, and Opportunities for Improvement

562 The results of this study are best understood with the data used and assumptions made

563 kept in mind. The sole source of biological data in this study came from the Genomes to Fields

564 Initiative (McFarland *et al.* 2020). The scale of this ambitious project increases the chances of

565 data being absent or compromised due to equipment malfunction, logistical or procedural

566 issues, and resource constraints. For example, many sites lack measurements for many soil

567 properties across the seasons considered here, and the timing of fertilizer applications was

568 absent in some cases. Our aim was to minimally filter the dataset while preventing missing or

569 distorted values (many of which are not missing at random) from altering model accuracy and

570 feature salience. We have aimed to reproducibly infer missing or aberrant values with relatively

571 simple methods (e.g., imputation using linear models, KNN, etc.) but more sophisticated

572 imputation techniques may have improved performance.

573 Alternatively, constraining the dataset to reduce the required imputation may have been

574 an effective strategy. We elected to minimally filter observations because machine learning

575 models, particularly deep learning, often benefit from having an abundance of data from which

24

576    to learn feature relationships. For models where this is not the case, restriction of observations

577    to the observations with the highest quality may be a preferable strategy. Note, however, that for

578    distortions that are not randomly distributed, filtering may bias the sample and result in a model

579    that appears to perform well but generalizes poorly (e.g., to sites similar to those with a

580    preponderance of observations excluded).

581        Beyond including as many distinct locations and seasons as we could, we approximately

582    balanced site by year groups through down sampling to avoid overfitting our DNNs to sites with

583    more observations or biasing the selection of hyperparameters. This reduces the size of the

584    dataset that can be used in training. Although outside the scope of this study, assessment of the

585    sensitivity of DNNs to unbalanced group sizes, or exploration of alternate means of balancing

586    groups (e.g., randomly *up sampling* small groups to equal the size of larger groups) would be

587    valuable. Indeed, if balance were not a concern, or if it could be effectively achieved without

588    discarding observations in some groups, one could potentially employ more strict data filtering

589    without producing a dataset too small to benefit from machine learning.

590        Substantial effort was devoted to producing testing, training, and validation sets that

591    would not lead to overconfidence in the accuracy of our models. To this end we kept

592    observations within site-by-year groups in the same partition of the data. In effect, this prevents

593    the model from being trained and tested on the same weather and management data.

594    Furthermore, except in cases where soil features are static from season to season, the model

595    will not be trained and tested on observations with identical soil features. Proceeding in this

596    manner rather than selecting observations at random for the testing set further reduces an

597    already small number of weather and management conditions. Incorporating historical data

598    (Washburn *et al.* 2021) or expanding the dataset to include data from other sources represent

599    two possible avenues to incorporate a greater diversity of weather and management conditions

600    without compromising the testing set.

601     Depending on the intended application of a model, one may be able to achieve higher

602     performance through altering some of the above decisions or replacing random assignment with

603     a targeted approach. For example, we assume that all group-by-year combinations are equally

604     likely to be of interest. However, if we assume that the distribution of sites collected match those

605     of interest for prediction (i.e., one is interested in predicting *any* future observation collected by

606     G2F and the number of observations per field site are representative of future number of

607     observations) then down sampling can be skipped, resulting a larger dataset. Similarly, with a

608     narrower aim, e.g., prediction of yield within a specific region, testing or validation sets could be

609     constrained to better select hyperparameters for or assess predictive accuracy of site-by-year

610     combinations within that region.

611     In summary, our decision to include as much data as possible and to limit the possibility

612     of overfitting to specific sites and seasons represent possible opportunities for improvement.

613     More sophisticated data imputation or more restrictive filtering, alternate means of balancing

614     groups, and the incorporation of other data sources have the potential to improve model

615     performance. Additionally, for more narrowly purposed models, non-random testing and training

616     sets may represent a more accurate metric of predictive power, and indeed may deviate

617     substantially from what we show here.

618     # Tradeoffs in Model Performance and Computational Resources

619     While the best performance was achieved with a deep neural network incorporating

620     genomic, soil, weather, and management data, simple linear models with fixed effects often

621     performed nearly as well (Figure 2). This is notable because tuning and training deep neural

622     networks requires significant computational resources and time. For example, hyperparameter

623     tuning in the machine learning models shown here took less than 24 hours to complete whereas

624  tuning a single DNN sub-network took up several days. In the case of the best performing model

625  this was repeated four times – once for each sub model.

626  By contrast, linear models, particularly those with only fixed effects, are quick to fit. They

627  also outperformed many of the machine learning models, despite not undergoing extensive

628  tuning for model structure. In cases where accuracy is not the sole factor under consideration,

629  or where time or computational resources are limiting, simpler models may be "good enough"

630  for the desired purpose.

# 631  Usefulness of Consecutive Optimization in Hyperparameter

# 632  Selection

633  We employed two strategies for hyperparameter optimization: consecutively optimizing

634  (CO) hyperparameters for distinct "modules" of the network and simultaneously optimizing (SO)

635  the network as a whole. CO reduces the range of possible combinations that are explored by

636  allowing only one module to vary at a time. However, if two features in different data sets have a

637  strong interaction effect (e.g., between genotype and weather patterns) then this approach will

638  not necessarily allow for optimization to better capture this interaction. SO represents the

639  reverse situation. With all features available, interactions between features in different tensors

640  can be leveraged, but the hyperparameter space to explore is larger as all the hyperparameters

641  are free to vary.

642  We find that the network resulting from CO substantially outperforms the one generated

643  through SO. This should not be taken as a problem with SO *per se*. In other applications, or with

644  a different optimization algorithm, it may prove to be a more efficient means of deriving a useful

645  architecture. Furthermore, it is conceivable that SO is effective but that additional trials were

646  required. The SO DNN architecture was selected based on 40 trials whereas the CO DNN

647  architecture was selected based on 40 trials *for each module* (160 trials across the whole

648    network) which confounds comparison. Selection of the training duration also warrants

649    consideration. The SO model is capable of performing comparably to the CO model, but overfits

650    more rapidly (Figure 1 B). Improved heuristics for selecting the training duration could increase

651    usefulness of the SO model while reducing computational demands as well.

652        As a pragmatic matter, CO benefits from the capacity to tuning multiple modules at once.

653    In our hands, total time spent tuning was driven more by modules with computationally intensive

654    components (e.g. convolution layers) rather than the number of modules to optimize. This

655    benefit is dependent on the tuning algorithm used. We used a bayesian optimization procedure

656    which aims to produce useful hyperparameter combinations in fewer cycles than a simpler

657    method such as grid approximation. However, because this method uses the performance of

658    previously evaluated hyperparameters in selecting the next set, it does not permit parallelization

659    in tuning a single network. If an optimization procedure that is conducive to parallelization were

660    used (e.g., Hyperband or grid approximation) with enough computational resources this benefit

661    would be non-existent.

662        Although we aimed to broaden the range of possible architectures relative to previous

663    modeling on G2F data (Washburn *et al.* 2021), we constrained the overall structure to

664    processing each tensor individually then allowing for interactions between the final layer of each

665    module. Other options might include, for example, allowing an interaction module use both the

666    first and final layers as input (instead of only the final one), or allow which layers were to be

667    used to be tuned.

668        An additional option that we did not explore is aiming to inform the structure of the

669    selected network based on known relationships between features. Similar to our decision to

670    minimally transform and filter the data, we elected to avoid "nudging" the architecture of the

671    network in any direction in order to allow the data to inform it instead. Informing the model

28

672    architecture based on known relationships, analogous to incorporating a prior, remains an

673    interesting and potentially fruitful avenue to pursue.

## Feature Importance

675          Similar to the results of previous modeling (Washburn *et al.* 2021), we find that no single

676    data grouping provides sufficient information to disregard all others. We note that weather and

677    management data does reduce error substantially relative to genetic and soil data, but the

678    variation in performance is large (Figure 2). Only after integration of all data types do we see a

679    relative reduction in error and consistency in this reduction.

680          Here we focus on salience in the weather and management data as it provided the best

681    average performance when used without other datasets. We find that the total water applied to

682    the field (including irrigation and rainfall, termed "WaterTotalInmm") is the most influential factor

683    for determining yield (Figure 3, Supplementary Figure 3 D). This is sensible from a biological

684    standpoint and is in agreement with previous models. Previous DNNs developed with a subset

685    of G2F data also identified precipitation as substantially influencing yield (Washburn *et al.*

686    2021). Linear modeling results find similar results and suggest a positive association between

687    precipitation early in development and yield (Rogers *et al.* 2021). Additionally, in a recent study

688    using a hybrid machine learning and crop growth model prediction system, the authors found

689    that water related features (e.g. average drought stress, average water table in season) were

690    important, although not as important as the trend in genetic and management improvements

691    over time (Shahhosseini *et al.* 2021). The daily average of solar radiation

692    ("SolarRadiationMean") is the next most salient feature of this dataset, followed by the

693    maximum temperature ("TempMax") and the average wind direction ("WindDirectionMean"). A

694    study employing a convolutional recurrent DNN to model county level data likewise found solar

695    radiation and maximum temperature as important features and note an apparent increase in the

696    importance of temperature near planting time (Khaki *et al.* 2020). A time dependent sensitivity

697    can be observed in our model as well (Figure 3).

698         The relationship driving the high average salience of the average wind direction is not

699    clear. This feature likely correlates with unrecorded variables. Assessment of the topology and

700    geographical surroundings of each field site to suggest what this measure may be linked to lies

701    outside the scope of this study.

702         With respect to management interventions, although addition of N, P, or K are not

703    among the most salient weather and management features, we observe that nitrogen does have

704    a mean salience comparable to relative humidity and photoperiod, while phosphorus and

705    potassium are far lower. As noted in previously (Washburn *et al.* 2021) limited salience of

706    fertilizers could be due to the quantities used being too low to exert a substantial effect, or

707    alternatively application of these elements may be insufficiently variable to reveal the effect.

## Importance of GEM Interactions Accuracy in Feature Salience

709         Incorporating interactions between genetic, environmental, and management factors

710    appears to have benefitted the accuracy of the resultant models. The CO DNN with interactions

711    performed best (nRMSE 14.554%, RMSE 0.948) with the next best models, CO weather model

712    and SO Models performing comparably (nRMSE 15.715%, RMSE 1.024). However, the two

713    DNNs with interactions have far lower dispersion in RMSE, with standard deviations of 0.013 in

714    the CO model and 0.035 in the SO model as compared with 0.074 in the CO weather model.

715         Interactions not only improve accuracy and model consistency across replicates, there

716    appear to be changes in the salience of individual features as well. This is most apparent in

717    considering weather and management features' salience (Figure 3 C). Relative to the sub

718    model, incorporating interactions appears to increase the salience of irrigation, although it is

719    highly salient in both models (relative to other time series factors). Additionally, several broadly

30

720    salient points in time, two of which are at the extreme end of the season, have diminished

721    salience with the incorporation of interactions. This reduction is not uniform across all highly

722    salient time points. A strong peak in salience shortly after planting is seen in both saliency maps

723    which agrees with previously reported results (Washburn *et al.* 2021).

## 724    Conclusions and Future Directions

725         The consecutively optimized deep neural network model developed here shows promise

726    for complementing existing models for crop selection and improvement, as it produces more

727    accurate estimates of yield than the other considered models. Of particular interest here is the

728    capacity of this and other convolutional neural networks to incorporate change in environmental

729    variables over time. This enables the generation of counterfactuals to examine the expected

730    effect of different planting times (shifting the planting date of a site relative to the true value),

731    planting in different sites, or planting under future possible climate scenarios. Additionally, the

732    ability to generate such estimates would enable breeders to consider not only the expected yield

733    of an individual cultivar but the expected *consistency* of yield as well.

734         For such a strategy to be adopted in genomic selection, further efforts are needed to

735    validate the predictions such a model produces. This will necessitate incorporating of and

736    validation on future data from the Genomes to Field Initiative (McFarland *et al.* 2020) or other

737    large-scale experiments. The Genomes to Field Initiative and other organizations sponsor

738    prediction competitions and other activities designed to advance this area of study.

739    Furthermore, applying the same model or the approach used to develop it to other crops would

740    be a valuable step towards assessing its' broad scale usefulness. This would also potentially

741    implicate groups of crops for which the same model may be used through transfer learning,

742    along with groups that require crop-specific models to be developed.

743   Additional improvements to accuracy that have the potential to transfer to modeling

744 efforts for other crops include improved heuristics for epoch selection and training set

745 construction. The simultaneously optimized model achieves a minimum error *lower* than our

746 selected model (see Figure 2) and does so in far fewer epochs, but overfits much faster as well.

747 If overfitting were preventable through a better heuristic for epoch selection than the one we

748 employed, simultaneous optimization would have produced a better performing model that was

749 simpler to generate. Training set construction is another opportunity for improvement with

750 transferable utility. Here we took an aggressive approach ensuring approximately balanced

751 groups, down sampling all groups with observations in excess of the smallest group in the test

752 set. Deep neural networks tend to perform better with an abundance of data, so alternate

753 approaches that retain more observations are of interest. In cases where there are few

754 observations or model development is heavily constrained by computational resources or model

755 development time, other models, especially linear regression models, may result in a model that

756 performs nearly as well as a deep neural network.

757   Deep learning models do not result in parameters which are as readily interpretable as

758 those of more standard statistical procedures and do not incorporate the physiology of the plant

759 as mechanistic crop growth models do. These represent ongoing challenges and limit the

760 scenarios in which a deep neural network may be useful. This can be partially addressed

761 through how the data is represented (e.g. using non-PC transformed data), which has been

762 explored for identification of genetic loci (Liu *et al.* 2019). Additionally, efforts to incorporate

763 known relationships into a deep learning model's structure have the potential to benefit

764 accuracy and interpretability. Improvements in the capacity to represent genetic or physiological

765 principles could allow for these methods to apply to a wider range of uses and address a

766 broader set of questions.

767

# Author Contributions

769   Genomes to Fields experiments were coordinated and designed by Natalia DeLeon, David Ertl,

770   Judith Kolkman, Dayane Cristina Lima, Danilo Moreta, James Schnable, and Maninder Singh.

771   Field experiments were conducted, and data was collected and curated by Barış Alaca, Tim

772   Bessinger, Natalia DeLeon, David Ertl, Sherry Flint-Garcia, Candice Hirsch, Joseph Knoll,

773   Judith Kolkman, Dayane Cristina Lima, Danilo Moreta, James Schnable, Maninder Singh, Jason

774   Wallace, Jacob Washburn, and Tecle Weldekidan. Joseph Gage provided novel genomic data.

775   David Ertl and Maninder Singh contributed funding to the Genomes to Fields Initiative.

776

777   The computational study was designed by Daniel Kick and Jacob Washburn. Additional data

778   cleaning and imputation was done by Daniel Kick, who developed the models and generated

779   the figures. The manuscript was written by Daniel Kick and edited by Jacob Washburn, Jason

780   Wallace, James Schnable, Judith Kolkman, and Daniel Kick.

781

# Funding

# 789 Acknowledgements

# References

Anaconda Software Distribution, 2021 Anaconda Documentation.

Bache, S. M., and H. Wickham, 2020 *magrittr: A Forward-Pipe Operator for R*.

Bates, D., M. Mächler, B. Bolker, and S. Walker, 2015 Fitting Linear Mixed-Effects Models Using lme4. Journal of Statistical Software 67: 1–48.

Bergstra, J., D. Yamins, and D. Cox, 2013 Making a Science of Model Search: Hyperparameter Optimization in Hundreds of Dimensions for Vision Architectures, pp. 115–123 in *Proceedings of the 30th International Conference on Machine Learning*, edited by S. Dasgupta and D. McAllester. Proceedings of Machine Learning Research, PMLR, Atlanta, Georgia, USA.

Bradbury, P. J., Z. Zhang, D. E. Kroon, T. M. Casstevens, Y. Ramdoss *et al.*, 2007 TASSEL: software for association mapping of complex traits in diverse samples. Bioinformatics 23: 2633–2635.

Buitinck, L., G. Louppe, M. Blondel, F. Pedregosa, A. Mueller *et al.*, 2013 API design for machine learning software: experiences from the scikit-learn project, pp. 108–122 in *ECML PKDD Workshop: Languages for Data Mining and Machine Learning,*.

Chollet, F. and others, 2015 Keras.

Couture-Beil, A., 2018 *rjson: JSON for R*.

Da Costa-Luis, C., S. K. Larroque, K. Altendorf, H. Mary, Richardsheridan *et al.*, 2022 *tqdm: A fast, Extensible Progress Bar for Python and CLI*. Zenodo.

fuzzywuzzy, 2017 SeatGeek.

Harris, C. R., K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen *et al.*, 2020 Array programming with NumPy. Nature 585: 357–362.

831    Hornik, K., M. Stinchcombe, and H. White, 1989 Multilayer feedforward networks are universal

832         approximators. Neural Networks 2: 359–366.

833    Hunter, J. D., 2007 Matplotlib: A 2D graphics environment. Computing in Science &

834         Engineering 9: 90–95.

835    Inc, P. T., 2015 Collaborative data science.

836    Izrailev, S., 2021 *tictoc: Functions for Timing R Scripts, as Well as Implementations of Stack*

837         *and List Structures*.

838    J. Liu, C. E. Goering, and L. Tian, 2001 A NEURAL NETWORK FOR SETTING TARGET

839         CORN YIELDS. Transactions of the ASAE 44:.

840    Jarquin, D., N. de Leon, C. Romay, M. Bohn, E. S. Buckler *et al.*, 2021 Utility of Climatic

841         Information via Combining Ability Models to Improve Genomic Prediction for Yield

842         Within the Genomes to Fields Maize Project. Front. Genet. 11: 592769.

843    Khaki, S., L. Wang, and S. V. Archontoulis, 2020 A CNN-RNN Framework for Crop Yield

844         Prediction. Front. Plant Sci. 10: 1750.

845    Kibirige, H., G. Lamp, J. Katins, Gdowding, Austin *et al.*, 2021 *has2k1/plotnine: v0.8.0*. Zenodo.

846    Kubota, Y., 2021 *tf-keras-vis*.

847    Kurtzer, G. M., V. Sochat, and M. W. Bauer, 2017 Singularity: Scientific containers for mobility

848         of compute (A. Gursoy, Ed.). PLoS ONE 12: e0177459.

849    Li, X., T. Guo, J. Wang, W. A. Bekele, S. Sukumaran *et al.*, 2021 An integrated framework

850         reinstating the environmental dimension for GWAS and genomic selection in crops.

851         Molecular Plant S167420522100085X.

852  Liu, Y., D. Wang, F. He, J. Wang, T. Joshi *et al.*, 2019 Phenotype Prediction and Genome-Wide

853      Association Study Using Deep Convolutional Neural Network of Soybean. Front. Genet.

854      10: 1091.

855  Martín Abadi, Ashish Agarwal, Paul Barham, Eugene Brevdo, Zhifeng Chen *et al.*, 2015

856      TensorFlow: Large-Scale Machine Learning on Heterogeneous Systems.

857  McFarland, B. A., N. AlKhalifah, M. Bohn, J. Bubert, E. S. Buckler *et al.*, 2020 Maize genomes

858      to fields (G2F): 2014–2017 field seasons: genotype, phenotype, climatic, soil, and inbred

859      ear image datasets. BMC Res Notes 13: 71.

860  Messina, C. D., F. Technow, T. Tang, R. Totir, C. Gho *et al.*, 2018 Leveraging biological insight

861      and environmental variation to improve phenotypic prediction: Integrating crop growth

862      models (CGM) with whole genome prediction (WGP). European Journal of Agronomy

863      100: 151–162.

864  Müller, K., 2020 *here: A Simpler Way to Find Your Files*.

865  O'Malley, T., E. Bursztein, J. Long, F. Chollet, H. Jin *et al.*, 2019 KerasTuner.

866  Pedersen, T. L., 2020 *patchwork: The Composer of Plots*.

867  Pedregosa, F., G. Varoquaux, A. Gramfort, V. Michel, B. Thirion *et al.*, 2011 Scikit-learn:

868      Machine Learning in Python. Journal of Machine Learning Research 12: 2825–2830.

869  R Core Team, 2021 *R: A Language and Environment for Statistical Computing*. R Foundation

870      for Statistical Computing, Vienna, Austria.

871  Richardson, N., I. Cook, N. Crane, J. Keane, R. François *et al.*, 2021 *arrow: Integration to*

872      *"Apache" "Arrow."*

873  Rogers, A. R., J. C. Dunne, C. Romay, M. Bohn, E. S. Buckler *et al.*, 2021 The importance of

874      dominance and genotype-by-environment interactions on grain yield variation in a large-

875        scale public cooperative maize experiment (E. Akhunov, Ed.). G3

876        Genes|Genomes|Genetics 11: jkaa050.

877    Rogers, A. R., and J. B. Holland, 2021 Environment-specific genomic prediction ability in maize

878        using environmental covariates depends on environmental similarity to training data (A.

879        Lipka, Ed.). G3 Genes|Genomes|Genetics jkab440.

880    Samek, W., T. Wiegand, and K.-R. Müller, 2017 Explainable Artificial Intelligence:

881        Understanding, Visualizing and Interpreting Deep Learning Models.

882    Seabold, S., and J. Perktold, 2010 statsmodels: Econometric and statistical modeling with

883        python, in *9th Python in Science Conference*,.

884    Shahhosseini, M., G. Hu, I. Huber, and S. V. Archontoulis, 2021 Coupling machine learning and

885        crop modeling improves crop yield prediction in the US Corn Belt. Sci Rep 11: 1606.

886    Simonyan, K., A. Vedaldi, and A. Zisserman, 2014 Deep Inside Convolutional Networks:

887        Visualising Image Classification Models and Saliency Maps.

888    SingularityCE Developers, 2021 *SingularityCE 3.8.3*. Zenodo.

889    Tavenard, R., J. Faouzi, G. Vandewiele, F. Divo, G. Androz *et al.*, 2020 Tslearn, A Machine

890        Learning Toolkit for Time Series Data. Journal of Machine Learning Research 21: 1–6.

891    team, T. pandas development, 2020 *pandas-dev/pandas: Pandas*. Zenodo.

892    Technow, F., C. D. Messina, L. R. Totir, and M. Cooper, 2015 Integrating Crop Growth Models

893        with Whole Genome Prediction through Approximate Bayesian Computation (I. De

894        Smet, Ed.). PLoS ONE 10: e0130855.

895    Techtonik, A., 2015 *wget 3.2*.

896    Thornton, M. M., R. Shrestha, Y. Wei, P. E. Thornton, S. Kao *et al.*, 2020 Daymet: Daily

897        Surface Weather Data on a 1-km Grid for North America, Version 4.

898    Van Rossum, G., and F. L. Drake, 2009 *Python 3 Reference Manual*. CreateSpace, Scotts Valley,

899        CA.

900    Virtanen, P., R. Gommers, T. E. Oliphant, M. Haberland, T. Reddy *et al.*, 2020 SciPy 1.0:

901        Fundamental Algorithms for Scientific Computing in Python. Nature Methods 17: 261–

902        272.

903    Washburn, J. D., E. Cimen, G. Ramstein, T. Reeves, P. O'Briant *et al.*, 2021 Predicting

904        phenotypes from genetic, environment, management, and historical data using CNNs.

905        Theor Appl Genet 134: 3997–4011.

906    Waskom, M. L., 2021 seaborn: statistical data visualization. Journal of Open Source Software 6:

907        3021.

908    Westhues, C. C., G. S. Mahone, S. da Silva, P. Thorwarth, M. Schmidt *et al.*, 2021 Prediction of

909        Maize Phenotypic Traits With Genomic and Environmental Predictors Using Gradient

910        Boosting Frameworks. Front. Plant Sci. 12: 699589.

911    Wickham, H., M. Averick, J. Bryan, W. Chang, L. D. McGowan *et al.*, 2019 Welcome to the

912        tidyverse. Journal of Open Source Software 4: 1686.

913    Zhou, D.-X., 2020 Universality of deep convolutional neural networks. Applied and

914        Computational Harmonic Analysis 48: 787–794.

915

916

# Figures

## Figure 1. Optimization strategy results in different network architectures and degree of overfitting in full model
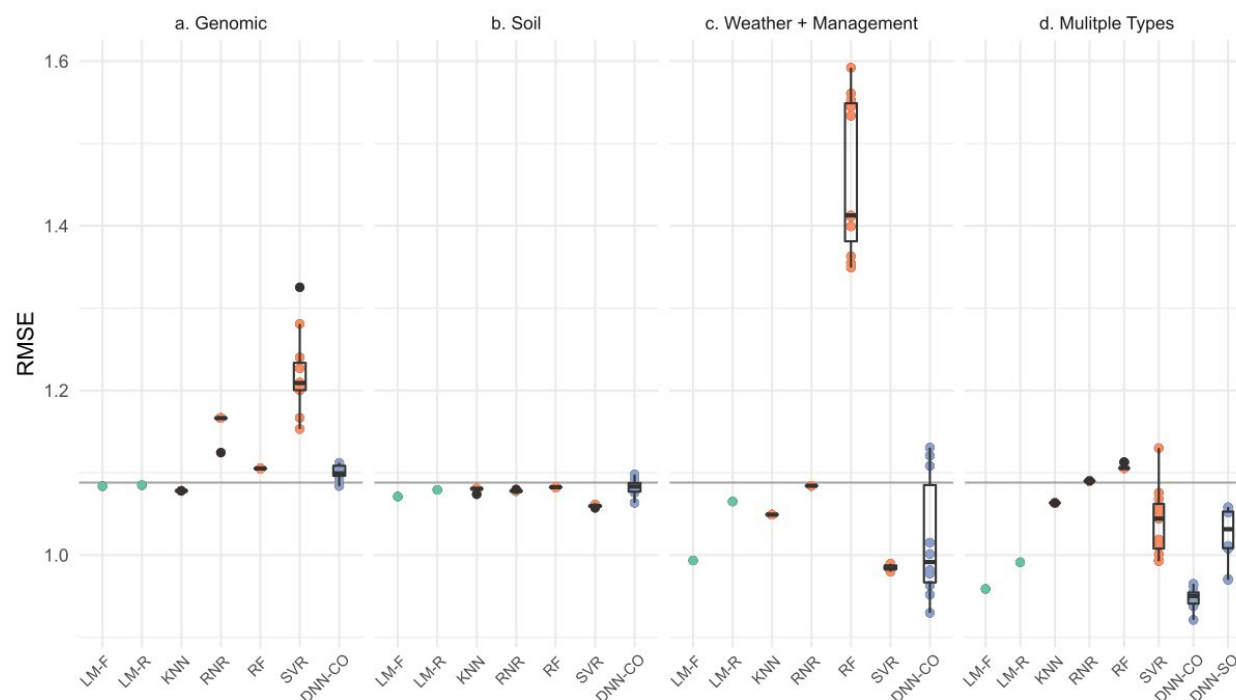


**A.** Hyperparameters for each optimization strategy are shown as a percent of the allowed

range. Data available for training is the same but the Consecutive Optimization (CO) and

Simultaneous Optimization (SO) strategy result in substantially different hyperparameter values

and thus network architecture. For exact values refer to Table 2 and Table 3. **B.** The average

RMSE of the test set (across 10 replicates to account for random initialization of weights) is

shown in black for each submodel (**a. – e.**). The horizontal dashed black line indicates the

minimum error achieved throughout the training duration. The vertical lines indicate the

difference in error and epochs of the minimum value and the values selected through minimizing

40

929     total validation error (red dashed line), the heuristic used in this study, and the mean plus

930     standard deviation of validation error (solid blue line), which was considered but not used. Both

931     strategies considered failed to select the epoch resulting in the minimum loss in the test set for

932     all submodels and resulted in apparent overfitting in the Weather and Management submodel

933     (**c.**) and the SO model (**e.**). For additional comparisons of heuristic performance see Table 4.

934

935 # Figure 2. Model Performance Across Methodologies and Data
936 Types

### A. Performance on Test Set in Root Mean Squared Error



937
938 **A.** The root mean squared error (RMSE) of the testing set is shown for each data grouping

939 (panels a - d) and class of model. Lower values indicate better model performance. The

940 horizontal gray line indicates the performance of an intercept model, i.e. using the mean of the

941 training set yield as the prediction for all observations in the test set. For models that depend on

942 a seed value the RMSE values for ten trials (evaluated on the same data) are shown and

943 standard Tukey box plots are provided. In deep learning models random initialization of weights

944 at the beginning of training result in different performance across trials. Three groups of models

945 are shown, linear models (green), machine learning models (orange), and deep learning models

946 (blue). Linear models are subdivided into those with exclusively fixed effects (LM-F) and those

947 with random effects (LM-R). The best performing linear model is shown. For LM-F and LM-R

948 respectively these are utilizing the first 8 PCs (explaining 31% of the variance) in **a**, utilizing all

949 soil variables (for both LM-F and LM-R) in **b**, utilizing the five most salient weather/management

42

950     factors (for both LM-F and LM-R) in **c**, and the first 8 genomic PCs, 5 most salient

951     weather/management variables, and interactions between the two (for both LM-F and LM-R) **d**.

952     Machine learning models used were K-Nearest Neighbors (KNN), Radius Neighbor Regression

953     (RNR), Random Forest (RF), and Support Vector Regression with a linear kernel (SVR). Deep

954     learning models are divided by whether they were part of the Consecutive optimization strategy

955     (DNN-CO) or the Simultaneous optimization strategy (DNN-SO). Note that DNN-SO requires all

956     data types and thus only appears in panel **d**.

957

958 # Figure 3. Influence of Interaction Effects on Feature Salience

959



960

961 **A.** Average salience across all weather and management factors for each day considered.

962 Interaction model values shown in black. **B.** The same values in **A.** are shown for the submodel

963 in blue. Salience peaks shortly after planting in both models. The submodel contains additional

964     peaks of salience prior to planting and near the end of the considered date range. **C.** Subtracted

965     salience values for the interaction model and the submodel. The interaction-containing model

966     appears to contain greater importance generally for certain features, e.g. irrigation and rainfall,

967     represented as "WaterTotalInmm". The difference between the two saliency maps indicates

968     additional times of sensitivity in the submodel (approximately -25, +180, +195) that the

969     interaction model is relatively insensitive to.

970 # Tables

971 ## Table 1. Hyperparameter Ranges: Deep Learning

| Category | Submodels | Hyperparameter | Range |
|---|---|---|---|
| Architecture | Genomic Only | Layers | 1-7 |
| | | Units | 4-256 |
| | | Dropout Fraction | 0-0.3 |
| | Soil Only | Layers | 1-7 |
| | | Units | 4-64 |
| | | Dropout Fraction | 0-0.3 |
| | Weather Only | Pooling Type | Max (1d), Ave. (1d) |
| | | Layer Repeats | 1-7 |
| | | Convolution Layers per Repeat | 1-4 |
| | | Filter Size | 4 - 512 |
| | Interactions | Layers | 1-7 |
| | | Units | 4 - 256 |
| | | Dropout Fraction | 0-0.3 |
| Training | Optimizer | Learning Rate | 0.1, 0.01, 0.001, 0.0001 |
| | | Beta 1 | 0.9 - 0.9999 |
| | | Beta 2 | 0.9 - 0.9999 |
| | Other | Batch size | 32-256, step=16 |
| | | Epoch | 1-1000 |

972

973

974    Table 2. Selected Deep Learning Hyperparameters: Architecture

| Submodel or Network | Hyperparameter | Specific Layer | Consecutive Optimization | Sequential Optimization |
|---|---|---|---|---|
| Genomic Only | Units | 1 | 83 | 196 |
| | | 2 | 133 | 47 |
| | Dropout Fraction | 1 | 0.163923177 | 0.15214 |
| | | 2 | 0.230663142 | 0.06061 |
| Soil Only | Units | 1 | 38 | 19 |
| | | 2 | 13 | 27 |
| | | 3 | 45 | |
| | | 4 | 29 | |
| | | 5 | 4 | |
| | | 6 | 4 | |
| | | 7 | 4 | |
| | Dropout | 1 | 0.148724301 | 0.21342 |
| | | 2 | 0.276340999 | 0.18589 |
| | | 3 | 0.005434164 | |
| | | 4 | 0.173380695 | |
| | | 5 | 0 | |
| | | 6 | 0 | |
| | | 7 | 0 | |
| Weather + Management Only | Pooling Type | N/A | Max | Max |
| | Convolution Layers per Repeat | N/A | 2 | 2 |
| | Filter Size | 1 | 433 | 370 |
| | | 2 | 436 | 303 |
| | | 3 | 52 | |
| | | 4 | 163 | |
| | | 5 | 400 | |
| | | 6 | 294 | |
| Interaction | Units | 1 | 152 | 10 |
| | | 2 | 207 | 25 |
| | | 3 | 206 | 126 |
| | | 4 | 188 | 204 |
| | | 5 | 44 | 45 |
| | | 6 | | 134 |
| | Dropout | 1 | 0.18658661 | 0.10201 |
| | | 2 | 0.289893588 | 0.14809 |
| | | 3 | 0.004841293 | 0.01536 |
| | | 4 | 0.198121953 | 0.15658 |
| | | 5 | 0.243027717 | 0.2428 |
| | | 6 | | 0.19048 |

975
976

47

977 ## Table 3. Selected Deep Learning Hyperparameters: Training

|  | Optimizer |  |  | Other |  |
|---|---|---|---|---|---|
| Network | learning_rate | beta1 | beta2 | numEpoch | batch_size |
| CO: Genomic Only | 0.0001 | 0.953368 | 0.985947 | 12 | 96 |
| CO: Soil Only | 0.01 | 0.928472 | 0.997516 | 199 | 176 |
| CO: Weather Only | 0.0001 | 0.903649 | 0.929582 | 629 | 240 |
| CO: Full Network | 0.01 | 0.98752 | 0.972311 | 364 | 112 |
| SO | 0.001 | 0.975893 | 0.994607 | 711 | 192 |

978
979

980   Table 4. Epoch Selection Underperformance

| Network | Epoch Selection By: | | | Average Test Loss: | | | Proportion of Minimum Loss | | |
| | Mean + Standard Dev. | Best Epoch | Sum of Losses | Mean + Standard Dev. | Best Epoch | Sum of Losses | Mean + Standard Dev. | Best Epoch | Sum of Losses |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |
| CO: Genomic Only | 10 | 10 | 12 | 1.209216 | 1.209216 | 1.21171 | 1 | 1.002063 | 0.002063 |
| CO: Soil Only | 161 | 14 | 199 | 1.178308 | 1.144793 | 1.172763 | 1.029276 | 1.024432 | -0.00484 |
| CO: Weather Only | 225 | 66 | 629 | 1.041072 | 0.945608 | 1.046358 | 1.100955 | 1.106545 | 0.00559 |
| CO: Full Network | 362 | 287 | 364 | 0.912883 | 0.893123 | 0.903884 | 1.022124 | 1.012048 | -0.01008 |
| SO | 796 | 14 | 711 | N/A | 0.86811 | 1.052109 | N/A | 1.211954 | N/A |

981

982

983 ## Table 5. Machine Learning Hyperparameter Optimization

| Model | Hyperparameter | Range | Genomic Only | Soil Only | Weather + Management Only | Multiple |
|---|---|---|---|---|---|---|
| KNN | Weight Metric | Uniform, Distance | Uniform | Distance | Uniform | Distance |
| | k | 1 - 250 | 237 | 248 | 248 | 49 |
| RNR | Weight Metric | Uniform, Distance | Distance | Distance | Uniform | Distance |
| | Radius | 0.01 - 2000 | 39.759518 | 3.406197 | 5.986679 | 40.375418 |
| SVR | Loss | Epsilon Insensitive, Squared Epsilon Insensitive | Epsilon Insensitive | Epsilon Insensitive | Epsilon Insensitive | Squared Epsilon Insensitive |
| | C | 1 - 5 (log uniform) | 2.772318 | 5.613996 | 4.623351 | 2.787589 |
| RF | Max Depth | 2 - 200, q = 1 (q uniform) | 64 | 10 | 102 | 7 |
| | Min Samples/Leaf | 1 - 200, q = 1 (q uniform) | 171 | 163 | 100 | 149 |

984

985

986

## Table 6. Performance Across Data Sets

| Data Set | Model | mean RMSE | Standard Dev. RMSE | Mean nRMSE | Standard Dev. nRMSE | Mean R2 | Standard Dev. R2 |
|---|---|---|---|---|---|---|---|
| a. Genomic | DNN-CO | 1.101 | 0.009 | 16.896% | 0.142% | -0.117 | 0.019 |
| | KNN | 1.078 | 0.000 | 16.548% | 0.000% | -0.072 | 0.000 |
| | LM-F | 1.084 | | 16.635% | | -0.083 | |
| | LM-R | 1.085 | | 16.653% | | -0.085 | |
| | Mean | 1.088 | | 16.701% | | -0.092 | |
| | RF | 1.105 | | 16.965% | | -0.126 | |
| | RNR | 1.163 | 0.013 | 17.846% | 0.194% | -0.246 | 0.027 |
| | SVR | 1.219 | 0.049 | 18.718% | 0.750% | -0.373 | 0.112 |
| b. Soil | DNN-CO | 1.083 | 0.010 | 16.622% | 0.152% | -0.081 | 0.020 |
| | KNN | 1.080 | 0.002 | 16.578% | 0.031% | -0.076 | 0.004 |
| | LM-F | 1.071 | | 16.441% | | -0.058 | |
| | LM-R | 1.079 | | 16.566% | | -0.074 | |
| | Mean | 1.088 | | 16.701% | | -0.092 | |
| | RF | 1.083 | | 16.616% | | -0.081 | |
| | RNR | 1.078 | 0.001 | 16.549% | 0.008% | -0.072 | 0.001 |
| | SVR | 1.059 | 0.001 | 16.262% | 0.017% | -0.035 | 0.002 |
| c. Weather + Management | DNN-CO | 1.018 | 0.074 | 15.627% | 1.141% | 0.040 | 0.142 |
| | KNN | 1.049 | | 16.105% | | -0.015 | |
| | LM-F | 0.993 | | 15.249% | | 0.090 | |
| | LM-R | 1.065 | | 16.349% | | -0.046 | |
| | Mean | 1.088 | | 16.701% | | -0.092 | |
| | RF | 1.461 | 0.095 | 22.430% | 1.455% | -0.977 | 0.256 |
| | RNR | 1.084 | | 16.643% | | -0.084 | |
| | SVR | 0.985 | 0.003 | 15.114% | 0.050% | 0.106 | 0.006 |
| d. Multiple Types | DNN-CO | 0.948 | 0.013 | 14.553% | 0.197% | 0.171 | 0.022 |
| | DNN-SO | 1.024 | 0.035 | 15.716% | 0.531% | 0.032 | 0.065 |
| | KNN | 1.063 | 0.000 | 16.322% | 0.000% | -0.043 | 0.000 |
| | LM-F | 0.959 | | 14.719% | | 0.152 | |
| | LM-R | 0.991 | | 15.217% | | 0.094 | |
| | Mean | 1.088 | | 16.701% | | -0.092 | |
| | RF | 1.107 | 0.003 | 16.991% | 0.045% | -0.130 | 0.006 |
| | RNR | 1.090 | 0.000 | 16.733% | | -0.096 | 0.000 |
| | SVR | 1.041 | 0.042 | 15.976% | 0.643% | 0.000 | 0.082 |

987