



**Shin-Han
Shiu**



**Gustavo de
los Campos**



**Andrew
McCarren**

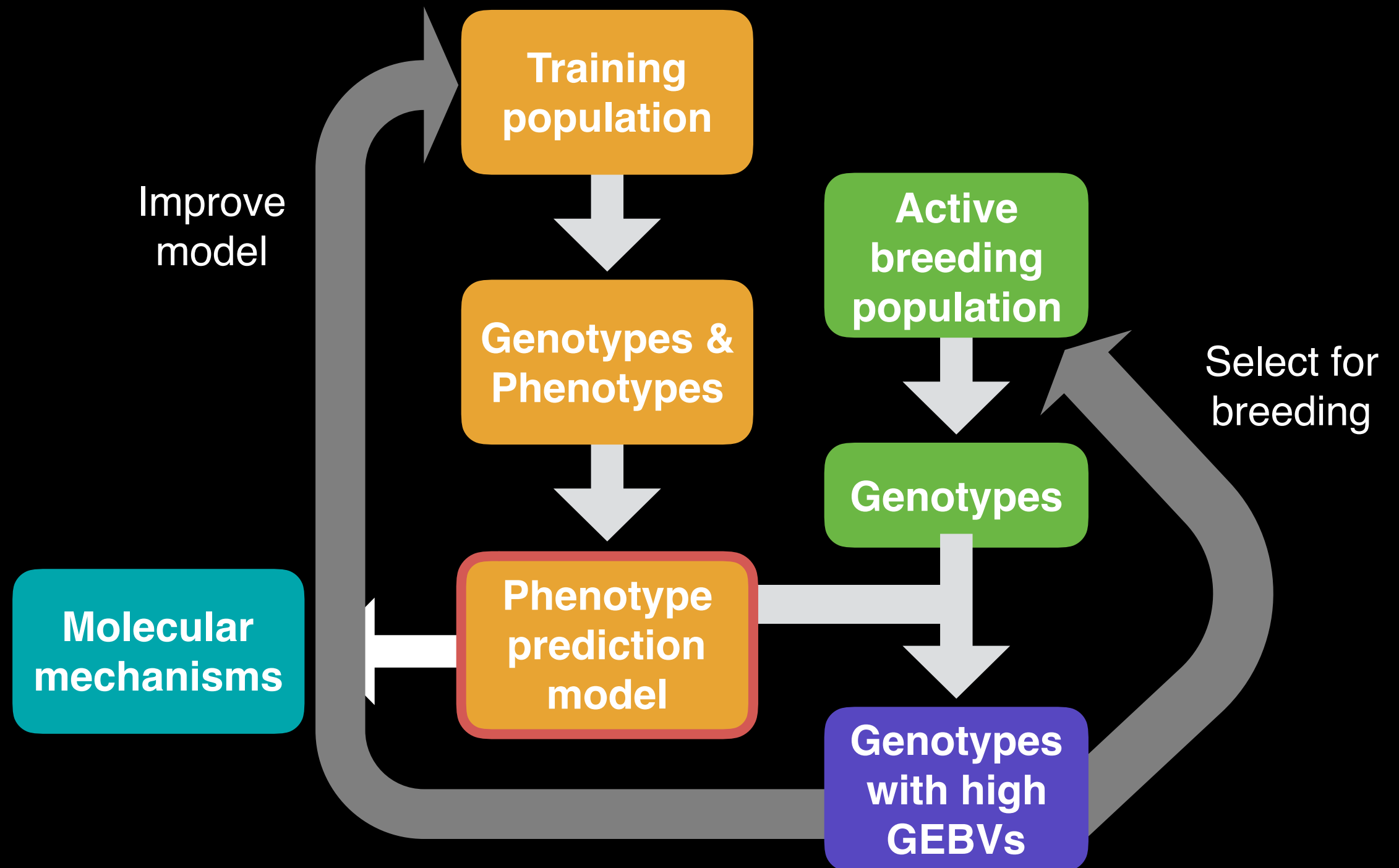


**Mark
Roantree**

Benchmarking algorithms for genomic prediction of complex traits in plants

Christina Azodi
Department of Plant Biology
February 26th, 2019
QuantGen Lab Meeting

Genomic Prediction: Breeding & Basic Science



GEBV: genomic estimated breeding value

Genomic Prediction Models

Quantitative
phenotype of
individual i

SNP genotype
of individual i

$$y_i = g(x_i) + e_i$$

Residual

Function relating
genotype to
phenotype

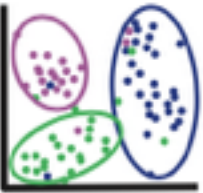
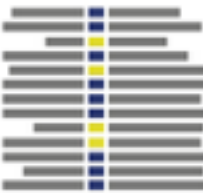


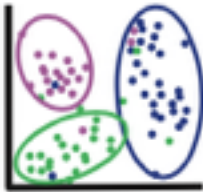

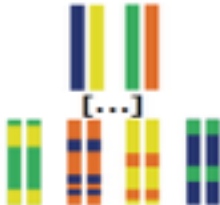
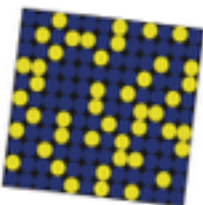
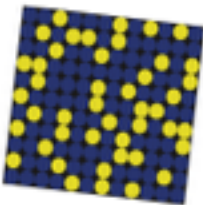

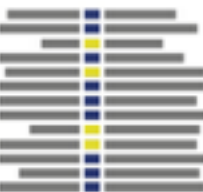
Performance
Metric (r)




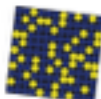
$$\text{cor}(y, \hat{y})$$

Outline

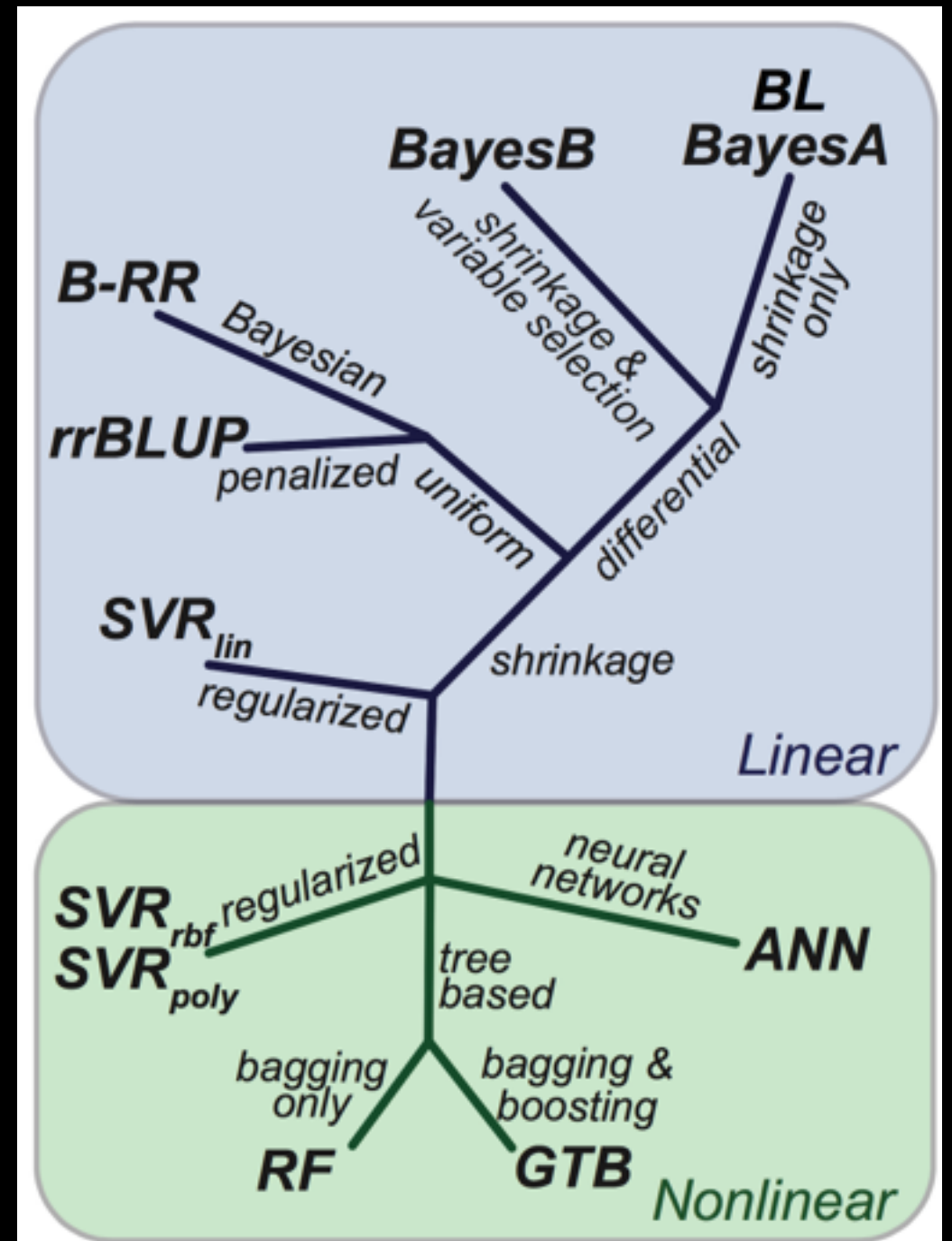
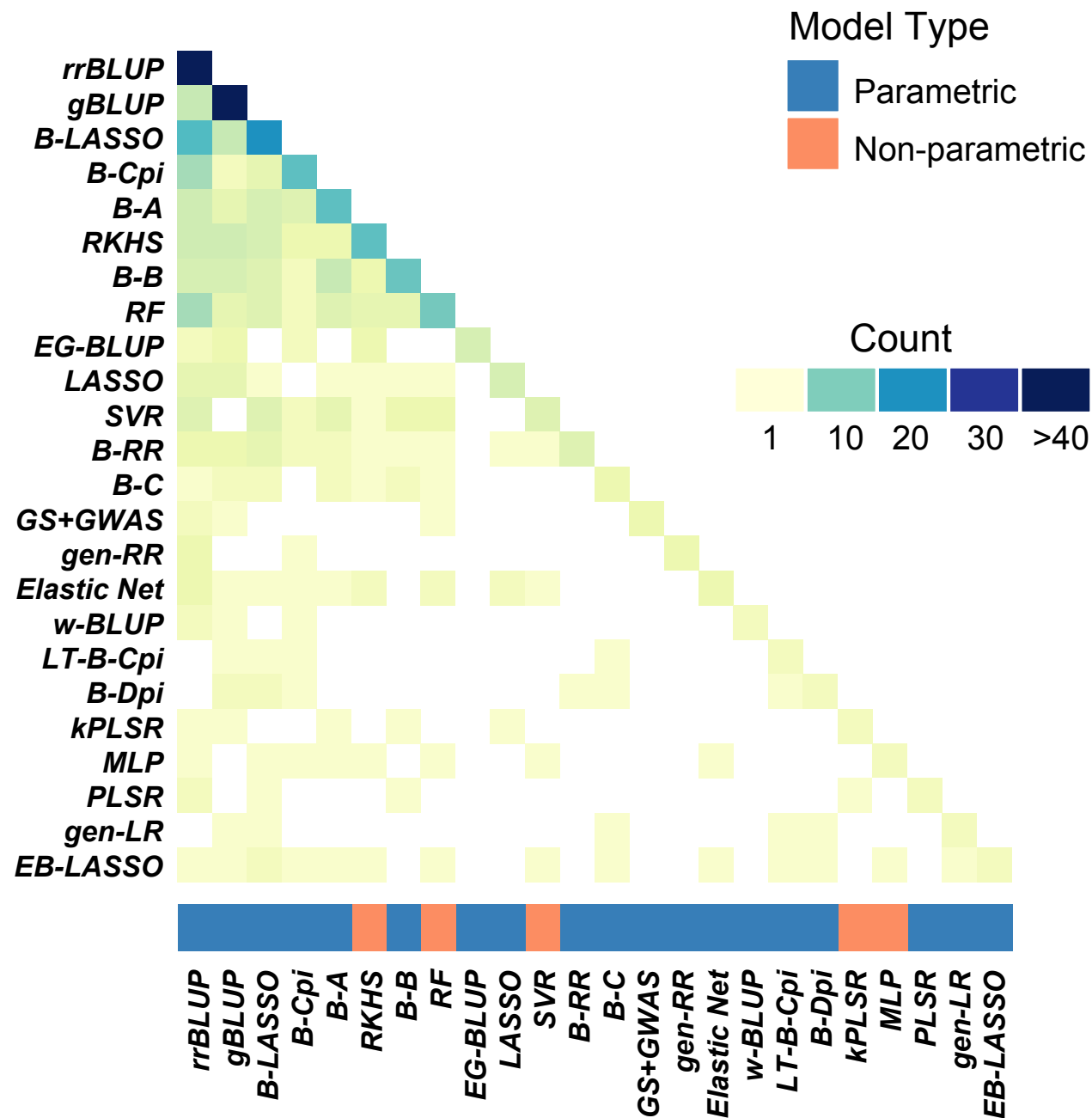
- ▶ Data and algorithms included in benchmark
- ▶ Intro to Neural Networks
- ▶ Performance on predicting height in 6 species
- ▶ Improvements to artificial neural networks
 - ▶ Feature selection and seeded starting weights
- ▶ Final results

Data Used

| Species | Population | Markers | Source |
|-------------|--|---|--|
| Maize |  391 |  332,178 | Hansey et. al 2011 Hirsch et. al 2014 |
| Rice |  327 |  73,147 | Spindel et. al 2015 |
| Sorghum |  451 |  58,961 | Fernandes et. al 2017 |
| Soy |  5,014 |  4,240 | Xavier et. al 2016 |
| Spruce | Partial DM 1,722 |  6,932 | Beaulieu et. al 2014 |
| Switchgrass |  514 |  218,528 | Lipka et. al 2014 Evans et. al 2017 |

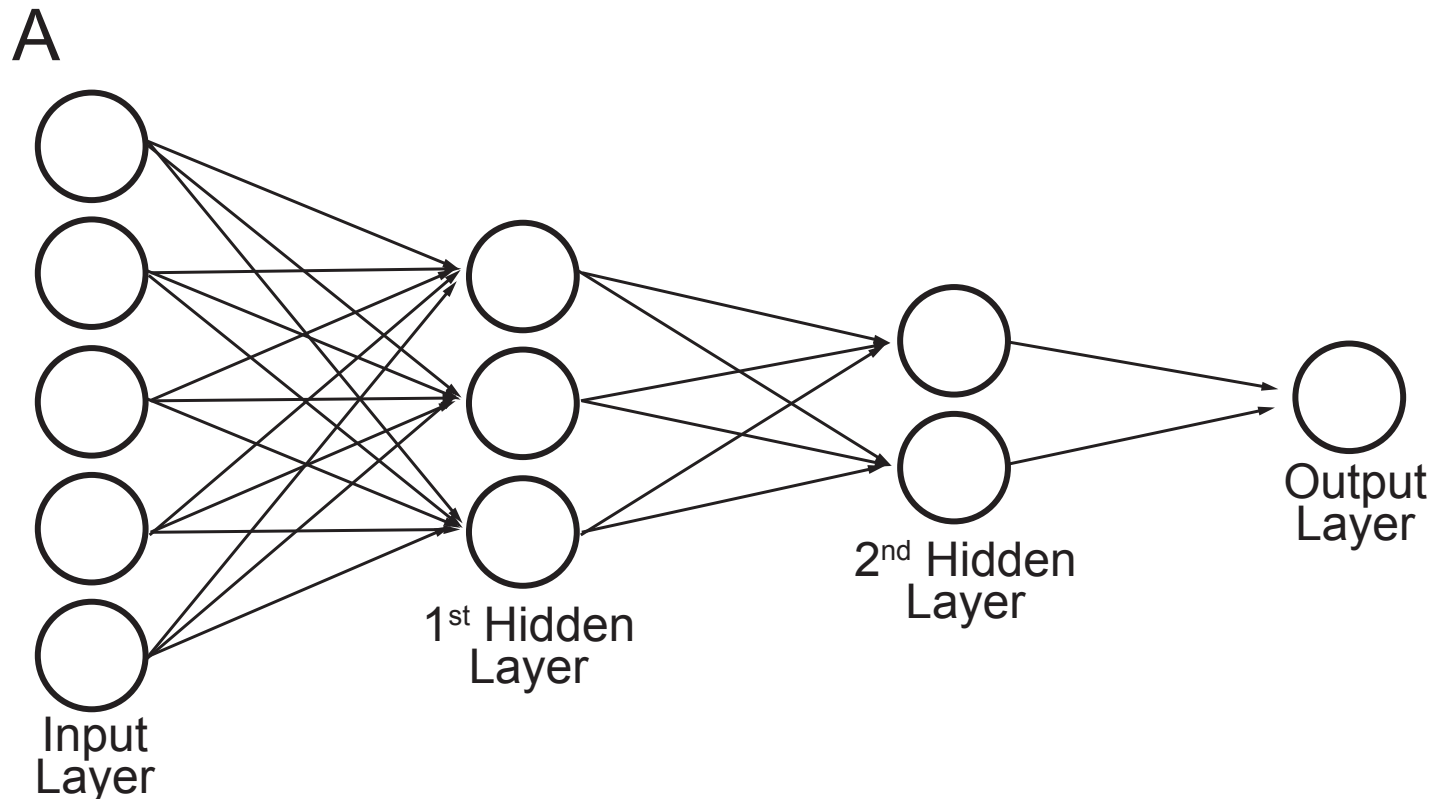
 Diversity Panel
  NAM
  GBS
  SNP-chip

GP Algorithm Comparisons



rrBLUP: ridge regression Best Linear Unbiased Predictor. BL: Bayesian Least Absolute Shrinkage and Selection Operator. SVR: Support Vector Regression. RF: Random Forest. ANN: Artificial Neural Network. GTB: Gradient Tree Boosting.

What is a Neural Network?



Difficulties

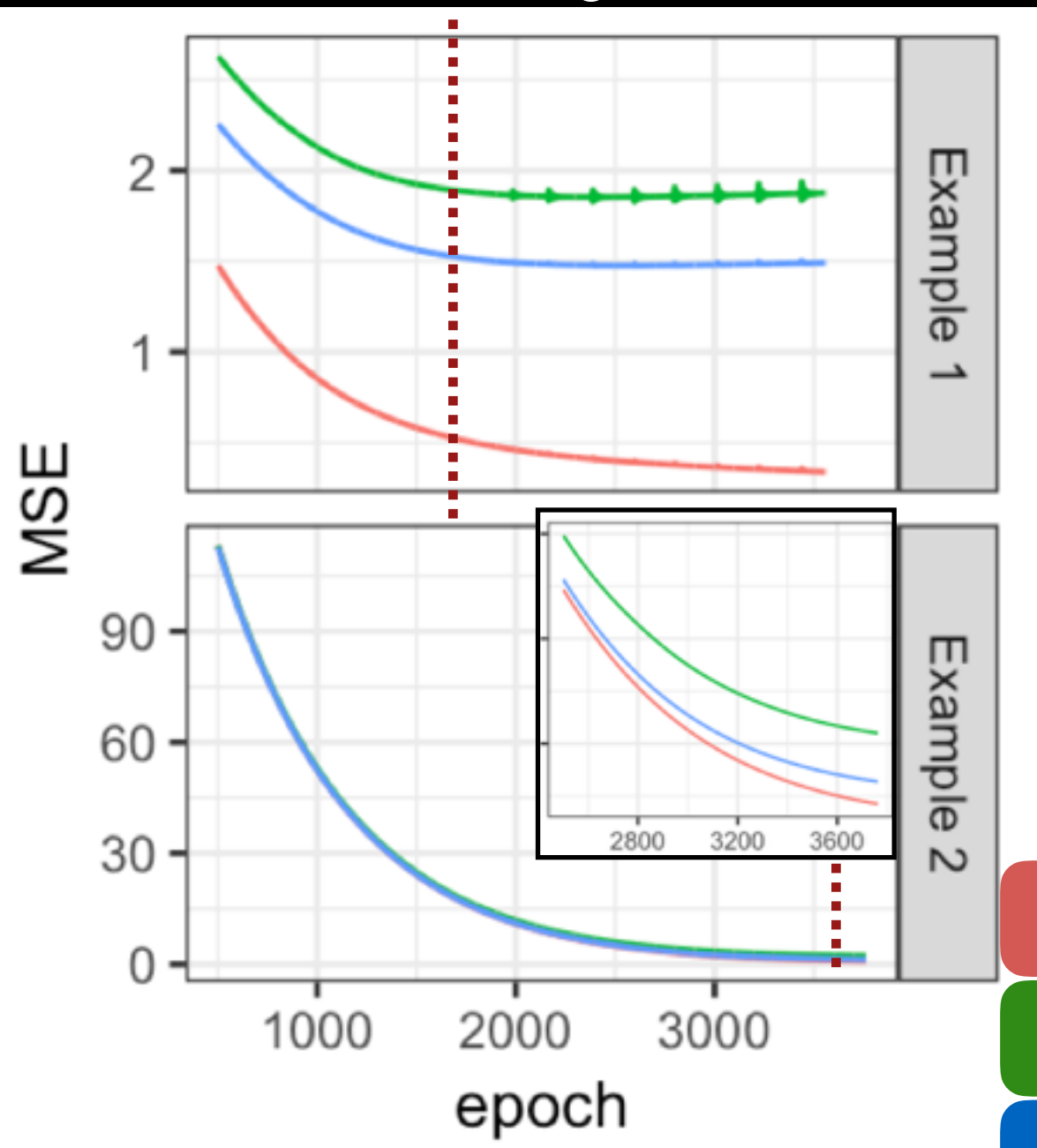
1. Large hyper-parameter space
2. When to stop training?

1. Grid Search

- Network Architecture
- Activation Function
- Learning Rate
- Type/degree of regularization

2. When to stop training?

Models can have very different ideal training times



Minimize validation error

Burn-in



$$\frac{\text{MSE}_{\text{valid Epoch X}}}{\text{MSE}_{\text{valid Epoch X+1}}} < 0.01 \times 10$$



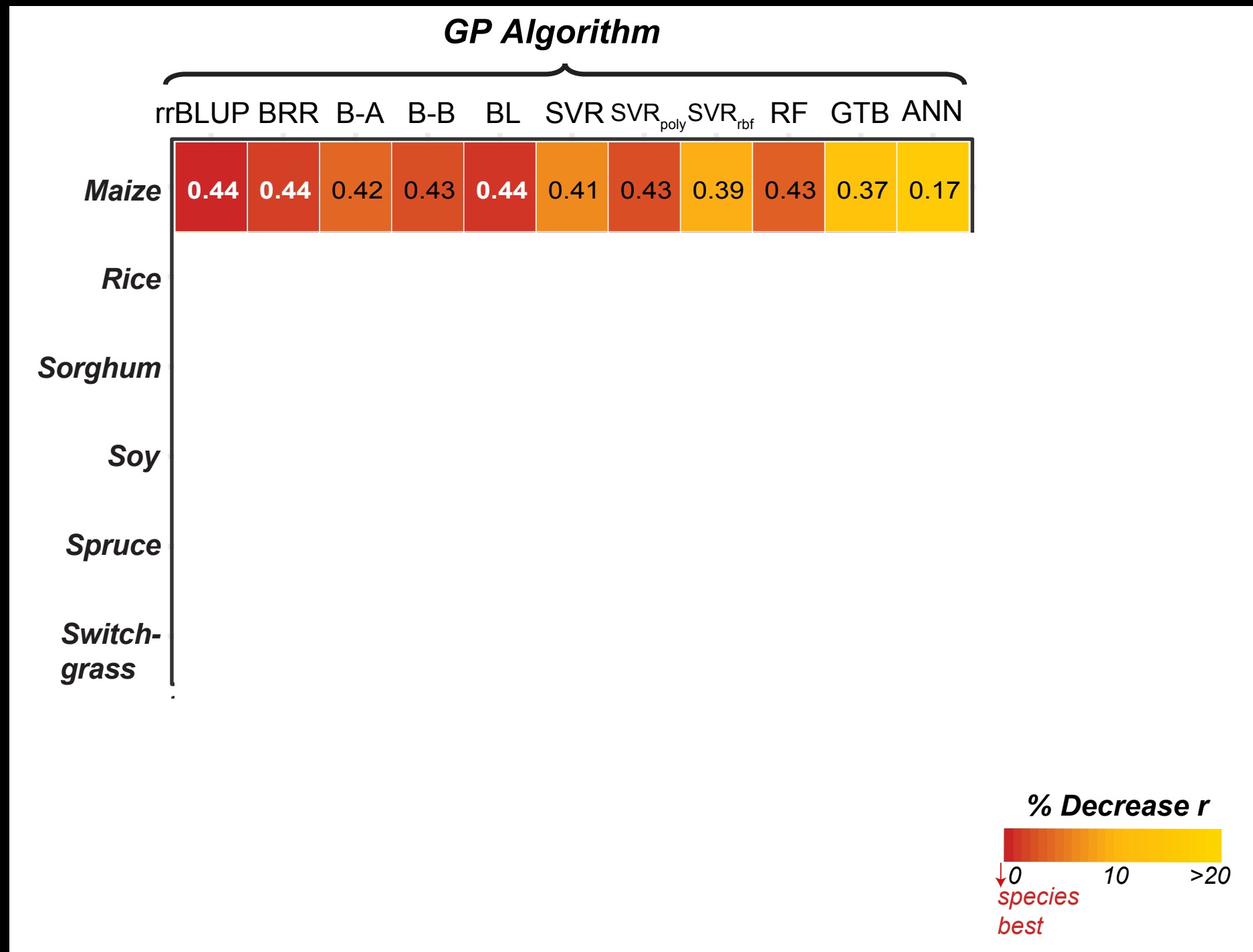
Stop
Training

Training

Validation

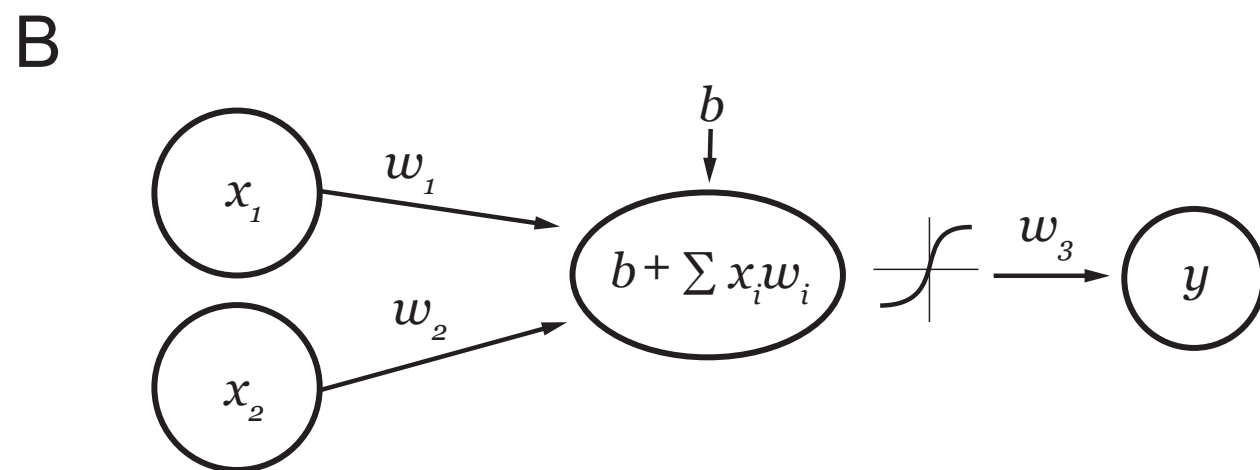
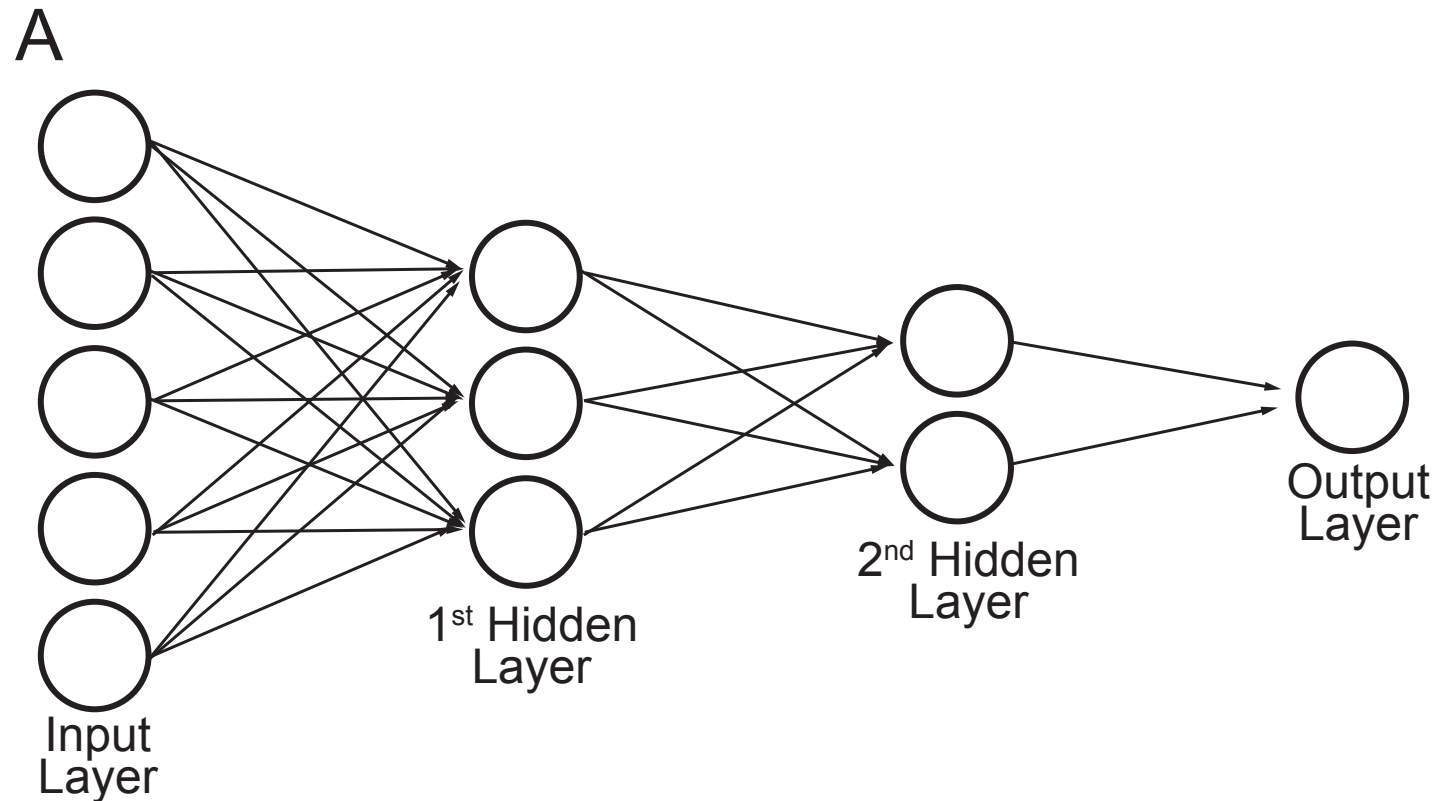
Testing

No one algorithm always performs best



model performance for predicting height in each species

What is a Neural Network?



Difficulties

1. Large hyper-parameter space
2. When to stop training?
3. Not robust when $p \gg n$
4. Large variation in model performance using same data and hyper-parameters.

3. Addressing $p \gg n$ with feature selection

1. Feature Selection

Training & Validation Data

Feature Selection

Random Forest
Bayes A
Elastic Net

2. Select hyperparameters & train model

Training Data

Grid Search

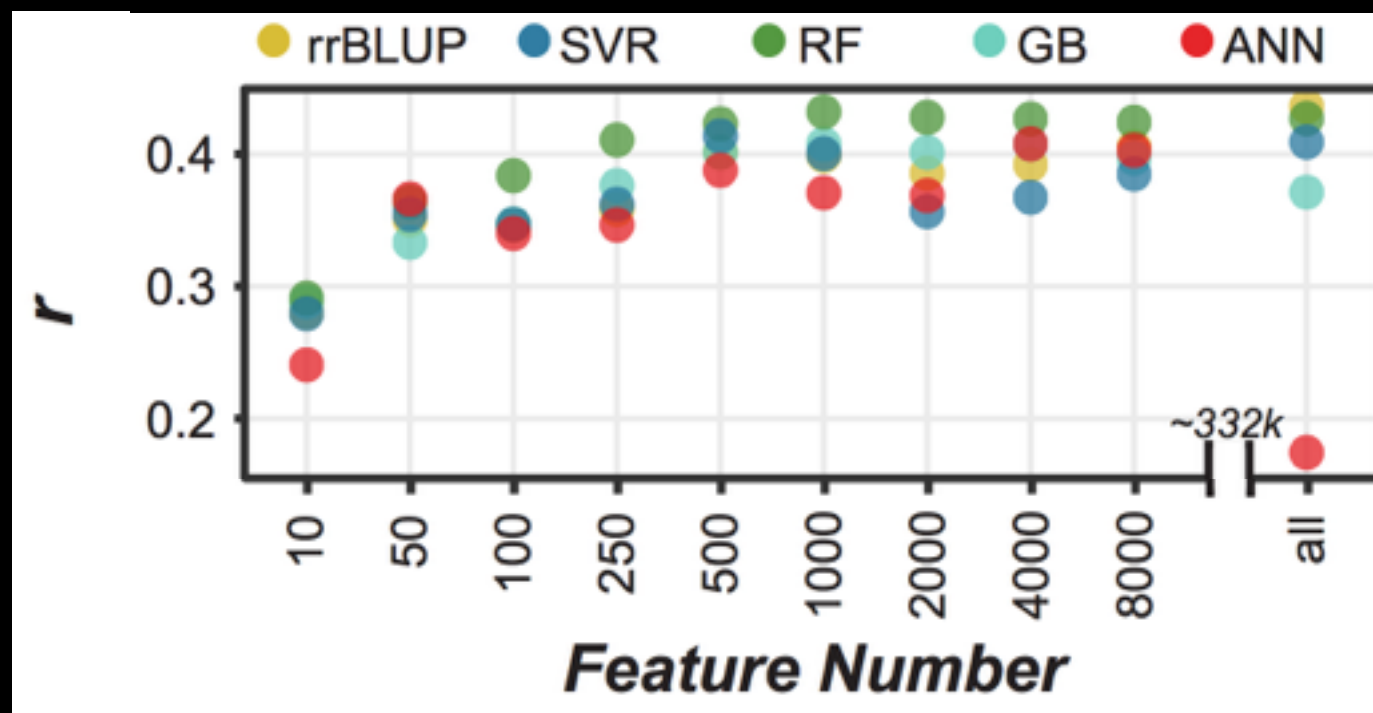
Validation Data

Train model

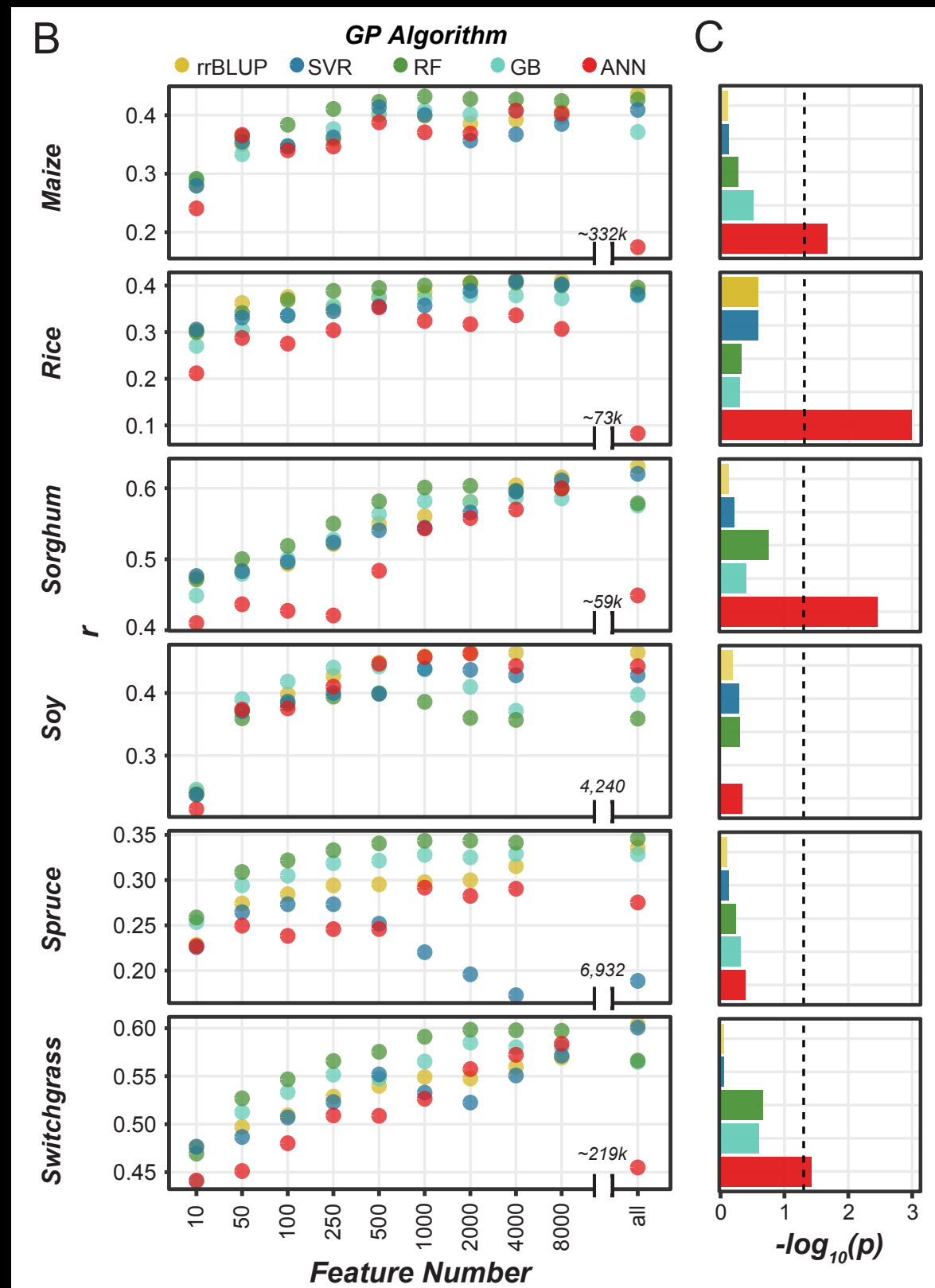
3. Performance

Testing Data

Height in maize



Dimension reduction improves ANN performance



p:n ratio

850

224

131

0.85

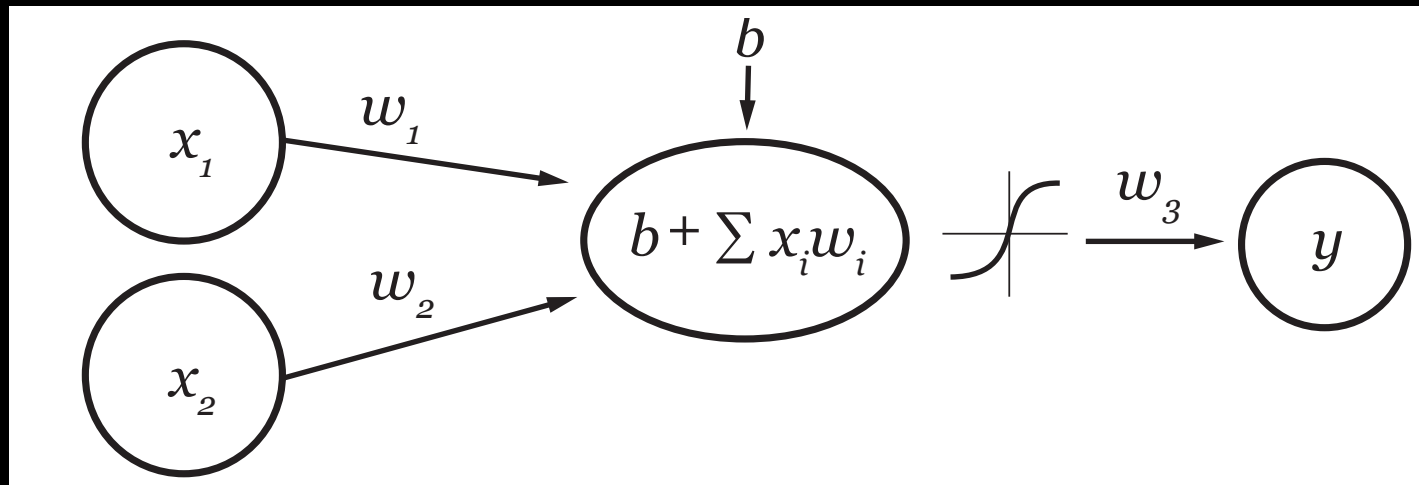
4

425

Difficulties

1. Large hyper-parameter space
2. When to stop training?
3. Not robust when $p \gg n$
4. Large variation in model performance using same data and hyper-parameters.

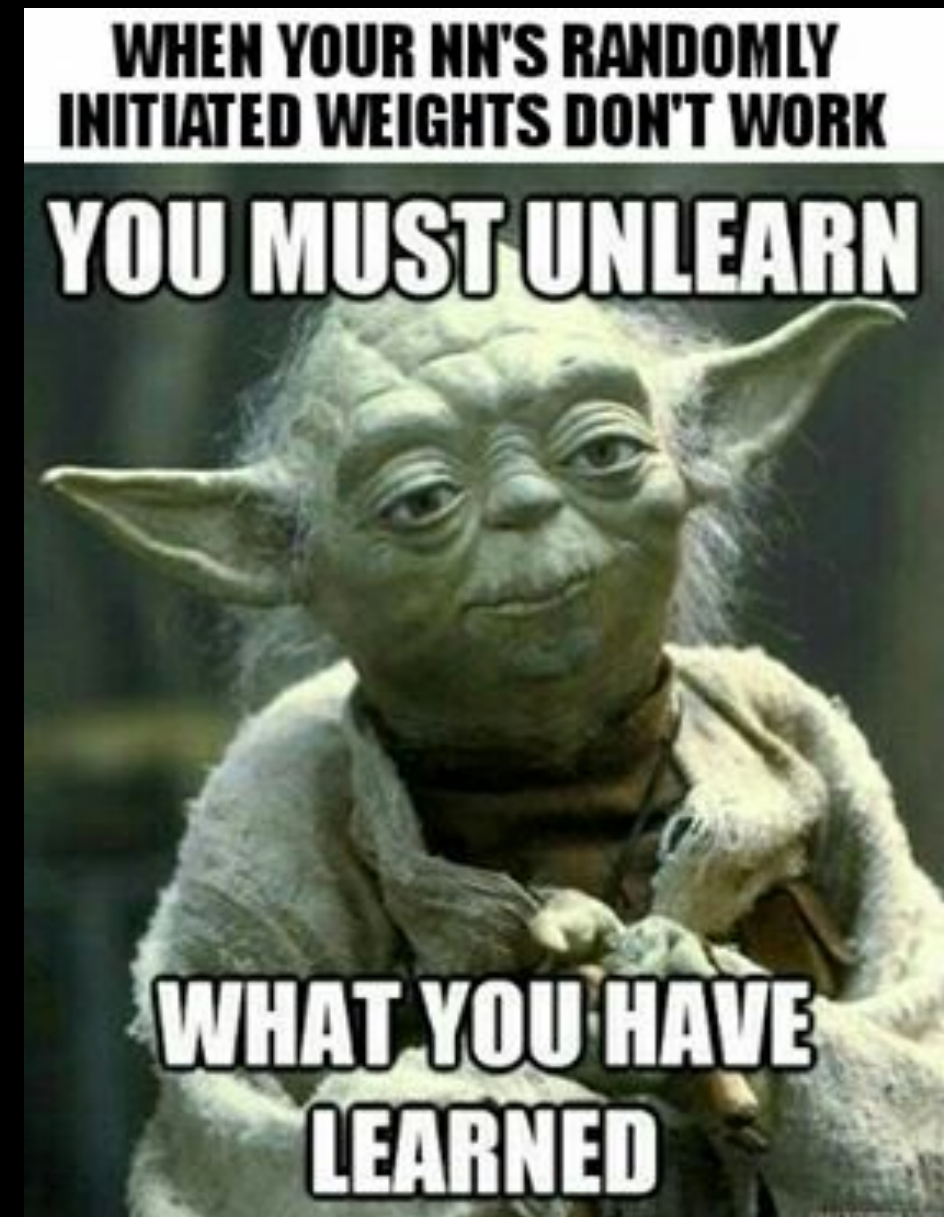
Random initialization of starting weights



Random Starting Weights

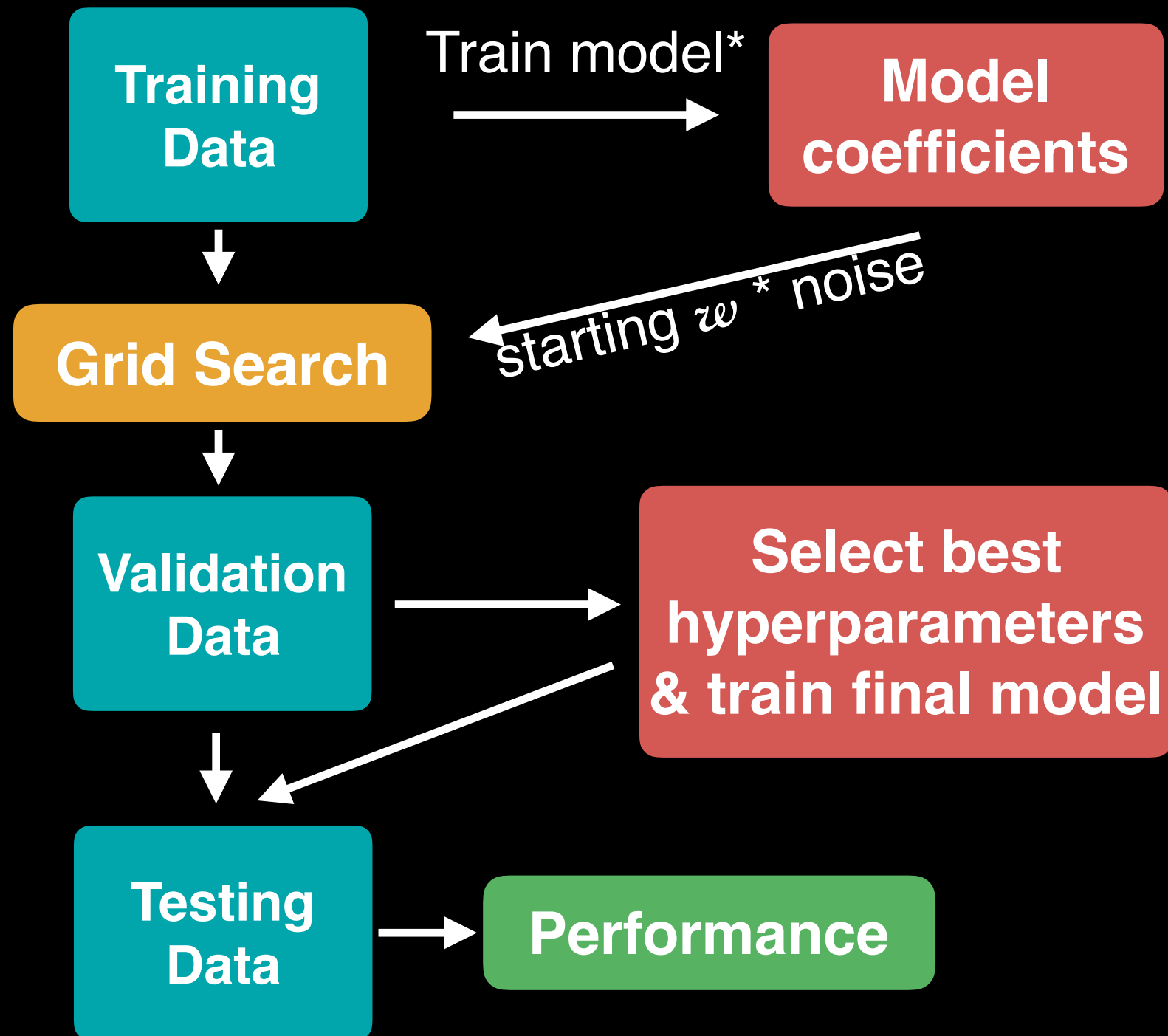
| | Node 1 | Node 2 | Node 3 |
|-------|--------|--------|--------|
| SNP A | -0.2 | 0.8 | -0.2 |
| SNP B | -0.4 | 1.2 | 0.6 |
| SNP C | 0.1 | 0.6 | 0.8 |
| SNP D | 0.4 | -0.1 | -0.5 |

****Reduces bias in the model****



Seeded starting weights approach

*rrBLUP, BayesB,
BL, and RF



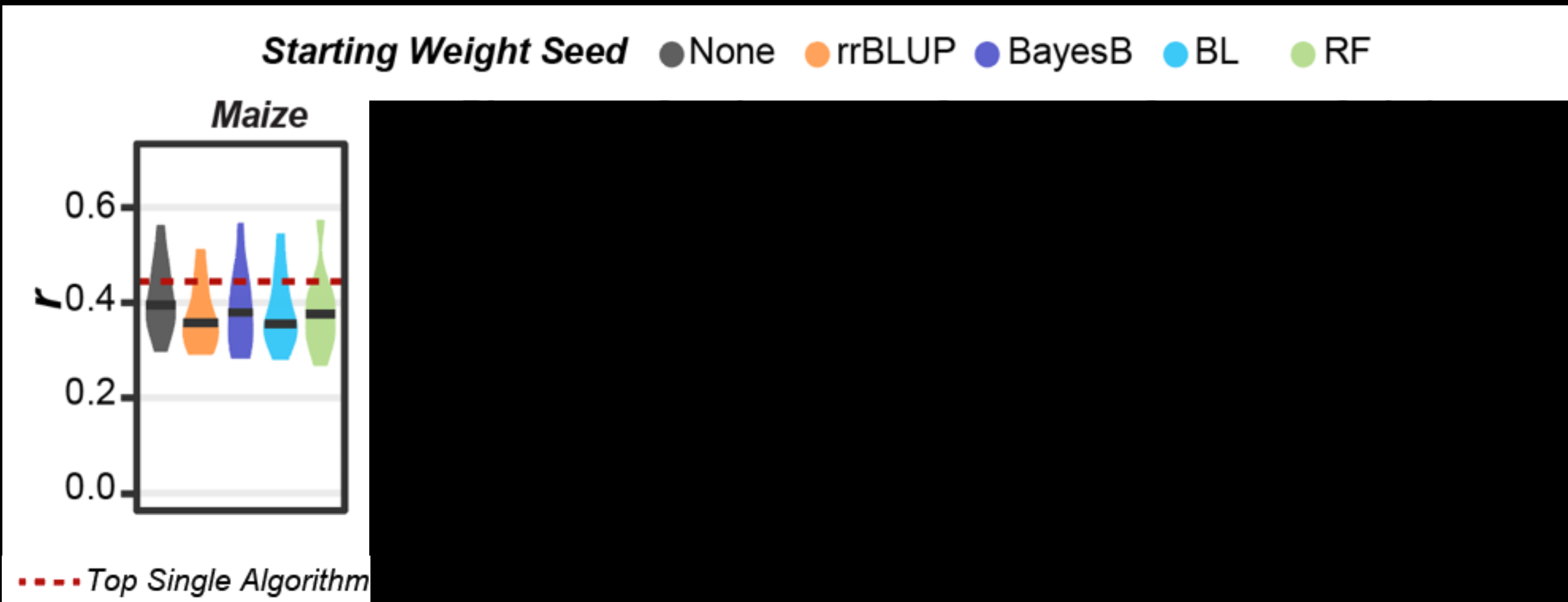
Seeded starting weights approach

25%
Seed Score +
Noise Infusion

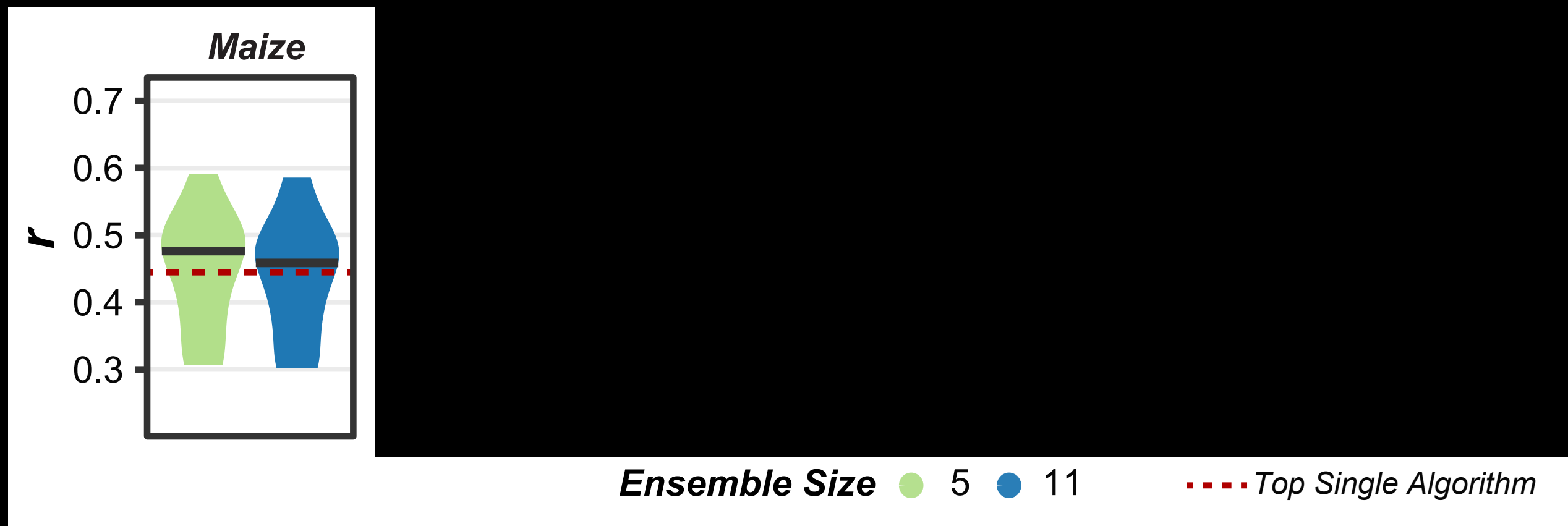
| | rrBLUP coef. |
|-------|-----------------|
| SNP A | -0.6 |
| SNP B | -0.4 |
| SNP C | 0.1 |
| SNP D | 0.4 |

| | Node 1 | Node 2 |
|-------|-----------|-----------|
| SNP A | -0.7 | -0.8 |
| SNP B | -0.5 | -0.4 |
| SNP C | 0.2 | 0.1 |
| SNP D | 0.4 | 0.5 |

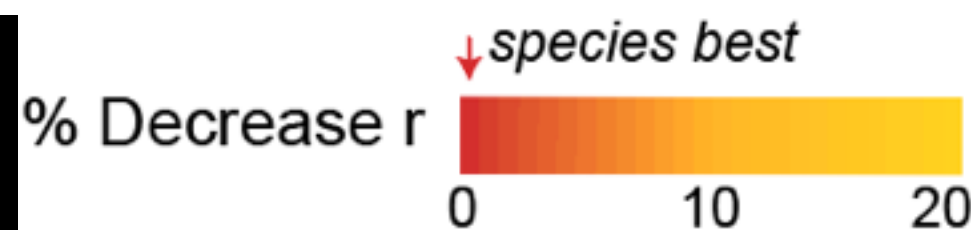
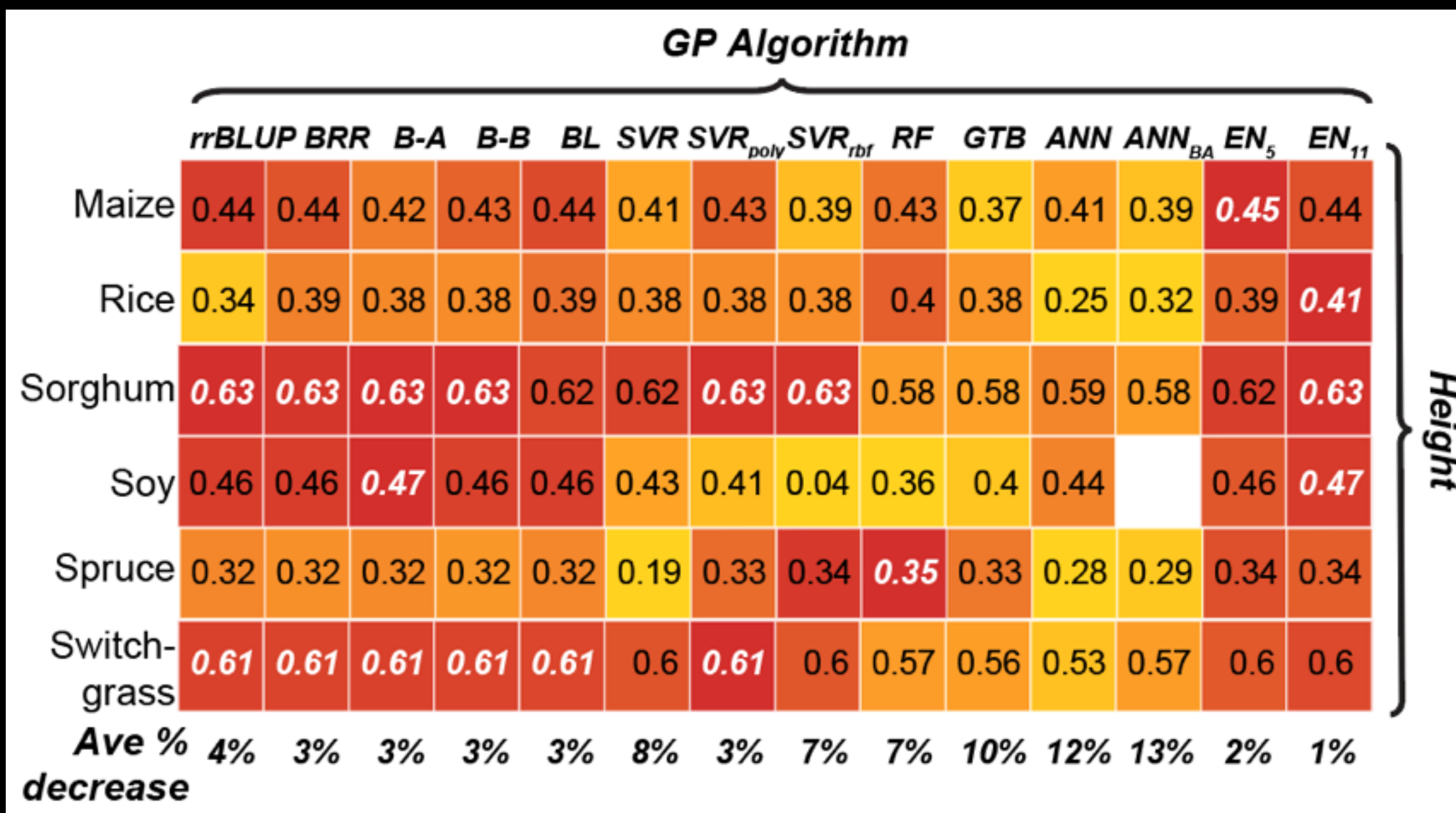
Seeded starting weights for top 8k maize markers sometimes improves performance



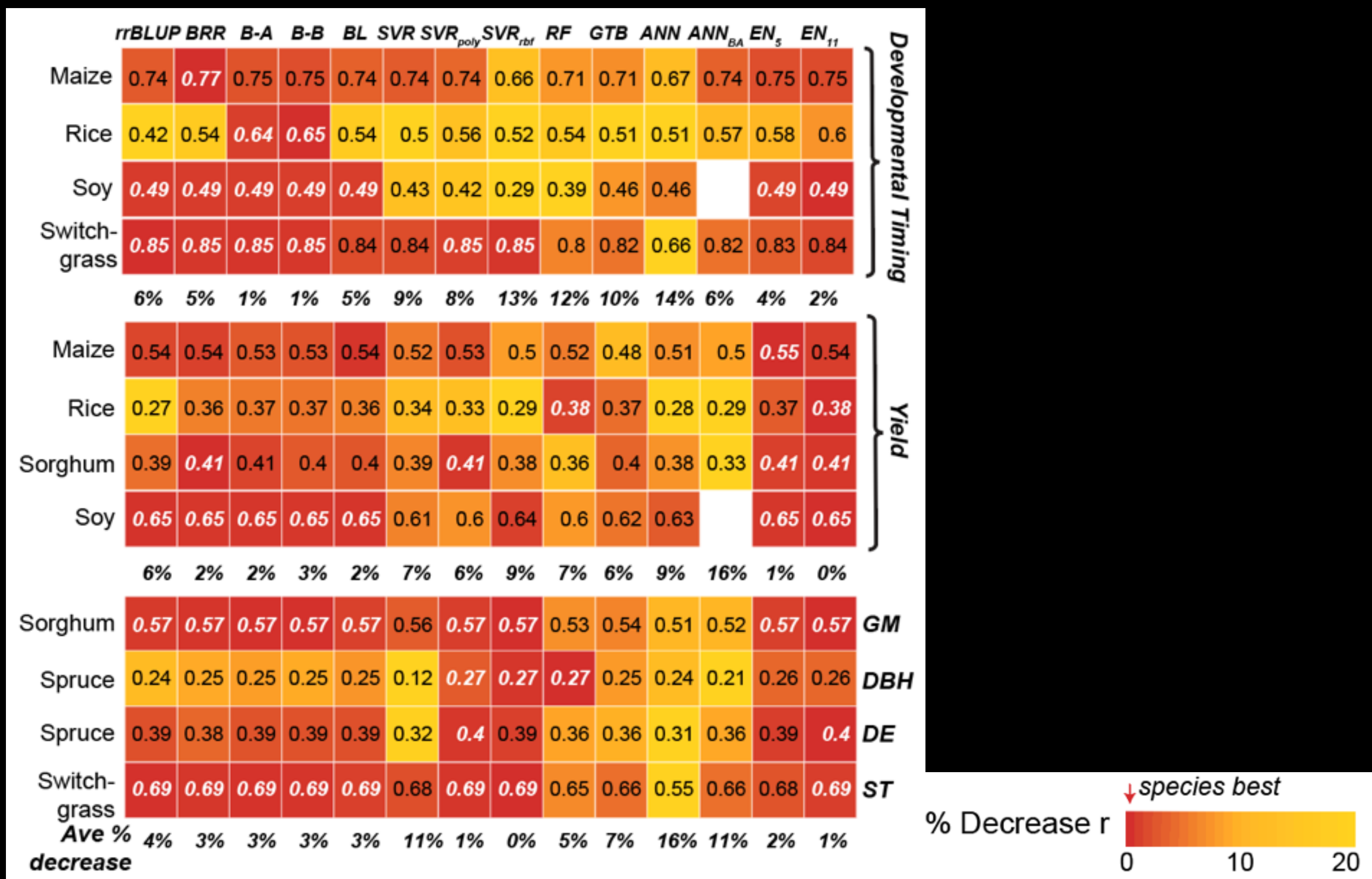
Ensemble Predictors: mean predicted trait value from 5/11 GP algorithms



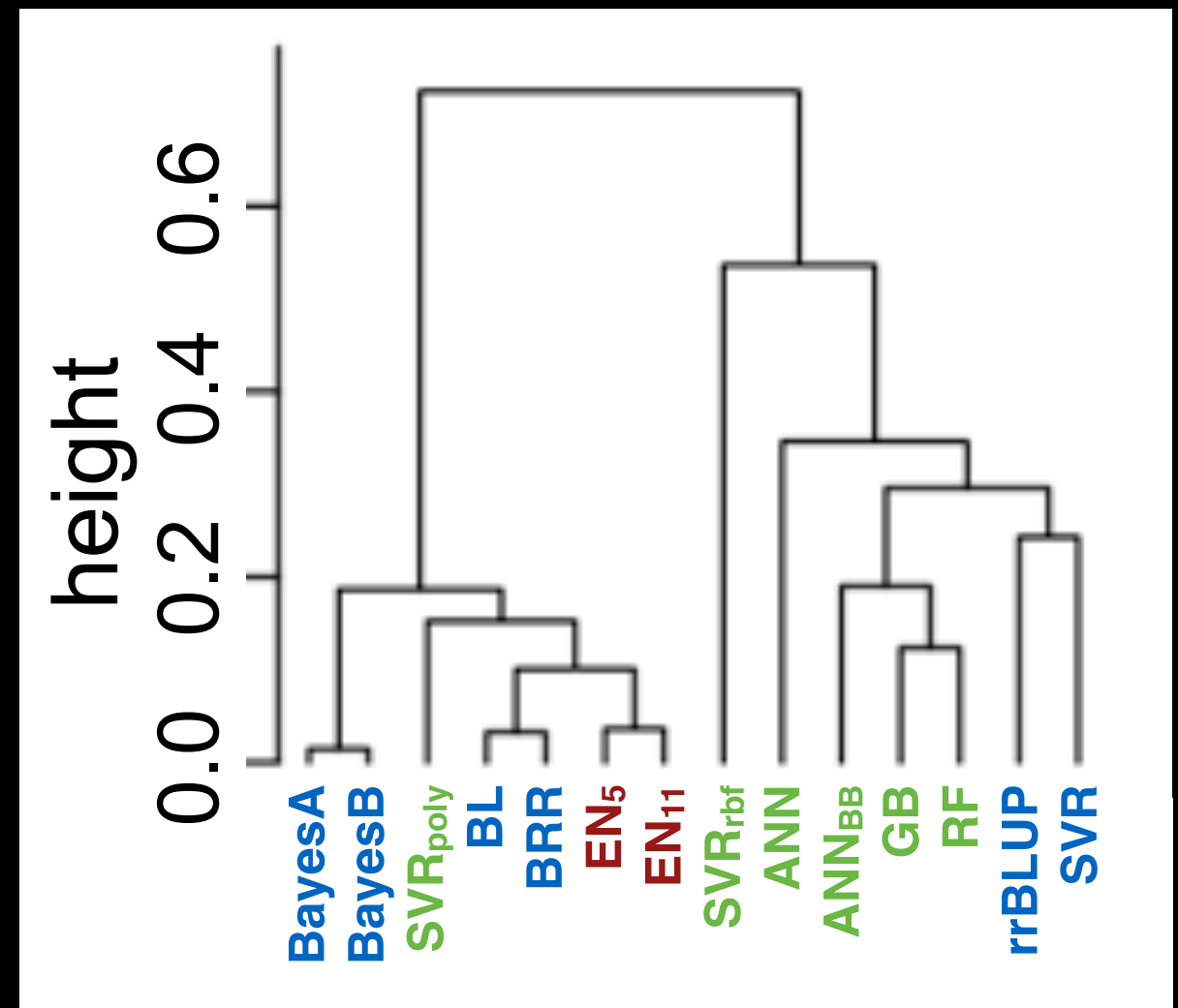
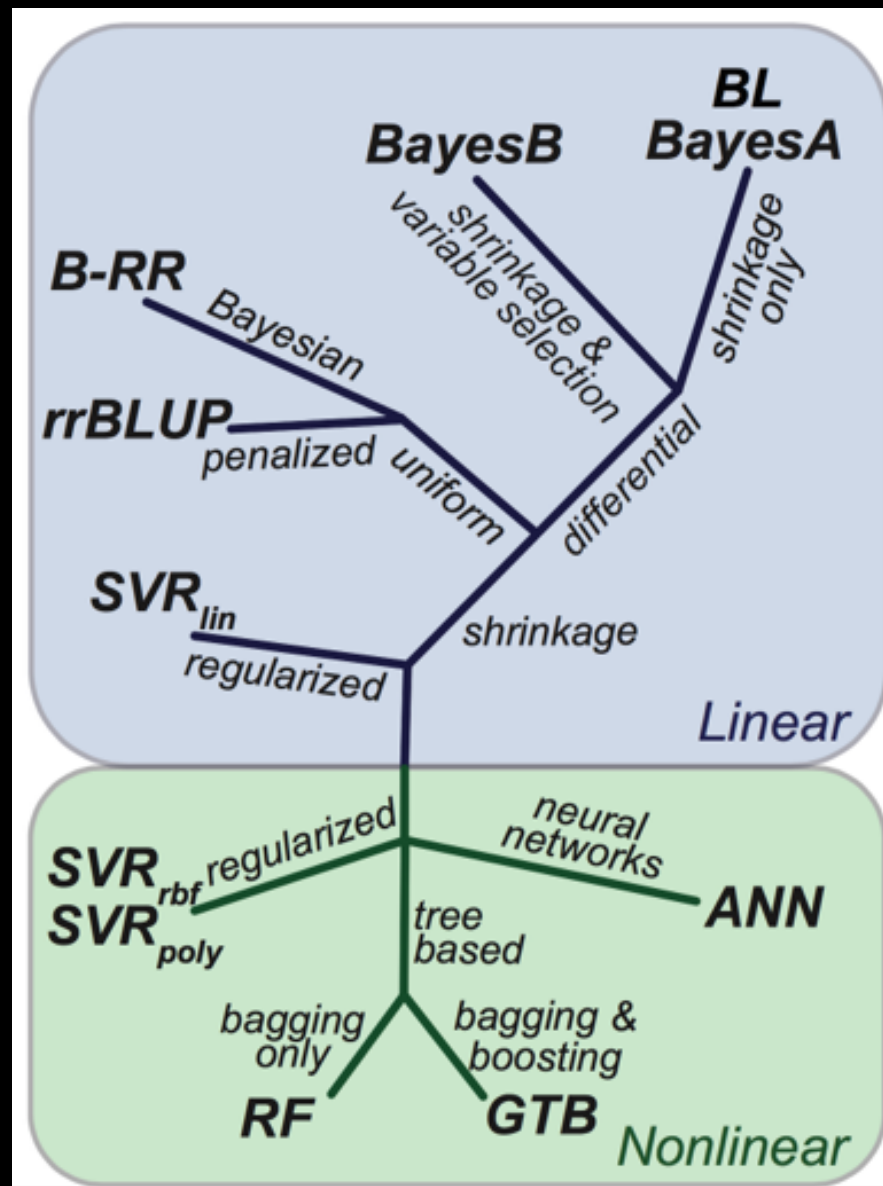
Final Benchmark Analysis Results...



Final Benchmark Analysis Results...



Hierarchical clustering of GP algorithms



Why do linear algorithms consistently perform well?

- ▶ H1: Complex biological systems generate signals that are consistent with linear, additive, genetic models (Hill, Goddard, & Visscher 2008).
- ▶ H2: The amount of training data available for most GP problems is insufficient for learning nonlinear interactions when $p \gg n$.

ML and Deep Learning Pipelines

azodichr / **ML-Pipeline**

<> Code

Issues 0

Pull requests 0

Projects 0

Wiki

Shiu Lab code base for Machine Learning implemented in SciKit-Learn

azodichr / **ANN_Pipeline**

<> Code

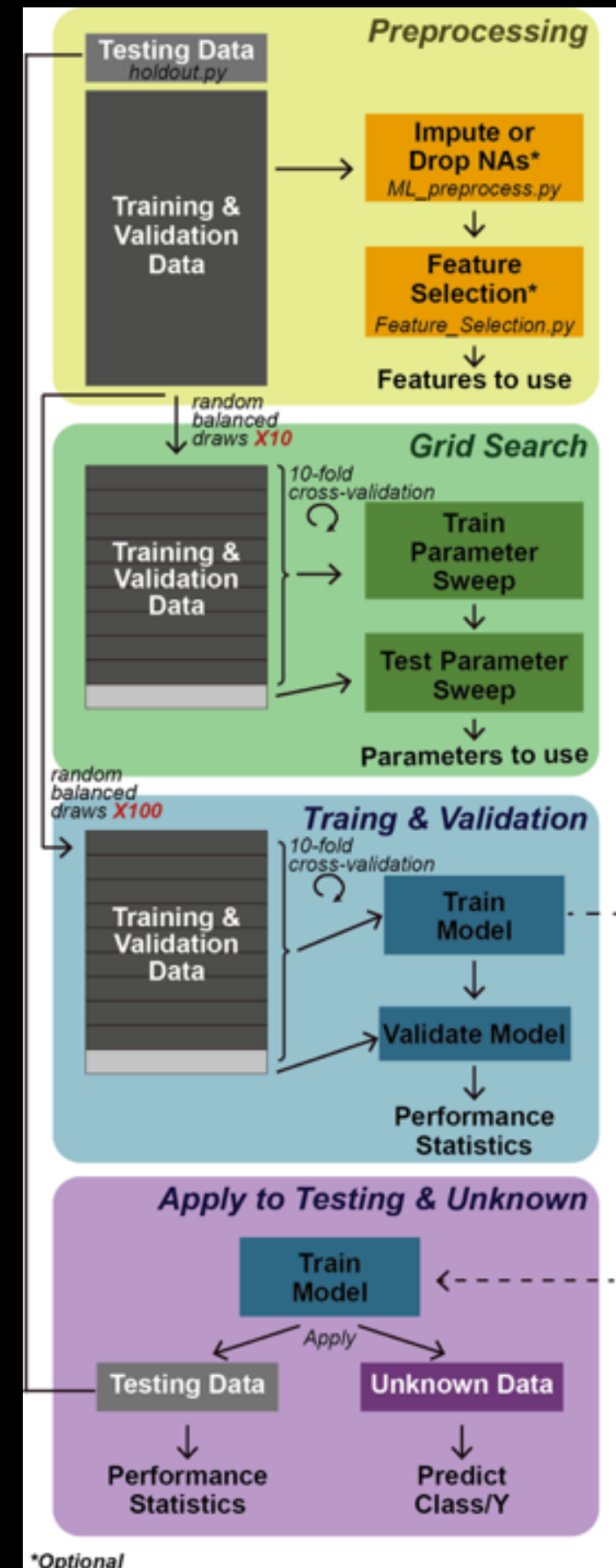
Issues 0

Pull requests 0

Projects 0

Wiki

Shiu Lab code base for Artificial Neural Networks implemented in Tensorflow

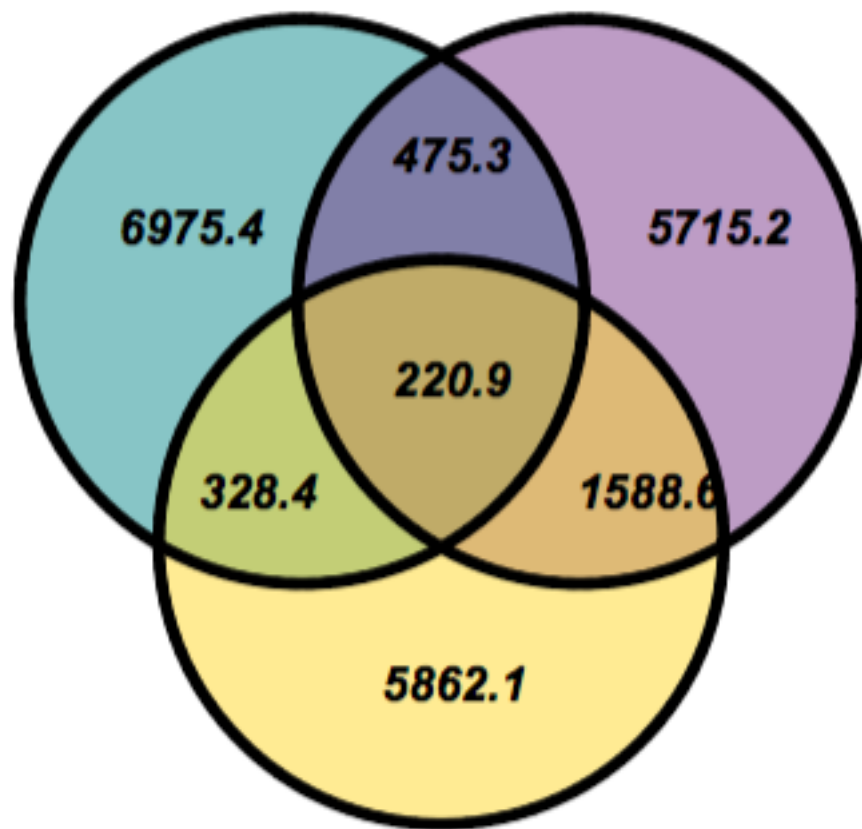


Future Directions

- ▶ Wrap up and submit GP algorithm comparison manuscript.
- ▶ Use gene expression levels as input into GP algorithms.
 - ▶ Are transcriptomes more predictive than genotype data?
 - ▶ Can we better predict traits by combining genotype & transcriptome data?
 - ▶ Is transcript data useful when dissecting GP models to learn about the genetic basis of a trait?
- ▶ Use “deep Connection Weight” (dCW) approach to interpret deep learning models to identify markers and interactions between markers important for a trait

Thanks!

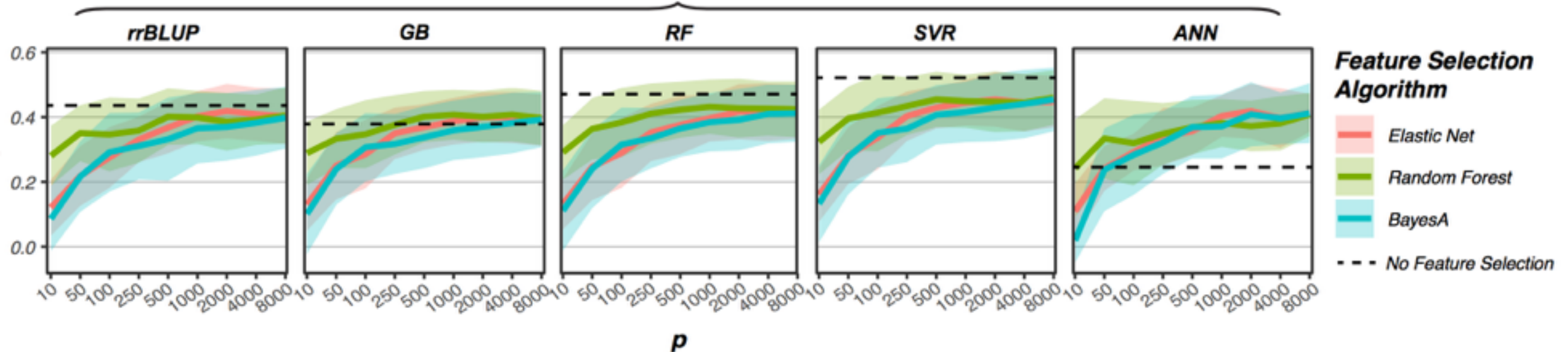
$p = 8000$



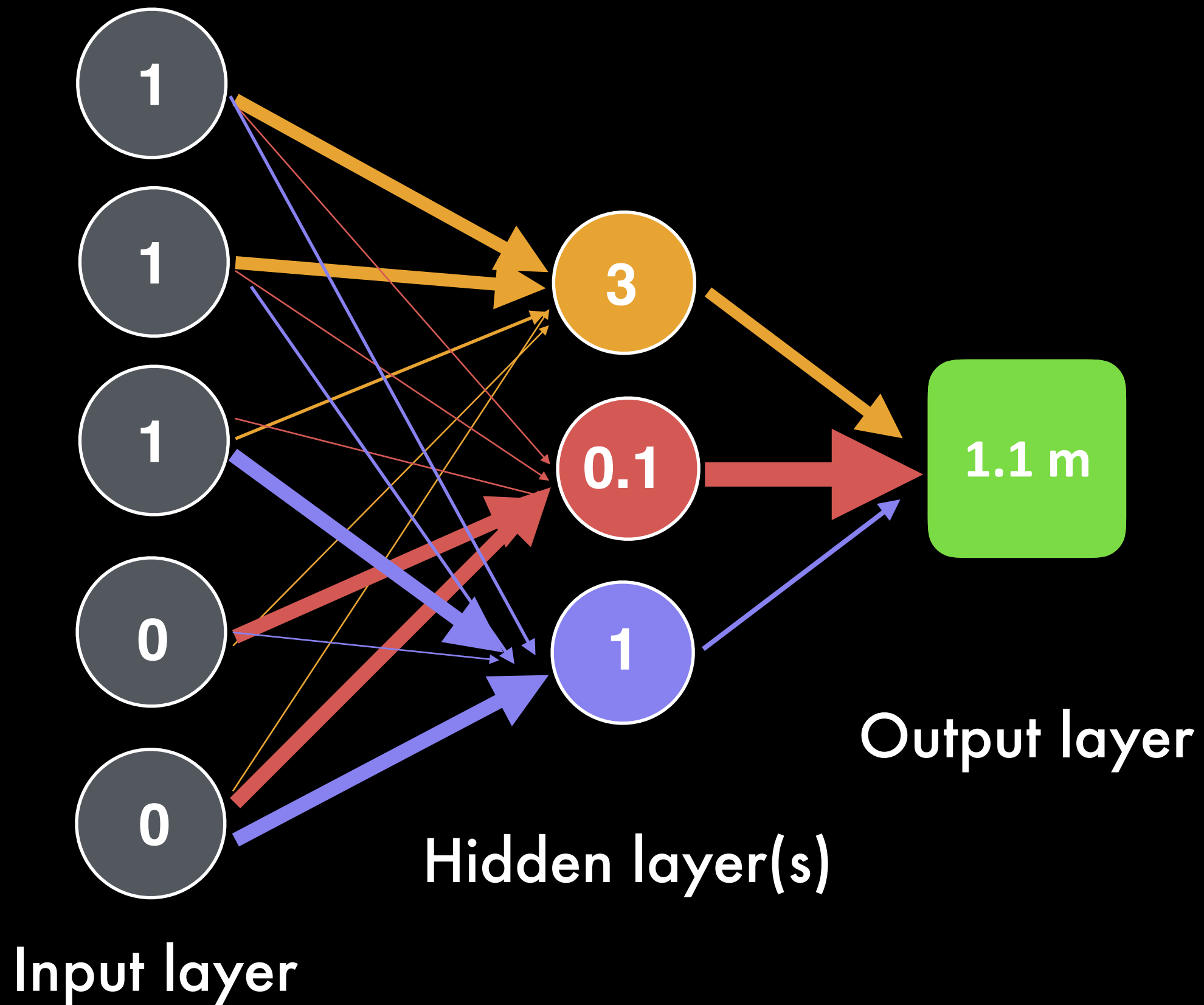
BayesA
 Elastic Net
 Random Forest

Random expectation
for triple overlap:
 ~ 10

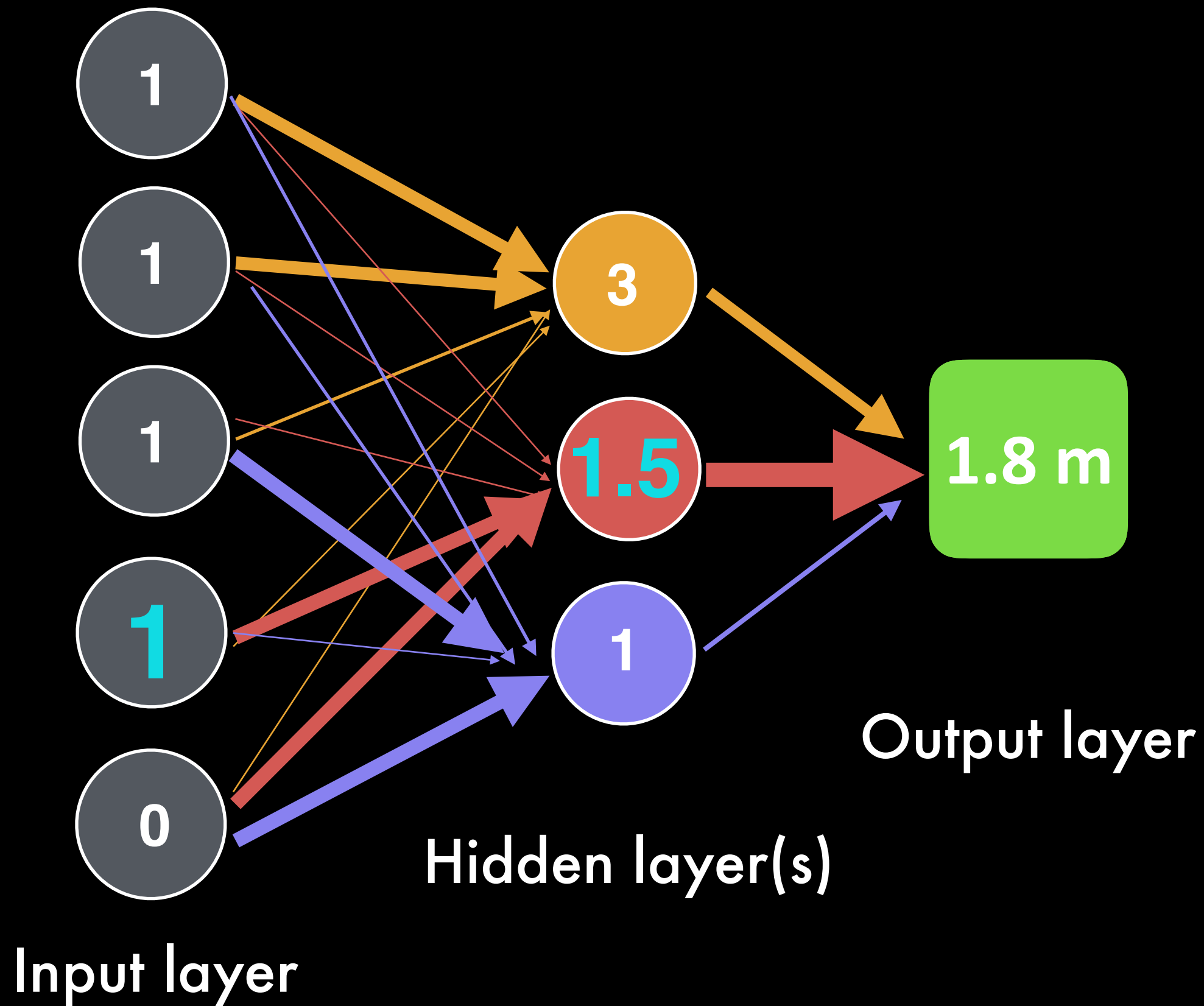
GP Algorithm



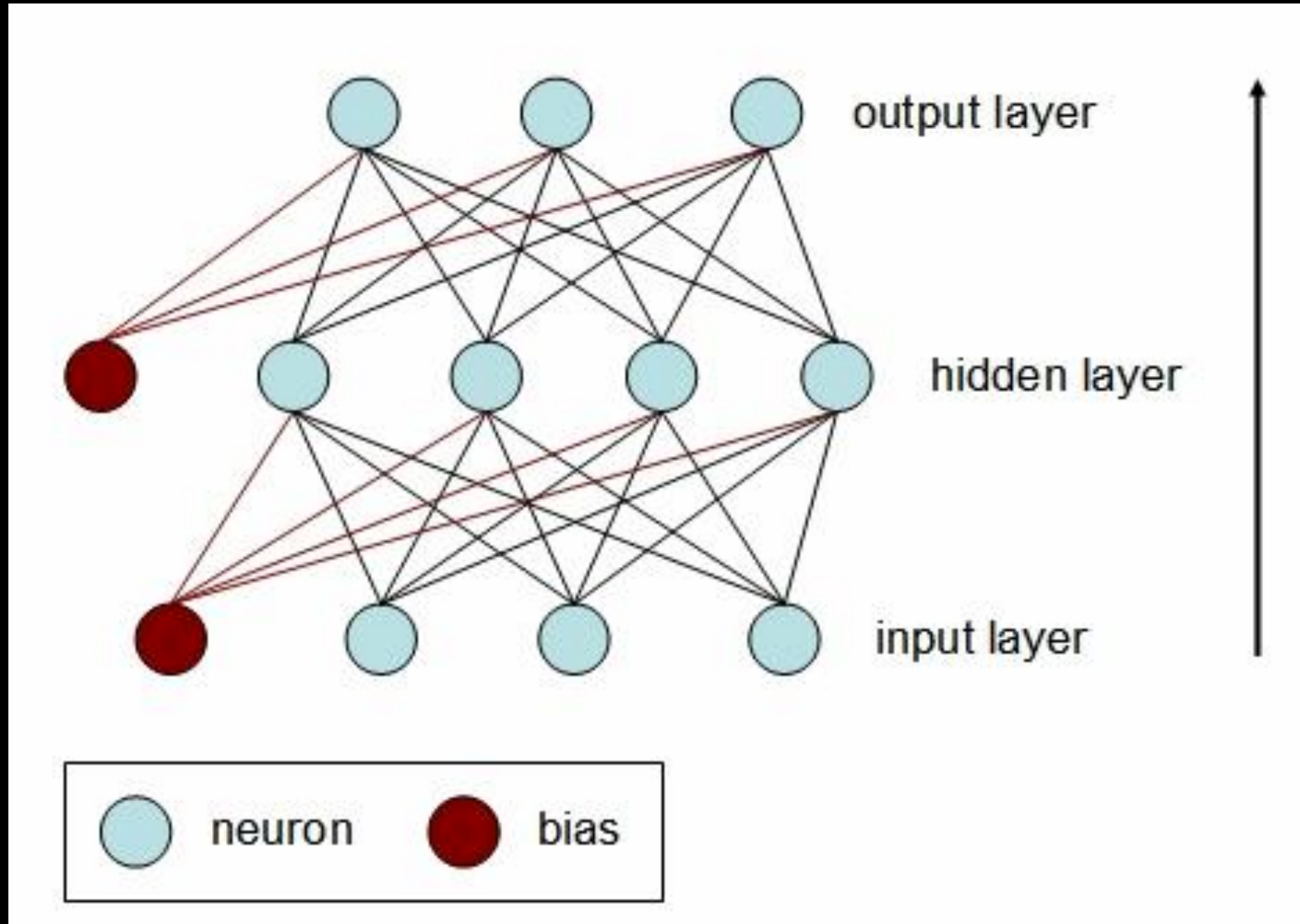
Genotype X



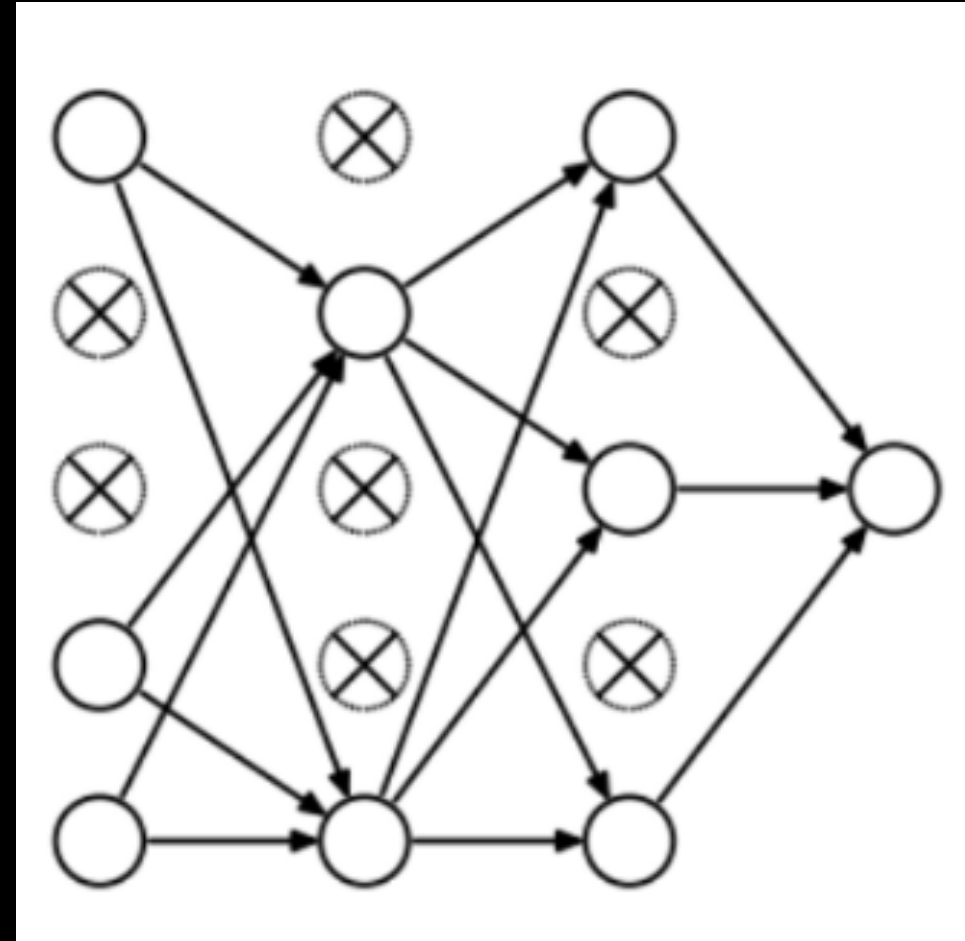
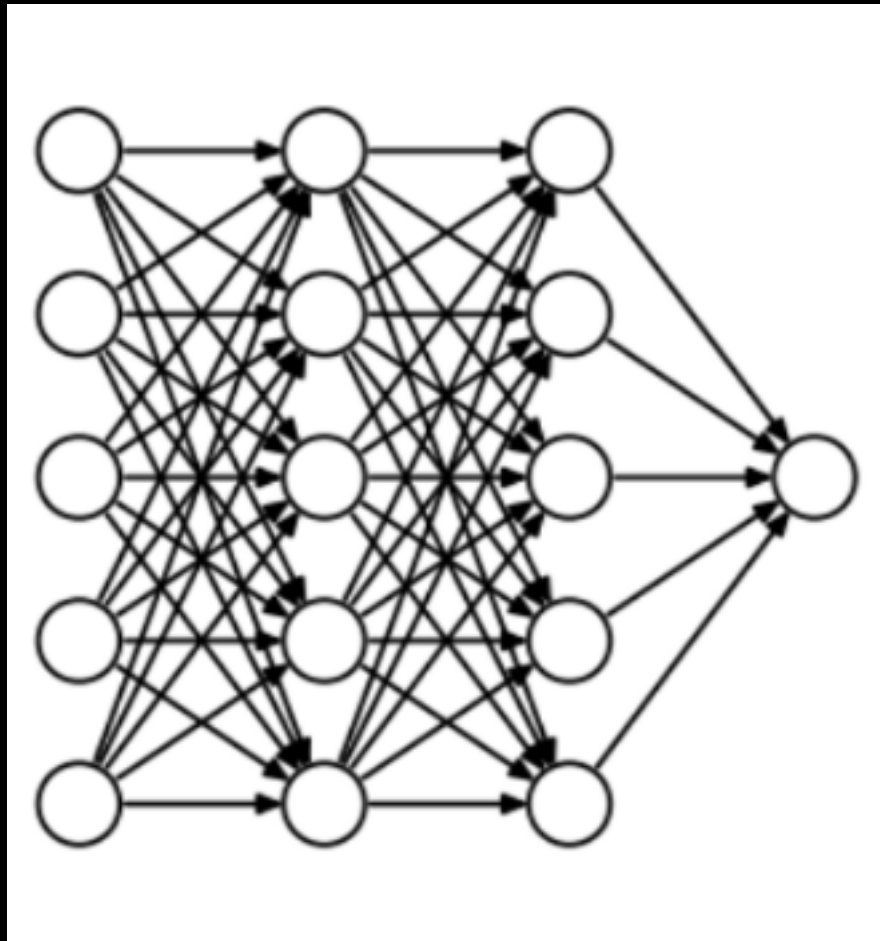
Genotype Y



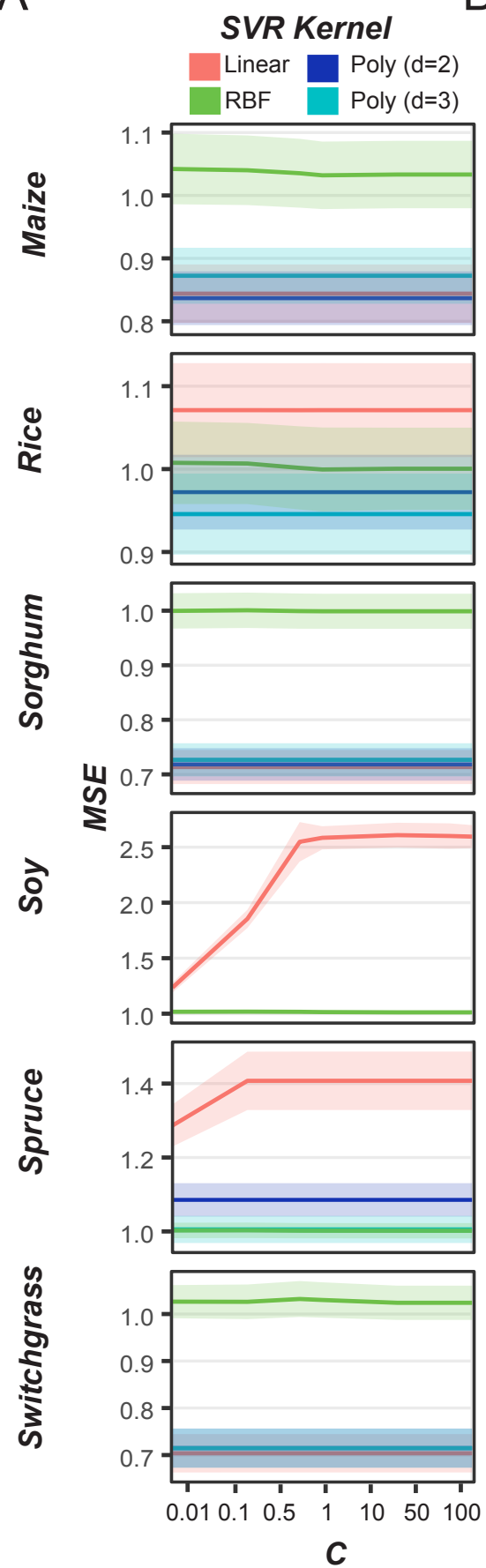
Role of the bias term....



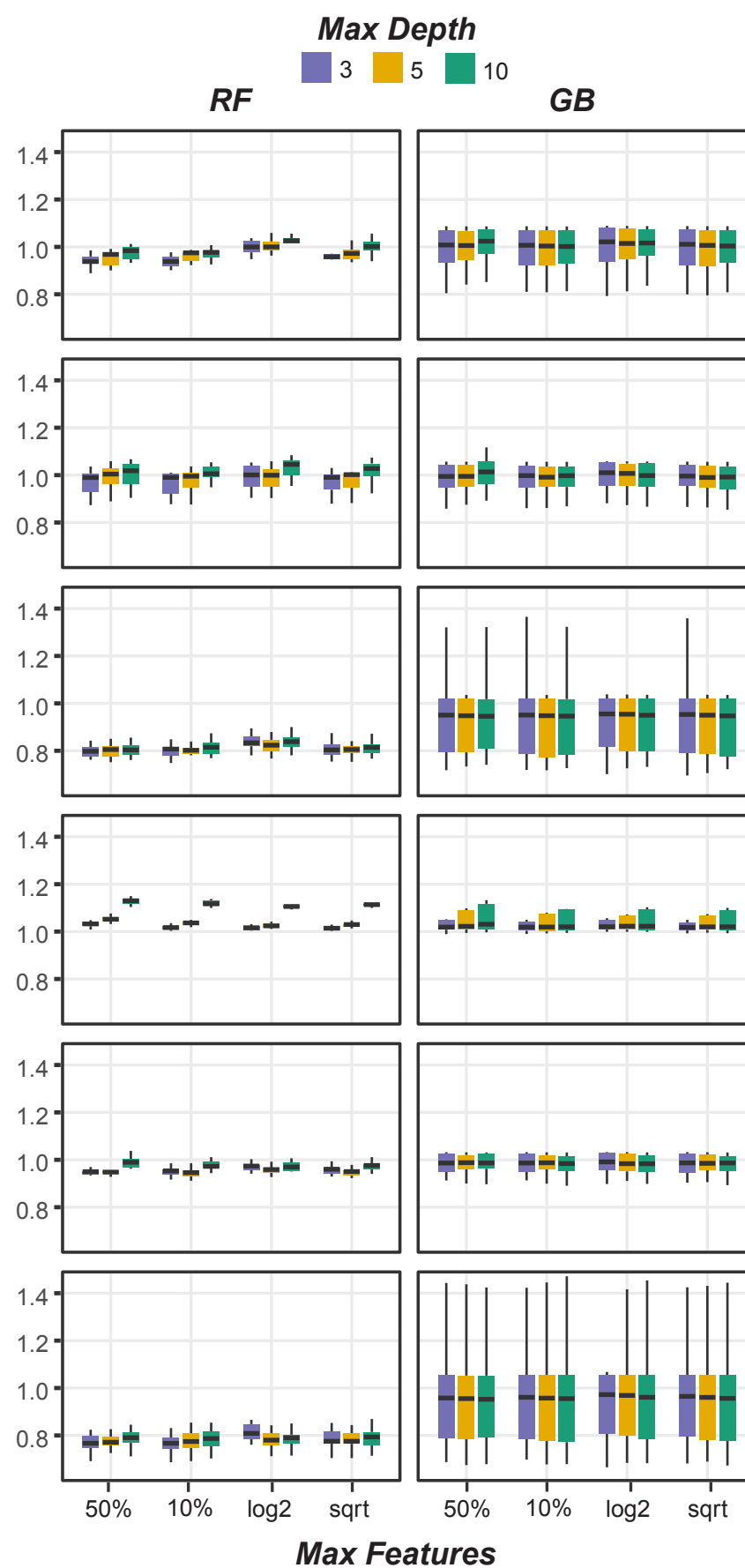
Drop out Explanation



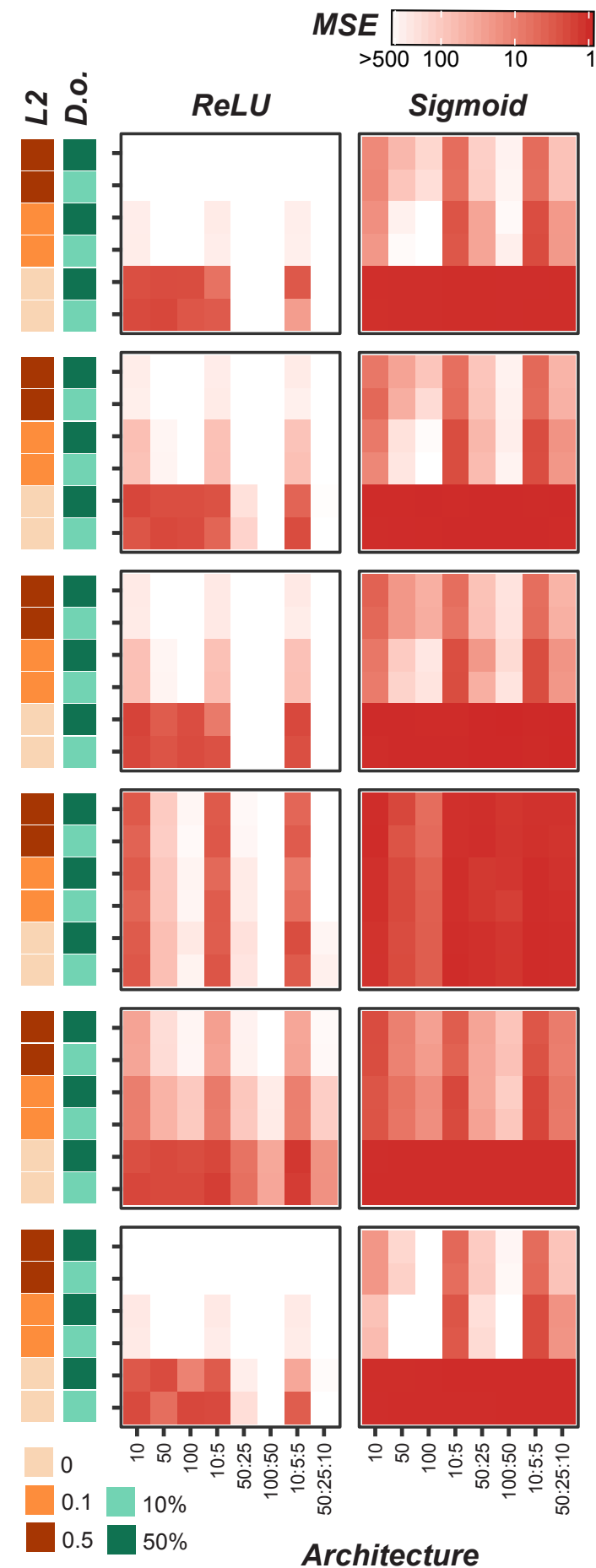
A

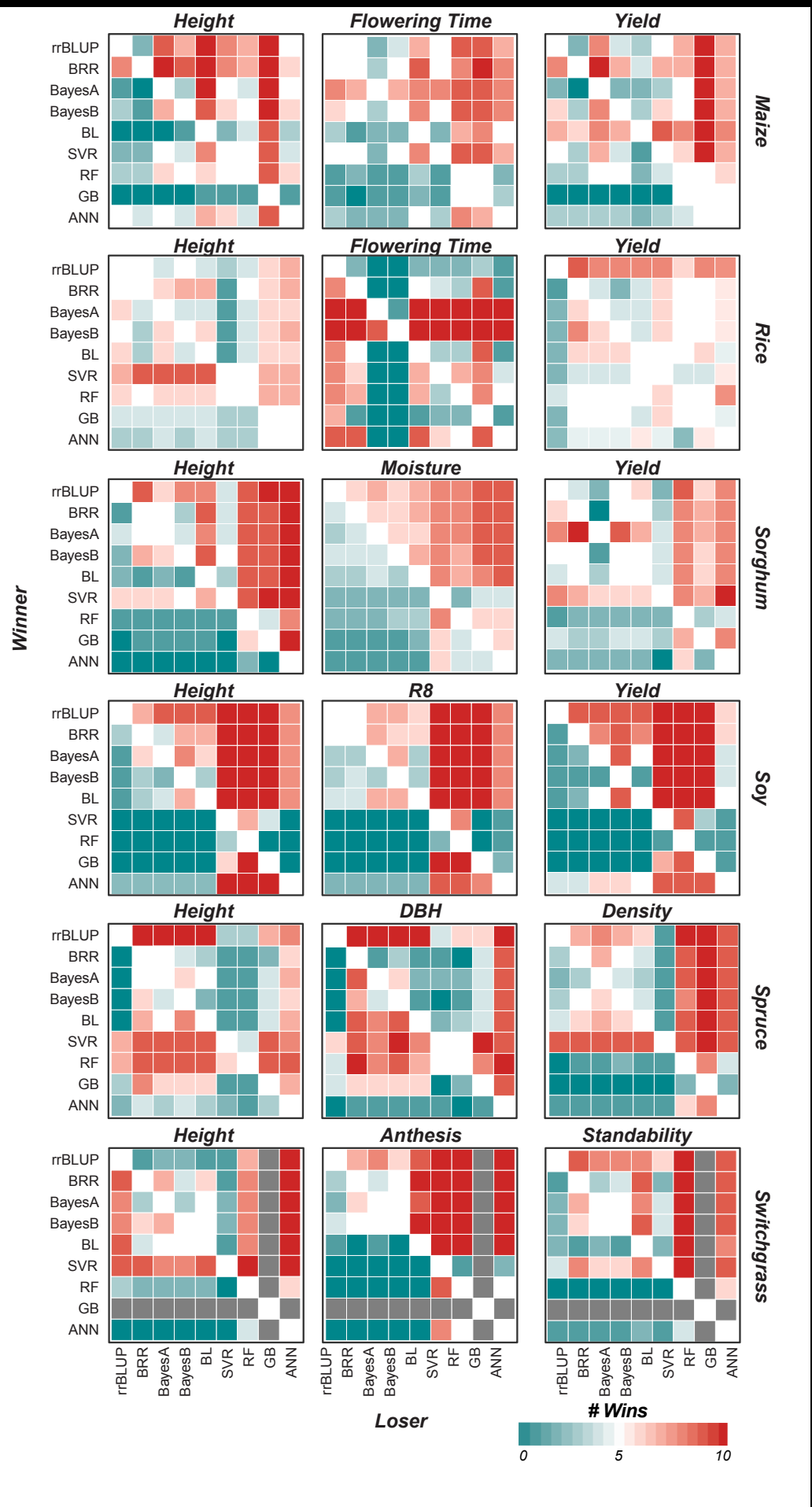


B



C





1st Hidden Layer

| | Node A1 | Node A2 | ... | Node An |
|-------|---------|---------|-----|---------|
| SNP A | -0.66 | -0.8 | ... | -0.21 |
| SNP B | -0.44 | 0.66 | ... | 0.64 |
| SNP C | 0.18 | 0.66 | ... | 0.84 |
| SNP D | 0.45 | 0.78 | ... | 0.52 |

2nd Hidden Layer

| | Node B1 | Node B2 | ... | Node Bn |
|---------|---------|---------|-----|---------|
| Node A1 | -0.44 | 0.18 | ... | -0.38 |
| Node A2 | 0.15 | 0.15 | ... | 0.24 |
| ... | ... | ... | ... | ... |
| Node An | 1.59 | 0.65 | ... | 0.37 |

Output Layer

| | Output Layer |
|---------|--------------|
| Node B1 | 0.00001 |
| Node B2 | 1.44 |
| ... | ... |
| Node Bn | 2.14 |