

**AI-based models and environmental data improves
genomic selection in canola**

Dr. Iulian Gabur

Department of Plant Breeding, Justus Liebig University Giessen,
Germany



Why A.I. in plant breeding?

Breeders equation

$$\Delta G = \frac{i * r * \delta_G}{L_g}$$

i = Selection intensity

r = Selection accuracy

δ_G = Genetic variance

L_g = Length of breeding cycle

Why A.I. in plant breeding?

Breeders equation

$$\Delta G = \frac{i * r * \delta_G}{L_g}$$

How to improve ΔG ?

i = Selection intensity

r = Selection accuracy

δ_G = Genetic variance

L_g = Length of breeding cycle

Select more intensely

Select more accurately

Introduce new variation

Do all faster!

Why A.I. in plant breeding?

Breeders equation

$$\Delta G = \frac{i * r * \delta_G}{L_g}$$

i = Selection intensity

r = Selection accuracy

δ_G = Genetic variance

L_g = Length of breeding cycle

How to improve ΔG ?

Select more intensely

Select more accurately

Introduce new variation

Do all faster!

ML models?

Digital phenotyping, drones (MLP, NN)

3D plant phenomics (FS, RF, CNN)

Conserve/increase diversity (GA, NLP)

Genomic selection (SVM, GBM, NN)

Why A.I. in plant breeding?

Breeders equation

$$\Delta G = \frac{i * r * \delta G}{L_g}$$

KEYNOTE TALK

#034 Omics-based optimisation of hybrid performance and heterosis in winter oilseed rape

Rod Snowdon, Sven Weber, HueyTyng Lee, Agnieszka Golicz, José Montero, Mauricio Orantes-Bonilla, Lennard Ehrig, Iulian Gabur, Eva Herzog, Amine Abbadi, Matthias Frisch

#093 Connecting rapeseed plant architecture to yield via 3D-imaging

Sven Weber, Andreas Eckert, Lennard Ehrig, Rod Snowdon, Andreas Stahl

#049 HaploMAGIC: phasing and accurately detecting recombination in multi-parental populations with genotyping errors

Jose Antonio Montero-Tena, Nayyer Abdollahi-Sisi, Amine Abbadi, Matthias Frisch, Rod J. Snowdon, Agnieszka Golicz

#074 Unravelling the *Brassica napus* epigenetic network with an integrated multi-omics approach

Silvia Zanini, Rod Snowdon, Agnieszka Golicz

#073 Phenomic selection: predicting quantitative traits in canola with NIRS-profiles

Lennard Ehrig, Sven Weber, Stefan Abel, Reinhard Hemker, Dirk Stulgies, Milka Malenica, Amine Abbadi, Benjamin Wittkop, Rod Snowdon, Andreas Stahl

Reducing dimensionality in big-data

Idea:

Identify the most important SNP for each Trait-SNP association (**feature selection**) and use them for training and predictions.

Avoid the **curse of dimensionality**

and

overfitting !



two wind turbines in two dimensional and three dimensional space

	Underfitting	Just right	Overfitting
Symptoms	<ul style="list-style-type: none">• High training error• Training error close to test error• High bias	<ul style="list-style-type: none">• Training error slightly lower than test error	<ul style="list-style-type: none">• Very low training error• Training error much lower than test error• High variance
Regression illustration			
Classification illustration			
Deep learning illustration			

Using linear and non-linear models in canola genomic selection

Idea:

Identify the most important SNP for each Trait-SNP association (**feature selection**) and use them for training and predictions.

Testing:

Used **three** linear and non-linear approaches for **feature selection**:

- **PCA** - data reduction (100 features)
- **RF-based** reduction (1000 features)
- Random selection of features (1000 features)

Statistical models used in a cross validation pipeline:

- Ridge-Regression Best Linear Unbiased Prediction (RR-BLUP)
- Least Absolute Shrinkage and Selection Operator (LASSO)
- "general" Linear Model (GLM)
- **Gradient Boosting Machine (GBM)**
- **Artificial Neural Network (ANN)**
- **Random forest (RF)**

Hybrid canola dataset

PREDICT data set – spring-type *B. napus* were evaluated at four different locations across Denmark, Germany, Poland and Estonia. Jan et al., 2016 PLoS ONE 11(1): e0147769.

Design:

two male sterile testers (MSL-T1 and MSL-T2, NPZ Lembke, Hohenlieth, Germany) and a diverse population of 475 spring-type “00” *B. napus* cultivars = 950 hybrids

SNP genotyping

Brassica 60K Illumina Infinium™ SNP genotyping array of parental lines

Population

950 F1 hybrids – *in silico* crossing scheme (0, 1, 2)

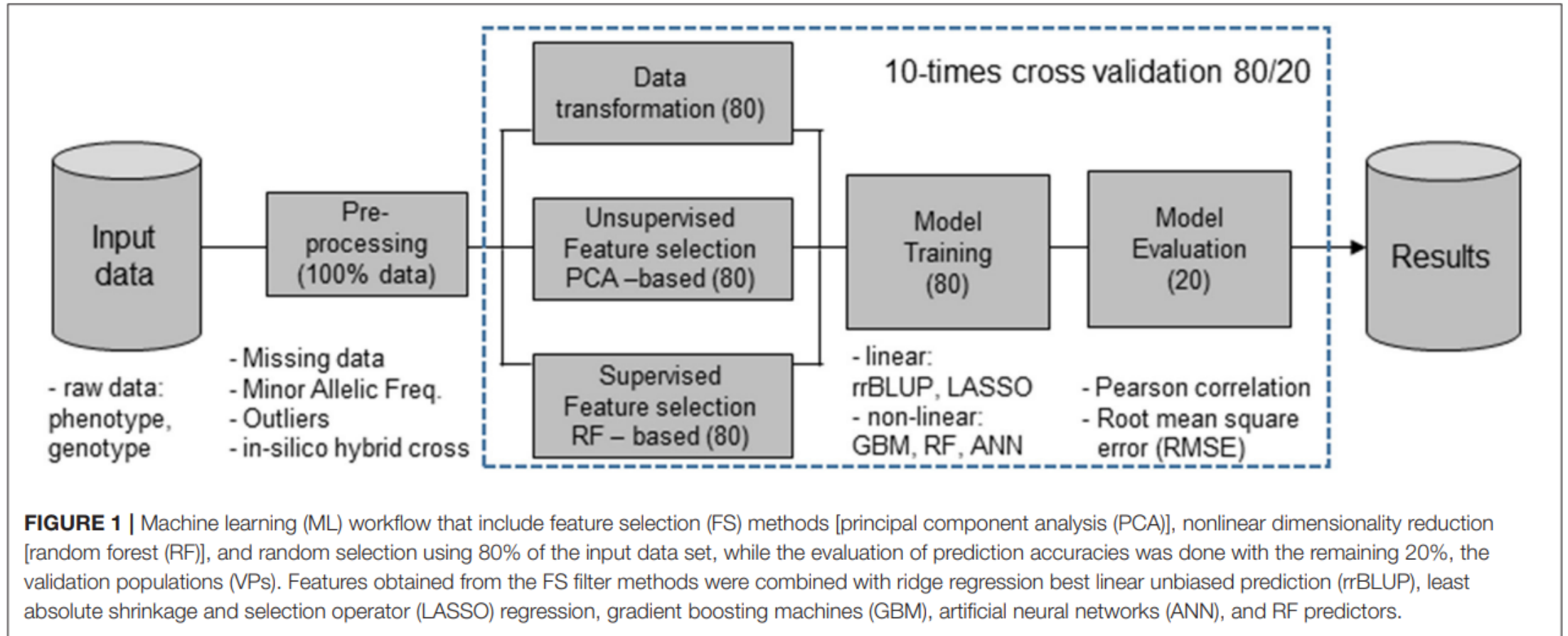
Werner et al. 2018 Theor Appl Genet 131, 299–317

Phenotype data of hybrids (14) :

Yield [dt/ha], Oil-Yield [dt/ha], Moisture [%], Oil [%], Protein Seed [%], Protein Meal [%]

GSL Meal ymol/g seed, GSL ymol/g seed, TSW [g], Days to flower, Straw length [cm]

Workflow



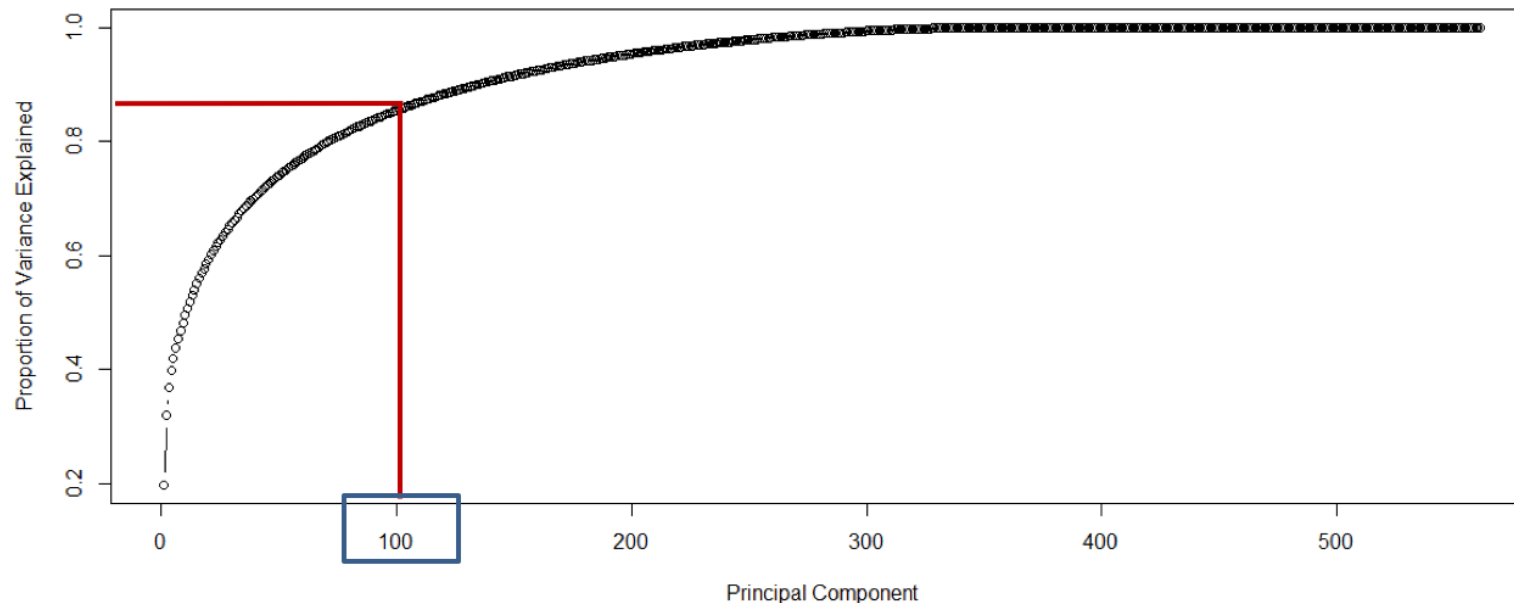
Top 100 PCs explain 87% of genetic variation in the canola dataset

$$\mathbf{x} = [x_1, x_2, \dots, x_d], \quad \mathbf{x} \in \mathbb{R}^d$$

$$\downarrow \mathbf{x}\mathbf{W}, \quad \mathbf{W} \in \mathbb{R}^{d \times k}$$

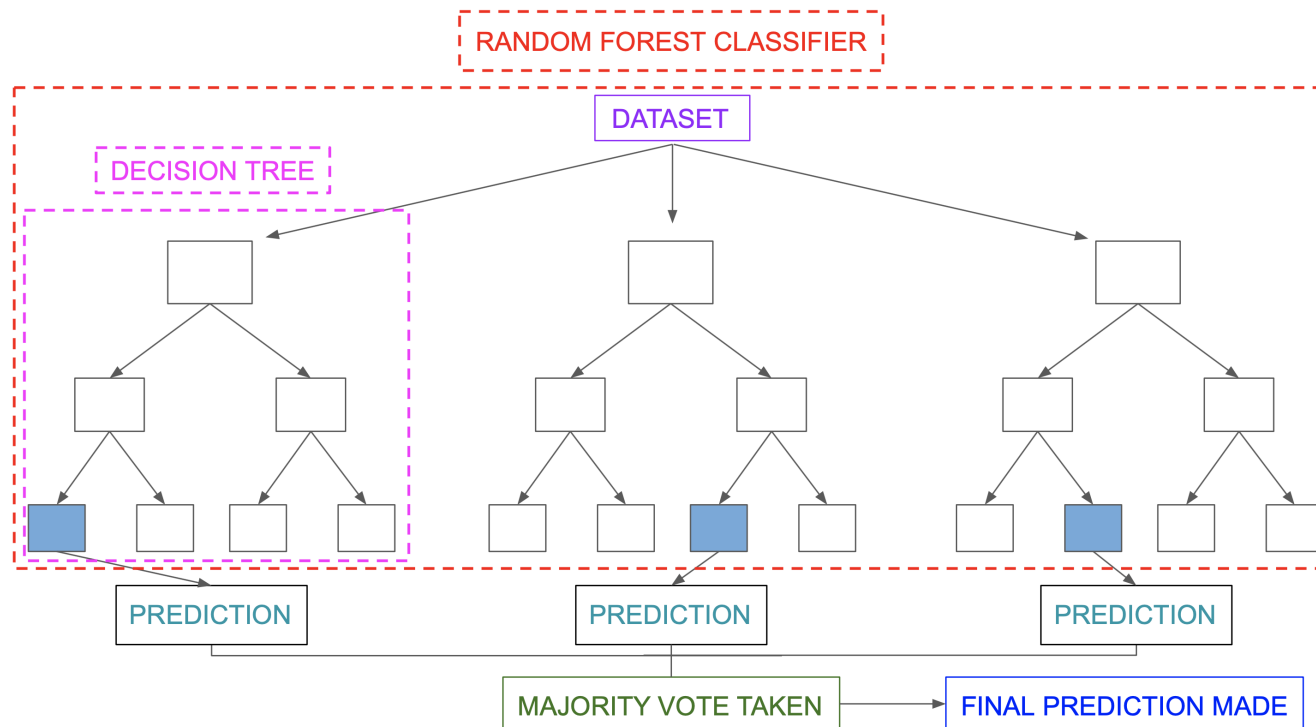
$$\mathbf{z} = [z_1, z_2, \dots, z_k], \quad \mathbf{z} \in \mathbb{R}^k$$

1. Select k eigenvectors for k largest eigenvalues ($k \leq d$).
2. Construct **matrix \mathbf{W}** from the “top” k eigenvectors.
3. Transform the d -dimensional input **dataset \mathbf{X}** using the projection **matrix \mathbf{W}** to obtain the new k -dimensional feature subspace.

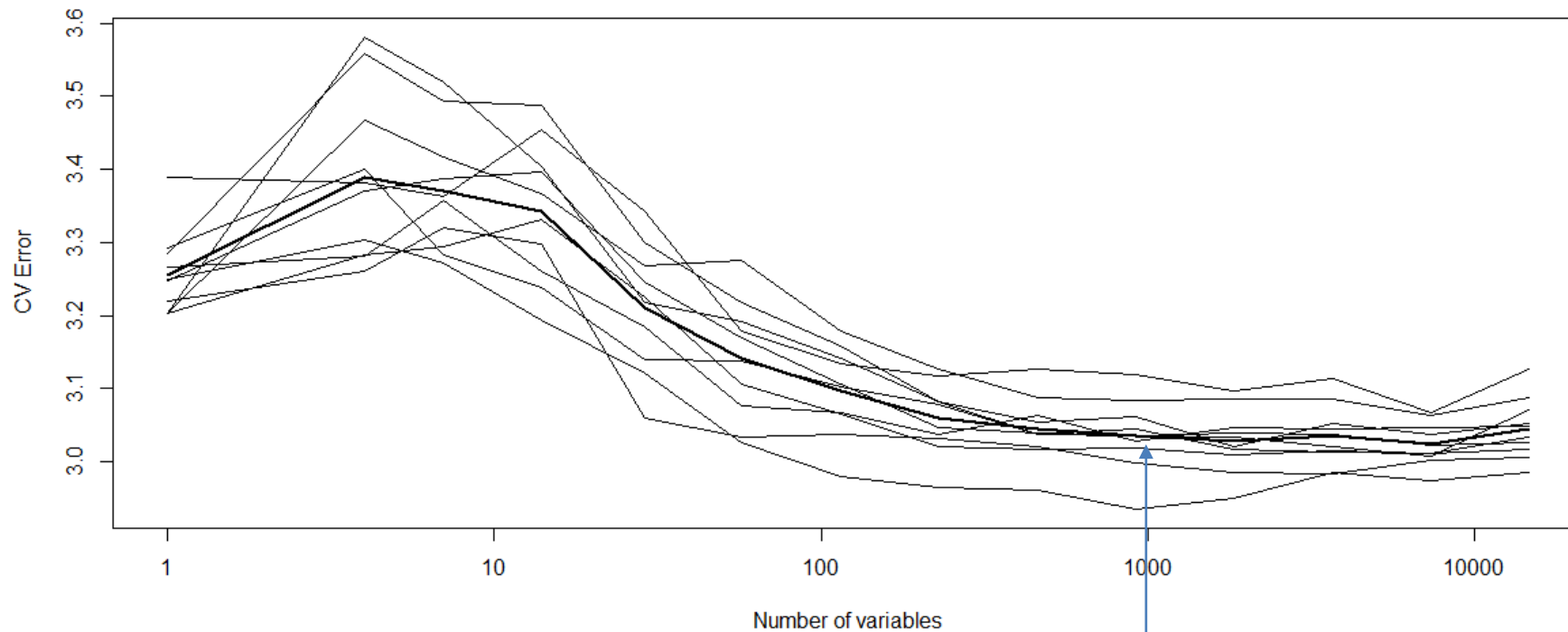


Random forest (RF) based data reduction for identification of top 1000 features

Identify the maximum variance in SNP matrix and calculate **feature importance** scores for each feature based on the 'Gini' criterion (RF trees).



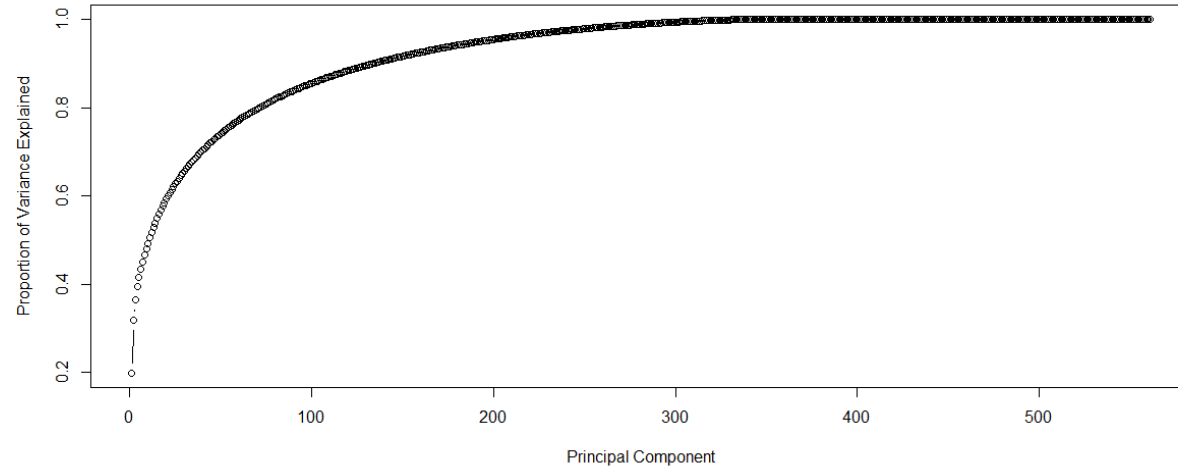
1000 trees reduce the error rate of Random Forest model



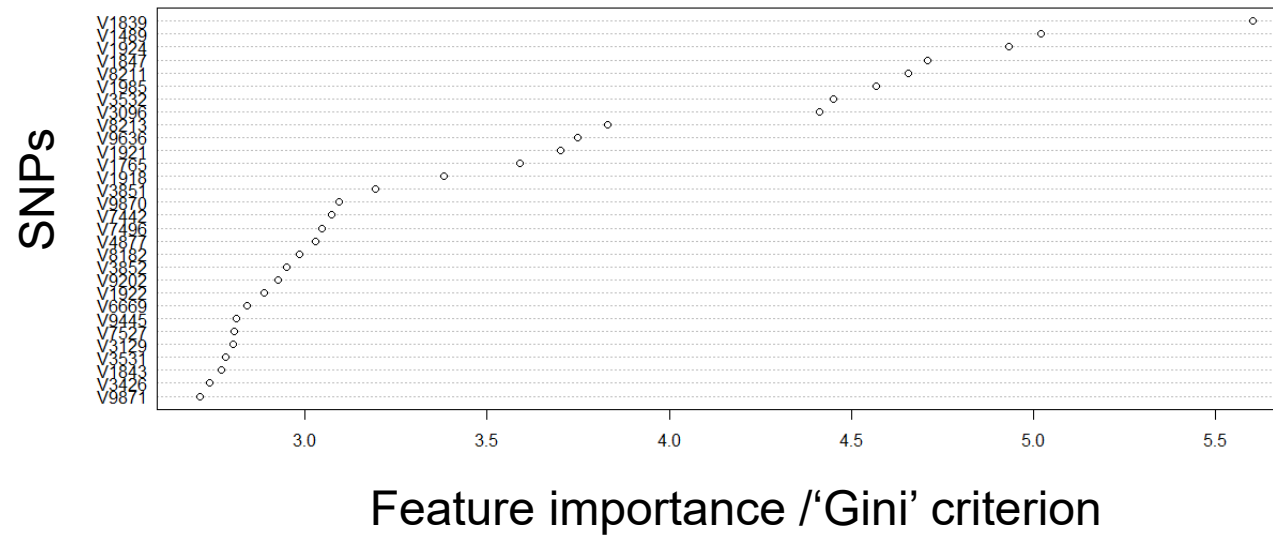
Select 1000 trees for predictions

Random Forest – based selection using features importance (top 1000)

PCA-based selection



RF selection



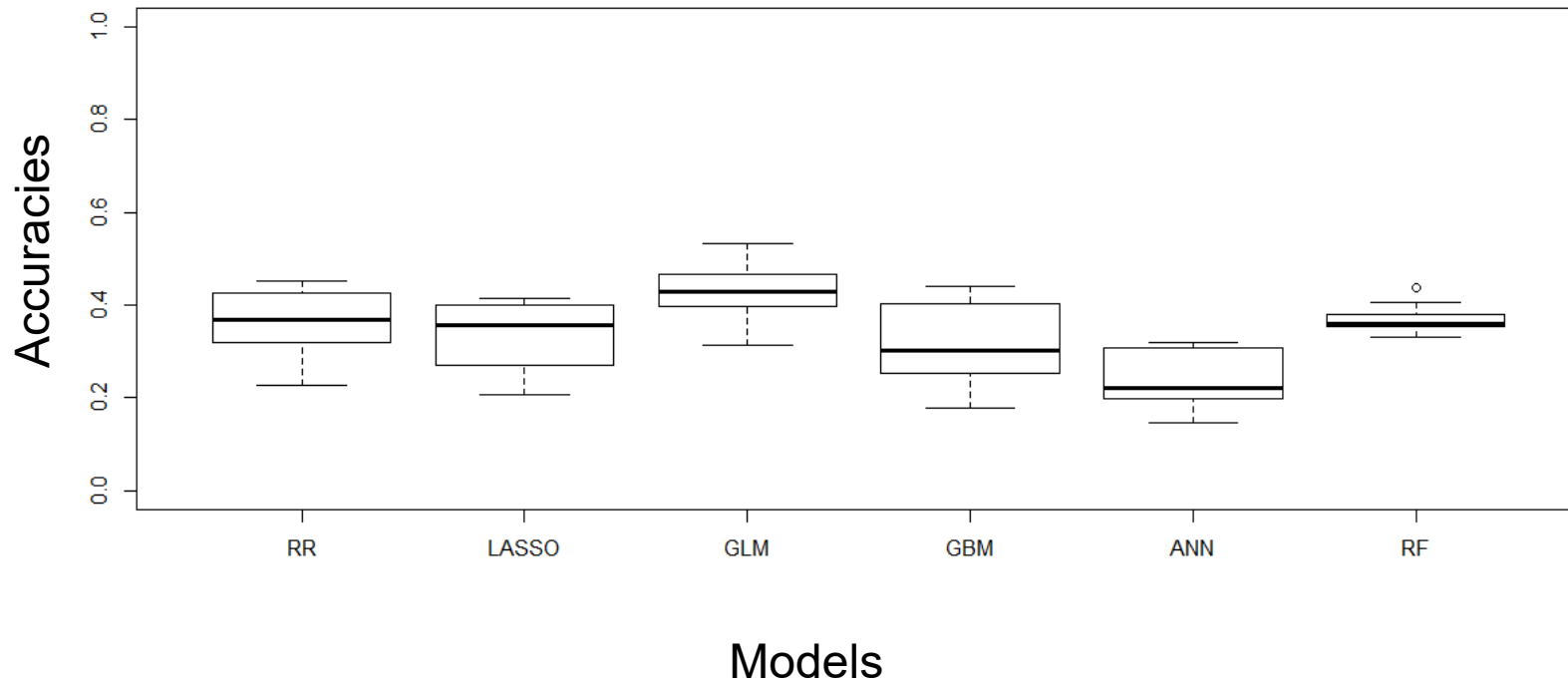
Predictions using the entire canola dataset are time consuming

Phenotypic trait : **hybrid yield**

Models: RR-BLUP, LASSO, GLM, GBM, ANN, RF

Cross-validation: 10x – 80%/20%

Entire set prediction – **approx. 1200 min/20h**



Predictions using the PCA-based FS are 200x faster !

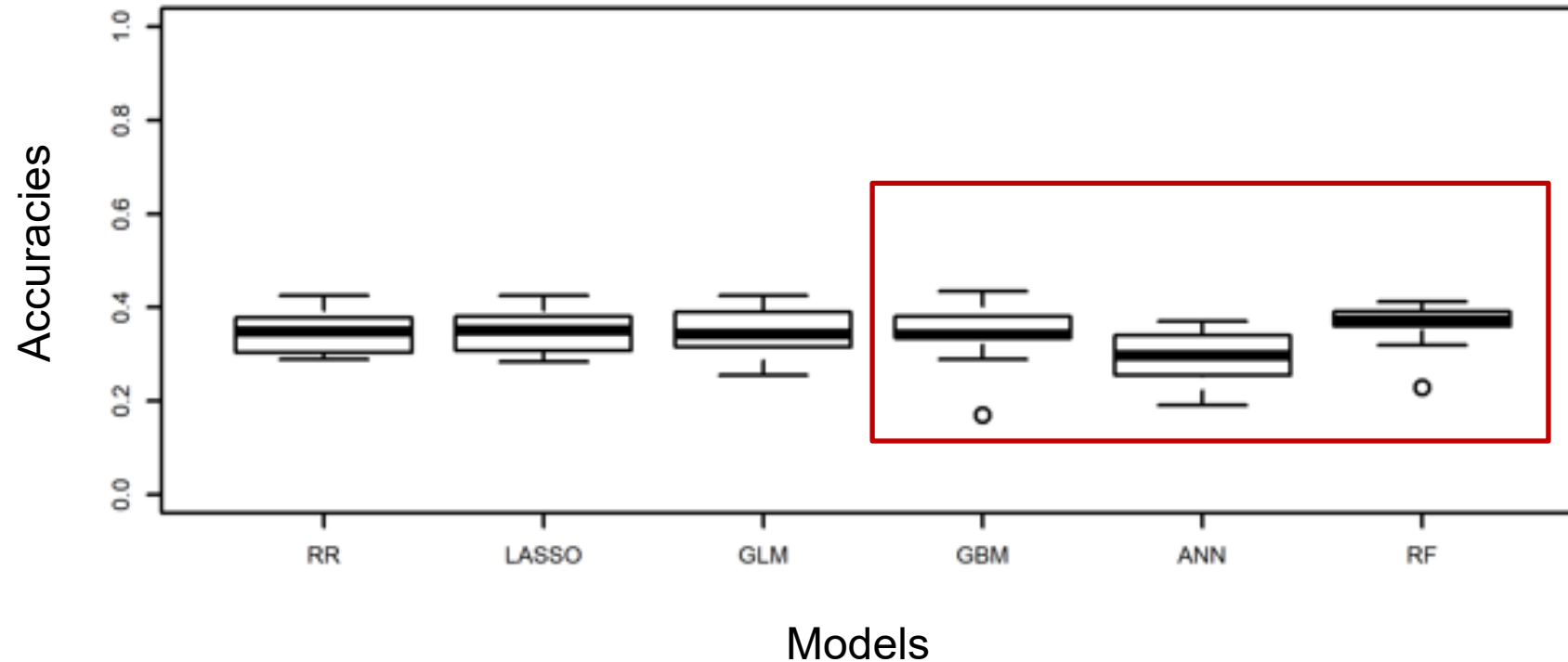
Phenotypic trait : **hybrid yield**

Models: RR-BLUP, LASSO, GLM, GBM, ANN, RF

Cross-validation: 10x – 80%/20%

Matrix : **first 100PCs**

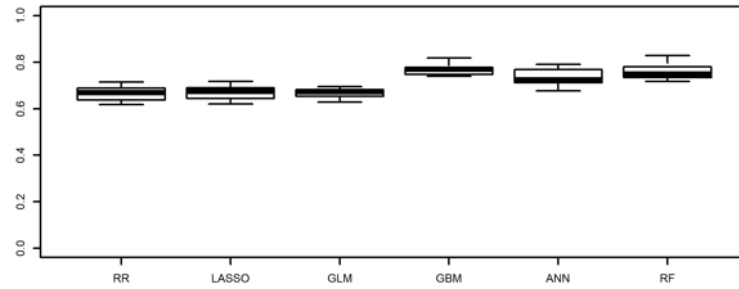
Entire set prediction – **approx. – 6.7 min** (200 x faster than the SNP matrix X)



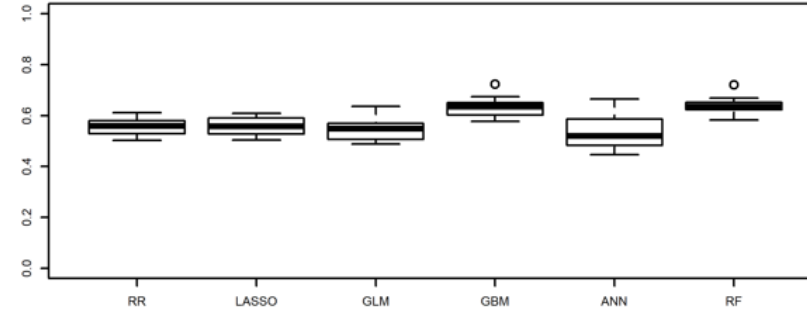
Similar results for all investigated traits using PCA-based FS

Accuracies

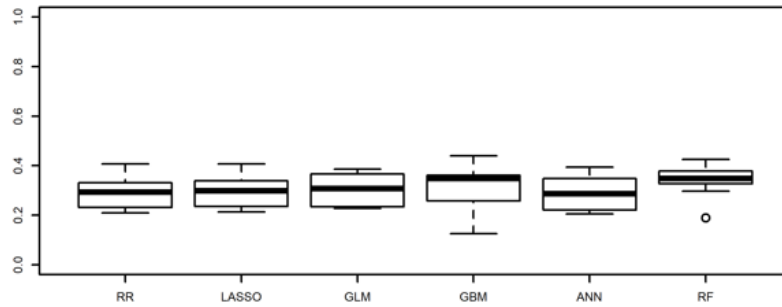
Days to flower



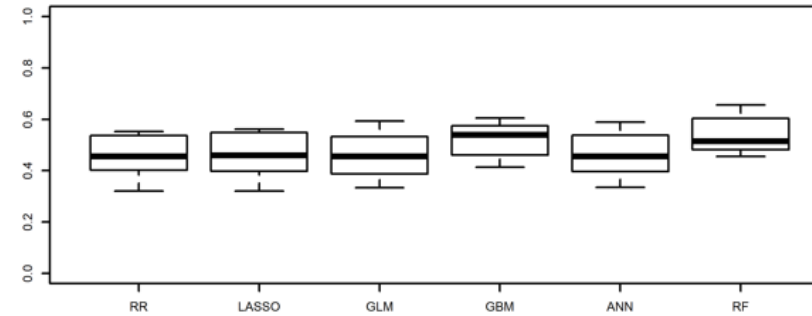
Oil content



Emergence



Protein content



Models

Predictions using the RF-based FS are 10x faster !

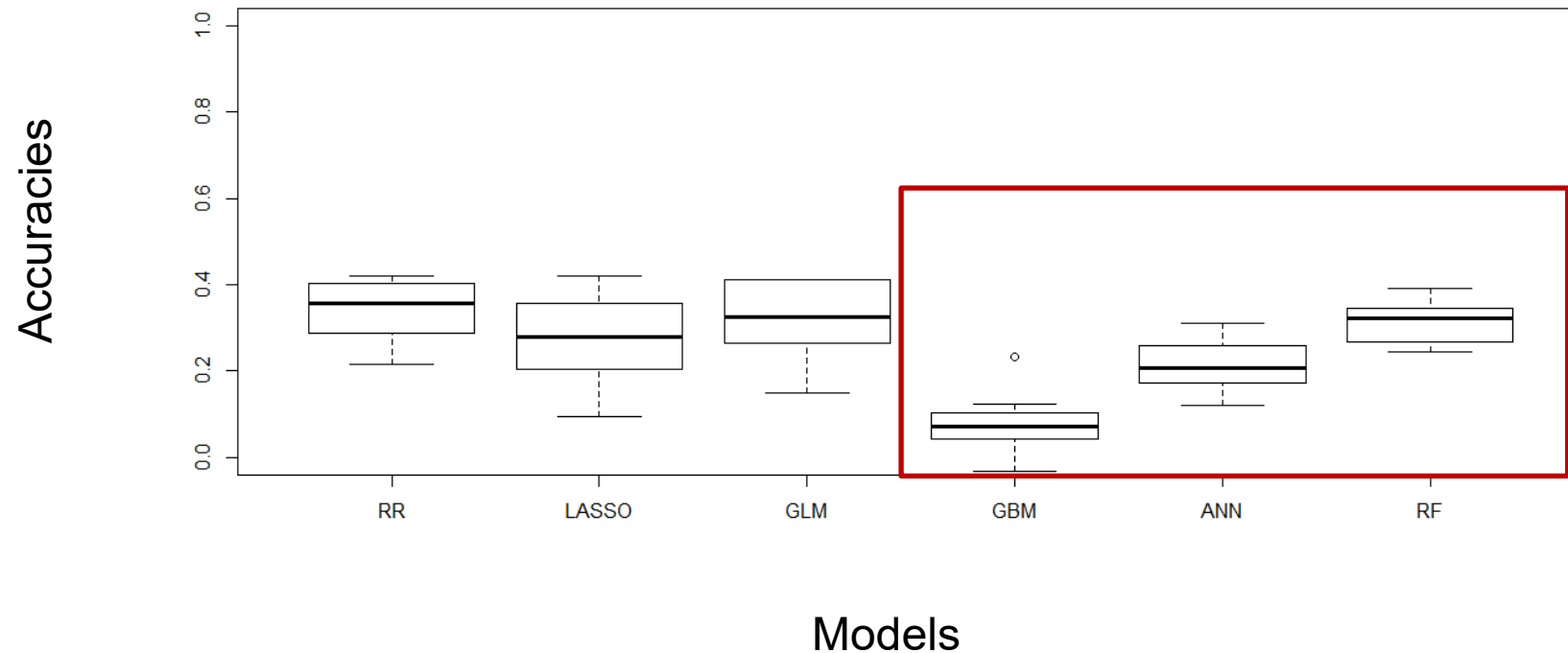
Phenotypic trait : **hybrid yield**

Models: RR-BLUP, LASSO, GLM, GBM, ANN, RF

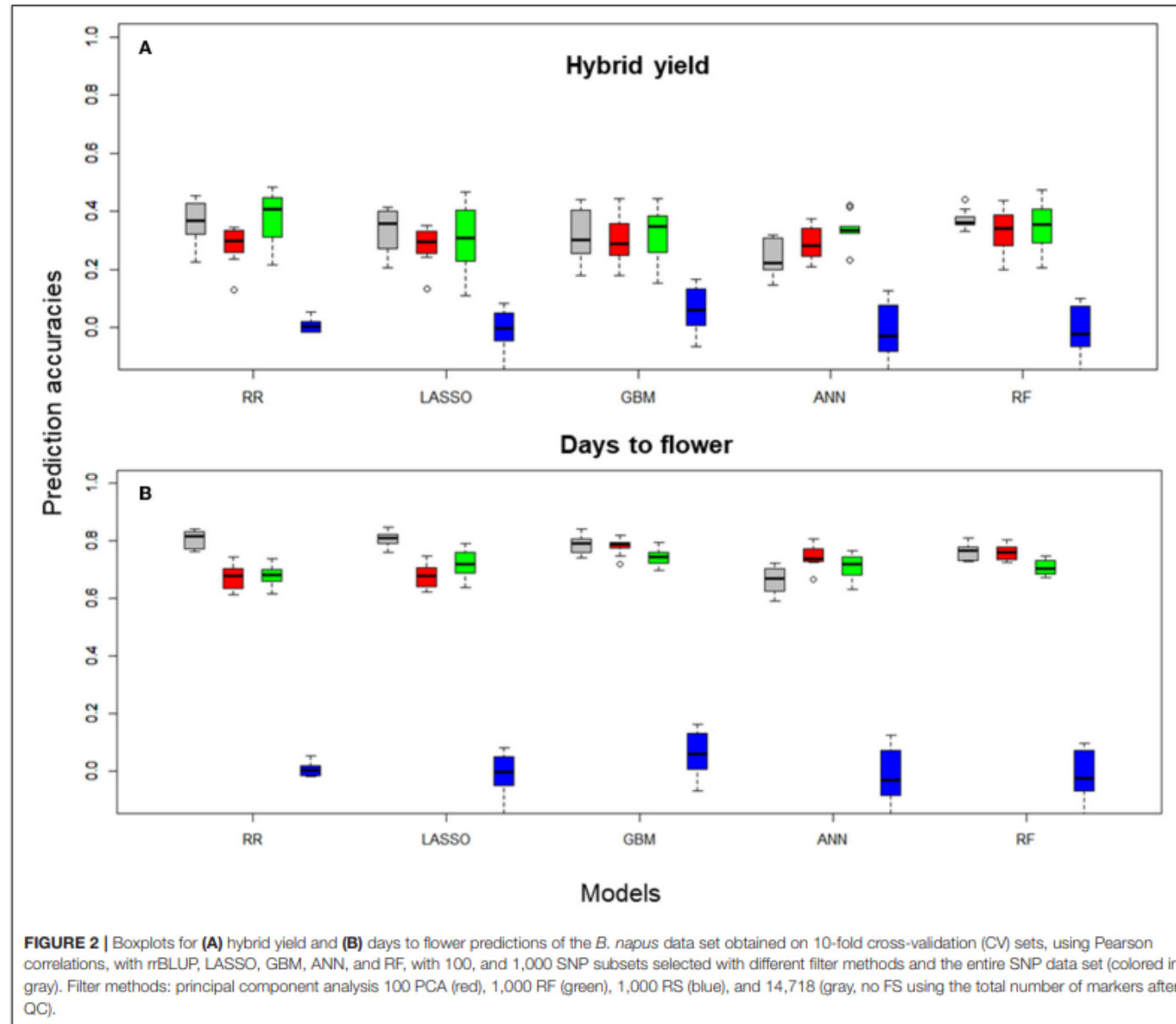
Cross-validation: 10x – 80%/20%

Matrix : **first 1000 RF features**

Entire set prediction – **approx. – 119.6 min** (10 x faster then the SNP matrix X)



Feature selection improves prediction accuracies of non linear models



Acknowledgements



Justus Liebig University, Germany:

Prof. Rod Snowdon
Prof. Mathias Frisch
Dr. Christian Obermeier
Lennard Ehrig
Sven Weber

Faculty of Computer Science, UAIC-Iasi, Romania:

Prof. dr. Dan Cristea
Asso. Prof. dr. Mihaela Elena Breaban
Lector dr. Ionut Pistol
Asso. Prof. dr. Madalina Raship

Department of Computer Science, Aalto University, Finland:

Institut of Bioinformatics, Kyoto University, Japan:

FiDiPro Prof. Hiroshi Mamitsuka

NPZ Innovation GmbH (NPZi), Germany:

Dr. Amine Abadi

All colleagues and technical assistants

