# Trends of predictive breeding

Changes in the plant breeding landscape driven by more, better and different data

**Alencar Xavier**
**Breeding Analyst at Corteva Agrisciences**
**Adjunct professor at Purdue University**

02/2024

1. **Introduction**
   - Rationale
   - Breeding with ML
2. **Modeling**
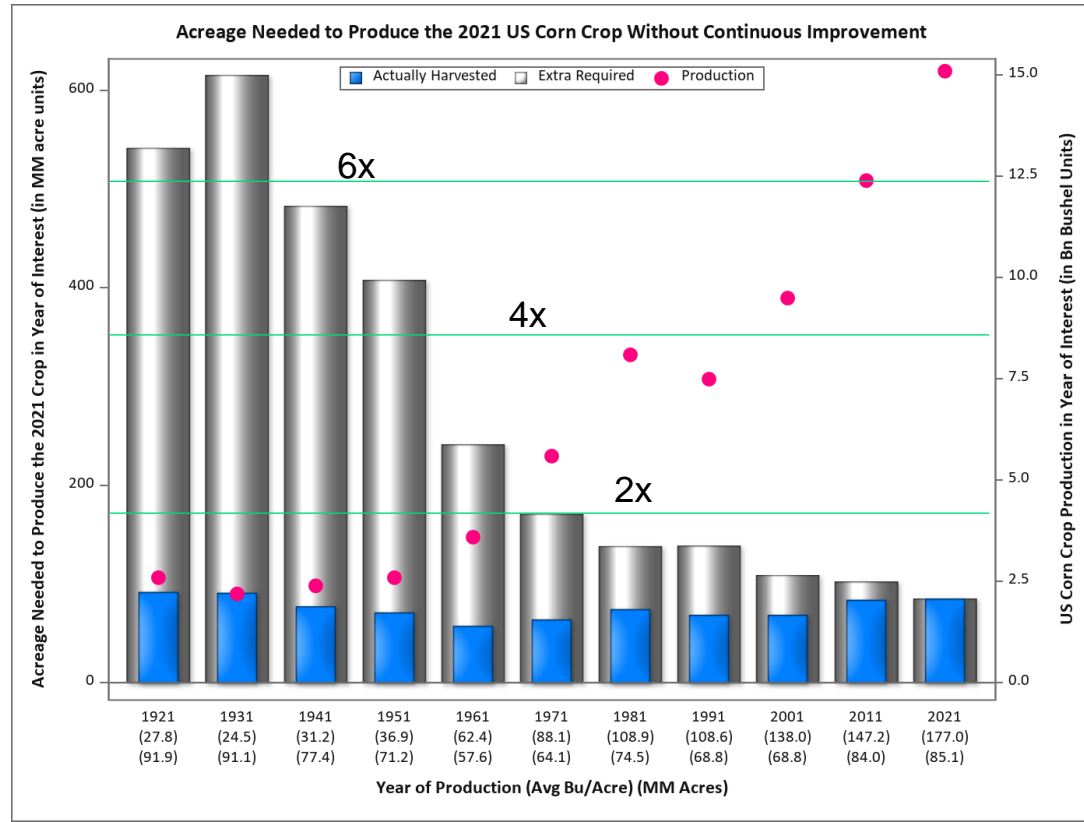   - Approaches
   - Correlated information
3. **Analytics**
   - Breeding objectives
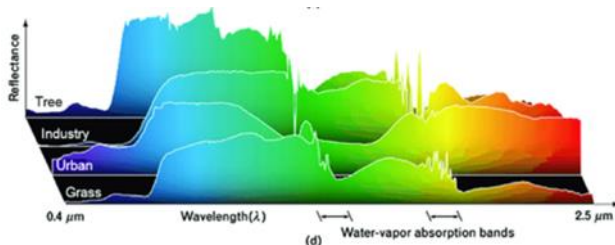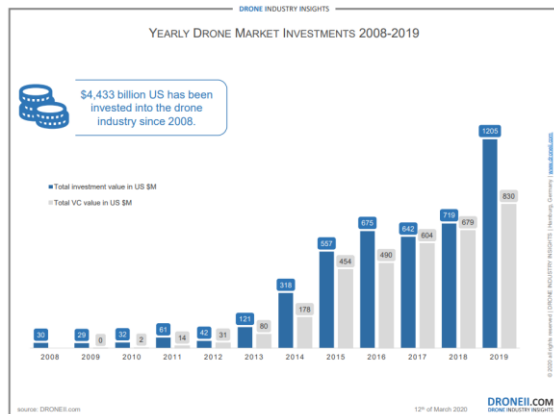   - Target G x E x M
   - Validation
4. **Conclusion**

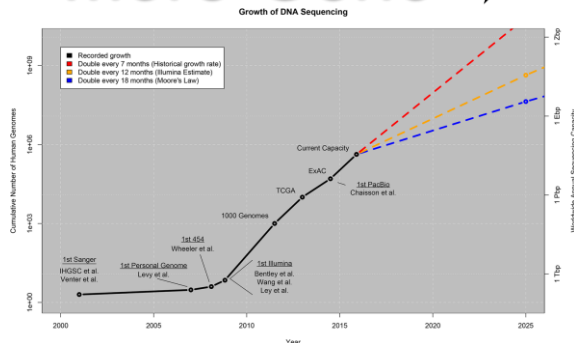# What are some implications of continuous Corn Improvement?



Source: Totir 2021, ASTA

*Based on 2021 USDA NASS data

Alencar.Xavier@Corteva.com
Quantitative Geneticist, Breeding Analyst LAAF

# More Pheno



YEARLY DRONE MARKET INVESTMENTS 2008-2019

$4,433 billion US has been invested into the drone industry since 2008.

https://www.mdpi.com/2076-3417/12/5/2570

# More Geno



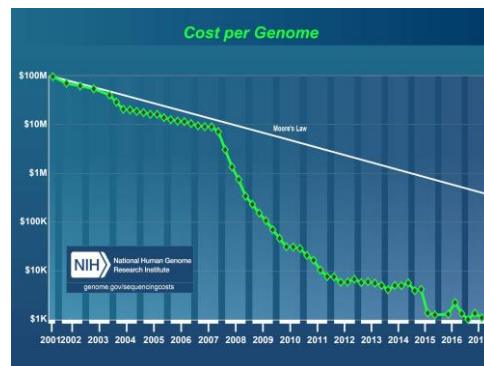The Cost of Sequencing a Human Genome. NIH.
https://www.genome.gov/27565109/the-cost-of-sequencing-a-human-genome/



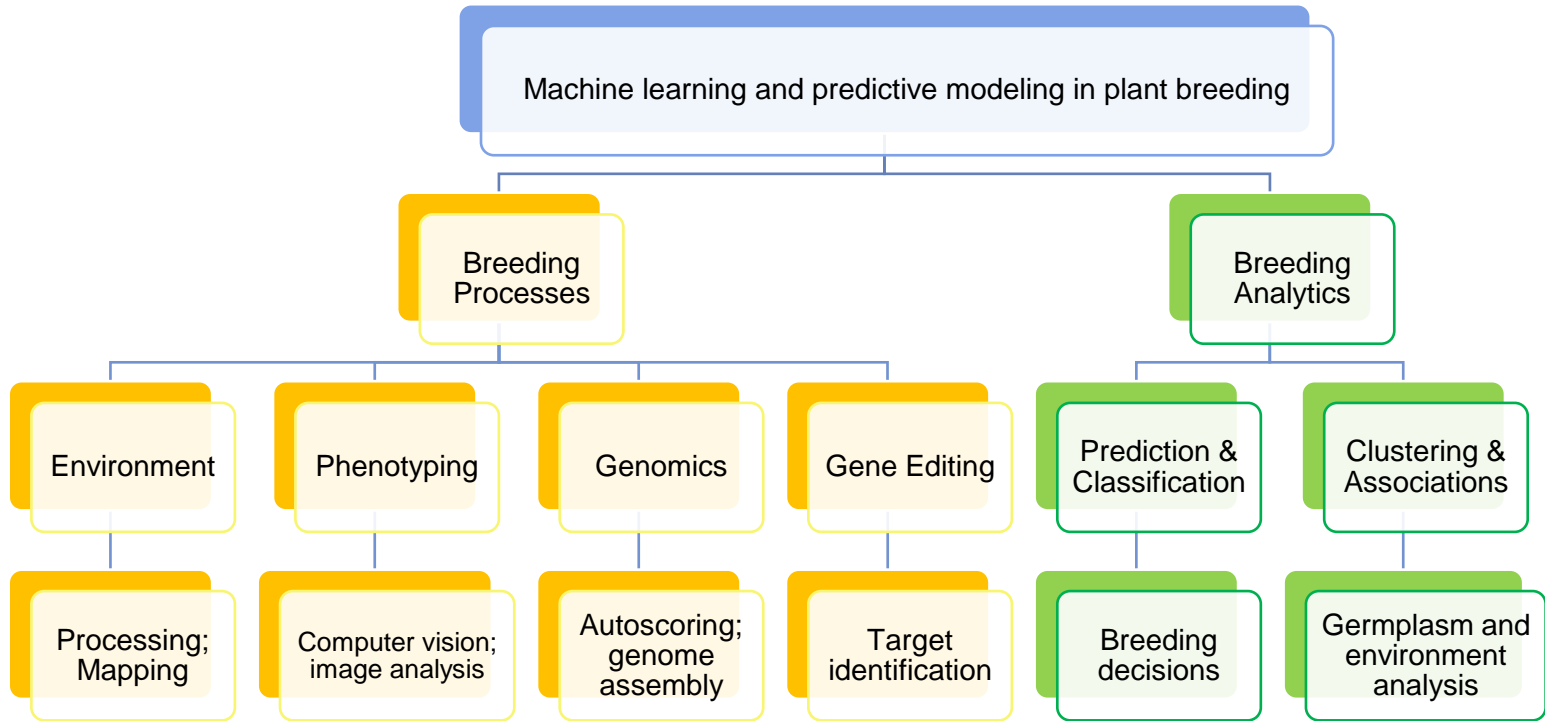Stephens, Z. D.et al. (2015). Big data: astronomical or genomical? *PLoS biology*, *13*(7), e1002195.

# More Env

- **UC Merced GridMET**
- **NWS NOAA**
- **NASA GISS, NASA power**
- **Harmonized SoilDB**
- **USDA SSURGO**

# More Computing

Alencar.Xavier@Corteva.com
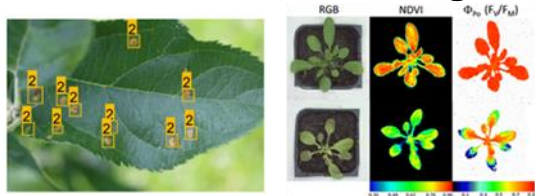Quantitative Geneticist, Breeding Analyst LAAF

CORTEVA agriscience

# Machine learning in breeding processes

**Enhancing databases, automating lab tasks field work**

## phenotyping

### Disease, stress scoring



https://www.mdpi.com/2673-2688/2/3/26
https://www.biomedcentral.com/collections/phenomics

### Phenotype automation
(e.g., plant height, identify new traits)



https://www.mdpi.com/2072-4292/8/12/1031
https://www.ncbi.nlm.nih.gov/pmc/articles/PMC7706325/

## environment

### Mapping / zoning



https://www.publish.csiro.au/cp/CP14007

### Latent weather, soil



https://doi.org/10.1093/bioinformatics/btaa971

## biotech

### SNP calls, genome assembly



https://doi.org/10.1186/1753-6561-3-s7-s58
https://www.nature.com/articles/s41467-022-29843-y

### Embryo rescue DH production



https://www.nature.com/articles/s41598-022-06336-y

### Gene editing targets



https://doi.org/10.1093/bioinformatics/btab268

# Challenges in Corn Improvement



Complex function: G * E * M

**G = Genetics; E = Environment; M = Management**

Source: Totir 2021, ASTA

**Alencar.Xavier@Corteva.com**
**Quantitative Geneticist, Breeding Analyst LAAF**

# Law of the minimum



Disease resistance

Standability

NUE

Moisture

Lodging score

Pest resistance

Yield potential

Drought resistance

Addressing the limiting factors

Realized yield in farmers' fields

# End goal of GxExM characterization



Phenotypic data
Genomic data
Environmental data
Management data
Business data

Modeling & Analytics
(Predictive, Prescriptive)

1) SELECTION
2) CHARACTERIZATION
3) PRODUCT PLACEMENT
4) SEED PRODUCTION
5) LOGISTICS OPTIMIZATION

Addressing yield limiting factors through breeding and agronomics

Addressing producibility and commercial questions

Alencar.Xavier@Corteva.com
Quantitative Geneticist, Breeding Analyst LAAF

CORTEVA
agriscience

# Linear models

# Crop models



Source: https://evolution.berkeley.edu/teach-resources/genes-environment-phenotype/

genes + environment = phenotype

https://www.nature.com/articles/s41477-019-0398-8

**Phenotype** ↔ **Secondary phenotypes**

Phenotype
- Environment
  - Local control
  - Management
  - Soil & Weather
- Genetics
  - Heritable
  - Non-heritable
- Interactions
  - Residuals
  - GxE, GxM, ExM, GxExM

**General framework, trait-crop agnostic**

**Biological knowledge, trait-crop specificity**

Phenotype
- Local control
- Gen | Env
  - Underlying physiological traits
    - Heritable genetics
    - Per se genetics
  - Growth conditions
    - Management
    - Soil & Weather
- Residuals

Walking through the statistical black boxes of plant breeding

Alencar Xavier[1] · William M. Muir[2] · Bruce Craig[3] · Katy Martin Rainey[1]

Integrating Crop Growth Models with Whole Genome Prediction through Approximate Bayesian Computation

Frank Technow[1]*, Carlos D. Messina[2], L. Radu Totir[1], Mark Cooper[2]

Alencar.Xavier@Corteva.com
Quantitative Geneticist, Breeding Analyst LAAF

CORTEVA agriscience

# Linearly correlated phenotypes

Using machine learning to infer connections from data

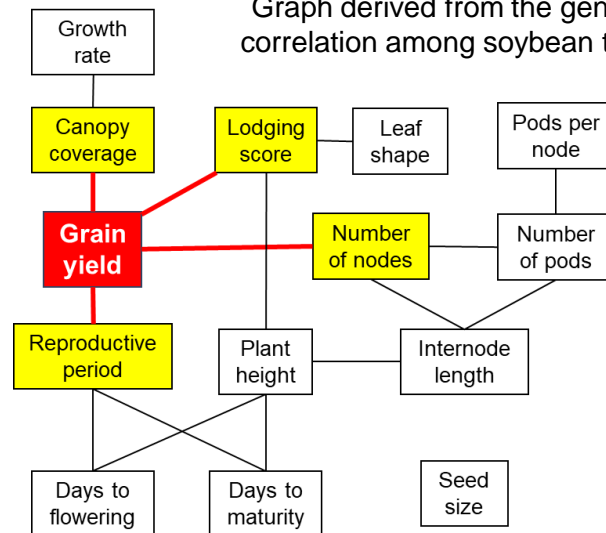Phenotype (y) = Gen|Trials(g) + Plot level noise(e)

$$\begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} g_1 \\ g_2 \end{bmatrix} + \begin{bmatrix} e_1 \\ e_2 \end{bmatrix}$$

$$\text{Variance} \begin{bmatrix} y_1 \\ y_2 \end{bmatrix} = \begin{bmatrix} \sigma_{g_1}^2 & \sigma_{g_{12}} \\ \sigma_{g_{12}} & \sigma_{g_2}^2 \end{bmatrix} + \begin{bmatrix} \sigma_{e_1}^2 & \sigma_{e_{12}} \\ \sigma_{e_{12}} & \sigma_{e_2}^2 \end{bmatrix}$$

**CORRELATED INFORMATION**

**Phenotypic network**

Graph derived from the genetic correlation among soybean traits



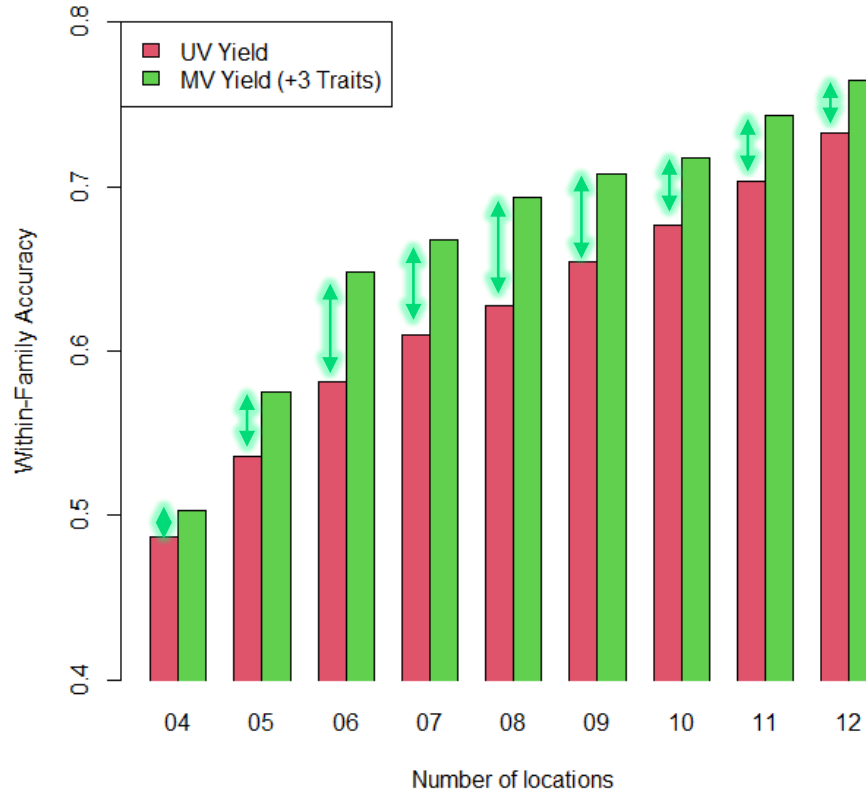A new approach fits multivariate genomic prediction models efficiently

Alencar Xavier[1,2*†] and David Habier[1*†]

Using unsupervised learning techniques to assess interactions among complex traits in soybeans

Alencar Xavier · Benjamin Hall · Shaun Casteel · William Muir · Katy Martin Rainey

Alencar.Xavier@Corteva.com
Quantitative Geneticist, Breeding Analyst LAAF

**NAWE2L1YRXS5T ( R1 experiment 2020 )**

# Leveraging information from <span style="color:red">secondary traits</span>

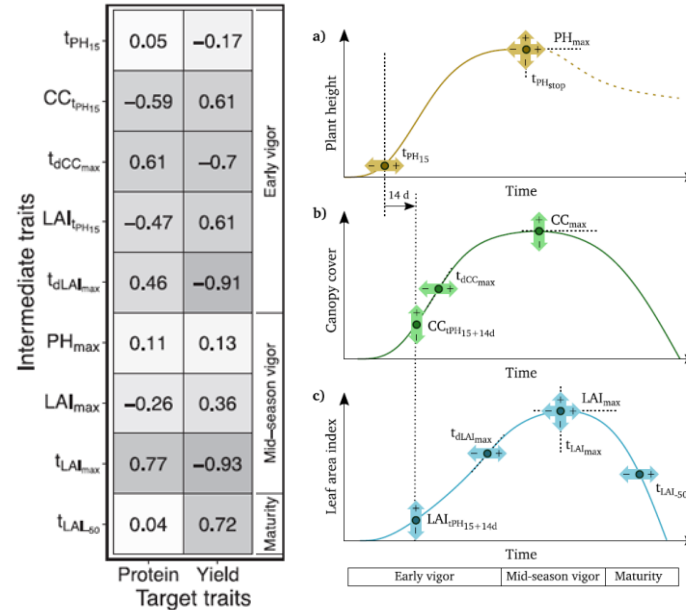Multi-trait analysis provided an average increase in accuracy of **0.03**

Equivalent to adding **~1.7 locations**
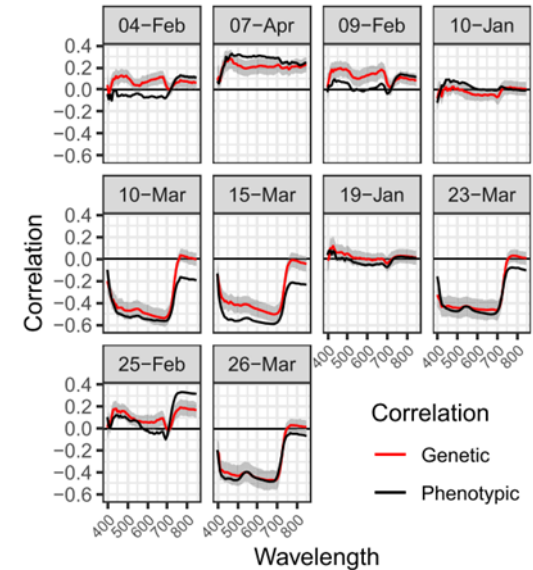
# High-throughput phenotyping

More data, more information, more complexity, more opportunities…

Correlation(HTP, main trait)

Correlation(HTP, main trait) with *large number of traits*



Singh A.K. et al. (2021)
High-Throughput Phenotyping in Soybean.
https://doi-org/10.1007/978-3-030-73734-4_7

Roth et al. (2022) High-throughput field phenotyping of soybean:
Spotting an ideotype. https://doi.org/10.1016/j.rse.2021.112797

Runcie et al. (2021) Mega-scale linear mixed models
for genomic predictions with thousands of traits.
https://doi.org/10.1186/s13059-021-02416-w

**Alencar.Xavier @Corteva.com**
**Quantitative Geneticist, Breeding Analyst LAAF**

**CORTEVA** agriscience

**Alencar.Xavier@Corteva.com**
**Quantitative Geneticist, Breeding Analyst LAAF**

# "Breeding objective"

$$f(\text{market segment}, \text{farming systems})$$

- Set of traits of interest (**TOI**) bred into a
  <mark>WHAT</mark>

| Yield, moisture, relative maturity, disease resistance, stability, **trait package, producibility** |
| --- |

- Target population of genetics (**TPG**) for a given
  <mark>WHO</mark>

| Corn 111-121, corn 122-130 white corn 118-123 |
| --- |

- Target population of environments (**TPE**) and management (**TPM**) practices   <mark>WHERE</mark>
  <mark>HOW, WHEN</mark>

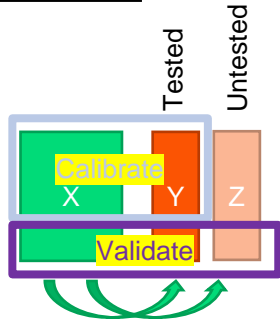| Drought, irrigation, early planting, varying levels of disease pressure, different soil types |
| --- |

$$\rho_{\text{GxExM}} = \rho_{\text{TPG}} \times \rho_{\text{TPE}} \times \rho_{\text{TPM}}$$

# Model testing and validation schemes

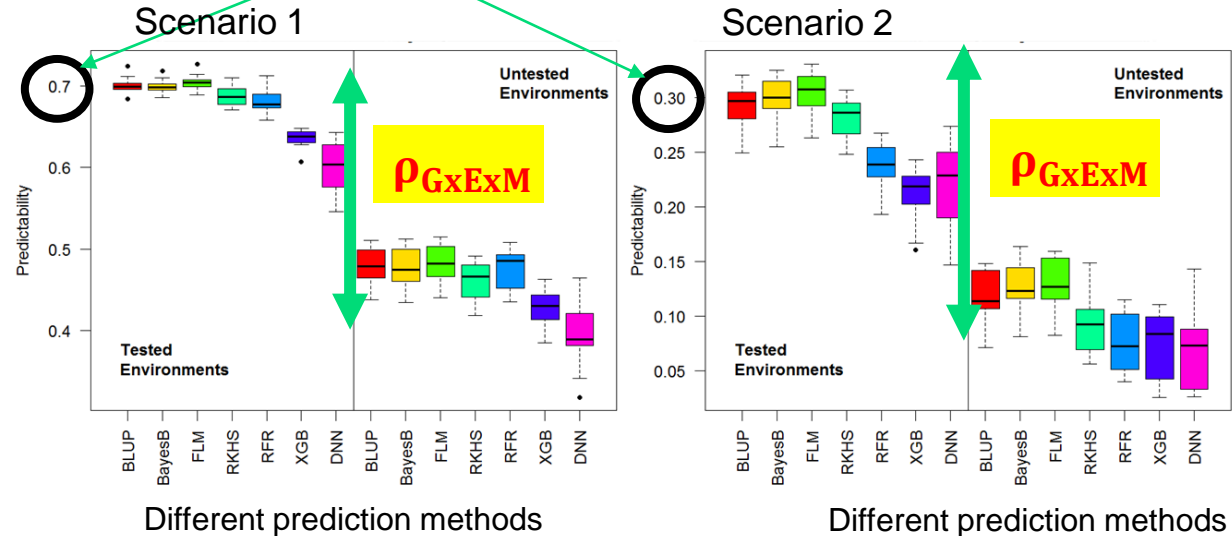$$\text{Prediction accuracy} \propto \sqrt{H^2} \times \rho_{\text{GxExM}}$$



Upper limit is determinate by the application $\cong \sqrt{H^2}$

CV scheme
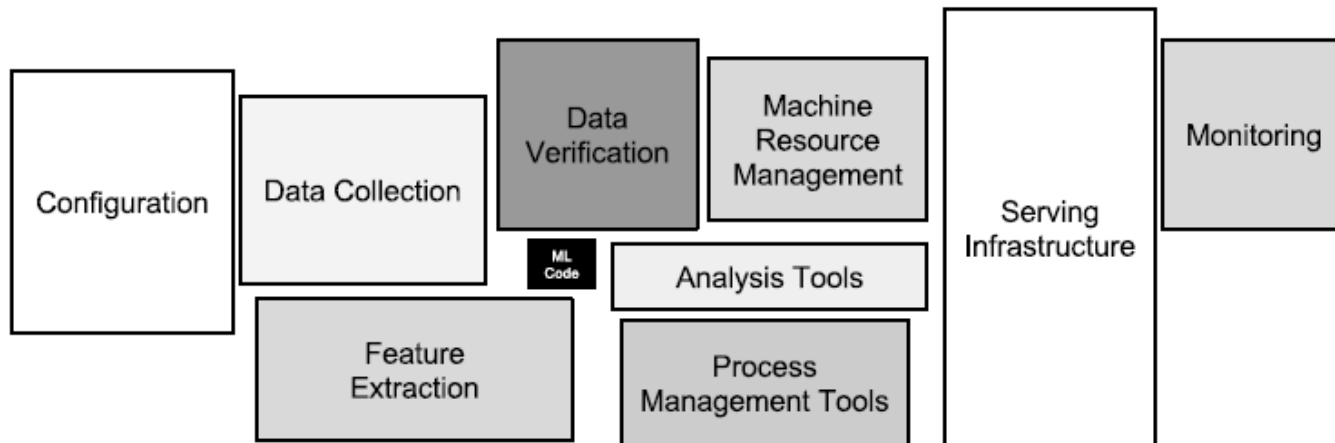
Technical nuances of machine learning: implementation and validation of supervised methods for genomic prediction in plant breeding

Alencar Xavier [1*]

Scenario 1

Scenario 2

Different prediction methods

Different prediction methods

# Multiple factors play a role on the implementation of <u>predictive breeding</u> and <u>automated systems</u> beyond proof of concepts using cross-validations
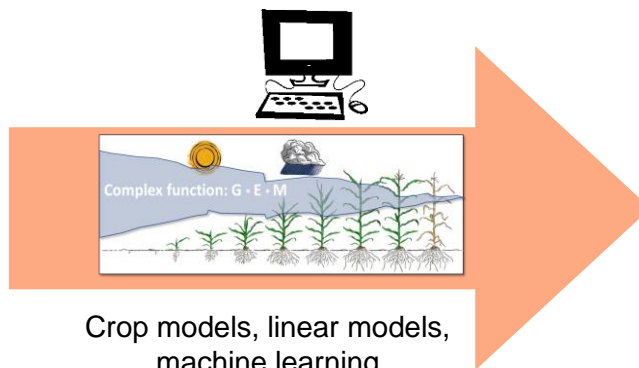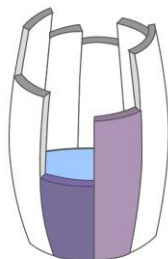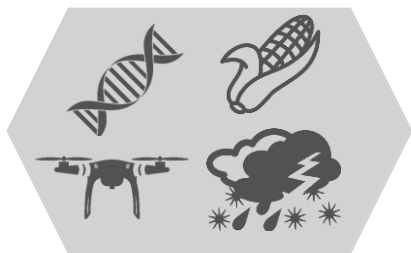
# With <u>more data and better analytics</u>, breeding can respond faster the new <u>farming challenges and trends</u>

- New & better management
- Changing environment
- New pests and diseases

Predictive breeding
(TPE, TPG, TPM)

Complex function: G · E · M

Crop models, linear models, machine learning

$$Acc \propto \mathbf{GxExM} \times \mathbf{h^2}$$

**Performance**

**Producibility**

**Robustness**

CORTEVA
agriscience

# Thank you for your attention!

**Final remarks**:

1) Plant breeding relies on analytics for multiple processes and analytics

2) Harnessing GxExM information benefits accuracy, business impact

3) Modeling is contingent to the target genetics, environments and management

# Questions??

*Alencar Xavier*

Alencar.Xavier@Corteva.com

**CORTEVA** agriscience