## Challenges and Opportunities in Connecting Genotype Data to Phenotype Data

#### The AgBioData Consortium

https://www.agbiodata.org/

#### Sushma Naithani

Dept. of Botany and Plant Pathology Oregon State University

sushma.naithani@oregonstate.edu



## **Genotype and Phenotype Working Group**





**Sushma Naithani** Oregon State Univ. Corvallis, OR, USA **Sook Jung**, Washington State Univ. Pullman, WA, USA Sunita Kumari Cold Spring Harbor Laboratory, NY, USA



Elsa H Quezada UNAM, Tizayuca Mexico



Irene Cobo-Simón Univ. of Connecticut, Storrs, CT, USA



Nicholas Gladman, Cold Spring Harbor Laboratory, NY, USA











Melanie Correll Univ. of Florida, Gainesville, FL, USA Maria Skrabisova Palacký Univ. The Olomouc, Czech for Republic I

Cecilia H. Deng The New Zealand Institute for Plant & Food Research Limited, New Zealand Wentao Zhang National Research Council Canada, Olusola Afuwape Univ. of Lagos, Lagos, Nigeria **Akeem Sikiru** Federal Univ. of Agri. Zuru, Nigeria **Ines Rebollo** Univ. de la República, Uruguay



#### • Data collection, sharing and integration issues :

- increased data size and complexity
- diverse data formats
- computer intensive analysis

- quality control
- o data formats, metadata annotation
- o data storage and access
- Limited data visualization and analysis tools: to address a biological question, output formats
- Data interoperability and integration: genotype and phenotype data
- Lack of Funding for Expert Biocuration: limited high
  quality knowledge synthesis

## The Genotype-Phenotype Working Group Aims & Goals

- To improve the data collection and data sharing
- To facilitate linking genotype and phenotype data
- To promote data interoperability and re-use

#### **Aim 1 To Improve the Data Collection**

Review of diverse data types, annotation, storage, and archiving in primary repositories



Data Generation, Annotation, & Archiving

#### AIM 2: Linking Genotype & Phenotype Data

![](_page_5_Figure_2.jpeg)

### Integration of Heterogenous Data Types Gramene-Ensembl Genome Browsers

![](_page_6_Figure_1.jpeg)

- Assembly structure and sequence
- Genes & transcripts
- Comparative alignments
- Baseline Annotation & Ontologies

- Genetic markers (SNPs, Indels, QTLs, SSRs)
- Repeats & transposable elements
- Regulatory & Epigenetic marks
- Transposon discovery

#### Integration of omics data in the Plant Reactome

![](_page_7_Figure_1.jpeg)

Links to the most

• Data extensions provided by collaborators at EBI Interactors (IntAct) and baseline expression data

#### Analysis of *Arabidopsis thaliana* Redox Gene Network Indicates Evolutionary Expansion of Class III Peroxidase in Plants

![](_page_8_Figure_1.jpeg)

#### Displaying gene-gene interaction data on Plant Reactome pathways

Synthesizing gene-gene interactions information

![](_page_9_Figure_2.jpeg)

#### AIM 2: Linking Genotype & Phenotype Data

#### Examples of existing resources-→secondary Knowledgebases

![](_page_10_Picture_2.jpeg)

#### Genomic variation in 3,010 diverse accessions of Asian cultivated rice

- ~29 million single nucleotide polymorphisms,
- ~2.4 million small indels
- ~ 90,000 structural variations that contribute to withinand between-population variation.
- ~10,000 novel full-length protein-coding genes that have a high number of presence–absence variations.
- Wensheng W. et al. (2018), Nature 557:43-49

![](_page_10_Picture_9.jpeg)

#### https://bar.utoronto.ca/eplant/

#### Data visualization tools for multiple levels of plant data.

![](_page_10_Figure_12.jpeg)

#### **AIM 3:** Limitations of data re-use and interoperability

![](_page_11_Figure_1.jpeg)

![](_page_12_Picture_0.jpeg)

# How does genotype relate to phenotype?

#### **Reverse Genetics Approach**

One Gene One Phenotype

#### **Mendelian Genetics**

Gray, W. M. (2004). "<u>Hormonal Regulation of Plant Growth and</u> <u>Development</u>". *PLoS Biology* **2** (9): e311. <u>DOI:10.1371/journal.pbio.0020311</u>

#### Extension to Mendelian Genetics Gene 1+ Gene 2- $\rightarrow$ Multiple Phenotypes **Example**: Gene A Fruit color in squash increases orange Gene B increases yellow Immature cukes are green Additional colors result from interaction of these two loci

#### How does genetic variation give rise to phenotypic variation? **One Gene with Multiple Phenotypes**

![](_page_14_Figure_1.jpeg)

Fambrini et al. (2014) Gene, 549(1), 198-207. https://doi.org/10.1016/j.gene.2014.07.018.

(turf)

2) [2]

-3) [1]

4) [3]

![](_page_15_Figure_0.jpeg)

#### Complexities in Relating Phenotype to Genotype (Exemption to Mendel's Law of one gene->one phenotype)

#### Perspectives

- Cell specific expression
- Stage specific expression
- Environmental effect
- Gene redundancy
- Multiple alleles → many phenotypes
- Penetrance or Expressivity
- Epistasis
- Lethal genes
- Linkage
- Pleiotropy

![](_page_16_Picture_12.jpeg)

Photo by Jing Yuan, Ph.D. student, Kessler Lab (https://ag.purdue.edu/stories/pollination-a-classic-taleof-romance-love-and-death/)

![](_page_16_Figure_14.jpeg)

#### Phenotypes are the outcomes of complex interactions (G(s) + E)

![](_page_17_Figure_1.jpeg)

#### **Phenotype -> Genotype**

#### Phenomics Approach & GWAS

For understanding

natural variability and biodiversity

Linking phenotype to genotype in genome-wide association studies (GWAS): Forward Genetics Approach

![](_page_19_Picture_1.jpeg)

Natural habitat of various accessions of *A. thaliana* 

Vegetative rosettes illustrating genetically determined variation in morphology among *A. thaliana* accessions

Weigel and Mott (2009): Genome Biol. 2009;10(5):107. doi: 10.1186/gb-2009-10-5-107

# Improving grain yield, stress resilience and quality of bread wheat using large-scale genomics

![](_page_20_Figure_1.jpeg)

Manhattan plots showing marker–trait associations for GY and agronomic traits from a genome-wide association mapping study. The 2AS chromosomal region was significantly associated with grain yield in several environments, whereas the 5BL and 2BS chromosomal regions were associated with days to heading and maturity in several environments.

Nature Genetics volume 51, pages1530–1539(2019)

#### Linking Phenotype with existing platforms

#### Plant Pathways would be very useful

![](_page_21_Figure_2.jpeg)

### Missing: Images of phenotypes

#### Linking genotype to Phenotype Data with gene pages of Ensembl and other crop databases like MaizeDB

![](_page_22_Figure_1.jpeg)

#### Aim 3: To Review Data Interoperability & Re-use

![](_page_23_Figure_1.jpeg)

Modified figure from Naithani et al. (2019), Database (Oxford), Volume 2019, bay146, https://doi.org/10.1093/database/bay146

#### Some Thoughts .....

- Linking genotype and phenotype data is one of the greatest challenges in research and development
- Why it is so difficult to map phenotypes to genotype data? Potential solutions:
  - Standardized data formats
  - Data annotations using ontologies, enriched metadata
  - Automated flows for data management and mappings

#### Input machine readable-- $\rightarrow$ output machine readable

#### Not useful for human mind: the way it works, thinks and creates new

# How to we integrate the diverse data so it is easy to understand ?

The Biocuration by experts could help in Knowledge synthesis

Additional investment in data visualization platforms will be of tremendous value to biologists

![](_page_25_Figure_3.jpeg)

Automation/AI

Illustration by **Dmitry Shevela** <u>https://www.instagram.com/scigrafik</u> **MaizeDIG**: provides high-quality data on Genotype—to-phenotype through value added expert curation, with easy-to-understand visual format in a model species rich in genotype and phenotype data.

https://maizedig.maizegdb.org (Cho et al. 2019: https://doi.org/10.3389/fpls.2019.01050)

![](_page_26_Picture_2.jpeg)

#### **Needed : Increased Support for Biocuration**

![](_page_27_Figure_1.jpeg)

## **Final Product: White Paper**

- Introduction
- Main Content
  - + Review the data types and method
  - + Primary repositories and secondary Knowledgebases
  - + Requirements of metadata and the minimum standards
  - + Examples of re-use and re-analysis
  - + Limitations: data re-use, FAIR policy, sustainability of public databases
- Recommendations: additional infrastructure
  & Biocuration needs

![](_page_28_Figure_9.jpeg)

![](_page_29_Picture_0.jpeg)

# The community is fundamental to advance we have deluge of Omics data: help is appreciated

![](_page_29_Figure_2.jpeg)

## Acknowledgements

#### **AgBioData SC members:**

Monica Poelchau Leonore Reiser Sunita Kumari Sook Jung Sushma Naithani Meg Staton Jacqueline Campbell Peter Harrison John P. McNamara

## Past AgBioData SC members:

Lisa Harper Eva Huala Marcela Tello-Ruiz Laurel Cooper Ethy Cannon

Past PC: Darwin Campbell

The AgBioData consortium

Award Abstract # 2126334

![](_page_30_Picture_8.jpeg)