# Key messages

1. ML in hands-on breeding: HTPs, BLUPs and GEBVs

2. Usage of the data relies on the nature of the signal

3. Breeding applications mostly on Gaussian process

4. Machines are important when signal is scarce

# Outline

Overview of ML

Topic 1 - Data processing

Topic 2 - Signal detection
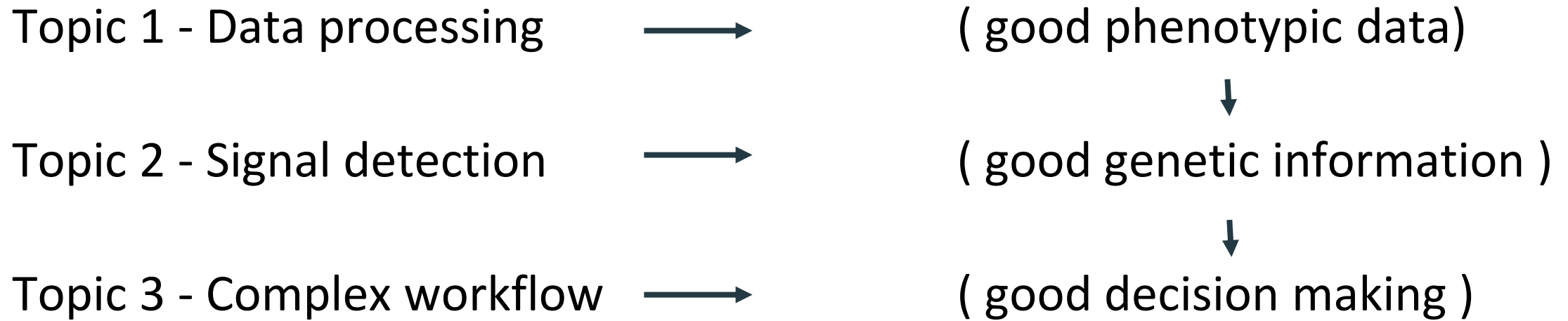
Topic 3 - Complex workflow

Concluding remarks

*"Begin at the beginning,"* the King said, *"and go on till you come to the end: then stop."*

**Lewis Carroll,** Alice in Wonderland

# Outline

Overview of ML

Topic 1 - Data processing $\longrightarrow$ ( good phenotypic data)

Topic 2 - Signal detection $\longrightarrow$ ( good genetic information )

Topic 3 - Complex workflow $\longrightarrow$ ( good decision making )

Concluding remarks

*"Begin at the beginning,"* the King said, *"and go on till you come to the end: then stop."*
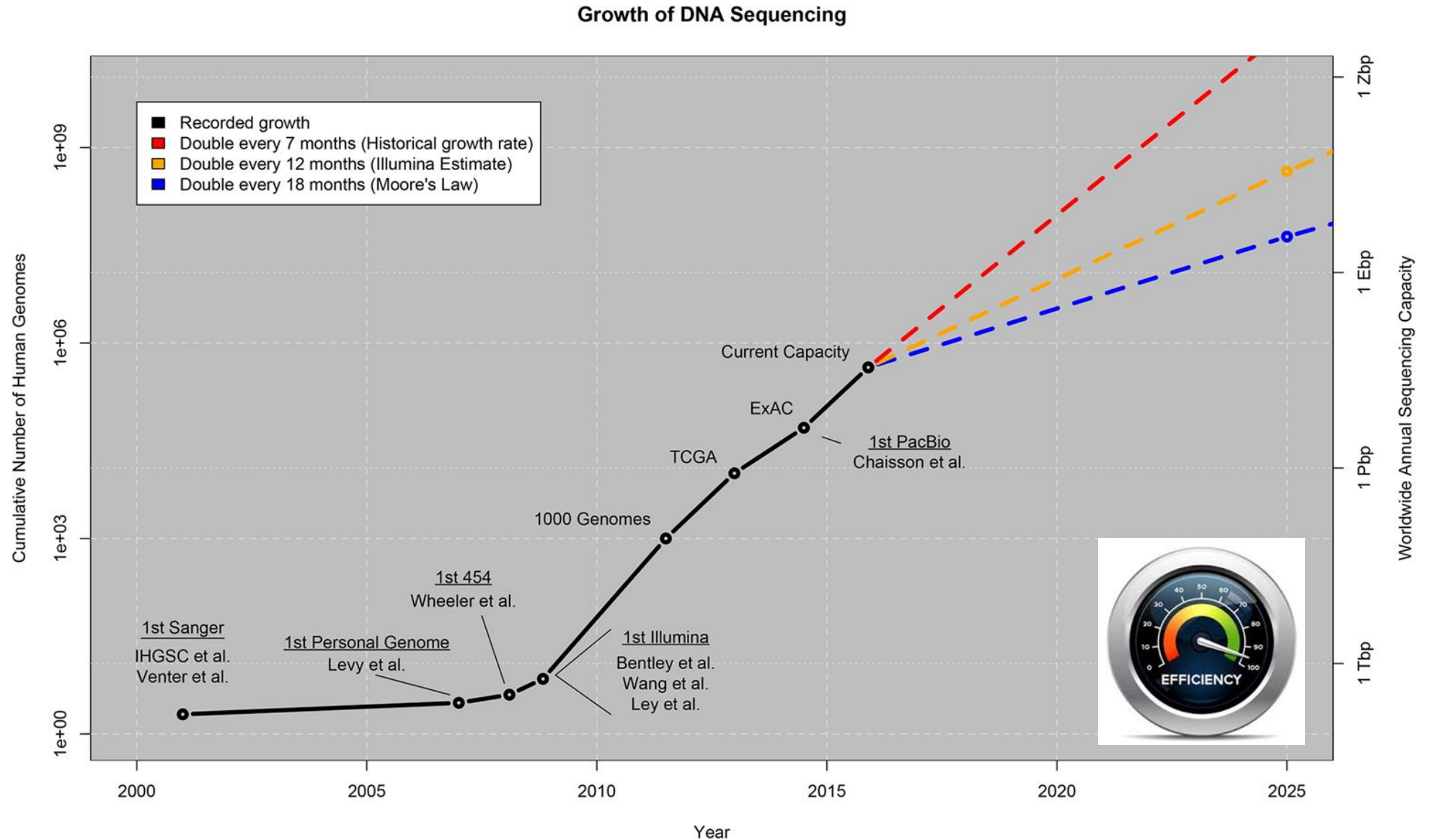
**Lewis Carroll,** Alice in Wonderland

Growth of DNA Sequencing

Stephens, Z. D. et al. (2015). Big data: astronomical or genomical?. *PLoS biology, 13*(7), e1002195.

- **Machine learning is** a major component of artificial intelligence (AI). ML is concerning with capturing specific patterns from data, often with the purpose of de-noising, classification and predictions for decision making.

**No machine can be optimally efficient in more than one task**

The example below postulates that one machine that climb stairs and make pancakes will be less efficient than two machines exclusively focus on climbing start and making pancakes, respectively.
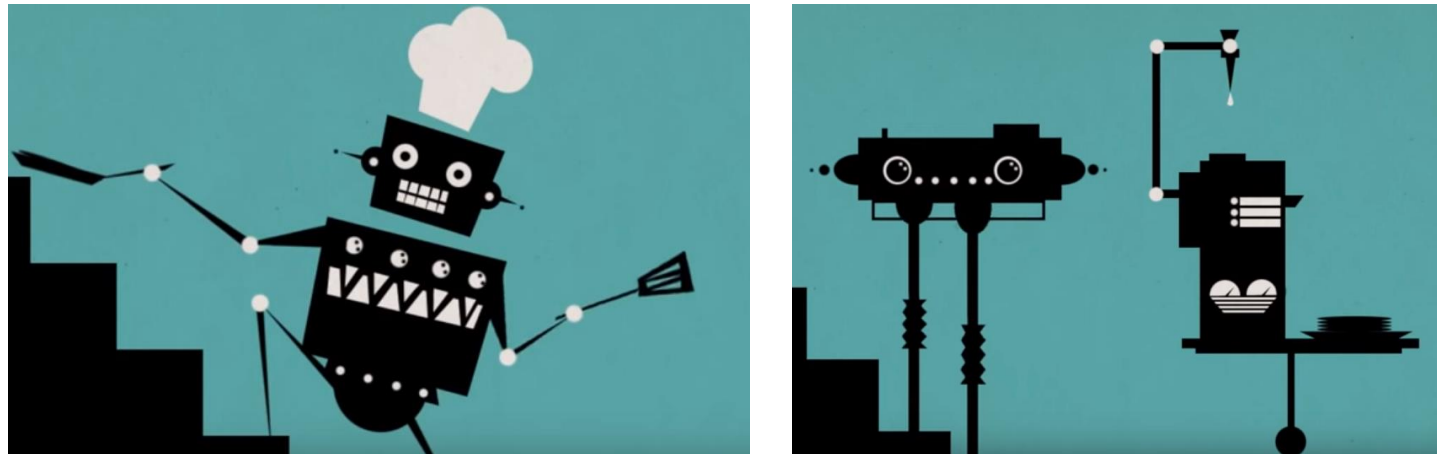


Figure source: https://www.youtube.com/watch?v=MPR3o6Hnf2g

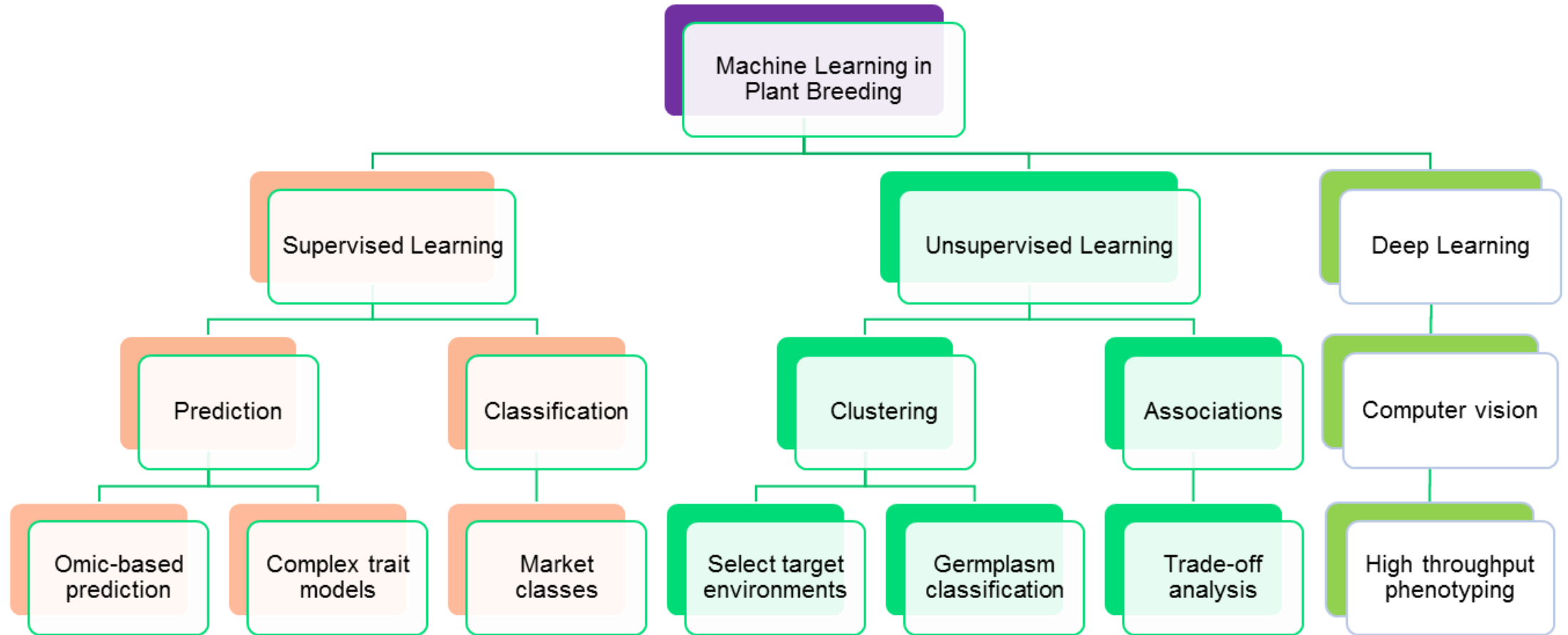# Where ML is being deployed on plant breeding

**On phenotyping and phenotypic analysis**

- Automated and precise scoring using drones and computer vision
- Discovery of (index and longitudinal) traits correlated to yield
- Phenotypic denoising: competition, environmental and spatial noise
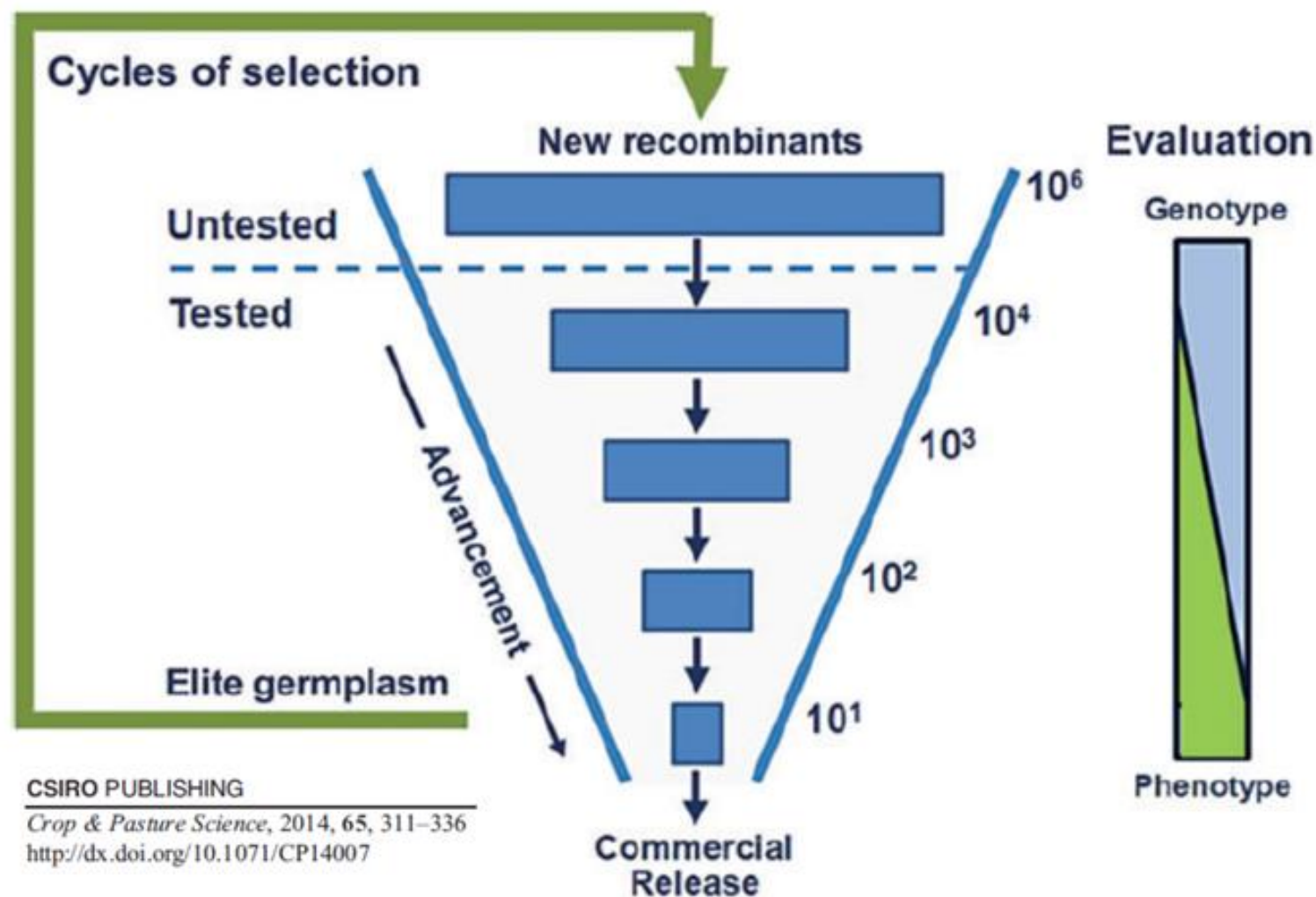
**On genotyping and genotypic analysis**

- Unsupervised analysis: Imputation and germplasm classifications
- Early generation selection and recycling: Additive genetic signal
- Advanced generations selection and stability: G and GxE signal

# Where is ML in plant breeding today?

# 1. Data processing
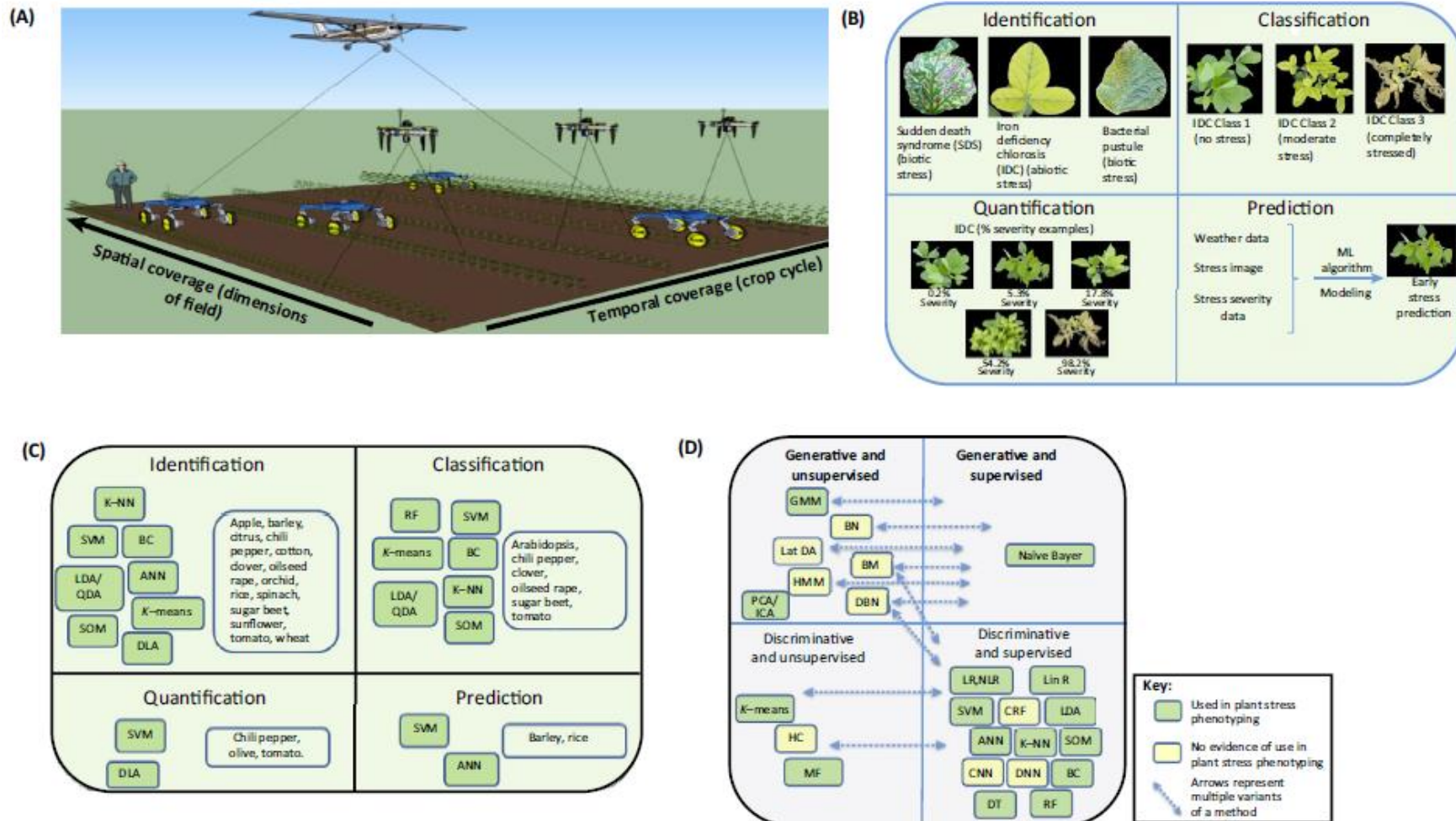
# Where does ML fit in the breeding pipeline?



Cycles of selection

New recombinants — $10^6$

Untested

Tested — $10^4$

$10^3$

$10^2$

Advancement

Elite germplasm

$10^1$

Evaluation

Genotype

Phenotype

Commercial Release

10

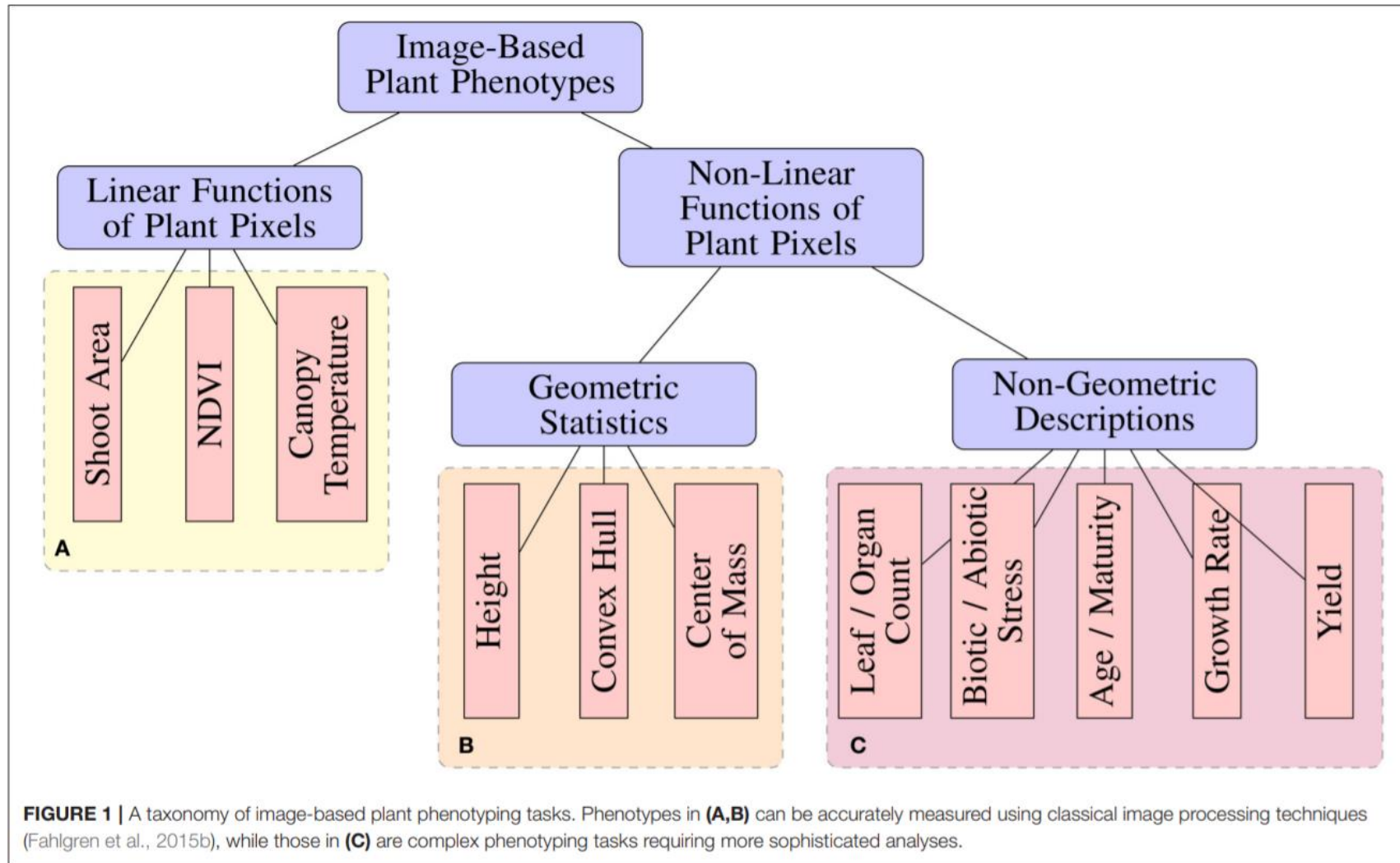# Pipeline complexity (highly computational)



**Key Figure** — Trends in Plant Science, February 2016, Vol. 21, No. 2   111
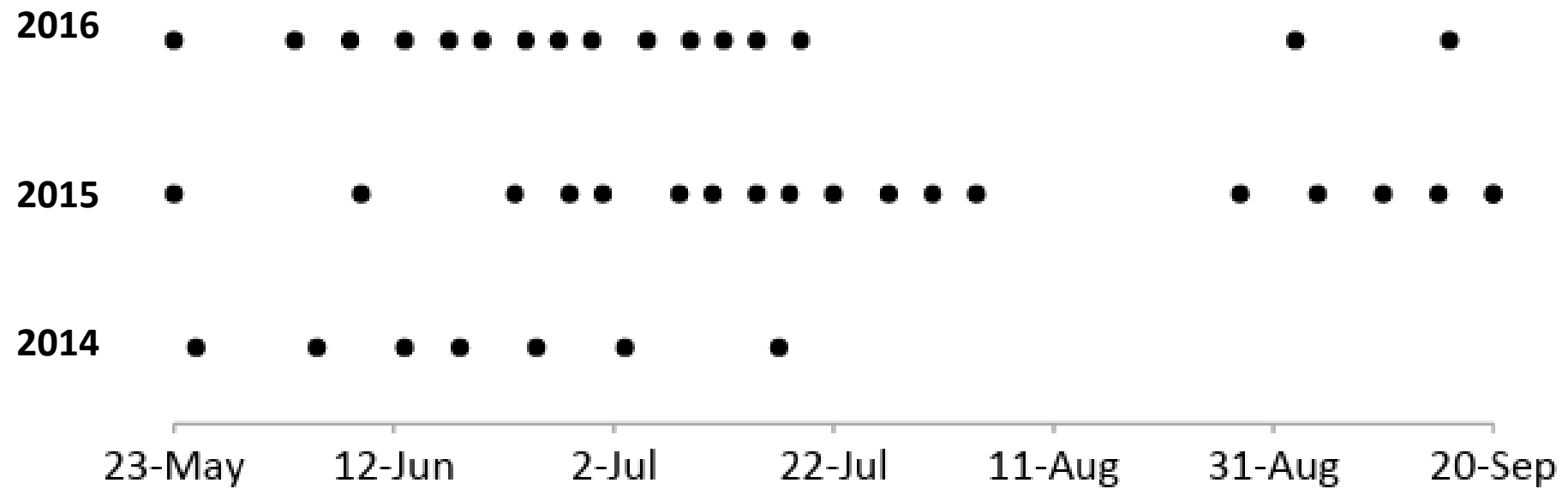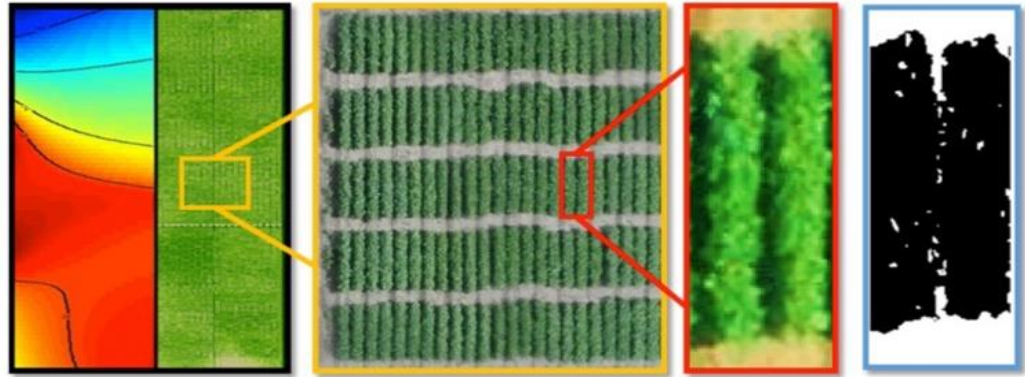
Machine Learning (ML) Tools for High-Throughput Stress Phenotyping

# Pipeline complexity (highly computational)



FIGURE 1 | A taxonomy of image-based plant phenotyping tasks. Phenotypes in (A,B) can be accurately measured using classical image processing techniques (Fahlgren et al., 2015b), while those in (C) are complex phenotyping tasks requiring more sophisticated analyses.
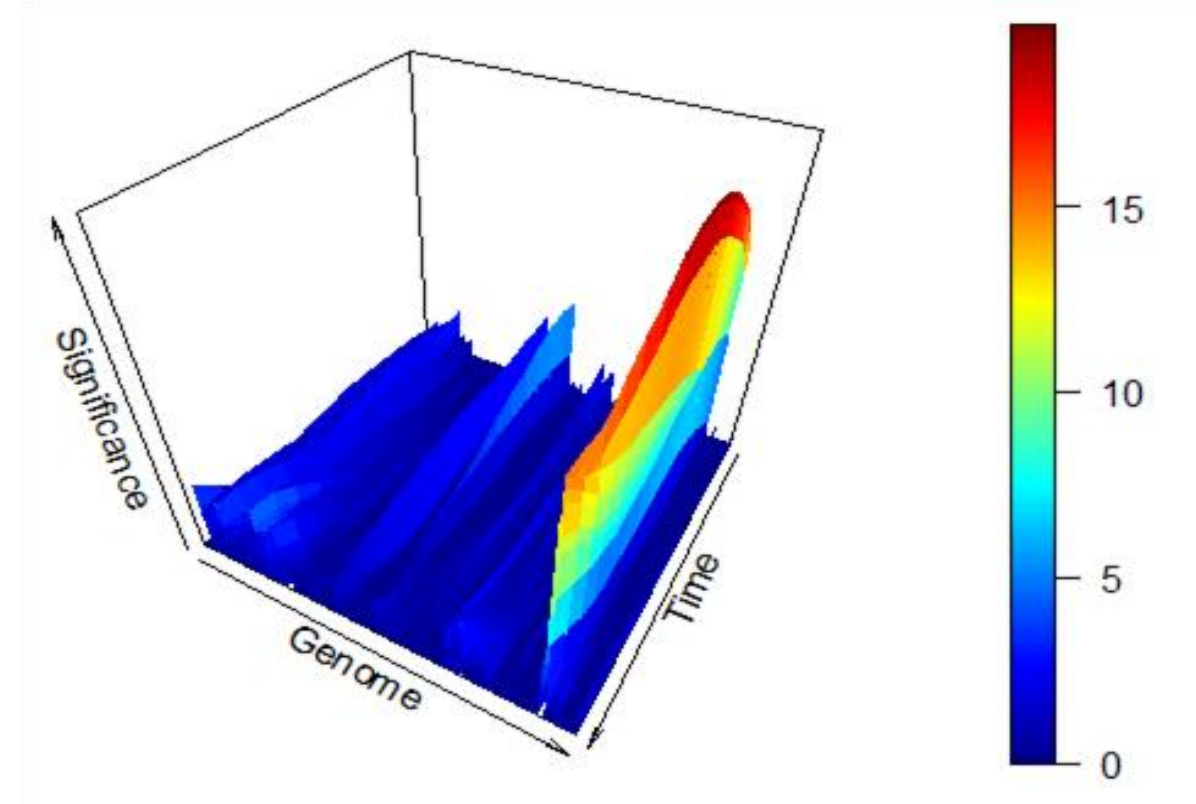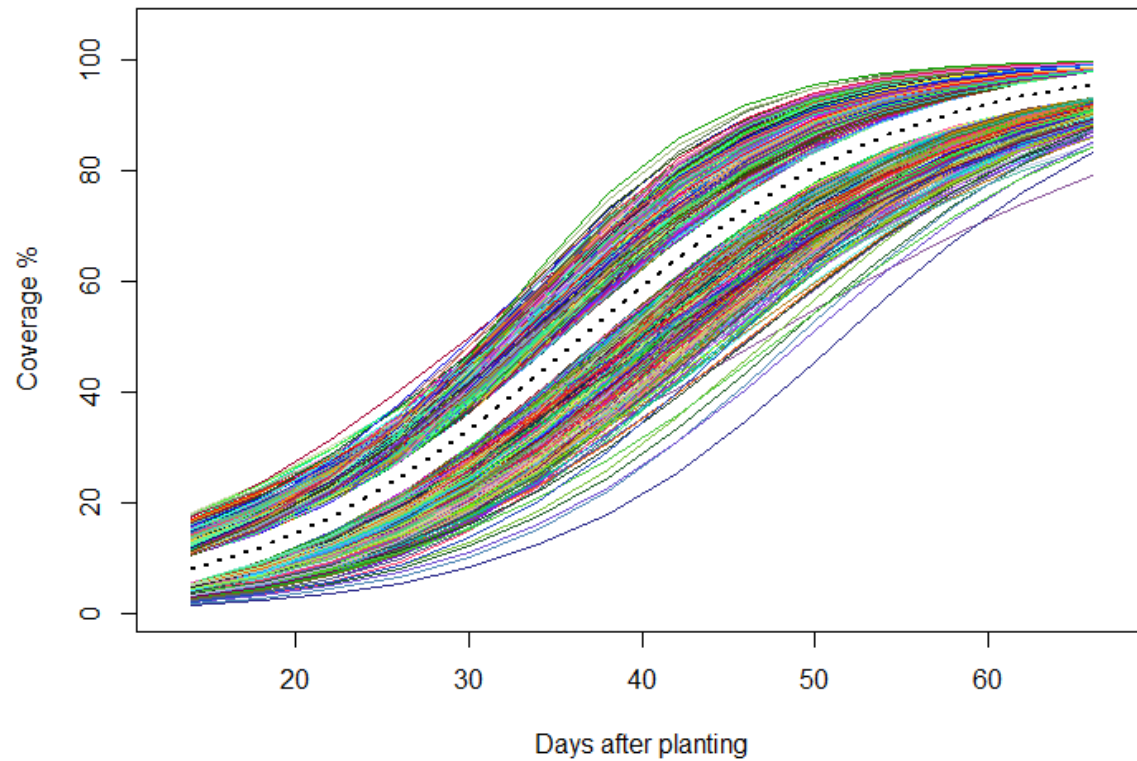
Ubbens, J. R., & Stavness, I. (2017). Deep plant phenomics: a deep learning platform for complex plant phenotyping tasks. Frontiers in plant science, 8, 1190.
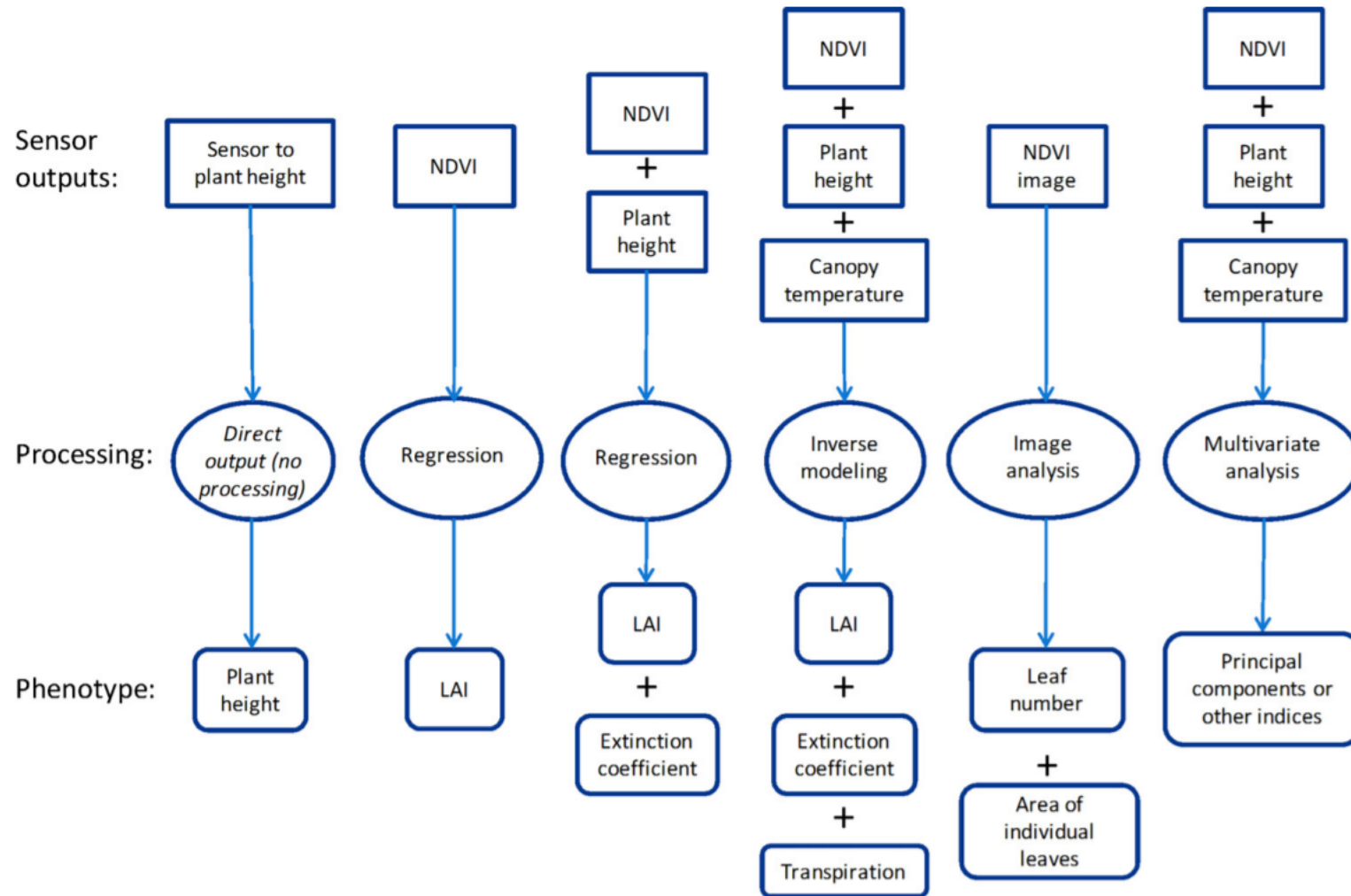
12

**2016**

**2015**

**2014**

23-May     12-Jun     2-Jul     22-Jul     11-Aug     31-Aug     20-Sep

# Different time-points, different genetic architecture



General behaviour of soybean canopy development

# Number and quality of phenotypes increase



**Fig. 6.** Examples of possible paths of data analysis whereby field measurements are processed to provide more biologically meaningful data. Field data usually would be recorded as time series, allowing estimation of growth or developmental rates.

# 2. Signal detection

**2.1 Robust machines**
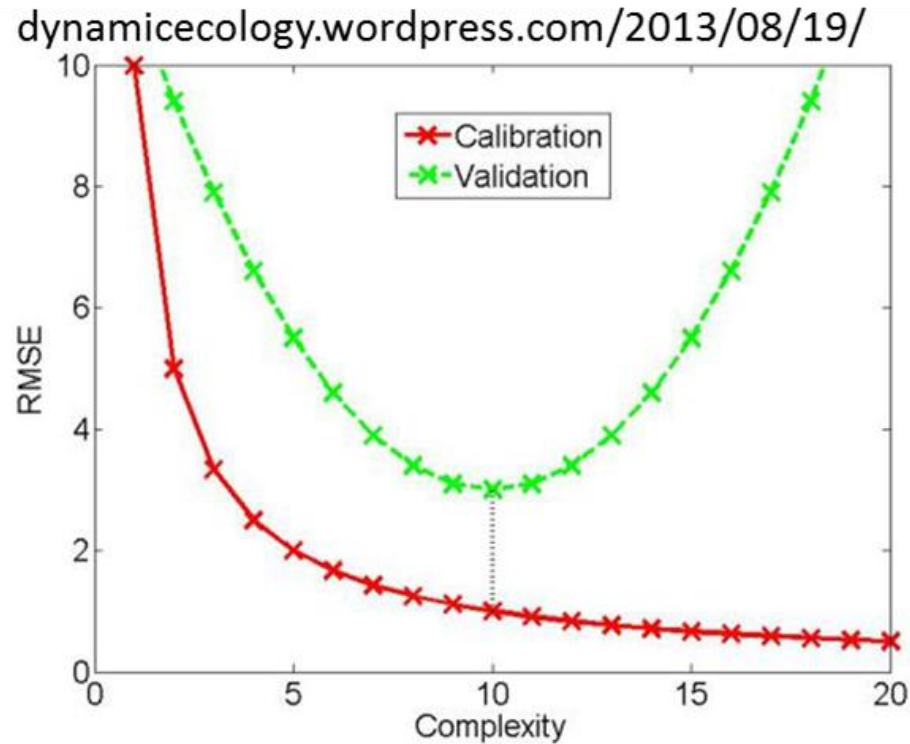
**2.2 Capturing signals**

# Robust machines

1. Regularization

2. Parsimony

3. Interactions

# 1) Overfitting: Complexity-Variance tradeoff



dynamicecology.wordpress.com/2013/08/19/
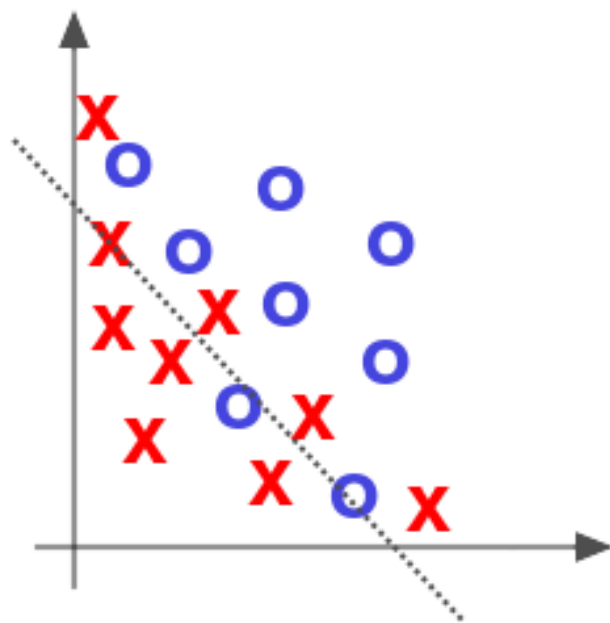
$\lambda \uparrow = Fitness\ (R^2)\ \downarrow = Number\ of\ SNPs\ \downarrow = Shrinkage \uparrow$

$\uparrow Complexity = \downarrow Variance = \downarrow Prediction$

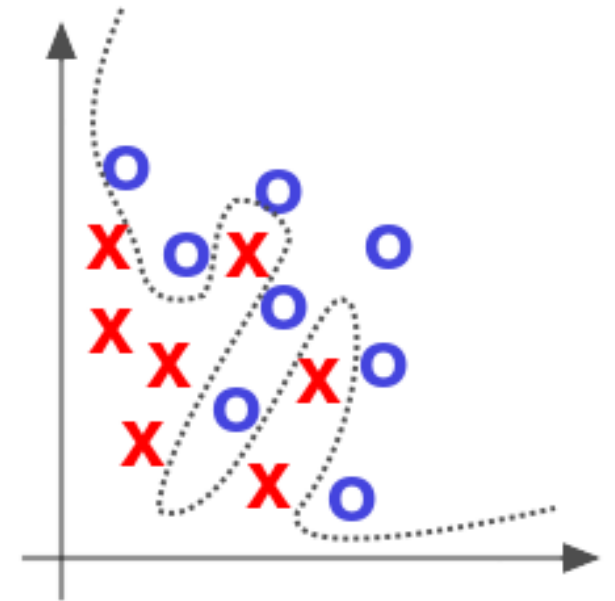# variables        Overfit        BAD!!

# Overfitting (kNN example)



Under Fit

Appropriate

Over Fit

https://www.safaribooksonline.com/library/view/deep-learning/9781491924570/ch01.html

# 2) Parsimony

- **Occam's razor** – A simpler explanation is better than a complex one
  - The less terms you have in your model, the better
  - An attempt to comprise most information with the least amount of factors

- **Example (model for genetic values). Check the two models:**
  1) <span style="color:red">**Yield = Block + Location + Year + Genotype + (Genotype x Year x Location) + …**</span>
  2) <span style="color:red">**Yield = Block + Genotype**</span>

- **Genomics: Models with too many parameters (p>>n)**
  - Most high-dimensional machines automatically perform variable selection

Occam's Razor: No more things should be presumed to exist than are absolutely necessary, i.e., the fewer assumptions an explanation of a phenomenon depends on, the better the explanation.

(William of Occam)

izquotes.com

# 3) Hierarchical principle

- **CLAIM**: Lower order effects more important than higher order effects
- Effects of same order equally important

**POWER IS AN ISSUE TO DETECT SIGNAL OF INTERACTIONS!!!!**

- Consider the type of the variable (continuous or categorical)
- This principle makes one wonder about the relevance of **Epistasis** and **GxE**
- Higher order terms are good to **run out of degrees of freedom**

# Capturing signals

1. Signal extraction

2. Supervised machines

3. Multiple signals

# Gaussian process: Signal and noise
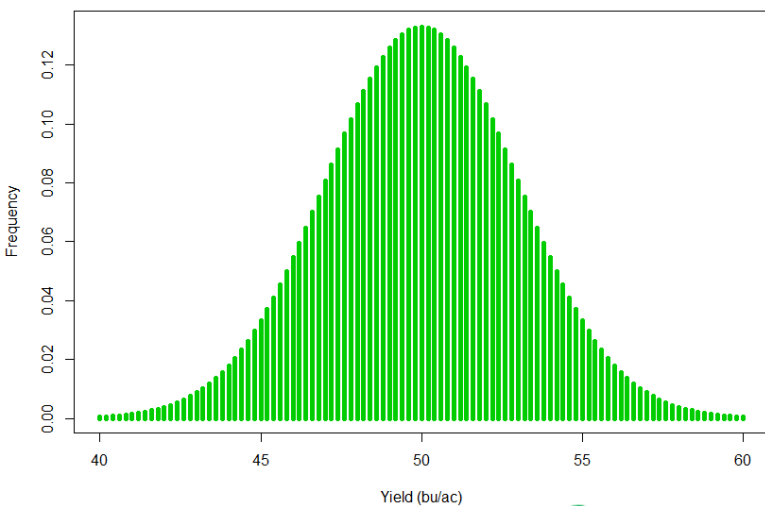
$$y = 1\mu + g + e$$

$$\sigma_y^2 = \sigma_g^2 + \sigma_e^2$$

$\mu = 50$
$\sigma_y^2 = 9$
$\sigma_g^2 = 4$
$\sigma_e^2 = 5$
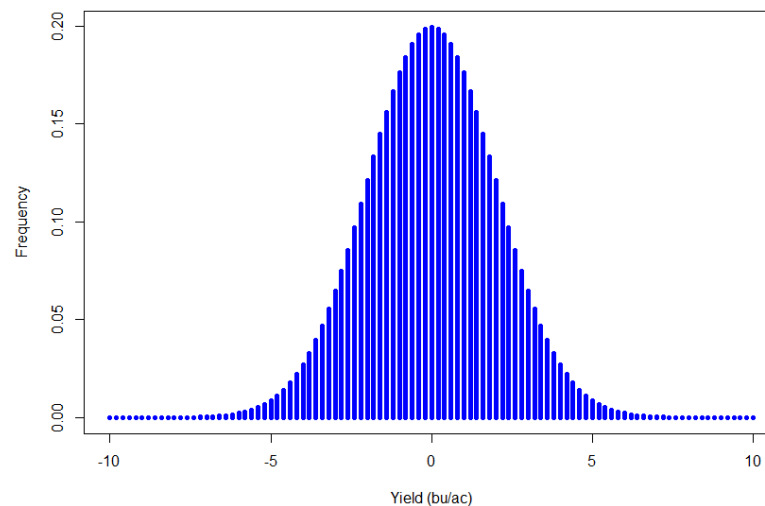$cov(g, e) = 0$



$$y \sim N(\mu, I\sigma_y^2)$$

$$u \sim N(0, I\sigma_g^2)$$

$$e \sim N(0, I\sigma_e^2)$$

# Distinction of signal and noise

$$y = \mu + g + e$$



Observation $y$ → Machine (eg. mixed model) $K$

$e$ Noise

$g$ Signal

# Most common supervised machines

## Regression trees
**(eg. simplified decision model)**



## Penalized regression
**(eg. random effects in mixed models)**



## Neural networks
**(eg. Google search engine)**



## Kernel regression
**(eg. GBLUP and GPS interpolation)**



26

# Multiple signals

# Multiple signals

**PHENOTYPE**

**ENVIRONMENT**

**GENETICS**

**INTERACTIONS**
Noise in early generation
Signal in advanced generation

**Controllable**
Important but Non-target
(eg experimental design)

**Additive**
Target signal:
"Breedable" variance

**Stochastic**
True noise

**Non-Additive**
Treated as noise

# Case of non-target signal (Lado et al. 2013)





**Figure 1** Diagram to calculate the covariable $x_i$. $Y_i$ is the phenotypic value in the plot. The neighboring plots are indicated with gray color.

■ **Table 3 Accuracy of predictions for each trial in 2011 using random training sets with 100 independent randomizations**

|  |  |  | IB | RC | RCB_MVNG | MVNG |
|---|---|---|---|---|---|---|
| SR_FI | GY | RR | 0.298 ± 0.117 | 0.296 ± 0.119 | 0.319 ± 0.114 | 0.319 ± 0.113 |
|  |  | GAUSS | 0.312 ± 0.117 | 0.310 ± 0.120 | 0.325 ± 0.117 | 0.326 ± 0.116 |
| SR_MWS | GY | RR | 0.236 ± 0.141 | 0.275 ± 0.147 | 0.231 ± 0.127 | 0.347 ± 0.134 |
|  |  | GAUSS | 0.231 ± 0.144 | 0.273 ± 0.150 | 0.260 ± 0.128 | 0.370 ± 0.132 |

IB, incomplete blocks, field design; RC, row by column model; RCB_MVNG, random complete block model with moving means as covariable; MVNG, linear regression model with moving means as covariable; SR_FI, Santa Rosa under full irrigation; GY, grain yield; RR, Ridge regression kernel; GAUSS, Gaussian kernel; TKW, thousand kernel weight; DH, days to heading; NKS, number of kernels per spike; SR_MWS, Santa Rosa under mild water stress.

# 3. Complex workflow

# Phenomics & genomics are usually tight together

**The phenomics wheel of fortune**

Drought-tolerant plant

6 Crop physiology

2 Physiology

7 Parameterisation

3 Validation

Models to compile traits in plant/crop

Identification of key traits

Germplasm enhancement

Dissect basis for tolerance

4 Reverse genetics

5 Breeding genetic modification

1 Forward genetics

Identify genetic basis for trait

*TRENDS in Plant Science*

**Figure 2.** Closing the gene to genotype loop with phenomics.

# Analysis of multiple traits

- In linear models: Covariate vs Multivariate
  - Covariate – on trait is used to predict the other. No strings attached to genetics.
  - Multivariate – modeling 2+ traits simultaneously. Genetics connected across traits.

- Computational burden increases exponentially $O(k)^7$ with the number of traits

- Multicollinearity
  - Many traits are nearly identical (eg. neighbor bandwidths from hyperspectral image)
  - OLS will not work due to singularities

- Modeling is often improved by accounting for time-space domains

# Multivariate models

## Core Ideas

- HTP platforms used to measure secondary traits across time
- Longitudinal data of secondary traits evaluated by SR, MT, and RR models, separately
- BLUPs of secondary traits used in the multivariate pedigree and genomic prediction
- Grain yield predictive ability was improved by 70%

Sun, J., Rutkoski, J. E., Poland, J. A., Crossa, J., Jannink, J. L., & Sorrells, M. E. (2017). Multitrait, Random Regression, or Simple Repeatability Model in High-Throughput Phenotyping Data Improve Genomic Prediction for Wheat Grain Yield. *The Plant Genome*.
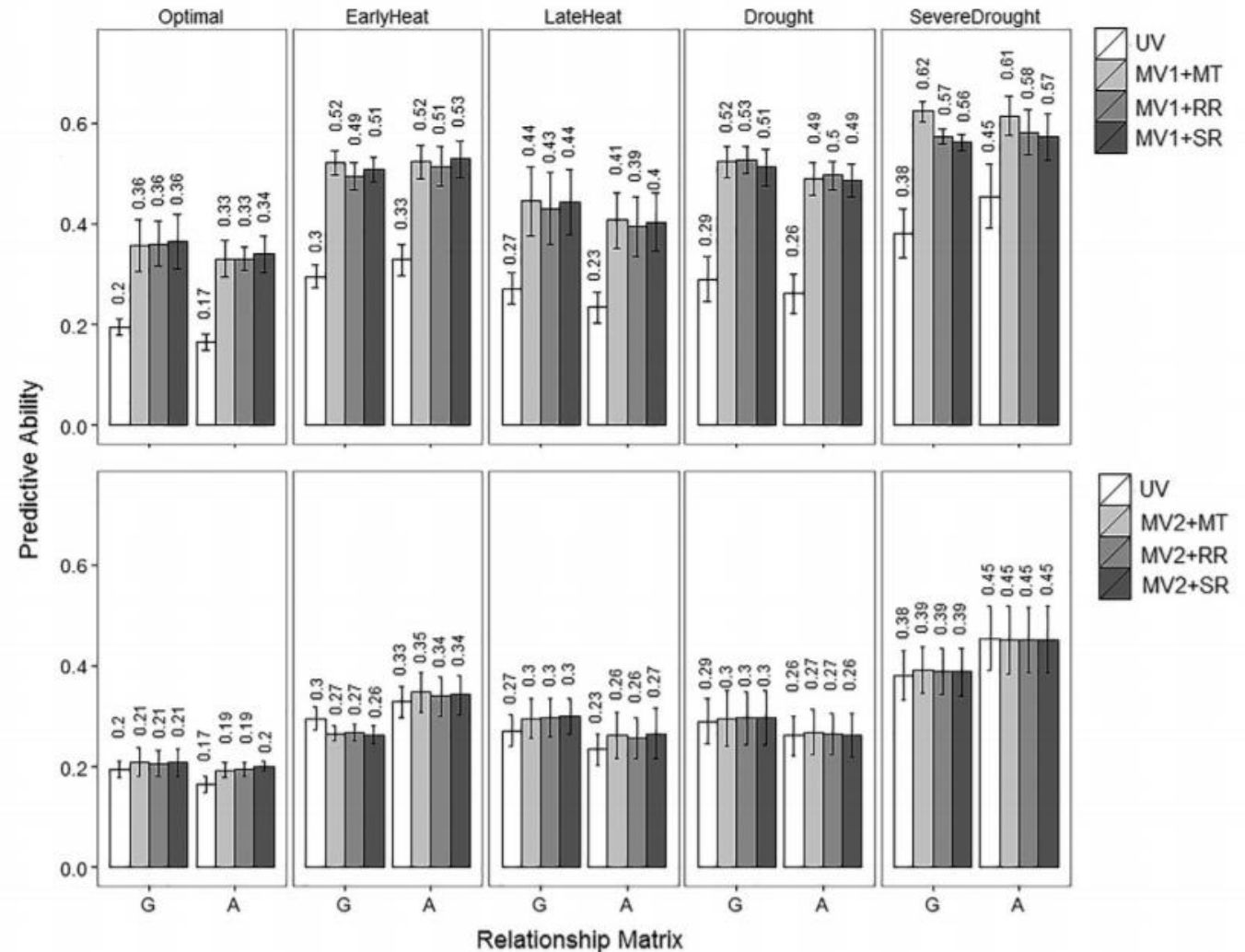


Fig. 1. Predictive ability comparison for grain yield between prediction models with secondary traits (MV1 and MV2) and without secondary trait (UV). MV1, multivariate prediction model with secondary traits in both training and testing populations; MV2, multivariate prediction model with secondary traits in training population only; UV, univariate prediction model with grain yield only; MT/RR/SR, multivariate prediction model MV1 or MV2 using best linear unbiased predictions (BLUPs) of secondary traits from multitrait (MT), random regression (RR), or simple repeatability (SR) model; G/A, genomic/pedigree relationship matrix.

# Alternatives to dense multivariate models:
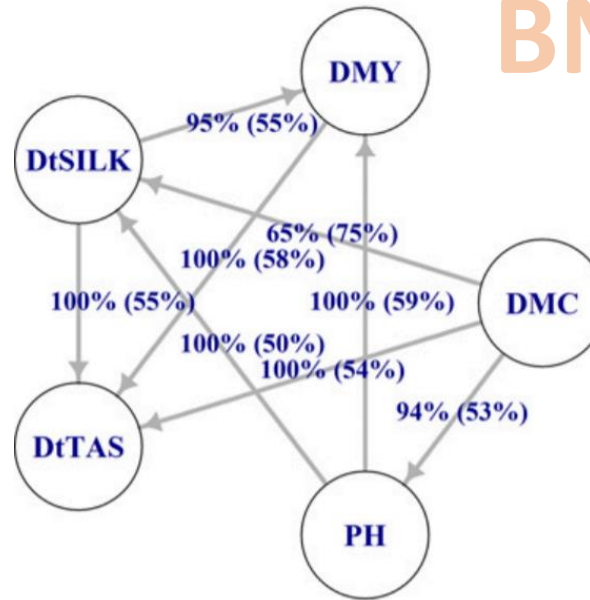## Structural Equation Models (SEM), Bayesian Networks (BN), and Markov Random Fields (MRF)



SEM

BN

MRF

Valente, B. D., Rosa, G. J., de los Campos, G., Gianola, D., & Silva, M. A. (2010). Searching for recursive causal structures in multivariate quantitative genetics mixed models. *Genetics*, *185*(2), 633-644.

Töpner, K., Rosa, G. J., Gianola, D., & Schön, C. C. (2017). Bayesian Networks Illustrate Genomic and Residual Trait Connections in Maize (Zea mays L.). *G3: Genes, Genomes, Genetics*, *7*(8), 2779-2789.

Xavier, A., Hall, B., Casteel, S., Muir, W., & Rainey, K. M. (2017). Using unsupervised learning techniques to assess interactions among complex traits in soybeans. *Euphytica*, *213*(8), 200.

# Better the signal, less important is the machine



Figure 1. Key parameters and changes during a breeding cycle, to consider in implementing genomic selection (GS). The triangles indicate increase or decrease of the quantity considered. QTL, quantitative trait loci.

Heslot, N., Jannink, J. L., & Sorrells, M. E. (2015). Perspectives for genomic selection applications and research in plants. *Crop Science*, *55*(1), 1-12.

Granier, C., & Vile, D. (2014). Phenotyping and beyond: modelling the relationships between traits. *Current opinion in plant biology*, *18*, 96-102.

# Concluding remarks

# Revisiting key messages

★ **ML in hands-on breeding: HTPs, BLUPs and GEBVs**
  ○ *Obtain better phenotypes and perform more accurate selections*

★ **Usage of the data relies on the nature of the signal**
  ○ *Best results come from simple and mindful model*

★ **Breeding applications mostly on Gaussian process**
  ○ *Most breeding problems can be tackled with mixed models*

★ **Machines are important when signal is scarce**
  ○ *Genomic breeding suits best modeling early generations*

# Acknowledgements

**Funding**
- USB (Genotyping, experiments, phenotyping), Corteva (experiment and phenotyping), USDA and NSF (students)

**Genetic resource and genotyping**
- USDA (Perry Cregan, Qijian Song), UNL (James Specht), UofIL (Brian Diers, Randy Nelson)

**Drone data**
- Purdue ABE (Keith Cherkauer, Anthony Hearst)

**Modeling**
- Purdue (William Muir), UC Riverside (Shizhong Xu), ISU (William Beavis, Vishnu Ramasubramanian), UNL (Diego Jarquin, Reka Howard), UMN (Aaron Lorenz)

**Phenotypes**
- Purdue (Katy Rainey, Ben Hall), UofIL (Randy Nelson, Brian Diers), UNL (James Specht, George Graeg), ISU (William Beavis), OSU (Leah McHale), KSU (William Schapaugh), MSU (Dechun Wang), UofMO (Grover Shannon), NCSU (Rouf Mian)
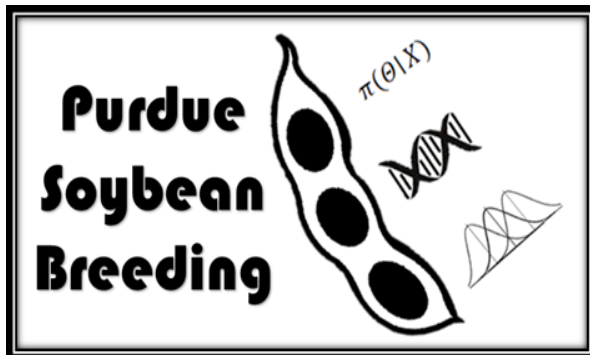
**Support**
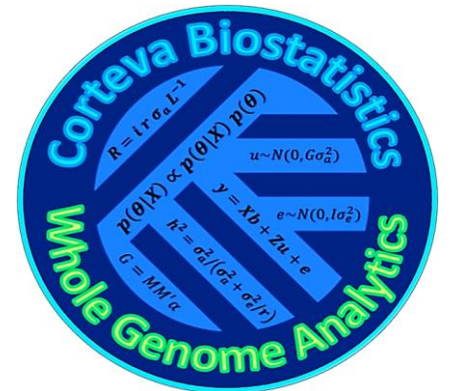- Corteva (Tabare Abadie, Radu Totir, Mak Geha, David Habier)

# That's all!

Thanks!

Questions?

**AX102418**

http://alenxav.wix.com/home

# Implementations in R

- Genetic signal based on linear (mixed) models: lme4, pedigreemm, SpATS
- Genomic signal based on supervised machine learning: glmnet, kernlab, ranger, pls, keras
- Genomic signal based on linear (mixed) models: BGLR, bWGR, rrBLUP, EMMREML

# Follow up readings

- Xu (2013). Mapping quantitative trait loci by controlling polygenic background effects. Genetics, genetics-113
- Morota and Gianola (2014). Kernel-based whole-genome prediction of complex traits: a review. Frontiers in genetics, 5, 363.
- Henryon et al. (2014). Animal-breeding schemes using genomic information need breeding plans designed to maximise long-term genetic gains. Livestock Science, 166, 38-47.
- Heslot, N., Jannink, J. L., & Sorrells, M. E. (2015). Perspectives for genomic selection applications and research in plants. Crop Science, 55(1), 1-12.
- Xavier et al. (2016). Walking through the statistical black boxes of plant breeding. Theoretical and applied genetics, 129(10), 1933-1949.
- Hickey et al. (2017). Genomic prediction unifies animal and plant breeding programs to form platforms for biological discovery. Nature genetics, 49(9), 1297.

# <u>*Tools for sceptical thinking*</u>

- Wherever possible there must be independent confirmation of the 'facts'.
- Encourage substantive debate on the evidence by knowledgeable proponents of all points of view.
- Arguments from authority carry little weight - 'authorities' have made mistakes in the past. They will do so again in the future. Perhaps a better way to say it is that in science there are no authorities; at most, there are experts.
- Spin more than one hypothesis. If there's something to be explained, think of all the different ways in which it *could* be explained. Then think of tests by which you might systematically disprove each of the alternatives. What survives, the hypothesis that resists disproof in this Darwinian selection among 'multiple working hypotheses', has a much better chance of being the right answer than if you had simply run with the first idea that caught your fancy.*
- Try not to get overly attached to a hypothesis just because it's yours. It's only a way-station in the pursuit of knowledge. Ask yourself why you like the idea. Compare it fairly with the alternatives. See if you can find reasons for rejecting it. If you don't, others will.
- Quantify. If whatever it is you're explaining has some measure, some numerical quantity attached to it, you'll be much better able to discriminate among competing hypotheses. What is vague and qualitative is open to many explanations. Of course there are truths to be sought in the many qualitative issues we are obliged to confront, but finding *them* is more challenging.
- If there's a chain of argument, *every* link in the chain must work (including the premise) - not just most of them.
- Occam's Razor. This convenient rule-of-thumb urges us when faced with two hypotheses that explain the data *equally well* to choose the simpler.
- Always ask whether the hypothesis can be, at least in principle, falsified. Propositions that are untestable, unfalsifiable are not worth much. Consider the grand idea that our Universe and everything in it is just an elementary particle - an electron, say - in a much bigger Cosmos. But if we can never acquire information from outside our Universe, is not the idea incapable of disproof? You must be able to check assertions out. Inveterate sceptics must be given the chance to follow your reasoning, to duplicate your experiments and see if they get the same result.

**Carl Sagan - The Demon Haunted World (p.197)**