



BIG data course: Dec 2018

Rabat, Morocco

Z.kehel@cigar.org

Data
Mining

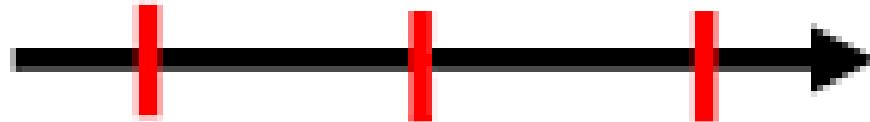






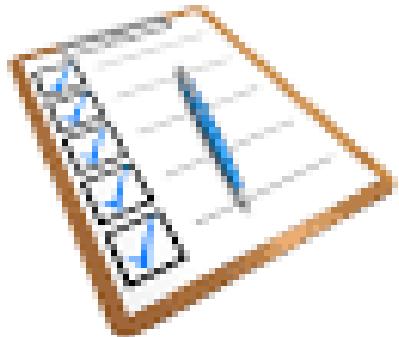
Where You
Are

Gap Analysis



Where You
Need to Be

*Skills Needed to
Get Where You're going*



Action Plan



Just to clarify!

Big data is a collection of data from ***traditional*** and ***digital*** sources ***inside*** and ***outside*** your company that represents a source for ***ongoing discovery and analysis***.

In defining big data, it's also important to understand the ***mix of unstructured and multi-structured data***.

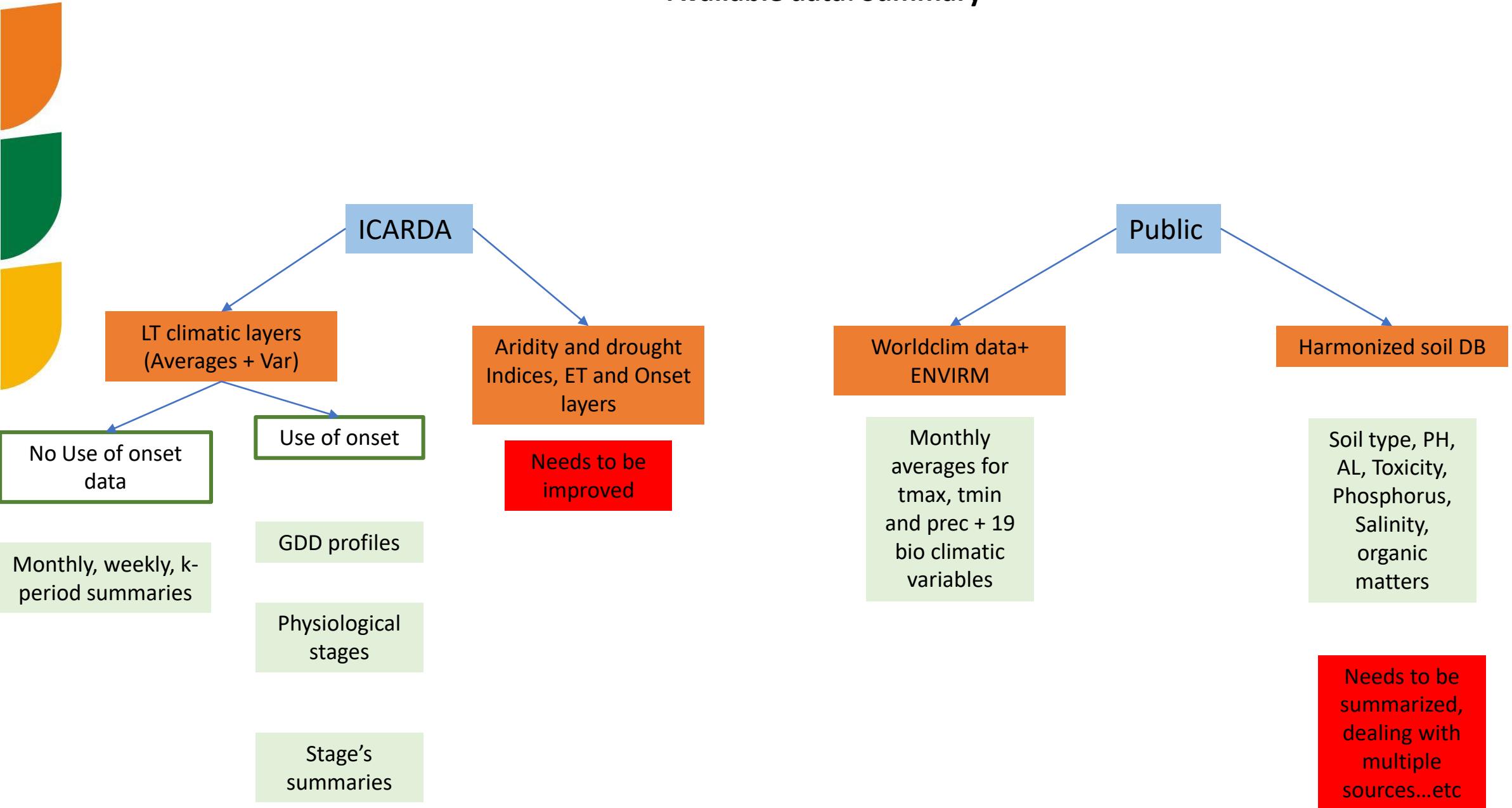
Every organization needs to ***fully understand big data*** (depending on the needs)– what it is to them, what it does for them, what it means to them –and the potential of ***data-driven decision making***. Don't wait. Waiting will only delay the inevitable and make it even more difficult to unravel the confusion.



Data Cycle at GRS (Passport, Eval, Charac)

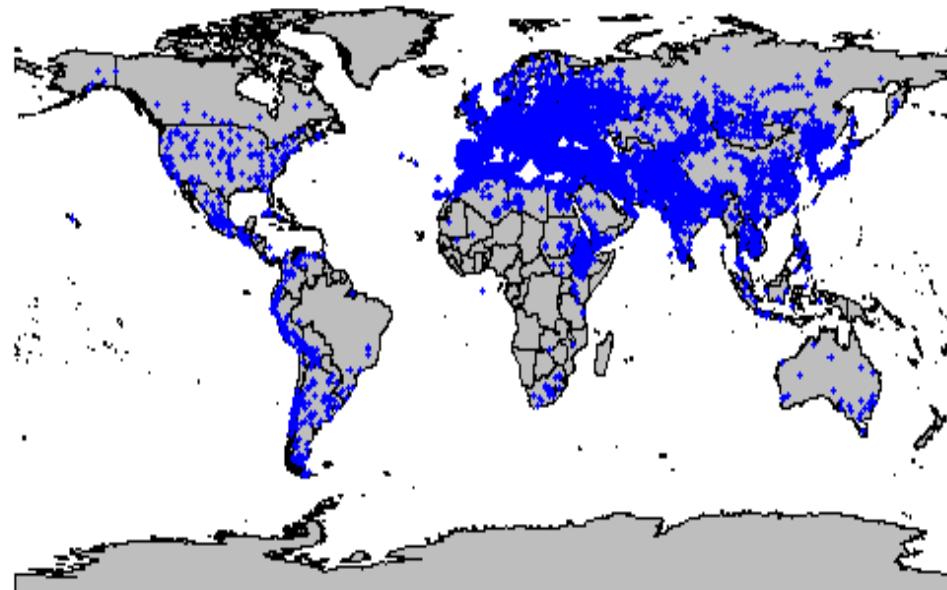
- Passport data and collecting information
- SMTA
- Characterization and evaluation
- Molecular data
- Reporting

Available data: Summary



Available data

- Passport and Daily long-term climatic surfaces; averages and variations over 30 years (GU)



- a. temp---temperature at 2 meter (degC)
- b. tmax---maximum temperature at 2 meter (degC)
- c. tmin---minimum temperature at 2 meter (degC)
- d. precip---precipitation (mm)
- e. ABSH---absolute humidity (kg/m^{**3} scaled by 10^{**6})
- f. RHY---relative humidity (%)
- g. PAR---photosynthesis active radiation ($\text{mol PPFD}=2.05*\text{MJ}$)
- h. uwind---wind at east-west direction(m/s)
- i. vwind---wind at north-south direction(m/s)
- j. VPD---vapor pressure deficit(Pa)

Available data

- Phenotypes + Markers

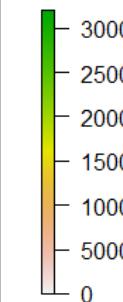
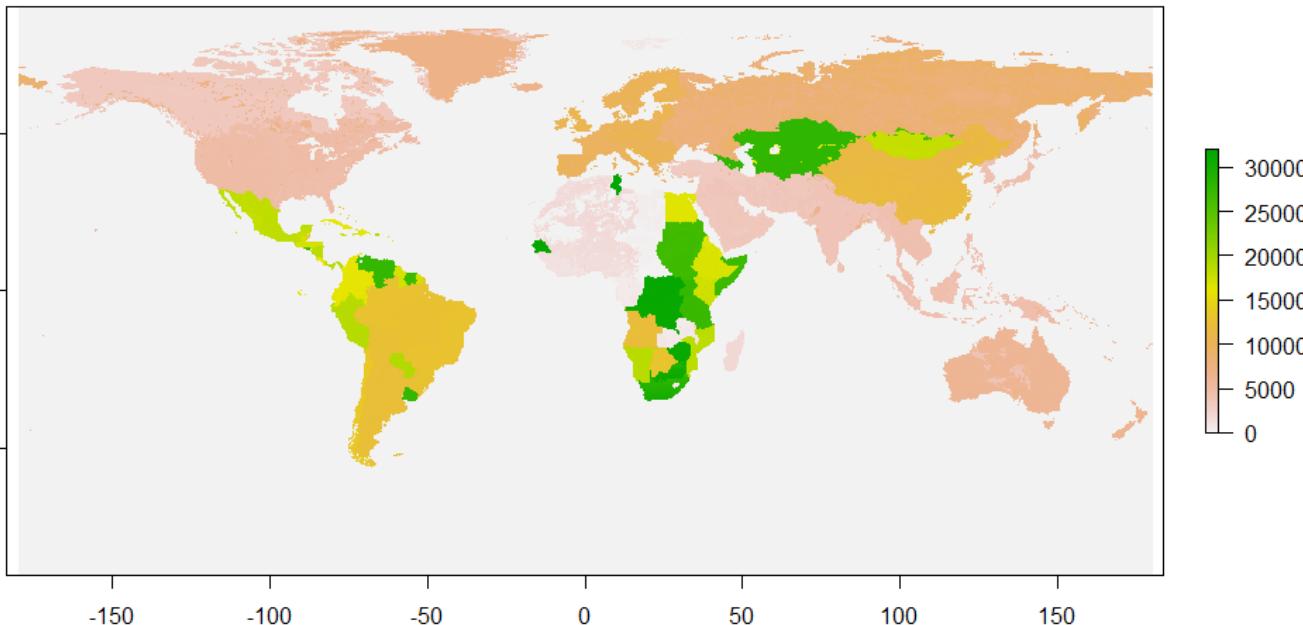
Morphology traits	Trait Type	No. of experiments	No. of accessions	No. of unique accessions
Growth habit	Categorical	79	73,912	29,699
Awnness	Categorical	72	74,512	33,276
Waxiness of plant	Categorical	73	74,417	33,305
Spike length	Quantitative	69	63,157	21,865
Spike density	Categorical	76	86,274	33,172
Seed color	Categorical	12	21,626	10,729
Stem solidness	Categorical	54	48,447	21,432
Number of spikelets groups per spike	Quantitative	91	94,256	33,745

Agronomy traits	Trait Type	No. of experiments	No. of accessions	No. of unique accessions
Early growth vigour	Categorical	11	19,089	9,014
Productive tillering capacity	Categorical	63	72,625	33,295
Plant height	Quantitative	81	88,094	33,379
Number of kernels per spike	Quantitative	70	82,007	31,191
1000 kernel weight	Quantitative	75	83,041	31,756
Seed protein content	Quantitative	3	22,118	11,028
Vitreousness	Quantitative	1	9,325	9,325
Grain yield per plot	Quantitative	57	53,663	19,382
Lodging resistance	Categorical	68	82,049	31,755
Agronomic score	Categorical	12	19,679	8,720

Phenology traits	Trait Type	No. of experiments	No. of accessions	No. of unique accessions
Days to heading	Quantitative	87	93,275	33,395
Days to maturity	Quantitative	86	89,031	33,339
Grain filling period	Quantitative	61	85,886	33,297
Stresses traits				
Yellow rust severity of infection	Categorical	11	15,529	7,685
Common bunt resistance	Categorical	1	13,984	7,685
Septoria resistance	Categorical	1	8,176	8,176
Yellow rust severity and infection type		26	19,288	11286
Yellow rust coefficient of infection	Quantitative	14	9,580	5,319
Yellow rust infection type	Categorical	24	15,453	8,360
Cold tolerance	Categorical	12	15,440	7,522

Available data

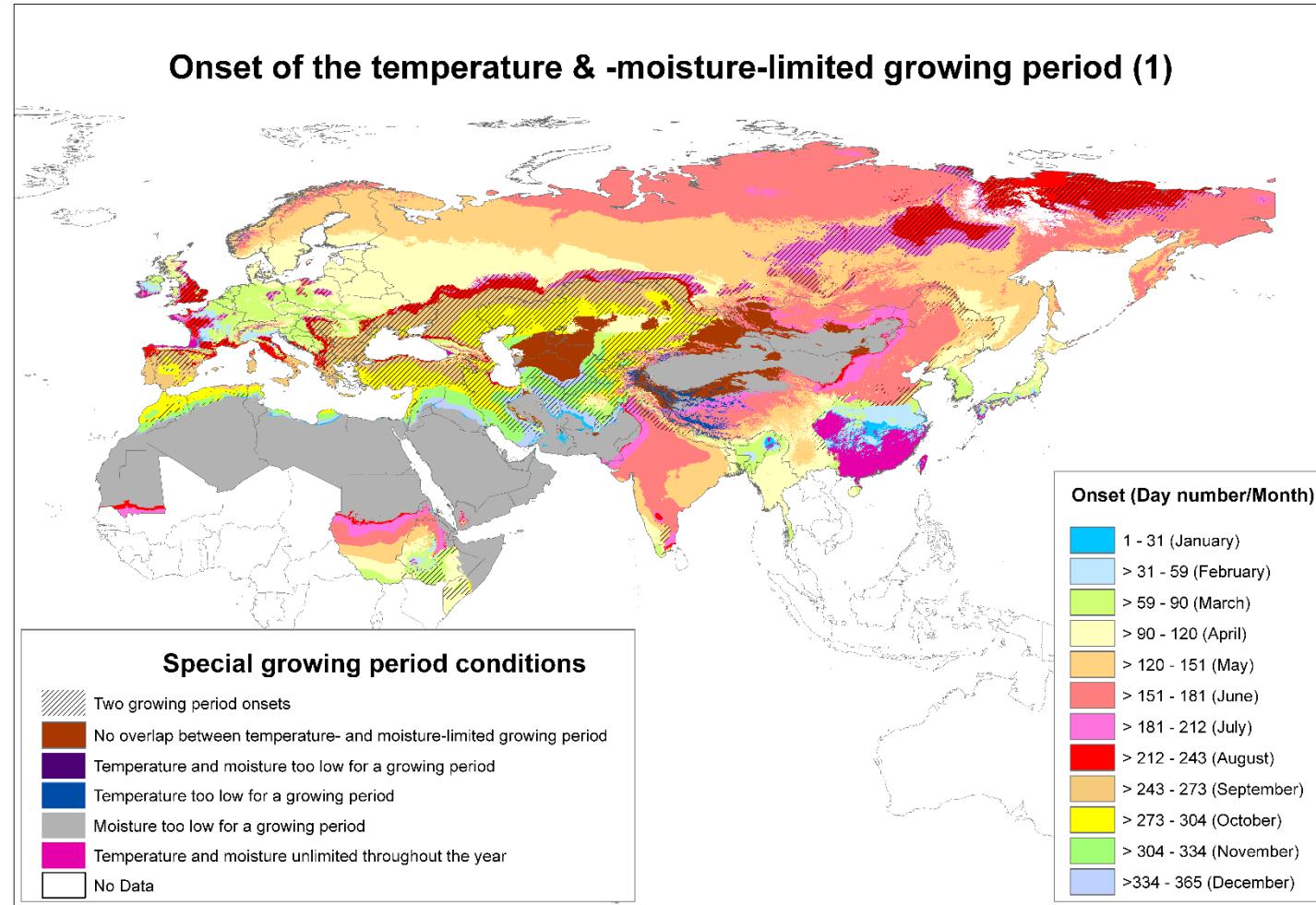
- Maps of soil characteristics from the harmonized world soil database HWSD



Variable	UNIT	DESCRIPTION
T_GRAVEL	%	Topsoil Gravel Content
T_SAND	%	Topsoil Sand Fraction
T_SILT	%	Topsoil Silt Fraction
T_CLAY	%	Topsoil Clay Fraction
T_USDA_TEX_CLASS		Topsoil USDA Texture Classification
T_REF_BULK_DENSITY	kg/dm3	Topsoil Reference Bulk Density
T_BULK_DENSITY	kg/dm3	Topsoil Bulk Density
T_OC	% weight	Topsoil Organic Carbon
T_PH_H2O		Topsoil pH (H2O)
T_CEC_CLAY	cmol/kg	Topsoil CEC (clay)
T_CEC_SOIL	cmol/kg	Topsoil CEC (soil)
T_BS	%	Topsoil Base Saturation
T_TEB	cmol/kg	Topsoil TEB
T_CACO3	% weight	Topsoil Calcium Carbonate
T_CASO4	% weight	Topsoil Gypsum
T_ESP	%	Topsoil Sodicity (ESP)
T_ECE	dS/m	Topsoil Salinity (ECE)
S_GRAVEL	%	Subsoil Gravel Content
S_SAND	%	Subsoil Sand Fraction
S_SILT	%	Subsoil Silt Fraction
S_CLAY	%	Subsoil Clay Fraction
S_USDA_TEX_CLASS		Subsoil USDA Texture Classification
S_REF_BULK_DENSITY	kg/dm3	Subsoil Reference Bulk Density
S_BULK_DENSITY	kg/dm3	Subsoil Bulk Density
S_OC	% weight	Subsoil Organic Carbon
S_PH_H2O		Subsoil pH (H2O)
S_CEC_CLAY	cmol/kg	Subsoil CEC (clay)
S_CEC_SOIL	cmol/kg	Subsoil CEC (soil)
S_BS	%	Subsoil Base Saturation
S_TEB	cmol/kg	Subsoil TEB
S_CACO3	% weight	Subsoil Calcium Carbonate
S_CASO4	% weight	Subsoil Gypsum
S_ESP	%	Subsoil Sodicity (ESP)
S_ECE	dS/m	Subsoil Salinity (ECE)

Available data

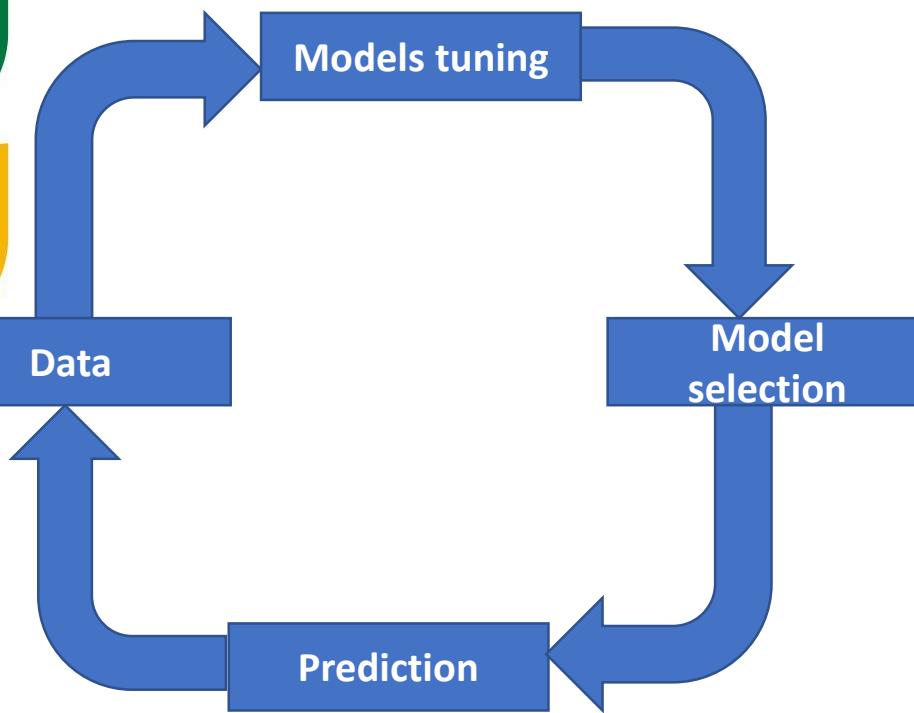
- Onset surfaces for ICARDA mandate crops, aridity and drought indices



Predictive characterization at ICARDA using FIGS



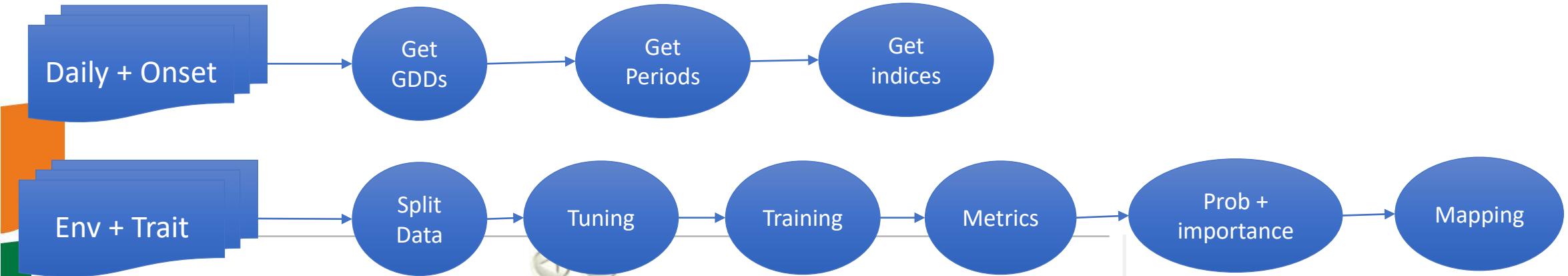
Durum wheat Grain filling period example



FIGS modeling pathway

Performance Measures	k-Nearest Neighbours	Random Forest	Support Vector Machine
Accuracy	0.834	0.838	0.817
95% CI	(0.799, 0.865)	(0.804, 0.868)	(0.781, 0.849)
No Information Rate	0.762	0.762	0.762
P-Value [Acc > NIR]	3.58E-05	1.37E-05	0.001423371
Kappa	0.563	0.557	0.467
Sensitivity	0.722	0.675	0.54
Specificity	0.869	0.889	0.903
True Positive	91	85	68
True Negative	351	359	365
False Positive	53	45	39
False Negative	35	41	58

High accuracy showing that there is a strong relationship between GFP and longterm climatic conditions



Documentation for package 'icardaFIGS' version 0.1.0

- [DESCRIPTION file](#).

Help Pages

durumDaily	DATASET_TITLE
durumWC	DATASET_TITLE
extractWCdata	Extracting WorldClim Data
getAccessions	Getting Accession Data from ICARDA's Genebank Documentation System
getCrops	Getting List of Crops Available in ICARDA's Genebank Documentation System
getDaily	Extracting Daily Values of Climatic Variables from ICARDA Data
getGrowthPeriod	Calculating Growing Degree Days and Lengths of Growth Stages for Various Crops Using Onset Data from ICARDA's Database
getMetrics	Obtaining Performance Measures from Confusion Matrix
getMetricsPCA	Obtaining Performance Measures from Confusion Matrix for algorithms with PCA pre-processing
getOnset	Extracting Daily Values of Climatic Variables from ICARDA Data Based on Onset of Planting
getTraits	Getting Traits Associated with Crops from ICARDA's Genebank Documentation System
getTraitsData	Getting Values of Accessions for Associated Traits
hello	Hello, World!
mapAccessions	Plotting Accessions on a World Map
septoriaDurumWC	DATASET_TITLE
splitData	Splitting Data into Train and Test Sets
tuneTrain	Splitting the Data, Tuning and Training the Data, and Making Predictions
varimpPred	Calculating Variable Importance and Making Predictions

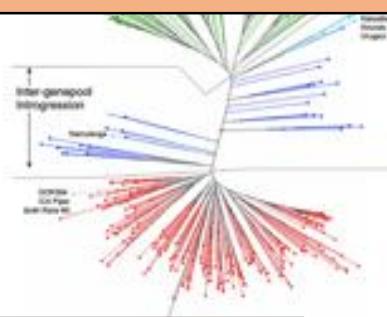
ICARDA-FIGS-R package

Landrace gap spatial analysis methodology

1

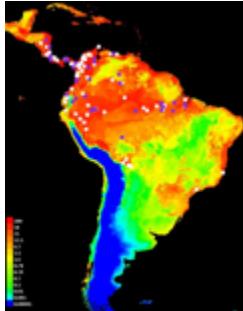
Review literature on landrace genetic structure and its relationship with ecogeography

Knowledge on crop genetic and ecogeographic variation



4

Model geographic distributions of each subgroup



Probability distribution of landraces in given sub-group

2

Gather passport data and other data relevant to genetic structure (e.g. characterization, names, molecular)

Database (passport, characterization, genetic)

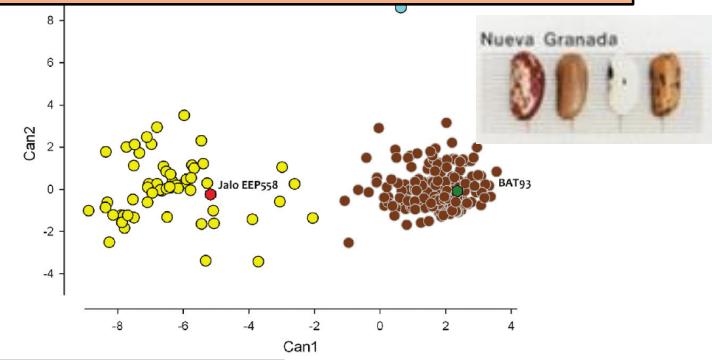


3

Classify dataset into its structural groups:

- 3.1. Define levels & sub-groups
- 3.2 Assess predictability
- 3.3 Classify at chosen level

Classified occurrence dataset



5

Calculate gap score for each subgroup:

- 5.1 Geographic score
- 5.2 Environmental score
- 5.3 Combined gap score

Maps of geographic, environmental and gaps

6

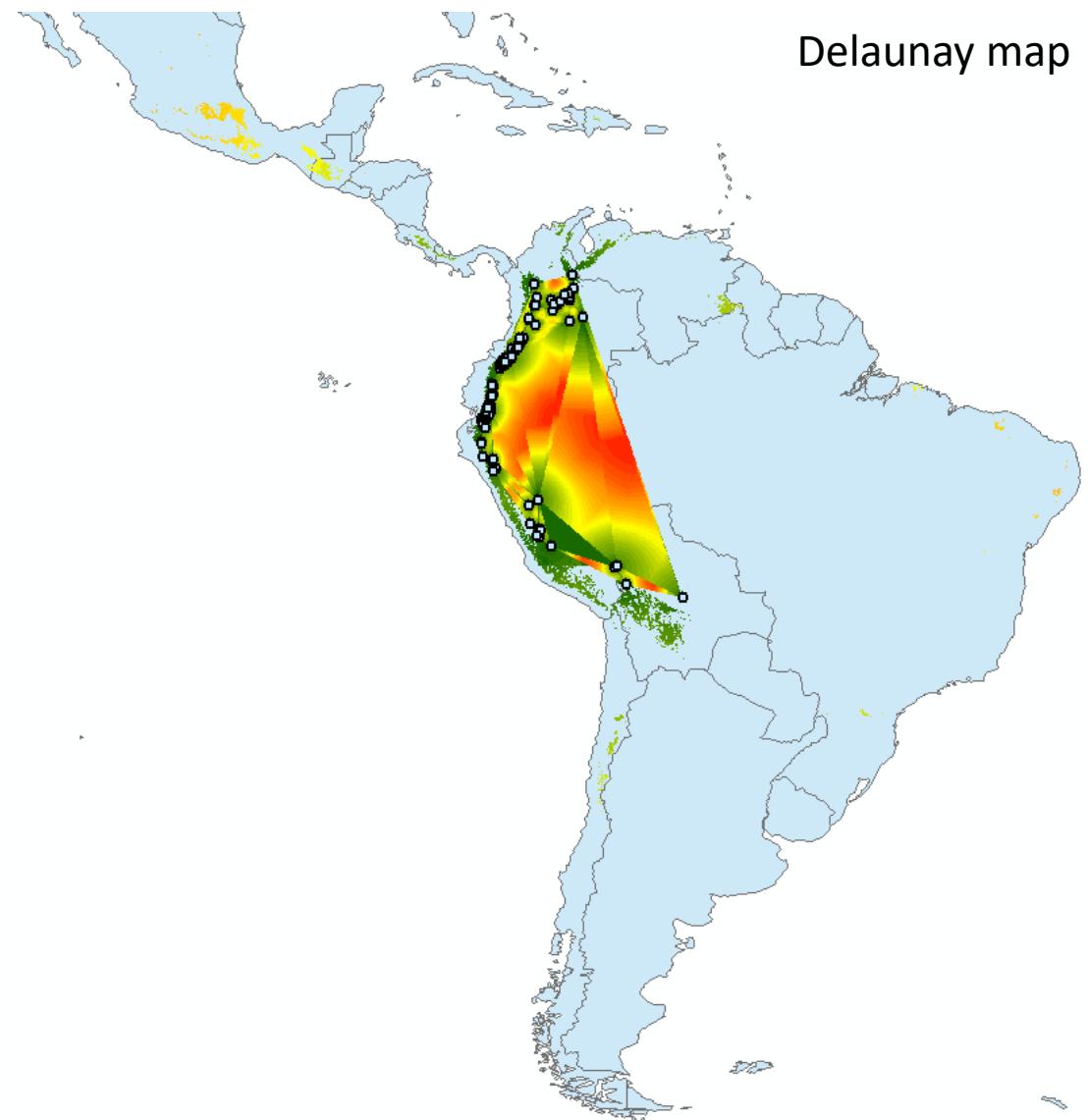
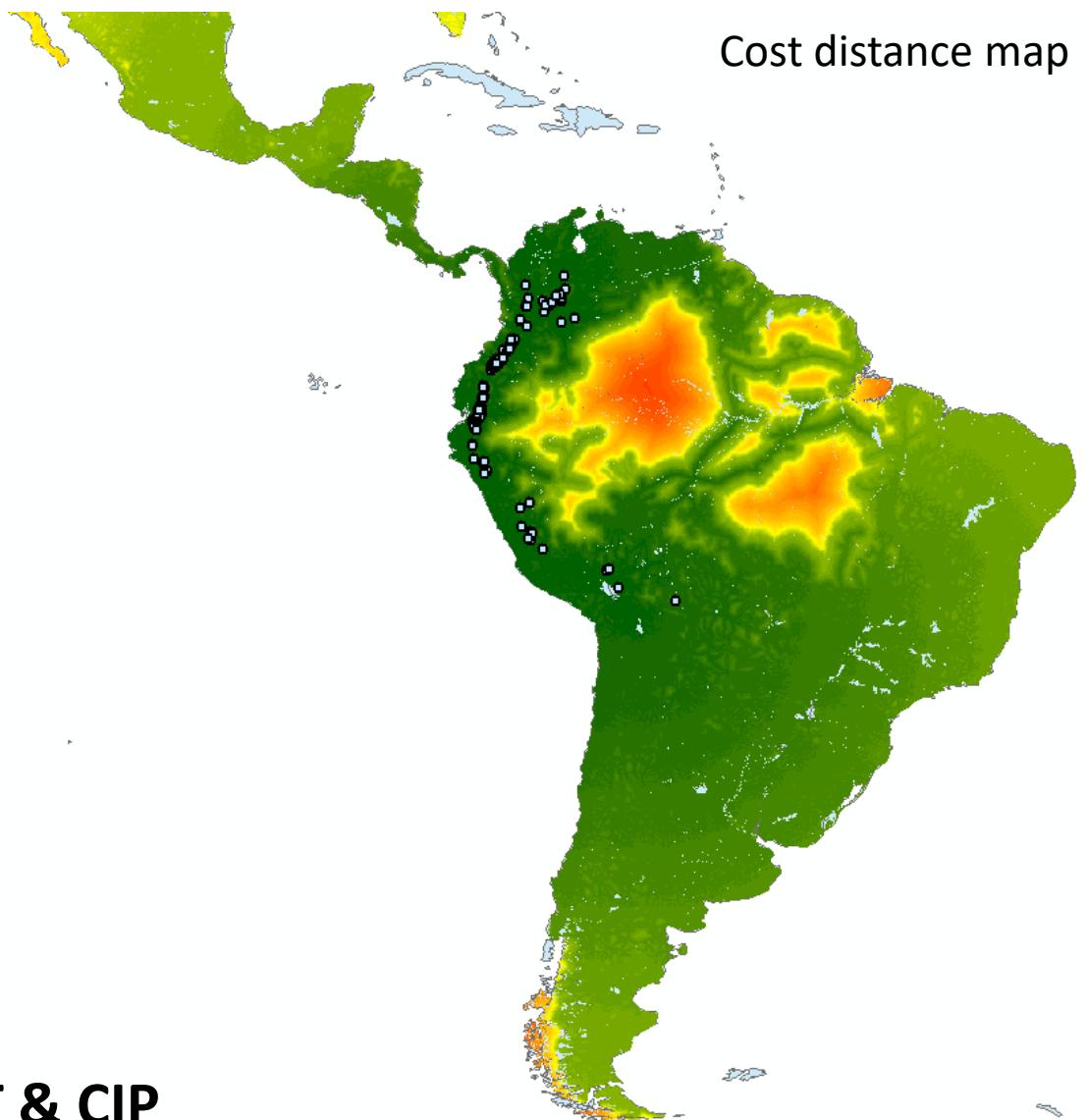
Expert evaluation of gap analysis results

Validated / updated results

LGA *S. phureja*: Species distribution model



LGA *S. phureja*: Geographic scores

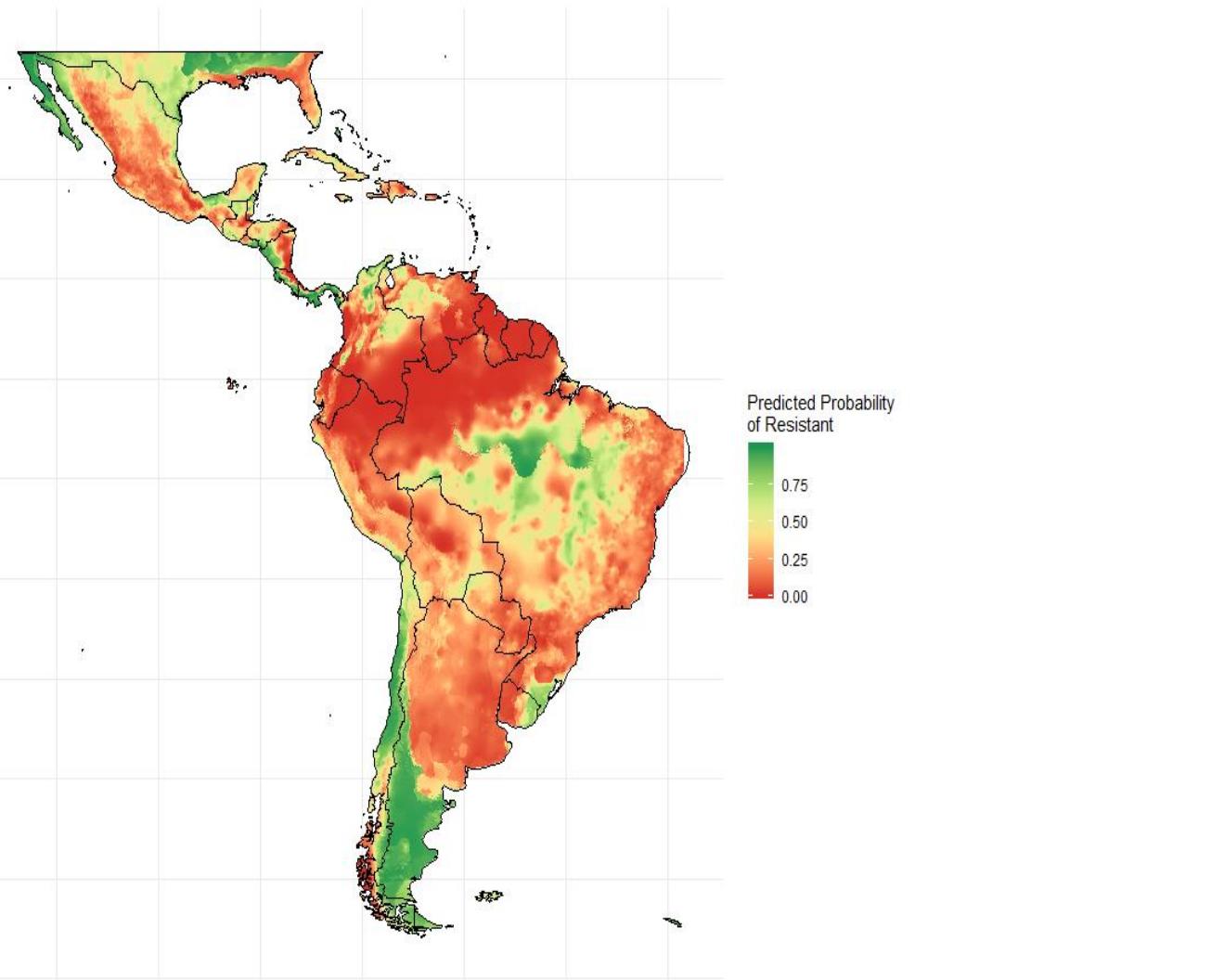


LGA *S. phureja*: Environmental score

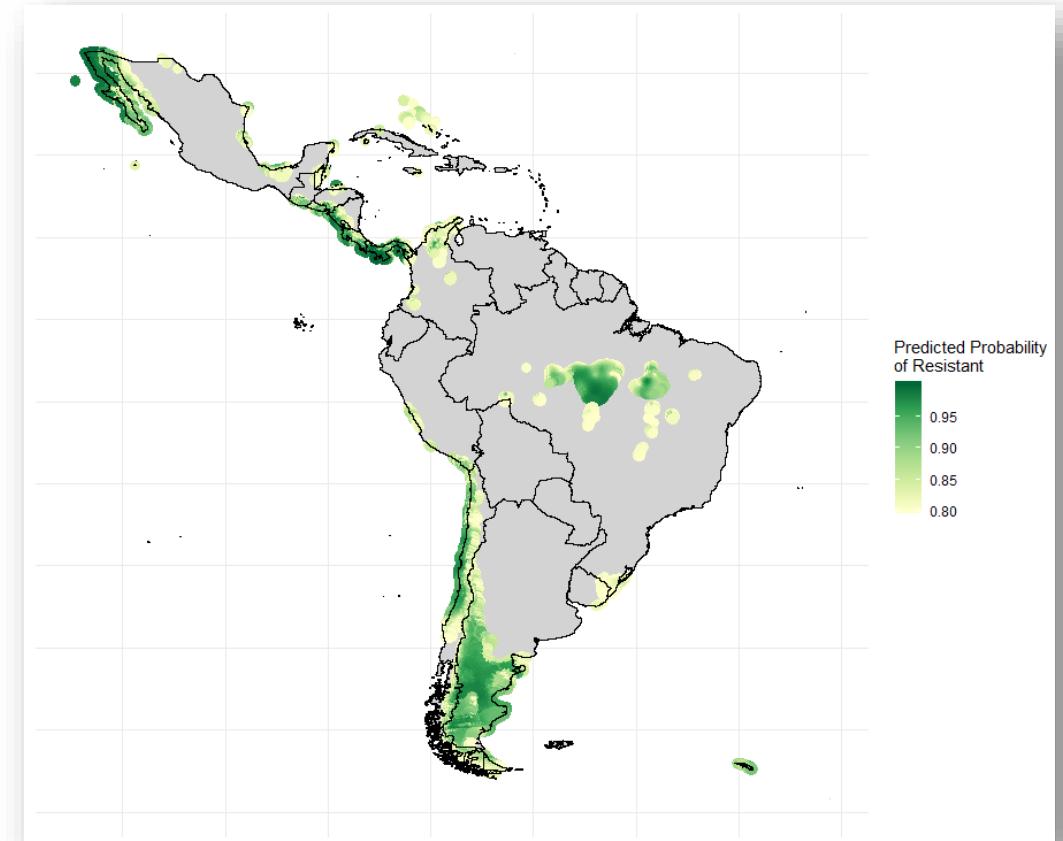


LGA *S. phureja*: Gap score





- Prob ≥ 0.8





Data Cycle at BCI (Pheno)

- Planning: which populations, how many, blocks
- Designing: experimental designs based on need, requirements and resources available
- Measuring: phenotypic evaluation in the field or controlled environments including diseases
- QC: spectrum of data, visualization in the field,
- Analyzing (QC as well): single field and MET analysis (outlier detection is done here not above): An outlier is a data point that has large studentized residuals from the model.
- Reporting: heritabilities, CVs, BLUEs and BLUPs, GxE correlations and latent regressions for stability

B Breeding Management X +

bms.icarda.org:48080/ibpworkbench/main

FABA BEAN BREEDING

PROGRAMS ?

PROGRAMS Click on a program to open it

PROGRAM NAME	CROP	LAUNCH
Winter Wheat	Wheat	▶
Durum Wheat	Wheat	▶
Winter Barley	Barley	▶
Spring Barley	Barley	▶
Grass Pea Breeding	Grasspea	▶
Grass Pea International Nurseries	Grasspea	▶
Spring Bread Wheat	Wheat	▶
Training Breeding Program	Tutorial1	▶
Wheat International Nurseries	Wheat	▶
Faba Bean International Nurseries	Faba	▶
Chickpea International Nurseries	Chickpea	▶

Phenotyping data management: Managing samples, germplasm lists, location lists and trial data (BrAPI)



Breeding Management × +

bms.icarda.org:48080/bpworkbench/main

BREEDING ACTIVITIES

- Manage Germplasm
- Manage Studies**
- Manage Samples

INFORMATION MANAGEMENT

STATISTICAL ANALYSIS

PROGRAM ADMINISTRATION

FABA BEAN BREEDING

MANAGE STUDIES

PYT Dis Low Latt Save Actions

BASIC DETAILS

Settings **Germplasm & Checks** **Environments** **Experimental Design** **Measurements**

Define Measurement Details Add

Name	Description	Input Variables
<input type="checkbox"/> Days_to_flowering	Flowering - count days after sowing (number)	
<input type="checkbox"/> Plant_height	Plant height - soil to tip at maturity (cm)	
<input type="checkbox"/> Days_to_maturity	Maturity - count days after sowing (number)	
<input type="checkbox"/> Grain_yield	Grain yield -dry and weigh (kg/ha)	

Select All Remove

Measurements

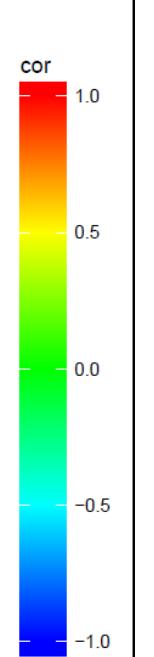
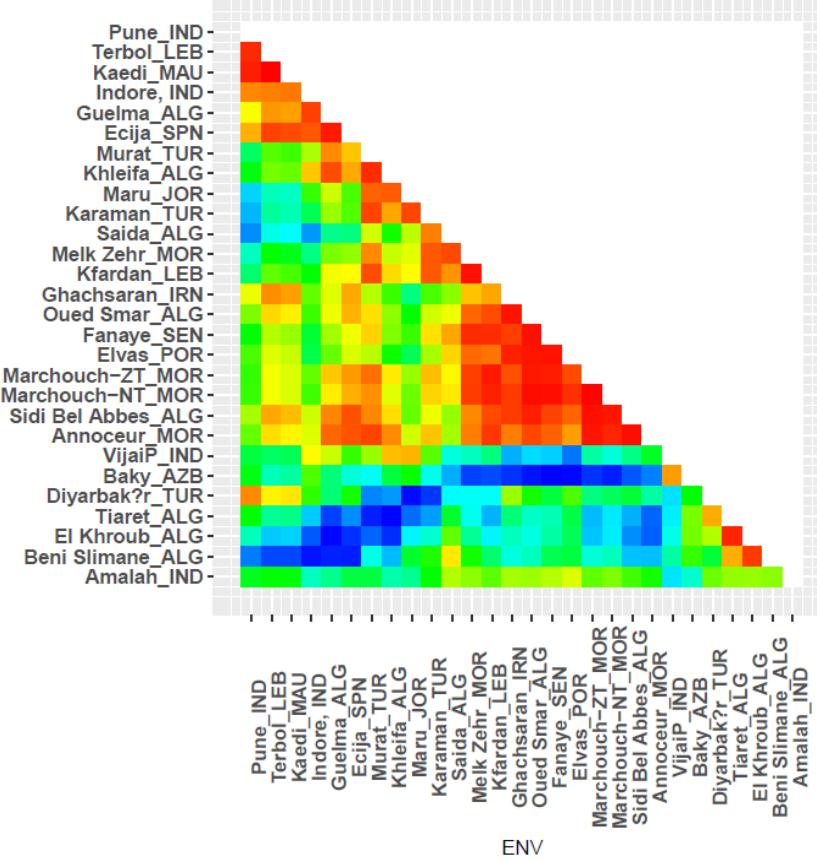
Select Environment: 1 - Unspecified Location Show Categorical Description

Records per page: 100 Showing 1 to 84 of 84 entries

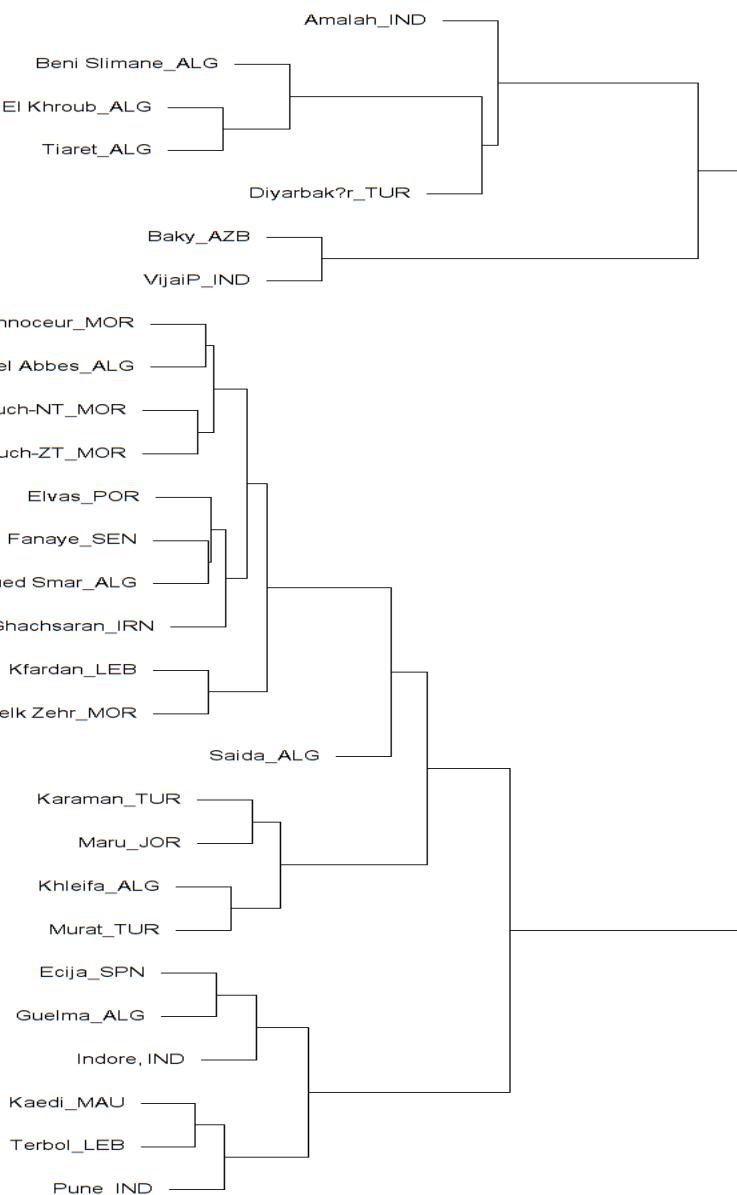
REP_NO	ENTRY_NO	PLOT_NO	GID	ENTRY_TYPE	BLOCK_NO	DESIGNATION	Days_to_flowering	Plant_height	Days_to_maturity	Grain_yield
1	9	21101	30072	Test entry	1	SelKFSH/013/237-1	109	77	168	1823.9
1	16	21102	30079	Test entry	1	SelKFSH/013/317-1	109	79	168	2473.79
1	5	21103	30069	Test entry	1	SelKFSH/013/878-1	112	74	168	2138.36
1	26	21104	30089	Test entry	1	SelKFSH/013/268-1	112	80	168	2430.56
1	1	21105	30065	Test entry	1	SelKFSH/013/680-1	109	93	168	2935.01
1	24	21106	30087	Test entry	1	SelKFSH/013/47-4	103	78	166	2243.19
1	21	21107	30085	Test entry	1	SelKFSH/013/988-1	106	75	166	2557.65
1	2	21108	40971	Test entry	2	SelKFSH/013/758-1	103	74	166	2222.22
1	19	21109	30082	Test entry	2	SelKFSH/013/622-4	112	79	168	2808.3
1	8	21110	40972	Test entry	2	Super Aguadulce	112	78	166	2809.22
1	14	21111	30077	Test entry	2	SelKFSH/013/203-2	112	73	168	1474.15
1	15	21112	30075	Test entry	2	SelKFSH/013/175-2	112	76	168	2002.14



Location genetic correlation matrix for 40

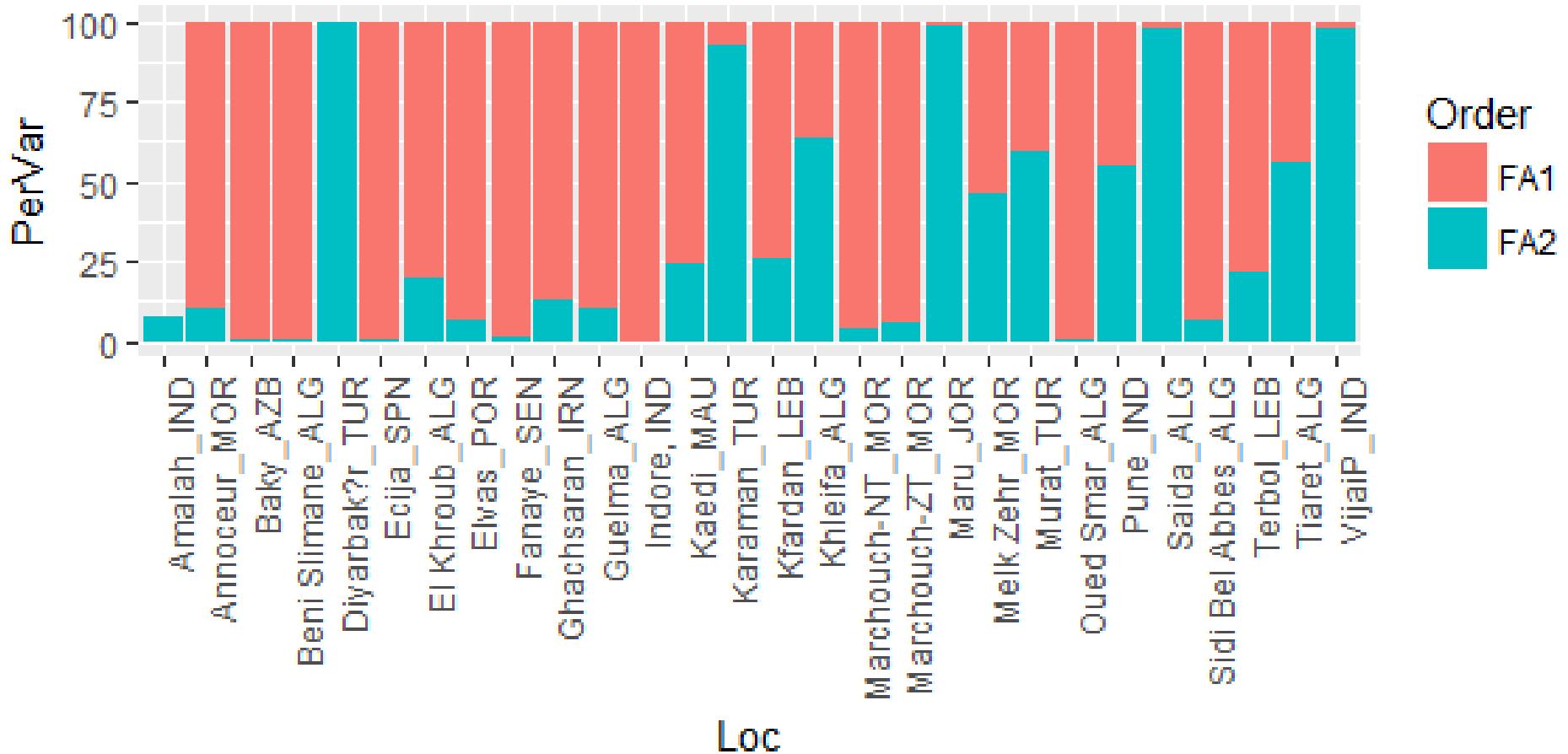


Agglomerative Coefficient = 0.98

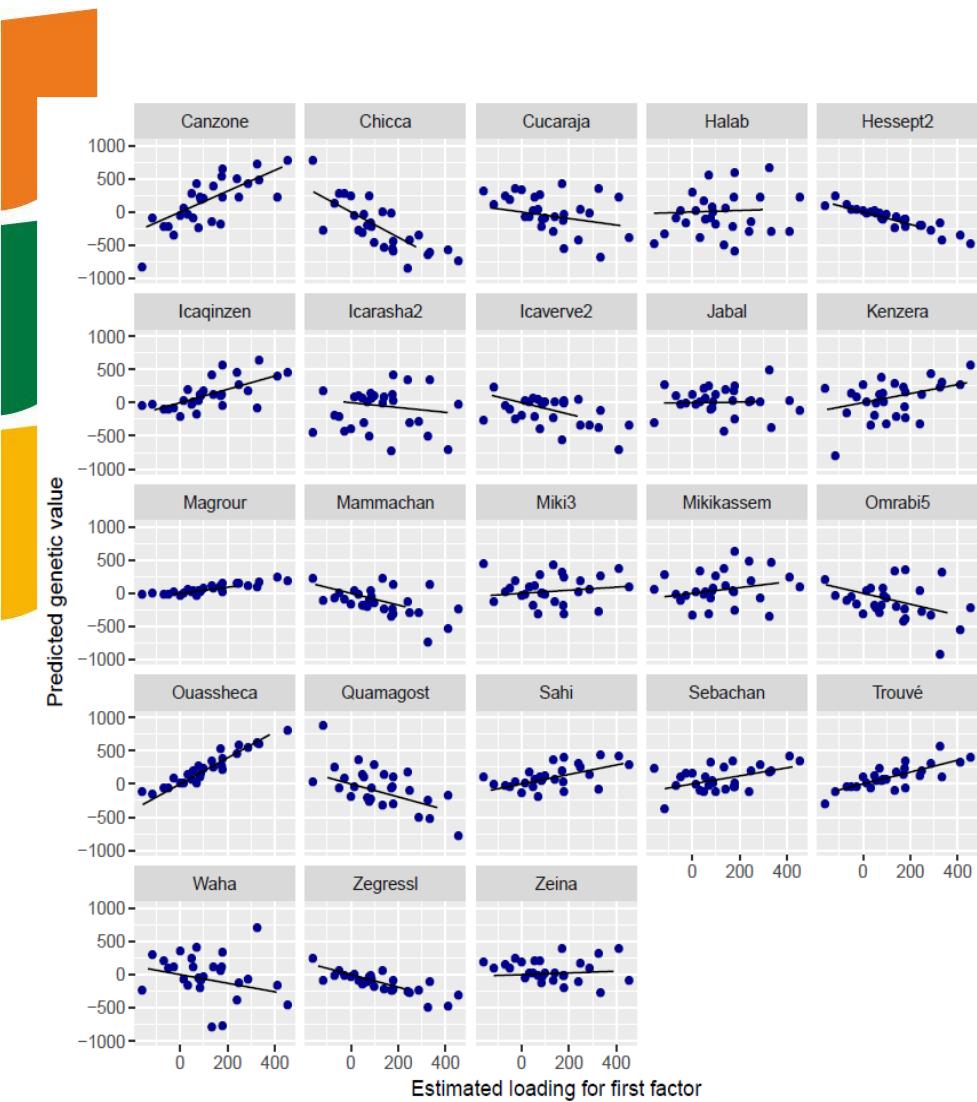


Cluster based on genetic correlations

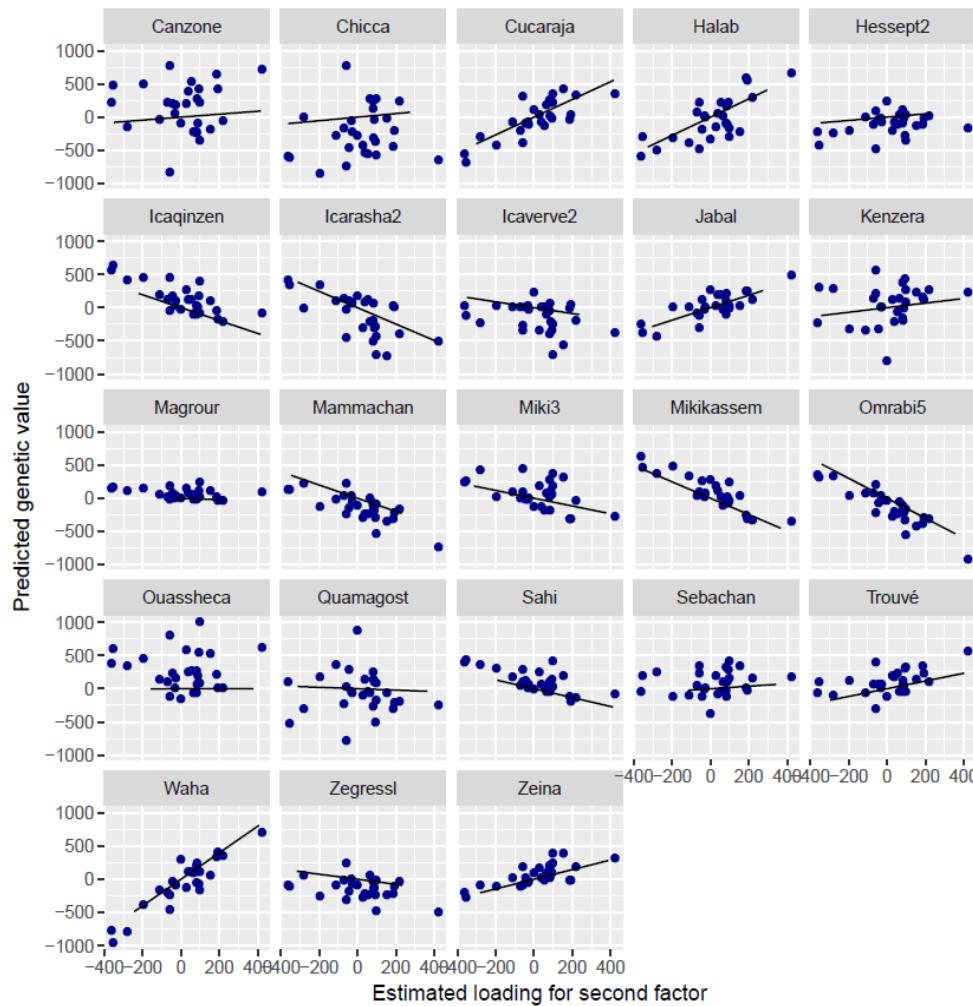
Variance explained by the two first FA model for 40



Latent variable to factor 1 for 40



Latent variable to factor 2 for 40

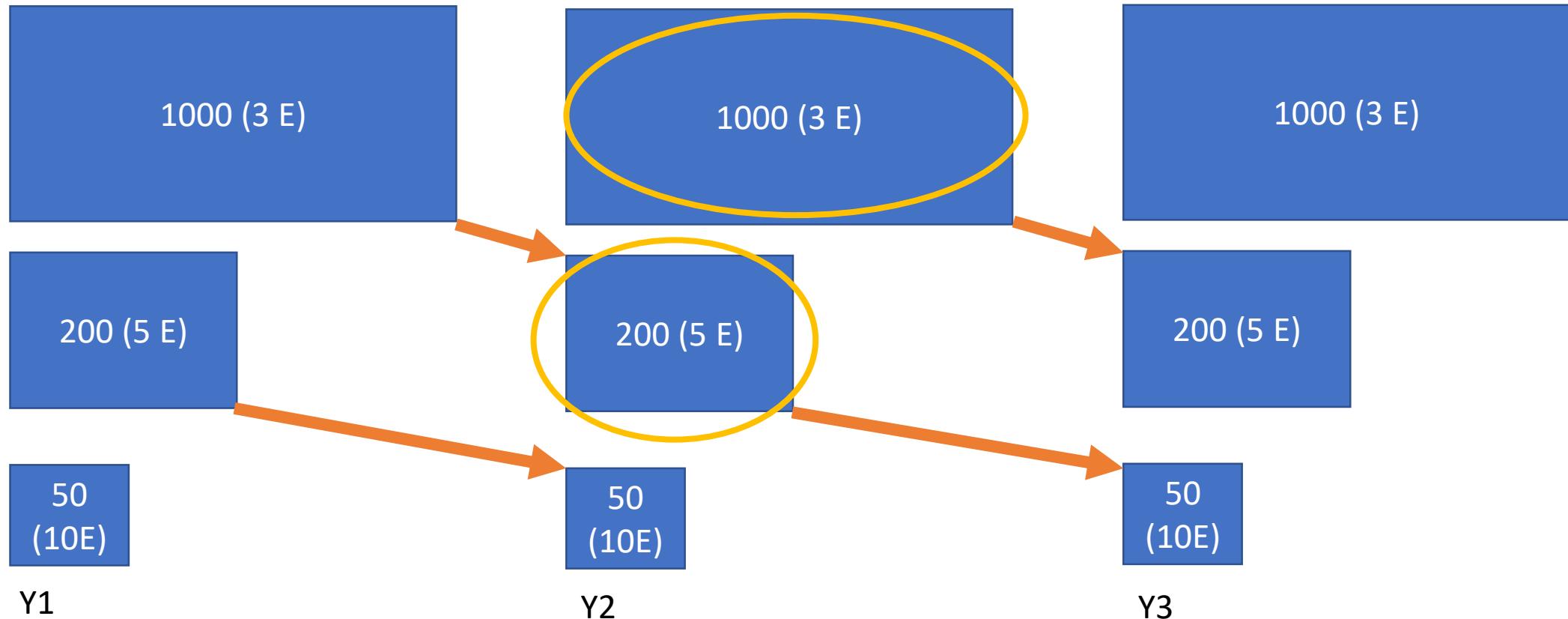


Rotated loading for envs

Env	fac_1	fac_2
Beni Slimane_ALG	-162.03	-58.46
VijaiP_IND	-118.01	-1.35
Baky_AZB	-69.98	81.86
Tiaret_ALG	-50.71	65.40
El Khroub_ALG	-28.10	97.20
Diyarbakr_TUR	-0.39	219.21
Kfardan_LEB	14.23	-31.52
Maru_JOR	31.97	-111.88
Indore_IND	48.41	82.49
Sidi Bel Abbes_ALG	54.48	87.70
Pune_IND	68.77	192.00
Amalah_IND	77.38	80.73
Fanaye_SEN	83.82	-71.75
Annoeur_MOR	86.22	-27.82
Khleifa_ALG	98.08	-44.16
Saida_ALG	133.14	-281.74
Guelma_ALG	139.93	37.35
Elvas_POR	169.94	152.45
Ecija_SPN	175.39	54.49
Kaedi_MAU	177.63	185.97
Karaman_TUR	178.81	-363.96
Murat_TUR	238.73	-196.15
Marchouch-ZT_MOR	248.02	26.75
Ghachsaran_IRN	285.80	93.59
Terbol_LEB	325.08	422.44
Melk Zehr_MOR	332.14	-354.42
Marchouch-NT_MOR	410.56	98.24
Oued Smar_ALG	453.77	-59.19

Opportunities from Big data

- Looking for patterns across time and region
- Most of scientist don't use data beyond the end of the project/experiment
- There are massive opportunities of enhancing analysis using historical data





Data Cycle at BCI (Geno)

- Planning: nature of genotyping based on need and objectives of the study (high versus low density)
- Sampling and tracking: get samples to go for genotyping from leaves (we send leaves or DNA). Every sample has an ID that is linked to the trial and GID (Sample tracking is mandatory)
- Genotyping: Send samples with their ID in the genotyping platform format
- Data management: That where the raw data is managed (GIGWA in ICARDA case)
- QC: Can also happen within the data management system
- Analyzing: Diversity, MT association, pop genetics, GS

	101032366_A_G_42	101032387_T_G_25	101032400_G_A_25	101032420_A_G_60	101032426_A_G_13
SEEDICAR10048	AA	GG	GA	AA	AA
SEEDICAR10052	AA	GG	AA	AA	NA
SEEDICAR10187	AA	GG	GG	NA	AA
SEEDICAR10189	AA	GG	GG	AA	AA
SEEDICAR10191	AA	GG	GA	AA	AA
SEEDICAR10197	AA	GG	NA	AA	AA
SEEDICAR10262	AA	GG	NA	AA	AA
SEEDICAR10438	AA	GG	GA	AA	AA

Breeding Management Syst Gigwa v2.0-alpha

bms.icarda.org:8080/gigwa/index.jsp

Gigwa v2.0 Barley Project LI_AM

Search Enable browse and export

1 - 100 / 6446

External tools IGV

Variant types Any

Sequences (8/8) Sequences

Position (bp)

Investigate genotypes on 1 group

Individuals (336/336) Group 1 Individuals

Max missing data 20 %

Minor allele frequency % %

Genotypes Any

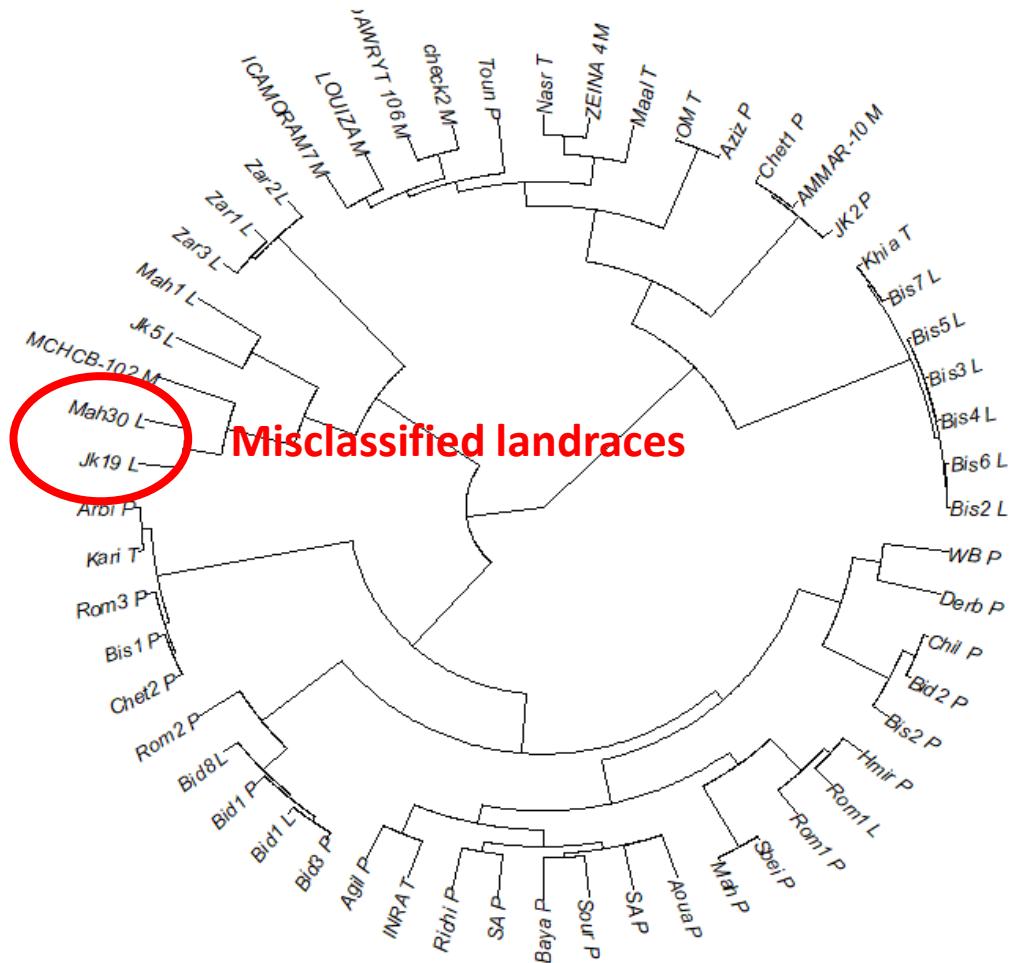
id	sequence	start	stop	alleles
11_21354	1	1	1	A T
12_10420	1	1	1	A T
12_30969	1	1	1	T A
SCRI_RS_120053	1	1	1	T A
SCRI_RS_120059	1	1	1	T A
11_10895	1	2	2	A T
12_30715	1	2	2	T A
11_20502	1	2	2	A T
11_21067	1	2	2	T A
SCRI_RS_113745	1	2	2	T A
12_31144	1	5	5	A T
11_10419	1	6	6	A T
12_10636	1	6	6	T A
12_11011	1	6	6	T A
SCRI_RS_60145	1	6	6	T A
SCRI_RS_60293	1	6	6	T A
SCRI_RS_66630	1	6	6	A T
SCRI_RS_82277	1	6	6	T A
SCRI_RS_194326	1	6	6	T A
SCRI_RS_232577	1	6	6	T A
12_30933	1	7	7	T A
11_21174	1	9	9	T A
11_21226	1	10	10	T A
SCRI_RS_162524	1	10	10	T A
SCRI_RS_184274	1	11	11	A T
SCRI_RS_1929	1	11	11	A T
SCRI_RS_224686	1	11	11	T A
12_30918	1	13	13	A T
12_30919	1	13	13	T A
12_30950	1	13	13	A T
12_30952	1	13	13	A T
SCRI_RS_110312	1	13	13	A T

Genotypic data management: Managing studies, snps, maps, QC and exporting (BrAPI)

Enhancing conservation and use of genebank accessions using molecular markers

Conservation

- Identification of duplicates
 - Identification of genetically close accessions
 - Identification of misclassified accessions at the taxa and name levels
 - Genetic structuring to inform on gaps in collections and collecting



Use

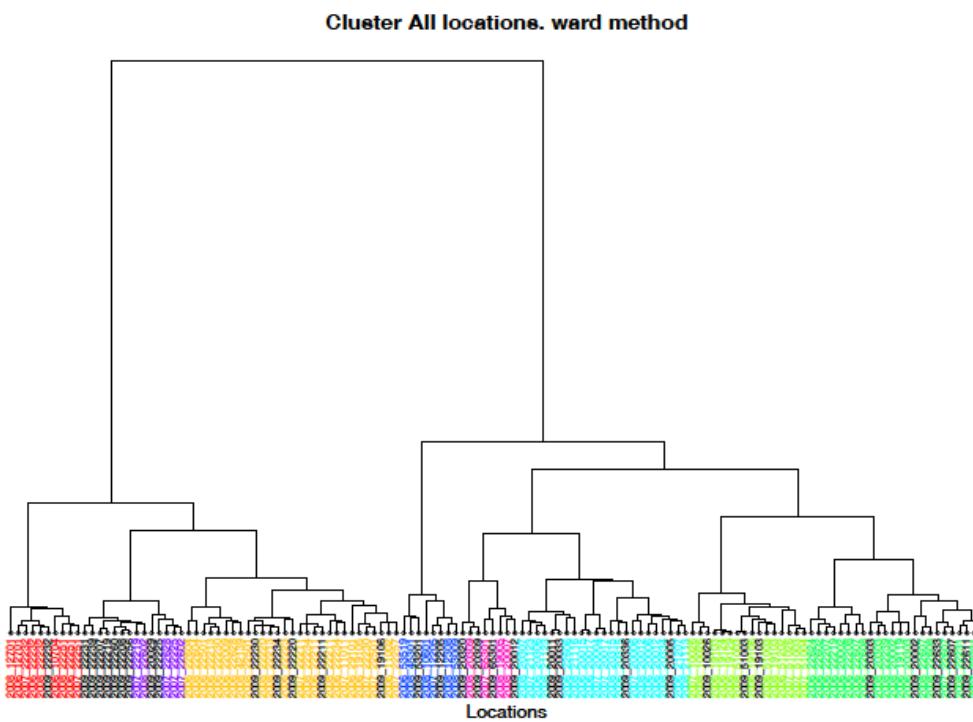
- Increasing FIGS efficiency to identify trait best-sets
 - Linking trait to accessions using genomic prediction
 - Trait-marker association (GWAS)
 - Adaptation mechanisms thru marker-environment association (EWAS)

Predictions of IN using GS

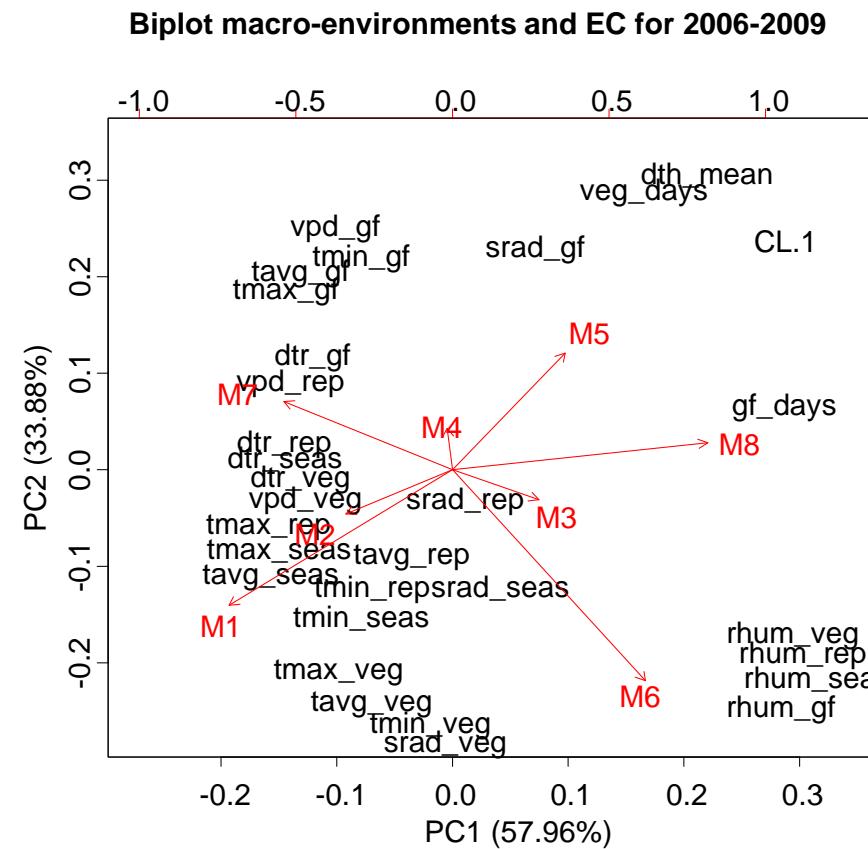
SAWYT - Correlation between the predicted values of 17 models trained with data from three years to predict the four different macro-environments in 2008, 2007, and 2006.

	Model	2008	2007	2006
0	ML	0.1540	0.0693	0.2268
1	MA	0.1608	0.0455	0.2289
2	MAW	0.3541	0.1863	0.3317
3	MAW+AW	0.2953	0.1377	0.3173
4	MAW+AW+MA	0.2327	0.1473	0.3383
5	MAS(W)	0.6140	0.5579	0.4854
6	MAS(W)+ML	0.5918	0.4981	0.5098
7	MAS(W)+MA	0.6495	0.5375	0.4115
8	MG	0.1250	0.0565	0.2230
9	MGW	0.2746	0.1186	0.3282
10	MGW+GW	0.1726	0.1521	0.2937
11	MGW+GW+MG	0.2351	0.0327	0.2937
12	MGS(W)	0.6669	0.5347	0.4123
13	MGS(W)+ML	0.2883	0.5136	0.4485
14	MGS(W)+MG	0.5898	0.5267	0.4396
15	MAGS(W)+ML	0.6343	0.5114	0.4931
16	MAGS(W)+MA+MG	0.6612	0.5261	0.4918

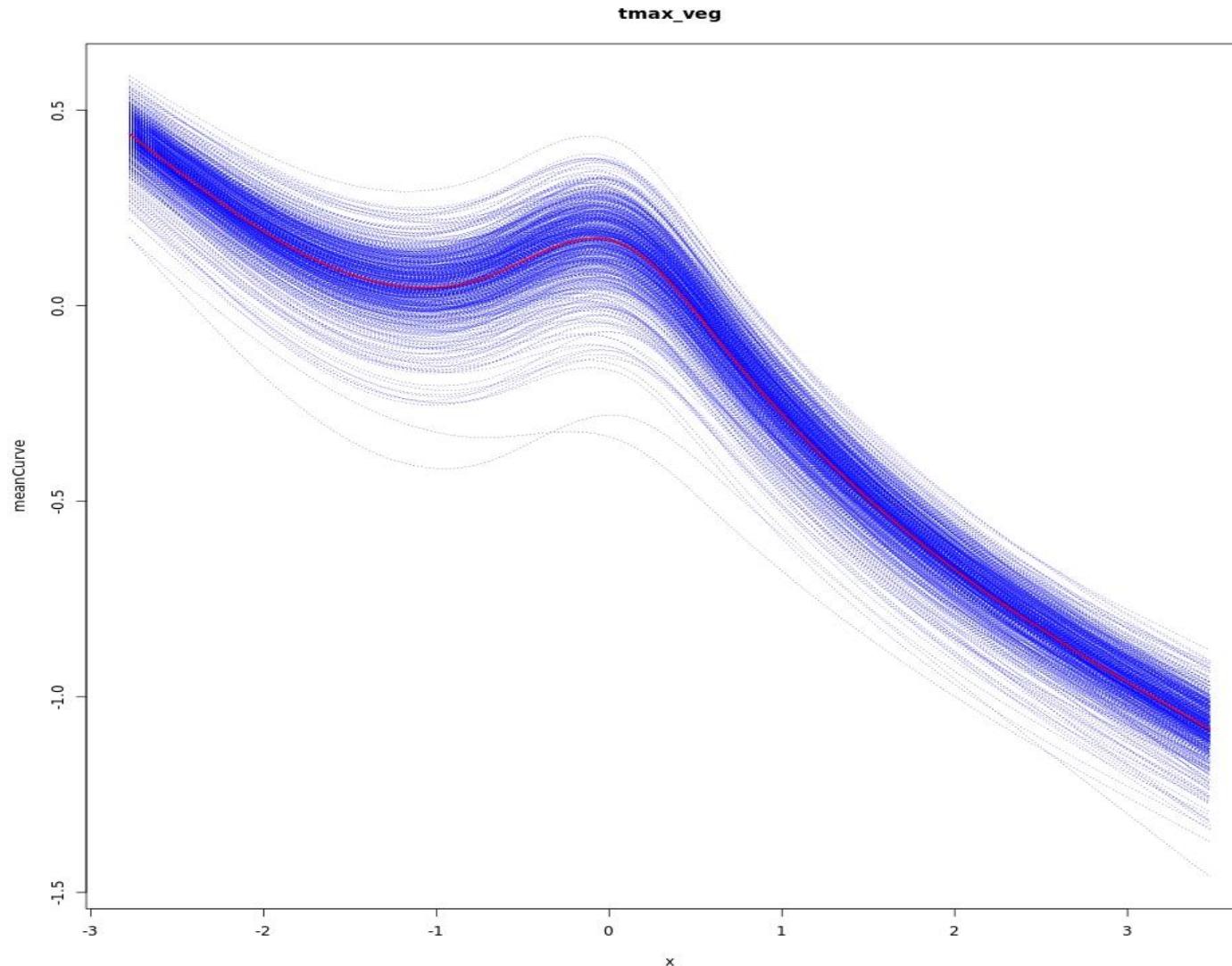
Understanding target environments in IN



8 macro environments



Using historical data to inform breeding strategies using GS





Thank you!

Questions and comments are welcome!