AgBioData SGV

# From Fragmented Data to Shared Knowledge: Adopting rsIDs in Agricultural Genomics

Marcela Karey Tello-Ruiz, PhD
AgBioData Standards for Genetic Variation Working Group

2026 PAG AgBioData Workshop

# AgBioData Standards for Genetic Variation WG



https://www.agbiodata.org/working_groups/sgv

AgBioData SGV

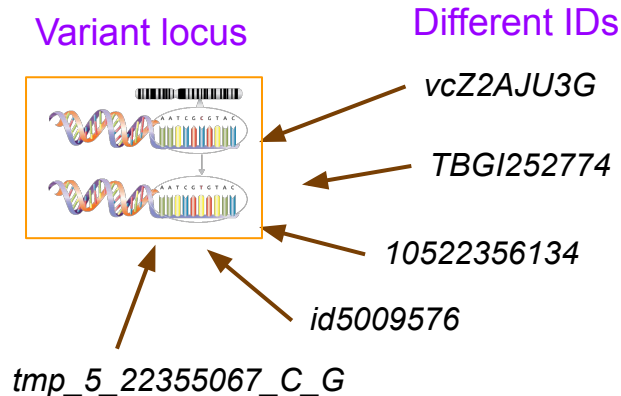# A Community-Led Path Toward FAIR Variant Data

AgBioData SGV

Goals:

- Assess existing standards

- Identify communities at different data-maturity stages

- Focus on promoting rsID adoption via EVA

- Provide practical FAIR guidance*

# Standardizing Genetic Variation for FAIR Agricultural Genomics

AgBioData SGV

- Genetic variation data are critical for **breeding, trait discovery & adaptation**

- **Agricultural genomics** generates vast amounts of SNP data that is hard to reuse & integrate

### *Why hard to reuse?*

Variant locus

Different IDs

*vcZ2AJU3G*

*TBGI252774*

*10522356134*

*id5009576*

*tmp_5_22355067_C_G*

SNPs often labeled with:

- Project- or array-specific identifiers
- Assembly-dependent coordinates

Data spread across:

- Community databases
- Supplementary tables
- Commercial genotyping platforms

# Challenges with (non)FAIR variation datasets

AgBioData SGV

Mandated by Funders. Not Enforced by Journals

PUBLISHED

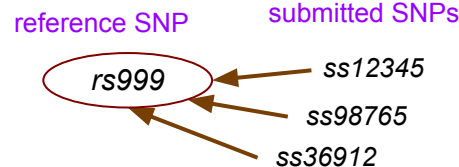| Only raw data is shared (no standard format) | Dataset is shared as supplementary data or local DB | Dataset is shared in generalist FAIR data repository | Dataset is shared in specialized FAIR data repository using standard identifiers & formats |

Not FAIR → FAIR

*How human genomics has tackled the SNP naming problem & why the problem is getting worse in agriculture?*
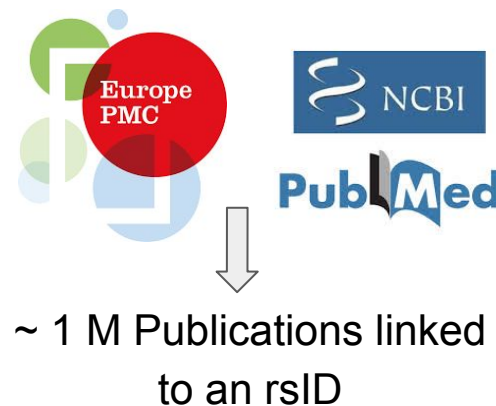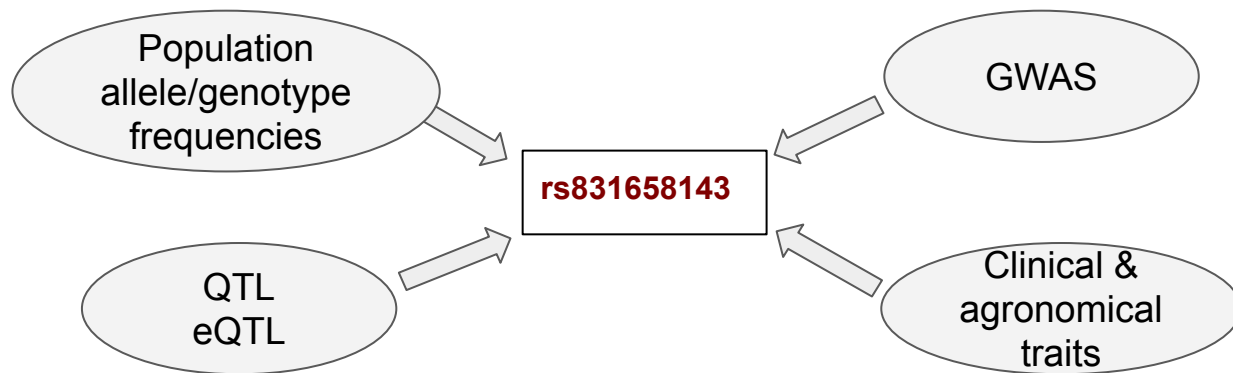
# The Power of Using rsIDs

rs999 ← ss12345
← ss98765
← ss36912

**AgBioData SGV**

*What is an rsID?*

rs831658143

- Reference SNP cluster ID
- Identifies a variable genomic locus
- Globally unique, persistent accession
- Stable across genome assembly versions & crop varieties

## Several data types aggregated around a marker

Population allele/genotype frequencies → **rs831658143** ← GWAS

QTL eQTL → **rs831658143** ← Clinical & agronomical traits

Europe PMC   NCBI   PubMed

~ 1 M Publications linked to an rsID

# Better Genomes, Bigger Integration Challenges

- **Long-read sequencing** → high-quality reference genomes

- **Pan-genomes** reveal:

  - Previously hidden variation
  - Reference-specific biases

Without stable identifiers:

- Variant IDs break across assemblies
- Data integration becomes harder

# Promoting use of rsIDs – SoyBase

**Comprehensive description of genomewide nucleotide and structural variation in short-season soya bean**

Davoud Torkamaneh[1,2], Jérôme Laroche[2], Aurélie Tardivel[1,2,3], Louise O'Donoughue[3], Elroy Cober[4], Istvan Rajcan[5] and François Belzile[1,2,*]
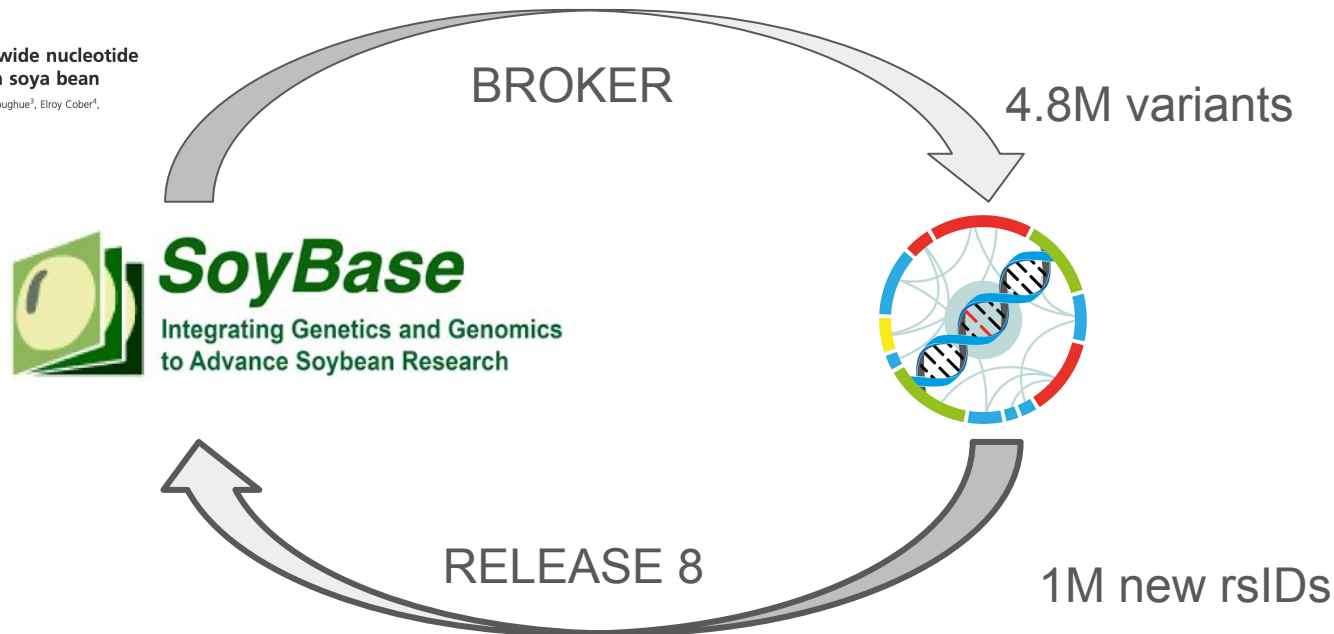
[1]Département de Phytologie, Université Laval, Quebec City, QC, Canada
[2]Institut de Biologie Intégrative et des Systèmes (IBIS), Université Laval, Quebec City, QC, Canada
[3]CÉROM, Centre de Recherche Sur Les Grains Inc., Saint-Mathieu de Beloeil, QC, Canada
[4]Agriculture and Agri-Food Canada, Ottawa, ON, Canada
[5]Department of Plant Agriculture, Crop Science Bldg., University of Guelph, Guelph, ON, Canada

BROKER

4.8M variants

**SoyBase**
Integrating Genetics and Genomics
to Advance Soybean Research

RELEASE 8

1M new rsIDs

*~90% concordance* across studies despite *differences in technology and SNP density*

Rex Nelson

2026 PAG AgBioData Workshop

9

# Promoting use of rsIDs – Gramene / SorghumBase

| Gramene PanGenome | Reference Crop | # rsIDs* (M) |
| --- | --- | --- |
| SORGHUMBASE. Genomic resources for the sorghum research community | Sorghum | 46 |
| MAIZE Gramene Comparative plant genomics focused on maize varieties | Maize | 79 |
| ORYZA Gramene Comparative plant genomics focused on rice varieties | Rice | 68 |
| GRAPEVINE Gramene Comparative plant genomics focused on grapevine varieties | Grape | 0.3 |
| Gramene | Arabidopsis | 26 |

Sharon Wei

Marcela K. Tello-Ruiz

# Promoting use of rsIDs – Community DB / Industry

Collaborative development of community marker panels with rsIDs:

- Sorghum 2.4K SNPs (AgriPlex panel)
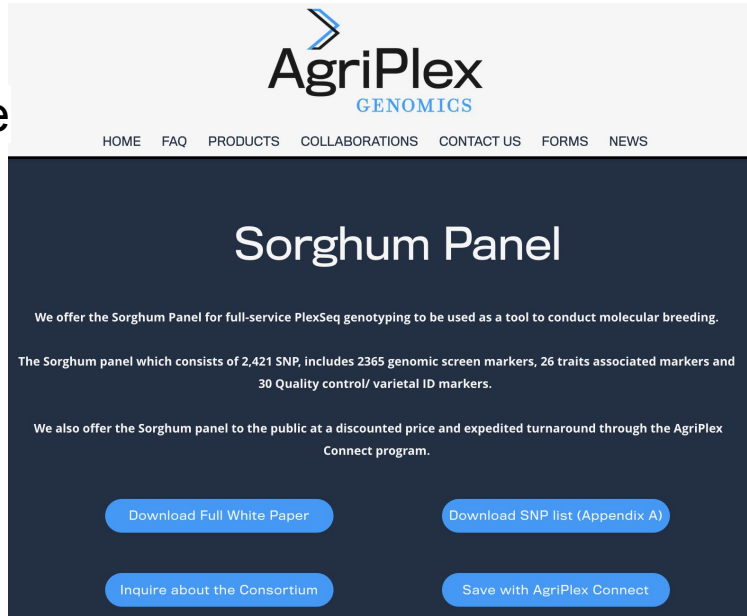- EVA assigned rsIDs to 26 markers without one



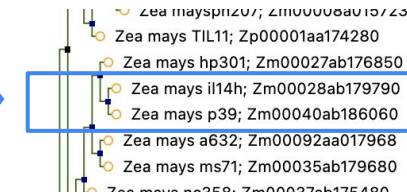Vivek Kumar    Marcela K. Tello

Tim Cezard

# rsIDs provide stable identifiers for transversing Pan-Genomes



**AgBioData SGV**

rs831658143 (T/C)

Tgg (**W**) → Cgg (**R**)
starchy        sweet

MAF <1%

Sweet corn

Zea mays B73 v5
*Sugary1 (Su1)* gene
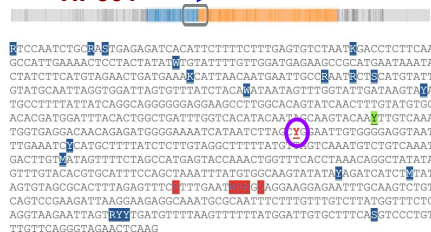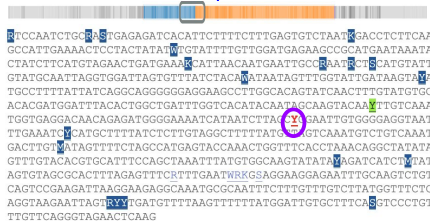
rs831658143 (T/C)

Starchy corn (B73, HP301)

**Variant: rs831658143**
more about rs831658143
Class: SNP
Location: 4:43436836
Alleles: T/C
Ambiguity code: Y
Amino acid: W/R
Consequence: missense variant
Sources: NV_PMID39611775_MaizeGDB, EVA5, HapMap2, EVA6
Codon change: Tgg/Cgg

Sweet corn (P39, IL14H)

**Variant: rs831658143**
more about rs831658143
Class: SNP
Location: 4:43555709
Alleles: C/T
Ambiguity code: Y
Amino acid: R/W
Consequence: missense variant
Source: EVA2024rsID_flanking_B73genes
Codon change: Cgg/Tgg

HP301

P39

IL14H

Tello-Ruiz et al, *unpublished*

* W578R previously reported 3X

12

# Data journey for GV datasets & recommended actions

**AgBioData SGV**

| | | | |
|---|---|---|---|
| Request assembly to be submitted to INSDC | Encourage or broker SNP submission to EVA; promote using standard IDs & formats | Adopt rsIDs, germplasm IDs & controlled vocabularies | Integrate with other standardized data types & link to other DBs & repos |

Other AgBio DBs provided overview & progress towards FAIRifying GV data for white paper

Not FAIR → FAIR

# A Tipping Point for Interoperable Agricultural Genomics

**Anchoring variants to rsIDs unlocks integration with**:

- Functional annotation tools
- GWAS, QTLs, and multi-omics data
- Breeding databases and germplasm repositories

**IMPROVES**

- Traceability
- Genotype-to-phenotype discovery
- Long-term reuse of genomic resources

**What's next?**

- Continue to raise community awareness of the benefits of rsIDs
- More AgBioData community databases adopting rsIDs
- More commercial SNP array providers using rsIDs
- rsID mapping to more pangenomes

# Thanks!