

## Galaxy for Genomics-enabled Breeding

### Star Yanxin Gao

### yg28@cornell.edu





## Introduction





Star Yanxin Gao, Ph.D., PMP

## Application Specialist, IT

- Enterprise Breeding System (2020+)
- GOBii (2015-Present)
- Breeding and Genetics
  - Corn, DAS (2008-2015)
  - Soybean, VT (2004-2008)
  - Vegetables, Cornell (1999-2004)



## Outline



- 1. Project
- 2. Products
- 3. Partnership







## **GS-Galaxy Project (1st** *P***)**

MAS

## **GOBii Mission**

Transform breeding by enabling genomic—assisted selection as **routine** breeding applications





KIRIBAT

## **GS-Galaxy Project (1st** *P***)** People



### **Cornell University**

- Angel Villahoz-Baleta (Former)
- Star Yanxin Gao
- Kelly Robbins
- Liz Jones
- Yaw Nti-addae

### CIMMYT

- Victor Ulat
- Susanne Dresigacker
- Mike Olsen
- Umesh Rosyara
- Xuecai Zhang
- Juan BURGUEÑO
- Fernado Toledo

### Collaborators

Alexis Dereeper

Greenland

RAZIL

- Michael Quinn
- Paulino Perez
- Jose Crossa
- Clay Sneller
- Kate Dreher
- Tom Hagen
- Yoseph Beyene
- Manje Gowda
- Nicholas Santantonio
- Isaak Tecle
- Milcah Kigoni
- Iain Milne
- Gordon Stephen
- Hiro Iwata
- Dave Clements

### IRRI

### • Venice Juanillas

- Ramil Mauleon (Former)
- Ken McNally
- Josh Cobb
- Carlos Ignacio
- Dmytro Chebotarov
- Nick Alexandrov Former)
- Jessica Rutkoski (Former)
- Juan Arbelaez (Former)

### ICRISAT

- Selvanayagam Siva
- Rajeev Varshney

Re Anuterdan

ALS CROTET

- Abhishek Rathor
- Manish Roorkiwal
- Hima Kudapa
- Santosh Deshpande



## **GS-Galaxy Project (1st** *P***)**

### Milestones



### **1st milestone:** Minimum set of GS tools installed in Galaxy- June 2018:

- ➤ GS workflow mapped for each crop and common minimum desirable tools and features identified.
- Pipeline developed with minimum desirable features.
- Basic and common tool components put in place.
- > Tested with well curated test datasets for each crop by product owners and testers

### **2nd milestone:** Production server with published GS workflows-**June 2019**

- Customized workflows for each crop.
- ▶ v1 for fully functioning GS pipeline.
- ➢ Working pipeline available to centers.
- > Training and workshop to users other than product owners and more users start using pipeline and tools

### **3rd milestone:** GOBII integrated as data source for Galaxy-June 2020

- Deliver a complete pipeline with no data manipulation required.
- > Access data extract for phenotypes and genotypes from Galaxy pipeline using remote Galaxy web servers.
- ➢ Found a solution for storing the output data long term so it can easily be accessed.
- Pipelines widely used within CG and outside CG in future.

·	
🔁 Galaxy / EiB-demo	
Tools	
search tools	
Get Data	
Genomic Selection	Geno
Data Format Conversion	<u>BL</u>
Marker Selection	(m
Imputation	Ha
Cluster Analysis	ba
GWAS tools	cri
Text Manipulation	Nu
Collection Operations	en
Filter and Sort	nu
Join, Subtract and Group	<u>Sa</u>
<u>Statistics</u>	sai
<u>Graph/Display Data</u>	an
VCF intersect	Pa
	Ma
HTPG TOOLS	ph
<u>Pre-genotyping</u>	Pa
<u>Post-genotyping</u>	Ma
TOOLS FROM TOOLSHED	ph
CNITDI AV	GE
SHIPLAT	
GOBII TOOLS	Cr
Flapjack Tools	aci

#### Workflows

All workflows

## **GS**-Galaxy **Products** (2nd **P**)

#### omic Selection

UP/BLUE calculator ultiple traits)

pmap filter Filter per SNPs sed on user-specified teria

imeric matrix encoder codes genotypes to meric

mple Matching compares mples between genotype d phenotype files

rallel Genotype-Phenotype atching genotype and enotype files

rallel Genotype-Phenotype atching genotype and enotype files

BV calculator

oss validation within and ross groups

#### **Cluster Analysis**

Calculate Distance Matrix from Genotype Data

SNPRelate GRM calculator

Get Distance Matrix from Genotype Data

Principal Component Analysis

Hybrid K-Means Clustering

K-Means Clustering

#### **GWAS tools**

Calculate Kinship from Genotype Data

Run PCA (Principal Component Analysis) on genetic data

Association using GLM (General Linear Model)

Association using MLM (Mixed Linear Model)

#### **Imputation**

Naive imputation using population mean or mode

Beagle imputation

Impute2 imputation



### http://galaxy-demo.excellenceinbreeding.org/

## **GS**-Galaxy **Products** (2nd **P**)

## Enable routine genomic selection analysis



(multiple traits) <u>Hapmap filter</u> Filter per SNPs based on user-specified criteria

BLUP/BLUE calculator

Genomic Selection

- <u>Numeric matrix encoder</u> encodes genotypes to numeric
- Sample Matching compares samples between genotype and phenotype files

<u>Parallel Genotype-Phenotype</u> <u>Matching</u> genotype and phenotype files

Parallel Genotype-Phenotype Matching genotype and phenotype files



\* 3

✤ 4

**\*** 6

GEBV calculator

<u>Cross validation</u> within and across groups

### Adopted Slide from Clay Sneller

### http://galaxy-demo.excellenceinbreeding.org/



## **GS-Galaxy Products (2nd** *P***<b>)**



## http://galaxy-demo.excellenceinbreeding.org/

Genomic Selection workflows



### Workflow1

**Prediction of genomic breeding values** Well defined training and prediction datasets, prediction within groups.

The prediction is based on the **GEBV** (Genomic Breeding Value Estimator) method implemented in BGLR R package.

**Input**: Genotyping matrix + Phenotyping file **Output**: Table of Predicted values for each trait

Access workflow

**GS Workflow** 

**GS** Workflow

Workflow 1: Predict GEBVs in untested lines Workflow 2: Clustering/population structure Workflow 3: Cross validation Workflow 4: Genome-wide association study

#### Examples of graphical outputs



## **GOBII GS-Galaxy Products (2nd** *P***<b>) GOBII-Galaxy Integration**

6	Galaxy - Mozilla Firefox		
Galaxy X	Siva   GOBii Project X GOBii_to_Galaxy/GOBii_ X +		
$\leftarrow \rightarrow C $	) localhost:8080/?job_id=b3fe17e79980f9fc&identifer=kt0isIno8r ***		II\
🔁 Galaxy	Analyze Data Workflow Visualize - Shared Data - Help - User -		Using 34.3 KB
Tools     1       search tools     3	GOBII Extractor Get genotype data (Galaxy Version 0.0.1)	History search datasets	2 * II 8
Get Data <u>Upload File</u> from your computer	Log-in options  Study name  KASB ED: 7 20 2010	Unnamed history 19 shown, 49 deleted 28.13 KB	<b>()</b>
<u>GOBII Extractor</u> Get genotype data <u>GOBII Extractor</u> Get tokens SNP-seek Get Data	Advanced options	68: GOBII output	• • ×
Dynamic Options list	forward-breeding-DS-7.29.2019		Sase: ?
Send Data Collection Operations	✓ Execute	1 2 3 marker_name Samplel Samp snp0S0019 NN TT	4 5 Dle2 Sample3 Samp TT CC
Lift-Over Text Manipulation Convert Formats		snp0S0054 GT TT snp0S0096 CC CC snp0S0002 GG GG	ТТ ТТ СС СС GG GG
Filter and Sort Join, Subtract and Group		error	• • • ×
		0 0 4 d 2 a	> 🚰 💷 🚫 💽 Right Ctrl

GALAXY



### **– Galaxy** PROJECT

## **Partnership Models**

### **Adopt products**

GALAXY

- Open-source and free
- Test with own datasets and use cases
- Download from toolshed

### **Community development**

- Implement new tools
- Develop crop-specific workflows
- Connect with databases





- Number of registered users: 214
- Number of different crops/programs: 3 (rice, wheat, maize)
- Number of analyses/jobs performed: 14200
- Total job file sizes stored: 148 Gb
- Average file size: 4.2982 Mb



## **BrAPI**





## Invite to Partner with GS-Galaxy

## **EiB Demo Server:**

## 1. Project <u>http://galaxy-demo.excellenceinbreeding.org/</u>

- 2. **Products** Publish tools to Galaxy main toolshed
- 3. Partnership
- Publish tools to Galaxy main toolshed Galaxy workshop at Cornell/BTI April, 2020













# Thanks







International Maize and Wheat Improvement Center



IRRI



