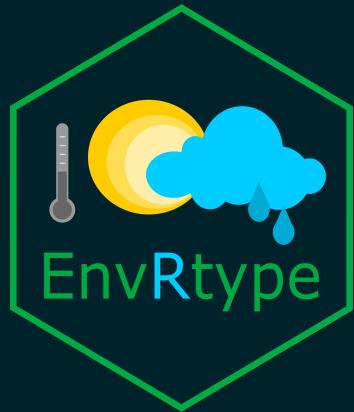


GEMS-R WEBINARS

Aug 28th, 2021

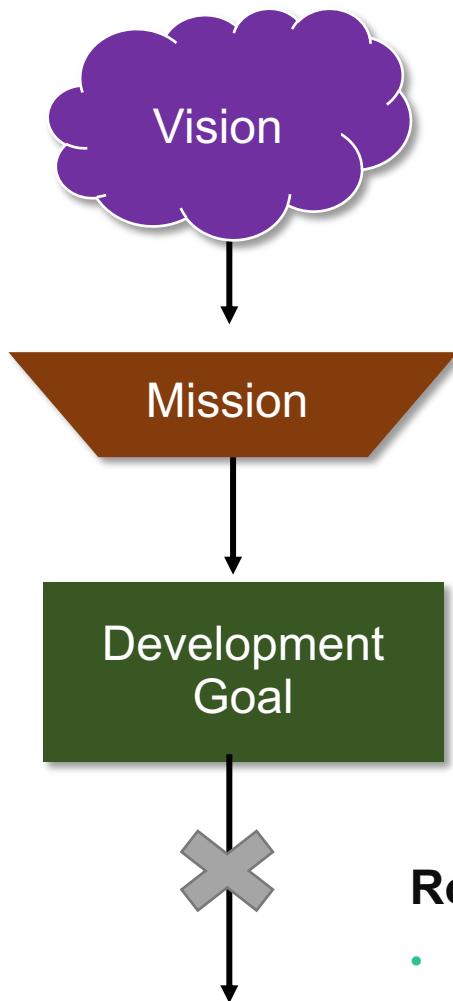


Germano Costa-Neto

Institute for Genomic Diversity – Cornell University
CIMMYT / BSU expert in envirotyping analytics



TRAINING COURSE
Basics for understanding
envirotyping analytics and starting
an enviromic pipeline in R



If $y = G+E$, there is an urgent need to develop ways to better understand E in quantitative genetics

Merge genomics and enviromics: interplay quantitative genomics and ecophysiology matters in the plant breeding context

A pipeline for collect, process and integrate environmental information in the current quantitative genetics analysis and genomic prediction models across diverse growing conditions

Roadblock

- Does enviromics really improves the prediction ability of multi-environment GP?
- Does the use of enviromics might fill the lack of phenotypic records across drastically sparse multi-environment phenotyping networks?
- Which is the best kernel method to model the environmental relatedness ?
- Can we create an open-source software and a pipeline for it?

Part I

Theoretical Concepts & Data Set

Core ideas

- Accounting for ecophysiology knowledge using environmental information in predictive breeding
- Interplay environment, GxE, reaction-norm and phenotypic plasticity in the quantitative genetics and genomic prediction context

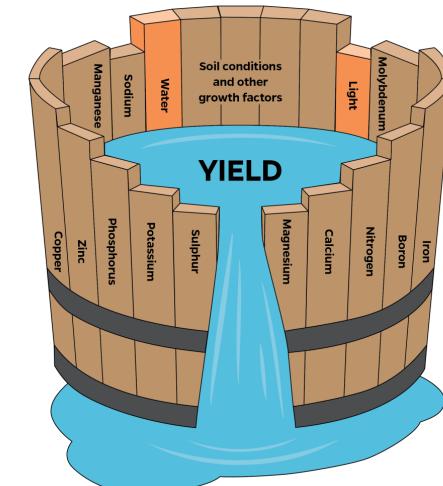
Environment

Check these references in the course repository

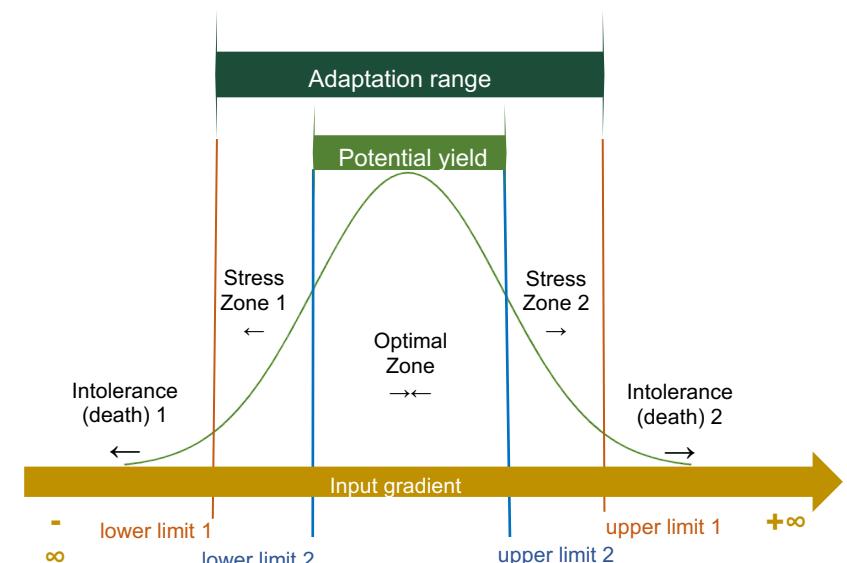
Cooper et al., (2014); Xu (2016); Resende et al., (2021); Costa-Neto et al. (2021a);
Crossa et al (2021); Costa-Neto & Fritsche-Neto (2021)

Liebig's Law of Minimum (1840)

Whitson and Walster (1912)



Shelford (1931, 1932)
Tolerance Limits and adaptation



Agronomic Standpoint

- Single-unit (planting date, location, management)
- Location + Planting date = Fixed + Random Factors
- Actual Management = Expected Management + biotic and abiotic interferences

Geneticist Standpoint (Plant Breeders)

- Near-Iso Environments (managed conditions)
- Drought-Stress Trials (screening)
- Nitrogen level conditions (screening)
- A non-genetic source of variation (to be controlled as a fixed factor)
- A sample of the theoretical **target population of environments (TPE)**

Ecophysiological Standpoint

- Relation between the **availability of inputs** and their resource absorption/alocation
- Availability as function of amount and frequency
- Core of events **linking** soil-plant-atmosphere dynamics (plant-env-pathogen)

Then, what is envirotyping? enviromics? envirome?

- **Envirotyping** (environmental + typing) gathers the steps of collecting, processing, and associating environmental data with phenotypic data to understand the typology of environments for some MET, which is conducted relating it to the target population of environments (TPE) of the breeding program.
- **Enviromics** is the large-scale envirotyping, based on the collections of environmental data across time and space to establish a global association between the *crops envirome* (the core of TPEs for a certain genotype or specie) and the phenotypic variation of key factors driving GxE.

Why is envirOMICS ? Just a wave?

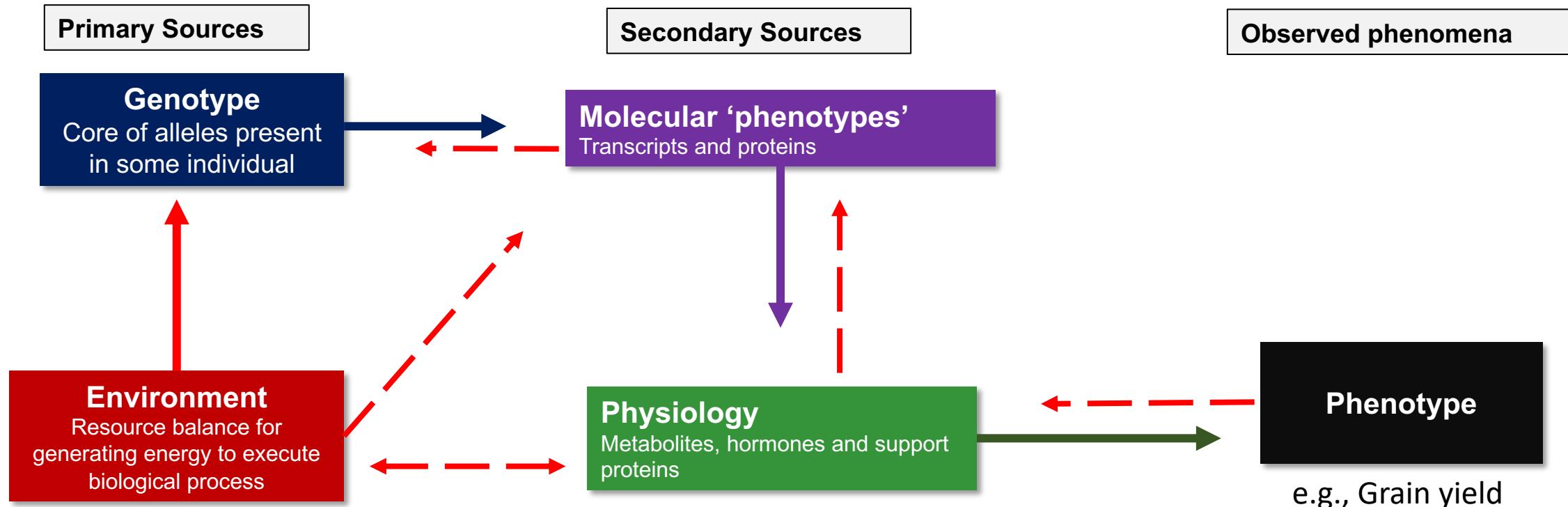
- **Because it can bridge different omics!**

Interplay Quantitative and Molecular Genetics in Predictive Breeding

- **Phenotype:** end-result of multiple interactions among expressed genes, epigenome and environmental signals
- Environmental acts at nuclear, cellular and physiologycal level (plant level)
- It is not a one-way road: compensatory effects, homeostasis, plasticity, pleiotropy, epigenetics...across time and space!

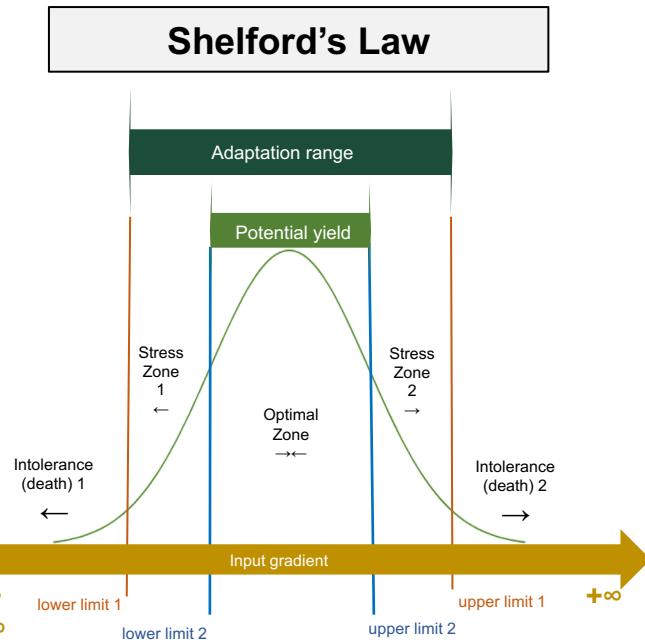
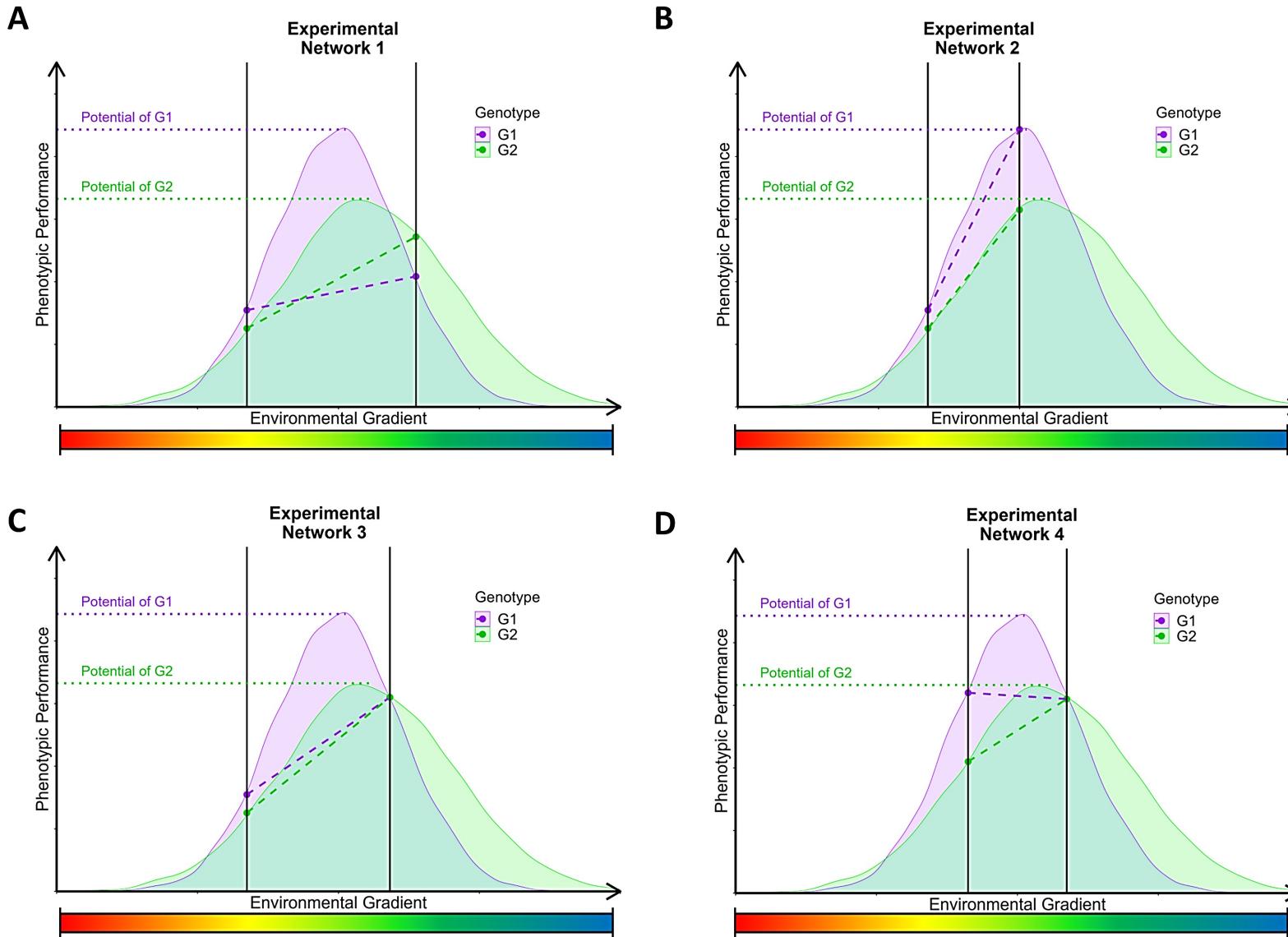
Solid lines: **Central Dogma**

Dashed lines: **Surrounding environment outside gene level**



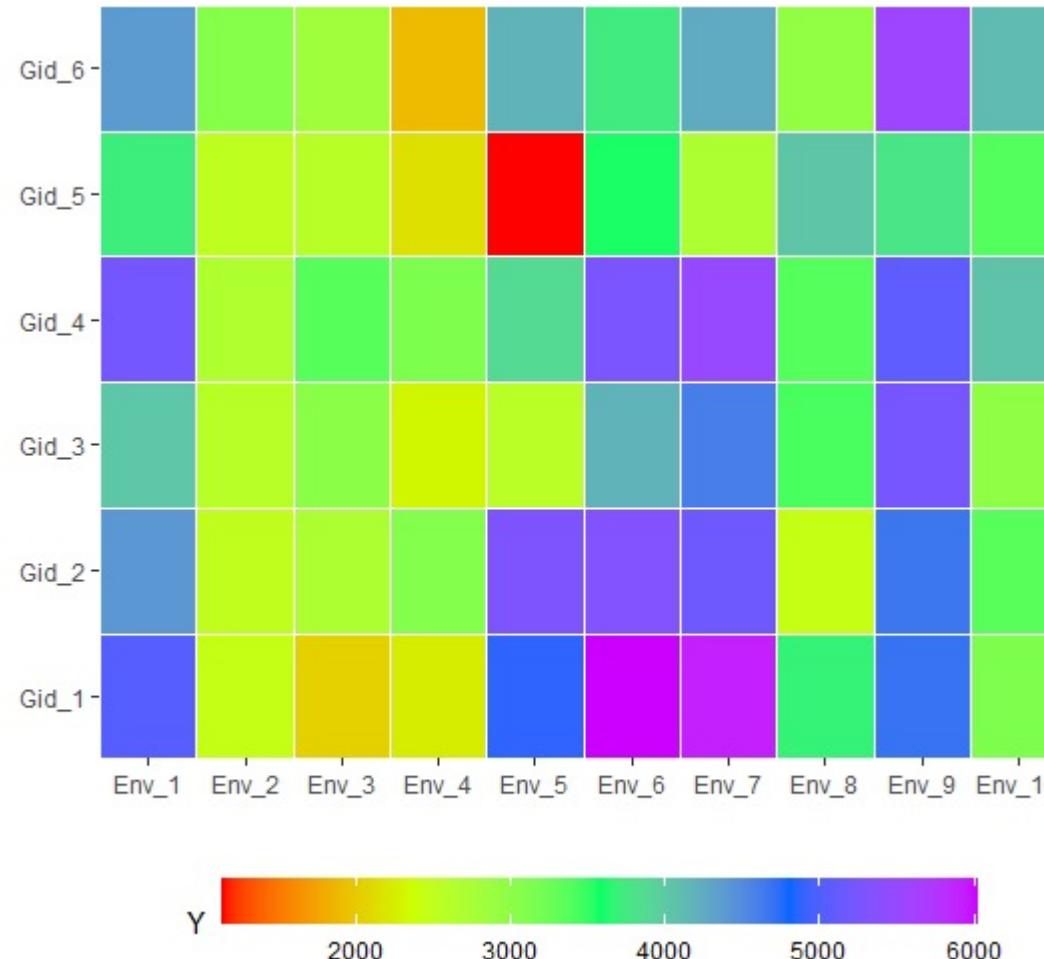
Genotype x Environment Interaction (GxE) as an emergent property

Product of reaction-norm, respecting the particular plasticity for each genotype, expressed across a given experimental gradient



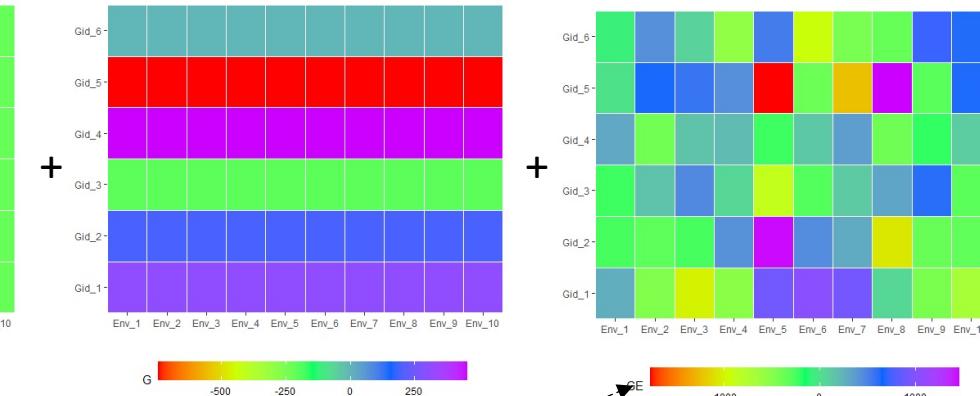
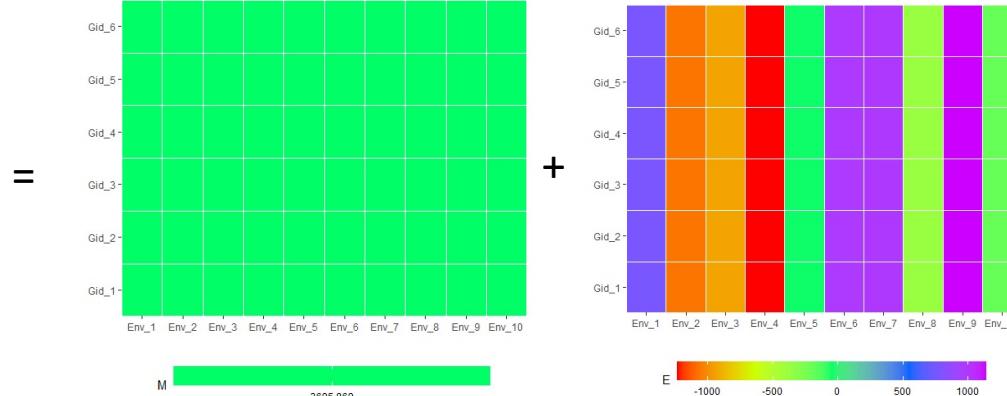
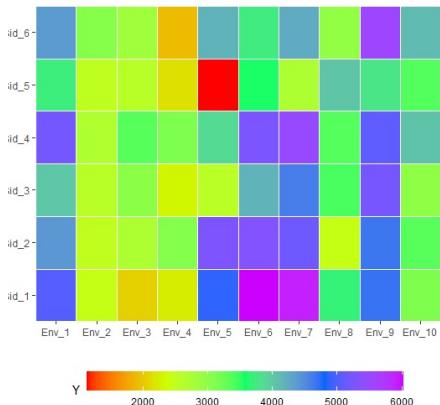
Two-way table (genotypes per environment)

$$\mathbf{Y}_{pq} = \begin{bmatrix} \hat{Y}_{11} & \cdots & \hat{Y}_{1q} \\ \vdots & \ddots & \vdots \\ \hat{Y}_{p1} & \cdots & \hat{Y}_{pq} \end{bmatrix} = (\mathbf{1}_{pq} \otimes \boldsymbol{\mu}) + (\mathbf{1}_p \otimes \mathbf{E}'_q) + (\mathbf{1}'_q \otimes \mathbf{G}_p) + \mathbf{G}\mathbf{E}_{pq}$$



Two-way table (genotypes per environment)

$$Y_{pq} = \begin{bmatrix} \hat{Y}_{11} & \dots & \hat{Y}_{1q} \\ \vdots & \ddots & \vdots \\ \hat{Y}_{p1} & \dots & \hat{Y}_{pq} \end{bmatrix} = (\mathbf{1}_{pq} \otimes \boldsymbol{\mu}) + (\mathbf{1}_p \otimes E'_q) + (\mathbf{1}'_q \otimes G_p) + GE_{pq}$$



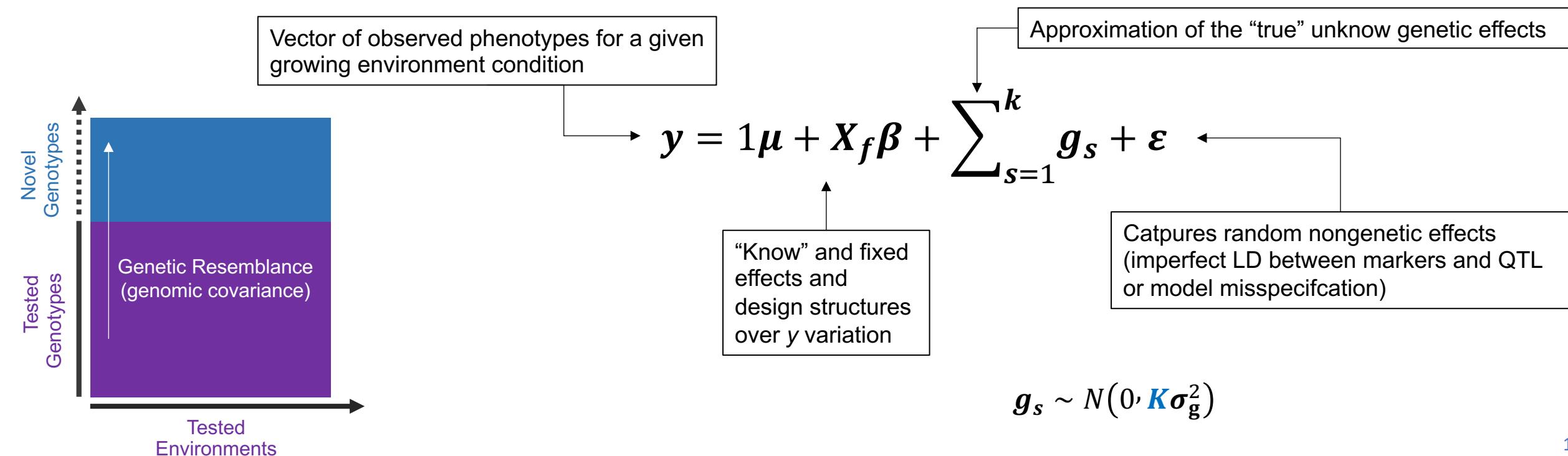
Enviromic (environmental covariables) realized relationship to 'estimate' reaction-norms or assembly environment-to-environment relatedness

Genomic (or Pedigree-based) realized relationship to explore genotype-genotype covariances, and a block-diagonal for mimic the genomic by environment effects

Genomic Prediction (GP)

Some key works:

- Introduced by Meeuwissen et al (2001) aimed for a whole-genome regression over some phenotype
- Burgueño et al (2012) proposed the marker x environment interaction model
- Jarquin et al (2014) and Lopez-Cruz et al (2015) presents the multi-environment GBLUP (using hadamard's product)
- Heslot et al (2014) and Ly et al (2018) presents a factorial regression with GP
- Bandeira e Sousa (2017) presents the use of gaussian kernel (GK) for modeling GxE in maize
- Morais-Júnior et al (2018) expands Jarquin's model to account for quantile covariables
- Crossa et al (2019) and Cuevas et al (2019) presents the use of Deep Kernel for modeling G effects

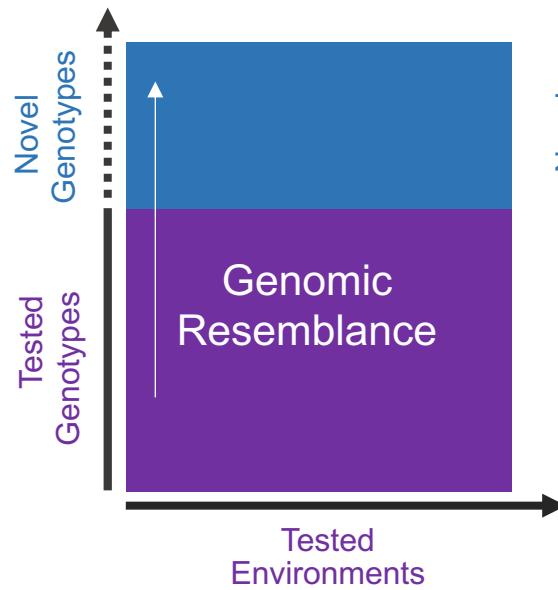


Why is *Envirotyping* useful for multi-environment genomic prediction?

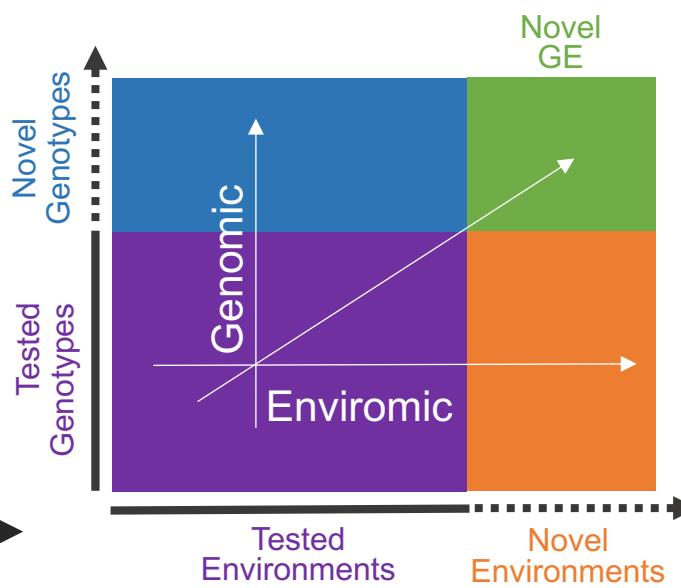
Try to deal with the environmental-phenotype (E2P) covariances in the phenotypic records in the context of MET

- Benchmark GP models were conceived for modeling genotype-phenotype covariances (G2P)
- Different field trials, different E2P; consequently, different G2P pattern (and GP accuracy)
- Multi-environment trials = more source of noise due to E2P?
- Implicit and Explicit covariates can shape E2P and translate some noise in exploitable pattern

MET Prediction (Pure Genomic)



MET Prediction (Enviromic + Genomic)



Adding value to Predictive Breeding

- Model genomic reaction-norms
- Realization of yet-to-be-seen GxE
- Optimize MET for training GP
- Orient cultivar targeting

Models for Multi-Environment Genomic Prediction (MET-GP)

$$\mathbf{y} = \mathbf{1}\boldsymbol{\mu} + \mathbf{Z}_E\boldsymbol{\beta} + \sum_{s=1}^k \mathbf{g}_s + \sum_{s=1}^k \mathbf{gE}_s + \sum_{r=1}^l \mathbf{W}_r + \sum_{s=1}^k \sum_{r=1}^l \mathbf{gW}_{sr} + \boldsymbol{\varepsilon}$$

$$\mathbf{y} = [\mathbf{y}_1, \mathbf{y}_2, \dots \mathbf{y}_q]^T$$

Model	Effect		
	Genetic	Environment	GxE
Main Genomic Effects (MM)	$\sum_{s=1}^p \mathbf{g}_s \neq 0$	$\sum_{r=1}^l \mathbf{W}_r = 0$	$\sum \mathbf{gE}_s = 0 \quad \sum \sum \mathbf{gW}_r = 0$
MM + single GxE deviation (MDs)	$\sum_{s=1}^p \mathbf{g}_s \neq 0$	$\sum_{r=1}^l \mathbf{W}_r = 0$	$\sum \mathbf{gE}_s \neq 0 \quad \sum \sum \mathbf{gW}_r = 0$
MM + main Enviromic effects (EMM)	$\sum_{s=1}^p \mathbf{g}_s \neq 0$	$\sum_{r=1}^l \mathbf{W}_r \neq 0$	$\sum \mathbf{gE}_s = 0 \quad \sum \sum \mathbf{gW}_r = 0$
EMM + single GxE deviation (EMDs)	$\sum_{s=1}^p \mathbf{g}_s \neq 0$	$\sum_{r=1}^l \mathbf{W}_r \neq 0$	$\sum \mathbf{gE}_s \neq 0 \quad \sum \sum \mathbf{gW}_r = 0$
EMM + reaction norm GxW (RNMM)	$\sum_{s=1}^p \mathbf{g}_s \neq 0$	$\sum_{r=1}^l \mathbf{W}_r \neq 0$	$\sum \mathbf{gE}_s = 0 \quad \sum \sum \mathbf{gW}_r \neq 0$
EMDs + reaction norm GxW (RNMDs)	$\sum_{s=1}^p \mathbf{g}_s \neq 0$	$\sum_{r=1}^l \mathbf{W}_r \neq 0$	$\sum \mathbf{gE}_s \neq 0 \quad \sum \sum \mathbf{gW}_r \neq 0$

Implementation

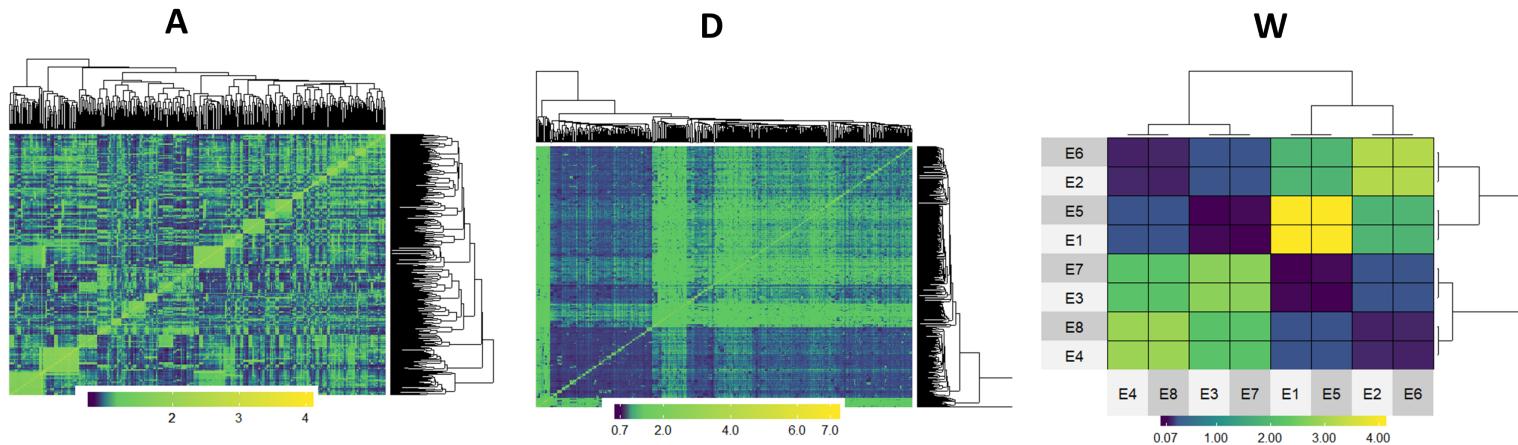
- Independent Residuals = I
- Hierarchical Bayesian Modeling
- BGGE (now within EnvRtype)

Phenology (example for maize)

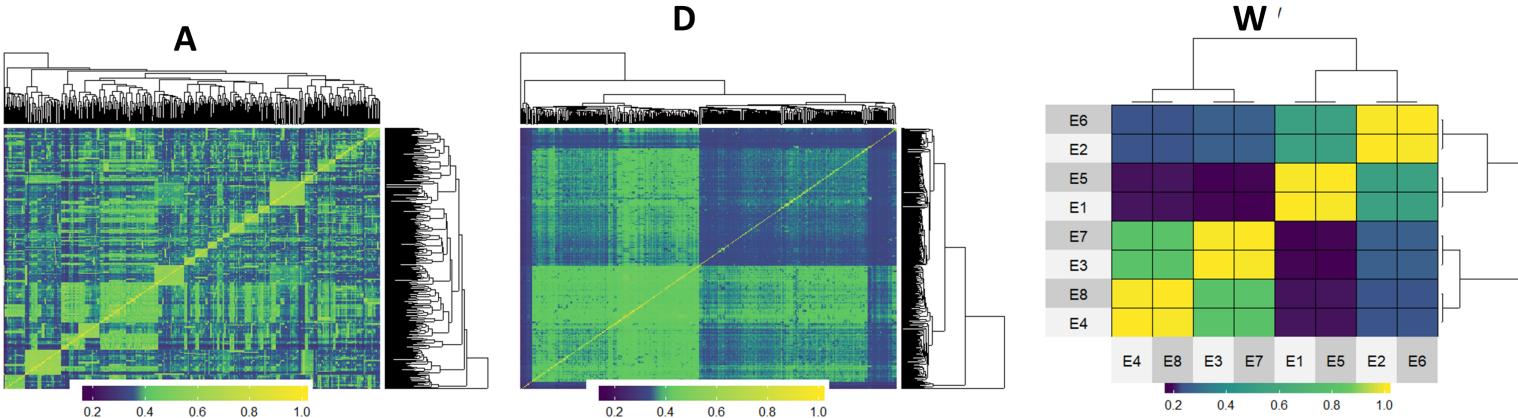
- emergence to the first leaf (V1, 14 DAE)
- V1 to the fourth leaf (V4, 35 DAE)
- V4 to the tasseling stage (VT, 65 DAE)
- VT to the kernel milk stage (R3, 90 DAE)
- R3 to physiological maturity (R6, 120 DAE)

Kernel Methods (GB, GK and DK)

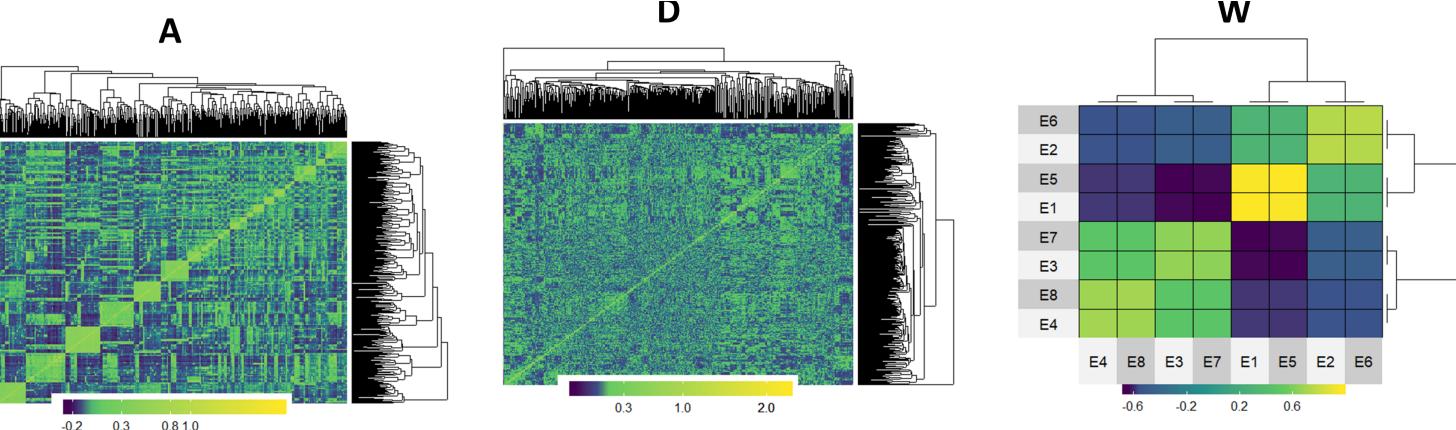
Benchmark GBLUP
(linear covariances)



Gaussian Kernel
(distances)

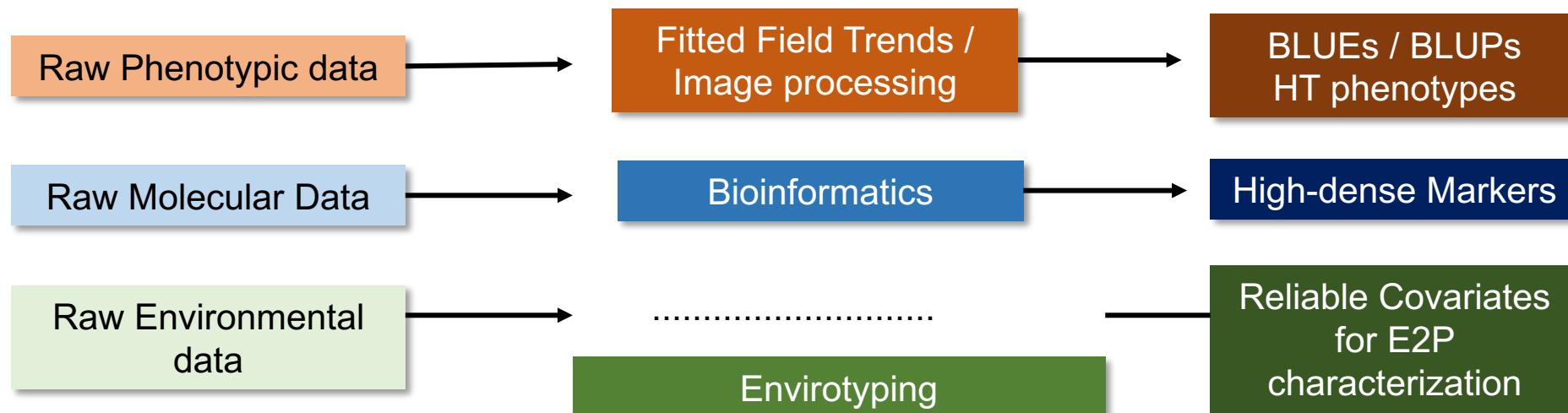


Deep Kernel based
on Arc-cosine
estimation (nonlinear)



Why EnvRtype?

- Availability vs Usefulness of Environmental Data
- Increase efficiency without increase costs of field trials
- Use any environmental information ?
- A clearly definition of how envirotyping can be exploited in predictive breeding
- From envirotyping we access the ‘envirome’ of the crop and better ‘explain’ some of the nongenetic effects

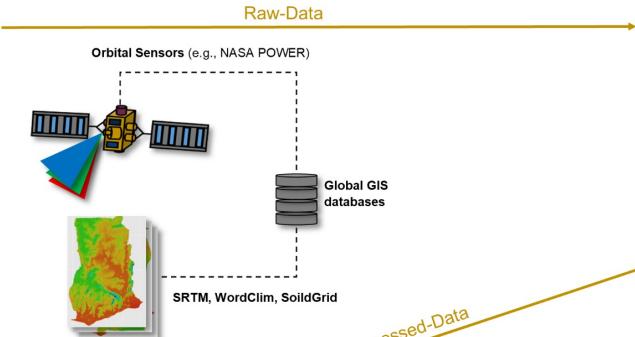


Module 1

Remote Data Collection

- `get_weather()`
- `extract_GIS()`

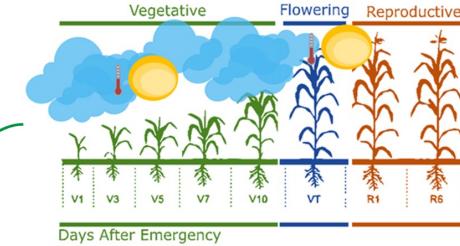
Raw-data from environmental sensors or `get_weather()` and `extract_GIS()` can be used in the further steps or processed



Raw-data Processing

- `SummaryWTH()`
- `ProcessWTH()`
- `param_temperature()`
- `param_radiation()`
- `param_atmospheric()`

Data processing involves the quality control and computation of additional variables

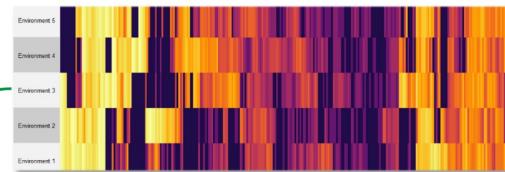


Module 2

Environmental Characterization

- `env_typering()`
- `W_matrix()`

Enviotype descriptors



Panel of enviotype descriptors or environmental covariabiles can be incorporated in predictive tools as environmental markers or to study similarity among environments

Environmental covariabiles

Module 3

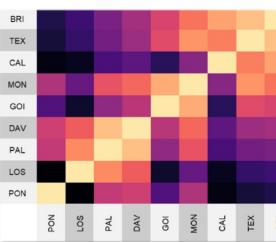
Environmental Similarity

- `env_kernel()`

Relatedness across environments derived from enviotyping data can be used to group environments with similar pattern

Enviromic Kernels

Kernel Models for phenotype prediction can be made involving genomic and enviromic-based sources derived from environmental similarity kernels



Enviromic-based Kernels Models

- `get_kernel()`
- `kernel_model()`

Models:

1. Genotypic Effects (MM and MDs)
2. Enviromic-enriched Main Effects (EMM and EMDs)
3. Enviromic-based Reaction-Norm (RNMM and RNMDs)

Genomic-based Predictions

Genomic-estimated Reaction-Norm Multi-environment Prediction

Outputs

Reaction-Norm

examples Factorial Regression GxE analysis

Environmental Grouping

Clustering (K-mean) TPE definition

Remote Sensing (Digital Envirotyping)



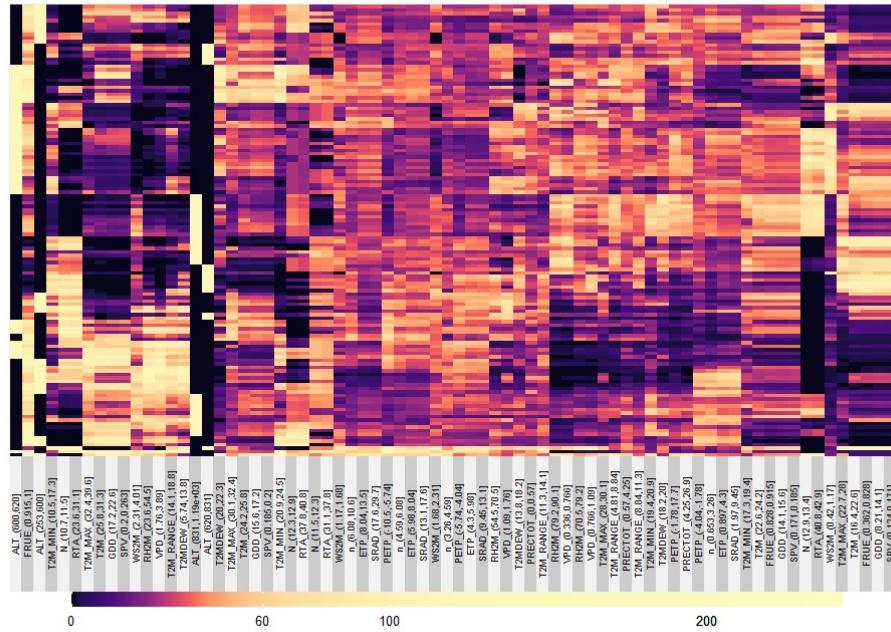
Source	ID	Environmental Factor	Unit
Nasa Power ¹	ALLSKY_SFC_SW_DWN	All sky insolation incident on a horizontal surface	MJ m ⁻² d ⁻¹
	ALLSKY_SFC_LW_DWN	Thermal infrared longwave radiative flux	MJ m ⁻² d ⁻¹
	WS2M	Wind speed at 10 m above the surface of the earth	m s ⁻¹
	T2M_MIN	Minimum air temperature at 2 m above the surface of the earth	°C d ⁻¹
	T2M_MAX	Maximum air temperature at 2 m above the surface of the earth	°C d ⁻¹
	T2MDEW	Dew-point temperature at 2 m above the surface of the earth	°C d ⁻¹
	RH2M	Relative air humidity at 2 m above the surface of the earth	%
	PRECTOT	Rainfall precipitation (P)	mm d ⁻¹
SRTM ²	ALT	Elevation (above sea level)	m
Computed ³	FRUE	Effect of Temperature on Radiation use Efficiency	-
	GDD	Growing Degree-days	°C d ⁻¹
	ETP	Evapotranspiration (ETP)	mm d ⁻¹
	PETP	Atmospheric water deficit P-ETP	mm d ⁻¹
	DVP	Deficit of vapor pressure	kPa d ⁻¹
	SVP	Slope of saturation vapor pressure curve	kPa C° d ⁻¹
	T2M_RANGE	Temperature Range	°C d ⁻¹
	RTA	Global Solar Radiation based on Latitude and Julian Day	MJ m ⁻² d ⁻¹

Cardinals for temperature (from literature)

Species	Suggested Cardinal Limit			
	Tbase1	Topt1	Topt2	Tbase2
Maize	8.0	30.0	37.0	45.0
Wheat	0.0	25.0	28.0	40.0
Rainfed Rice	8.0	30.0	37.0	45.0
Irrigated Rice (only vegetative stage)	8.0	28.0	40.0	45.0
Irrigated Rice (only reproductive stage)	15.0	25.0	35.0	45.0
Sorghum	8.0	30.0	37.0	45.0
Soybean	8.0	30.0	35.0	45.0
Peanut	8.0	30.0	35.0	45.0
Canola	0.0	25.0	28.0	40.0
Sunflower	8.0	30.0	34.0	45.0
Dry Bean	8.0	30.0	35.0	45.0
Chickpea	0.0	25.0	30.0	40.0
Barley	0.0	25.0	28.0	40.0
Sugarcane	5.0	22.5	35.0	40.0

Environmental Matrix (W_{matrix} or `env_typing`)

Panel of Environmental Types



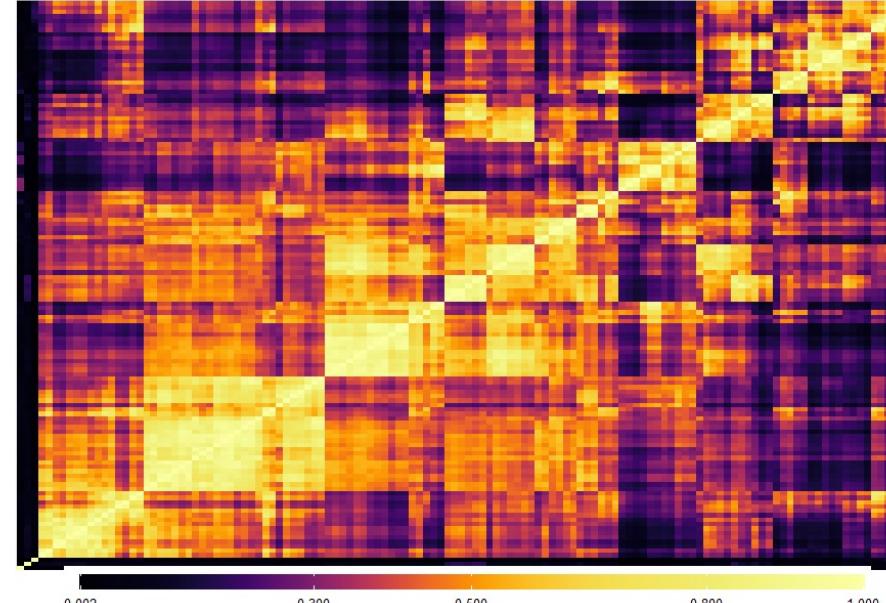
Approaching Similarity
(*in silico*)

Used as Explicit Covariates

Kernelization
(`env_kernel`)

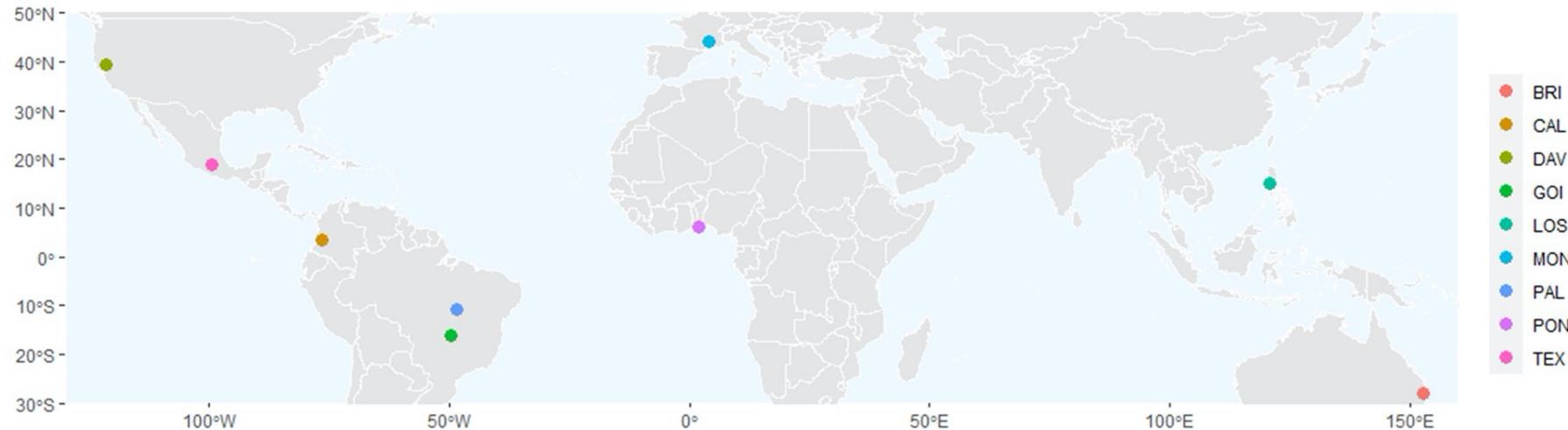
Environmental Relatedness

1. Linear Kernel (same as GBLUP)
2. Nonlinear Gaussian Kernel (GK)
3. Nonlinear Deep Kernel (Arc-cosine approach)

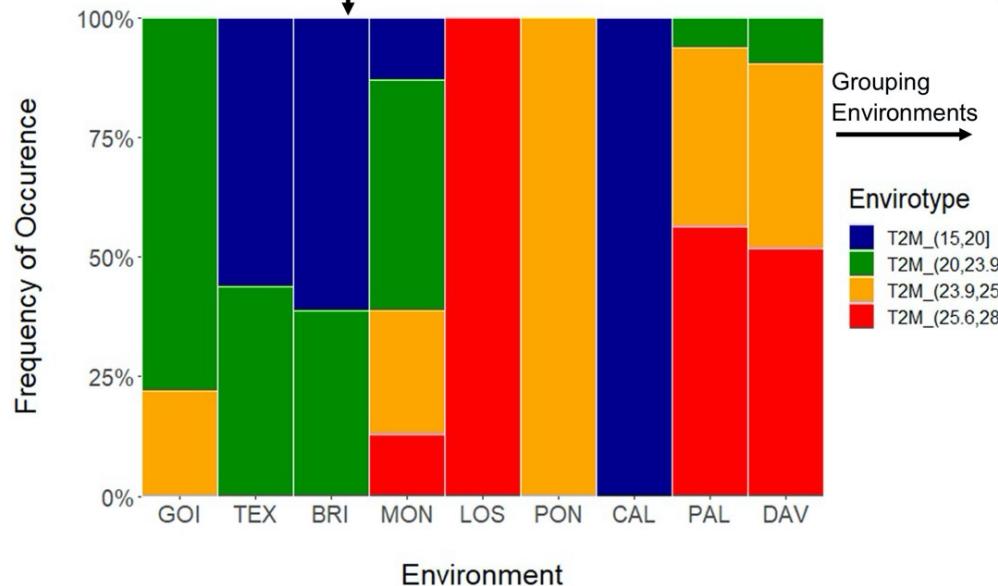


Global Envirotyping Network

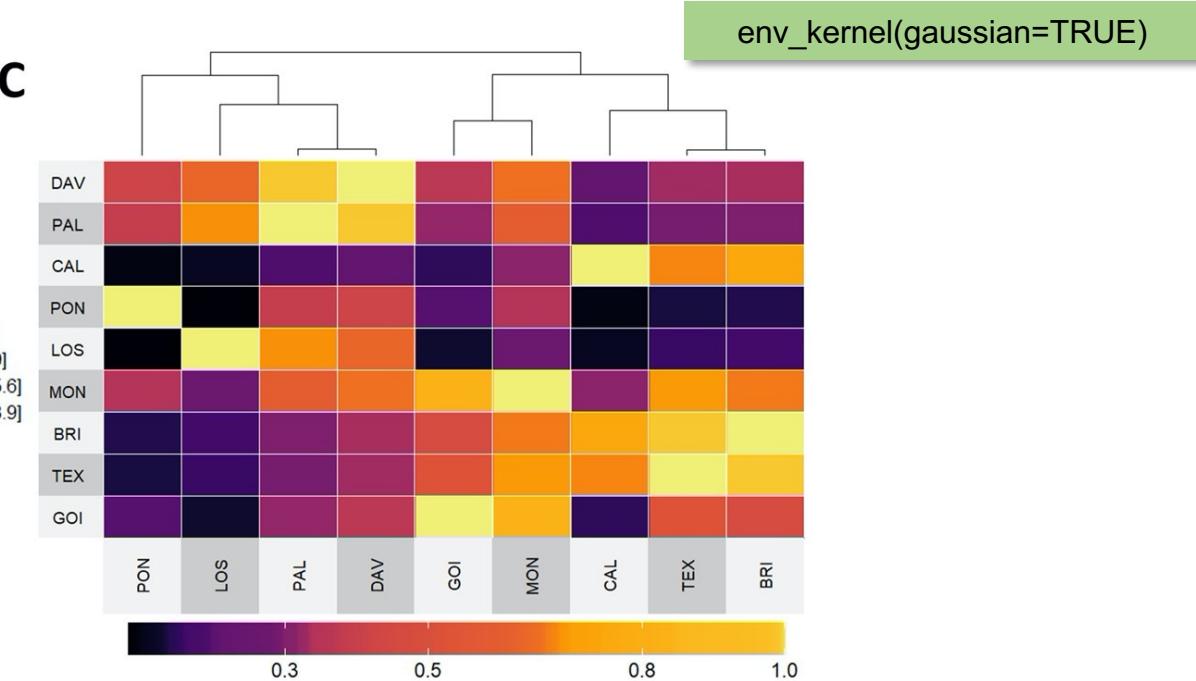
A



B



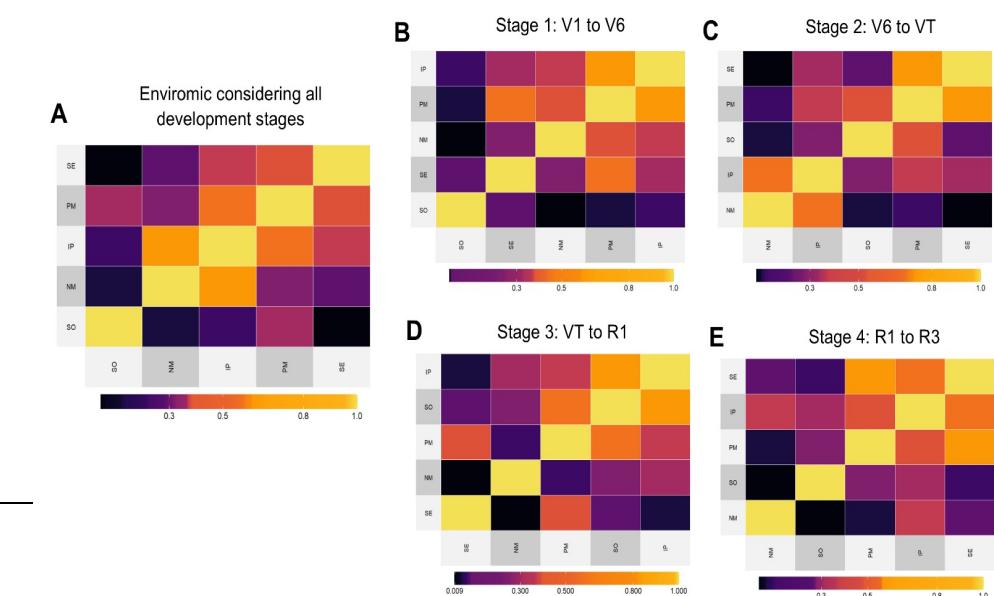
C



“New” reaction-norm approaches

GK+different model structures

Random Effect	Model		
	M1	M2	M3
Genomic (G)	0.426 [0.389;0.470]	0.509 [0.464;0.562]	0.555 [0.506;0.612]
Environment (E)	-	2.686 [2.470;2.505]	-
Stage 1 (S ₁ : V1 to V6)	-	-	3.507 [3.227;3.827]
Stage 2 (S ₂ : V6 to VT)	-	-	2.711 [2.494;2.958]
Stage 3 (S ₃ : VT to R1)	-	-	3.940 [3.626;4.300]
Stage 4 (S ₄ : R1 to R3)	-	-	4.018 [3.697;4.385]
GxE*	0.353 [0.322;0.390]	0.269 [0.246;0.297]	-
GxS ₁	-	-	0.308 [0.269;0.326]
GxS ₂	-	-	0.295 [0.278;0.337]
GxS ₃	-	-	0.306 [0.279;0.336]
GxS ₄	-	-	0.304 [0.280;0.339]
Residual	0.848 [0.773;0.936]	0.269 [0.245;0.296]	0.262 [0.238;0.289]



Model	Prediction Scenario	
	CV1	CV00
M1 (Baseline Genomic × Environment)	0.130 ± 0.047	0.102 ± 0.045
M2 (Benchmark Reaction-Norm)	0.762 ± 0.024	0.485 ± 0.211
M3 (Reaction-Norm for Each Development Stage)	0.760 ± 0.028	0.504 ± 0.194

“New” reaction-norm approaches

GK+specific or joint covariates effects

$$\mathbf{y} = \mathbf{1}\boldsymbol{\mu} + \mathbf{X}_f\boldsymbol{\beta} + \mathbf{g} + \mathbf{gE} + \mathbf{EC} + \mathbf{gEC} + \boldsymbol{\varepsilon}$$

Envirotyping Level	Model	Random Effect			
		Environment (E)	Genotype (G)	G×E	Residual
No Envirotyping	M0	-	0.435	0.329	0.837
Envirotyping by environment	M1	4.117	0.425	0.764	0.849
	(EC1 = FRUE)	[3.789;4.493]	[0.387;0.468]	[0.696;0.843]	[0.773;0.936]
	M2	3.440	0.384	0.786	0.726
	(EC 2 = PETP)	[3.165;3.754]	[0.350;0.423]	[0.716;0.867]	[0.662;0.801]
	M3	4.279	0.497	0.664	0.456
	(EC3=FRUE+PETP)	[3.938;4.670]	[0.453;0.548]	[0.605;0.733]	[0.416;0.503]
Envirotyping by development stage at each environment	M4	8.802	0.522	0.484	0.266
	(EC4 = FRUE)	[8.099;9.605]	[0.476;0.576]	[0.441;0.534]	[0.243;0.294]
	M5	3.514	0.548	0.425	0.267
	(EC5 = PETP)	[3.233;3.835]	[0.500;0.604]	[0.388;0.469]	[0.243;0.295]
	M6	1.595	0.514	0.464	0.262
	(EC6 = FRUE+PETP)	[1.468;1.740]	[0.468;0.566]	[0.423;0.512]	[0.238;0.289]

Final remarks

"[...] Enviromics is like a chocolate box: you never know what you gonna get"

