# Standardizing Biocuration of Genetic Variation Data to Promote FAIRification

**Standards for Genetic Variation Working Group**
AgBioData Consortium

Marcela Karey Tello-Ruiz, PhD
Cold Spring Harbor Laboratory

# Outline

1. The SGV Working Group

2. Standards for Genetic Variation & Interoperability

3. Data Submission to the European Variation Archive

4. Challenges

5. Progress towards FAIRifying plant data sets

# AgBioData SGV Working Group Goals

**AgBioData SGV**

- Support the harmonization and adoption of standards for genetic variation (GV) data from various platforms in Plants & Animals

- Bring together a community of data providers, biocurators & computer scientists to promote interoperability and access to GV datasets

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

# Standards for Genetic Variation Working Group

**AgBioData SGV**

- **Specific objectives**:
    - Enable sharing of GV data to support agriculture
    - Identify existing GV and technical barriers for data exchange
    - Review technical standards for GV to support adoption
    - Review GV workflows
    - Engage community to support ingestion and usability of GV data into community and archival resources

- **Activities**:
    - Regular monthly meetings (engagement with Education & Sci. Literature WGs)
    - Participation at AgBioData annual workshop & PAG workshop
    - Data biocuration & coordination across participant resources
    - Promoting FAIRification of GV data & recruiting members at relevant events
    - Merging with Public Genetics Resources WG
    - Reporting to funders

# AgBioData Standards for Genetic Variation WG

**AgBioData SGV**

**Co-Chairs:**

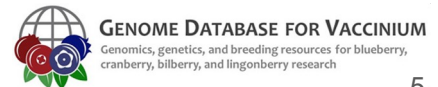Marcela K. Tello-Ruiz
Timothee Cezard

**Most active members:**

- Nahla Bassil
- Sebastian Beier
- Irene Cobo
- Sarah Dyer
- Osman Gutierrez
- Melanie Harrison

- Jodi Humann
- Rex Nelson
- Mazdak Salavati
- Moira Sheen
- Doreen Ware
- Sharon Wei

Full list at https://www.agbiodata.org/working_groups/sgv

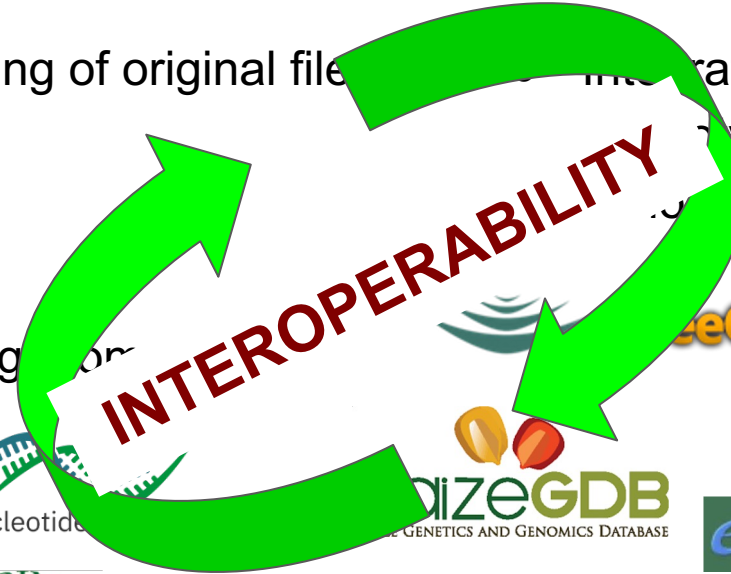# Different databases are serving different purposes

**Archival DBs**

- Long-term archiving of original file
- Accessioning
  - Study
  - Samples
  - Variants
- Update to newer g

**Community/species DBs**

- Integration between genotypes and phenotypes
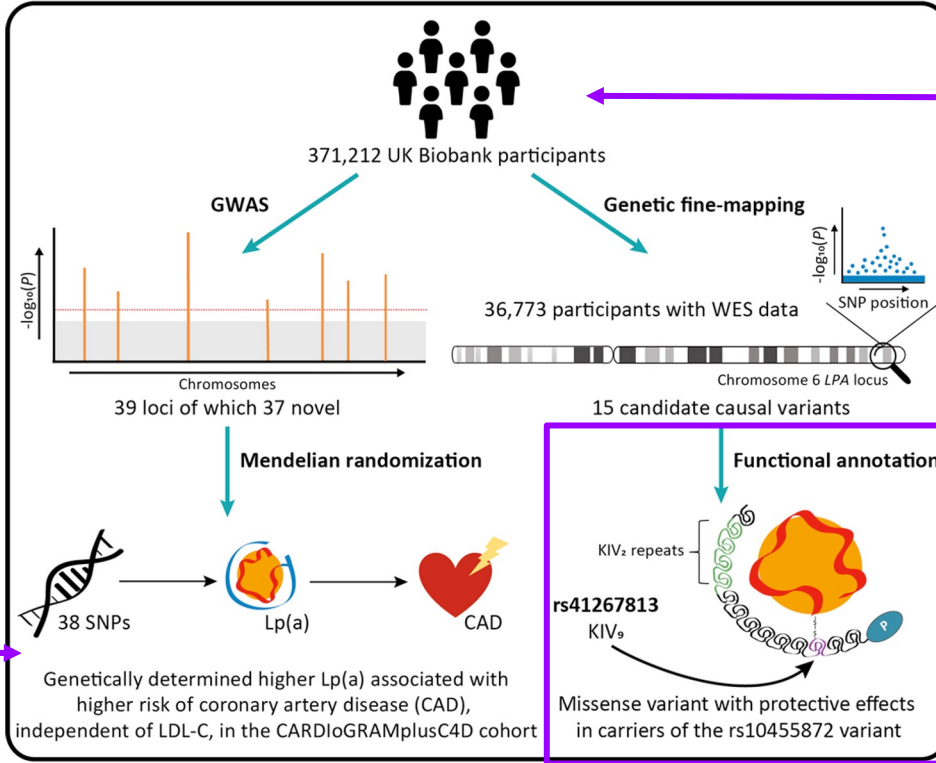- feature/toolsets

INTEROPERABILITY

# Lessons from human genetics

Biosamples

Trait & disease associations

Functional annotations

Genetic variants

rsID (missense) associated with protective trait in disease

Image taken from doi: 10.1161/ATVBAHA.120.315300

# Standards for Genetic Variation – Interoperability

**AgBioData SGV**

rs1234
rs4567
rs7890
…

## rsIDs

Reference clu...

Stable/unique for a...

EVA provides 'ss' ...

and 'rs' (ref) ids for ...

variant...

## Biosample IDs

- BioSample ID (EVA ...ment)
- ...plasm ID (genebanks ...T: IS 12661, GRIN: PI

**VCF**

```
##fileformat=VCFv4.2
##contig=<ID=2,length=51304566>
##INFO=<ID=AC,Number=A,Type=Integer,Description="Allele count in genotypes">
##INFO=<ID=AN,Number=1,Type=Integer,Description="Total number of alleles in called genotypes">
##FORMAT=<ID=GT,Number=1,Type=String,Description="Genotype">
##FORMAT=<ID=DP,Number=1,Type=Integer,Description="Read Depth">
##FORMAT=<ID=GQ,Number=1,Type=Integer,Description="Genotype Quality">
#CHROM  POS    ID  REF ALT QUAL FILTER  INFO       FORMAT     SAMPLE1     SAMPLE2     SAMPLE3     SAMPLE4
2       81170  .   C   T   .    .       AC=9;AN=7424  GT:DP:GQ  0/0:4:12    0/0:3:9     0/1:1:3     0/1:9:24
2       81171  .   G   A   .    .       AC=6;AN=7446  GT:DP:GQ  0/1:4:12    0/0:3:9     0/0:1:3     0/0:9:24
2       81182  .   A   G   .    .       AC=5;AN=7506  GT:DP:GQ  0/0:5:15    0/0:4:12    0/0:5:15    0/0:9:24
2       81204  .   T   G   .    .       AC=2;AN=7542  GT:DP:GQ  1/0:5:15    0/0:9:27    0/0:10:30   0/0:15:39
```
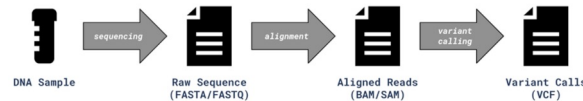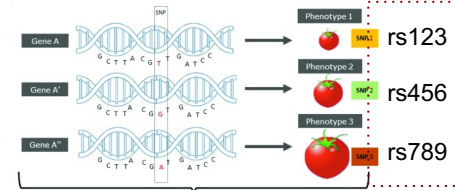
## Traits
Controlled vocabularies
for GWAS, QTLs, etc.

Phenotype 1 — rs123

Phenotype 2 — rs456

Phenotype 3 — rs789

## VCF

Variant Call Format

Text file format with meta-info and data for a

variant position in a *genome sequence assembly*

at INSDC

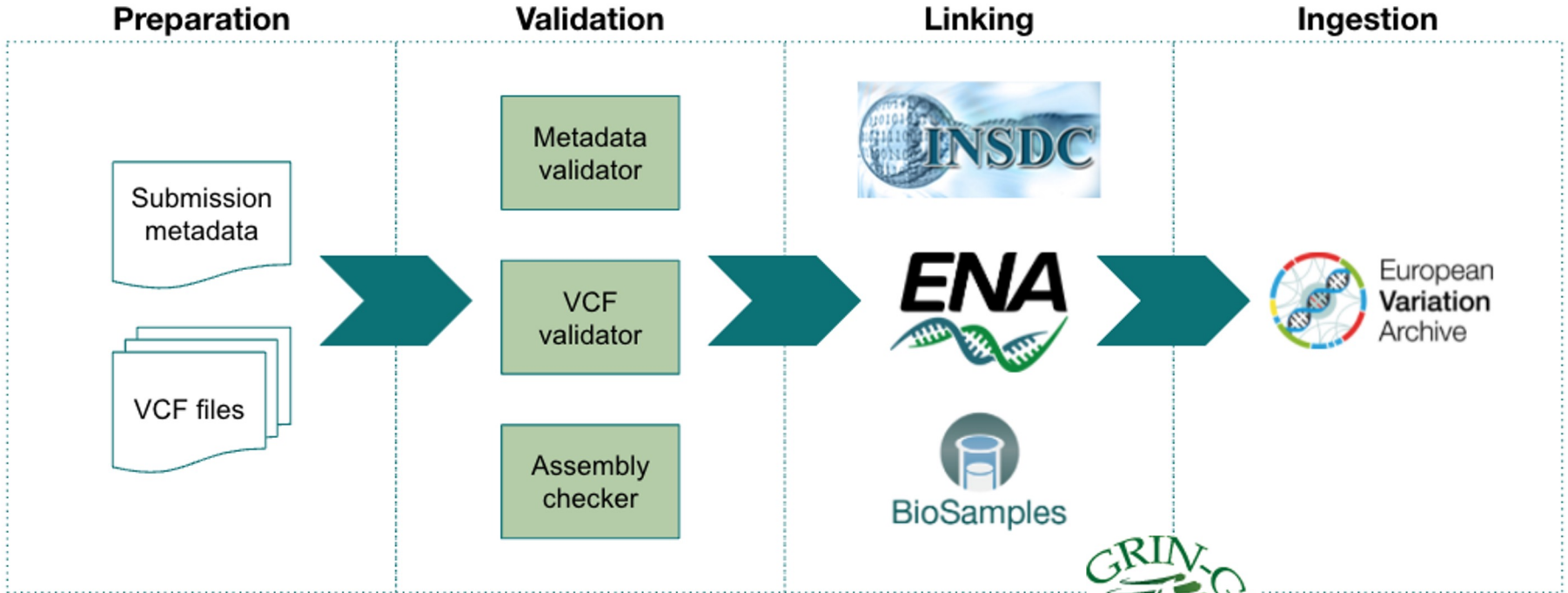DNA Sample → sequencing → Raw Sequence (FASTA/FASTQ) → alignment → Aligned Reads (BAM/SAM) → variant calling → Variant Calls (VCF)

# Data submission to the European Variation Archive

PAG31

9

# Major challenges associated with genetic variation

**AgBioData SGV**

- All data could be reused

- Remapping to a newer assembly may result in reduced precision & data loss

- Moving from a single reference to a PanGenome

- Improvements in assays and algorithms to determine GV (GBS, WGS, etc.)

- Converting from SSRs to SNPs

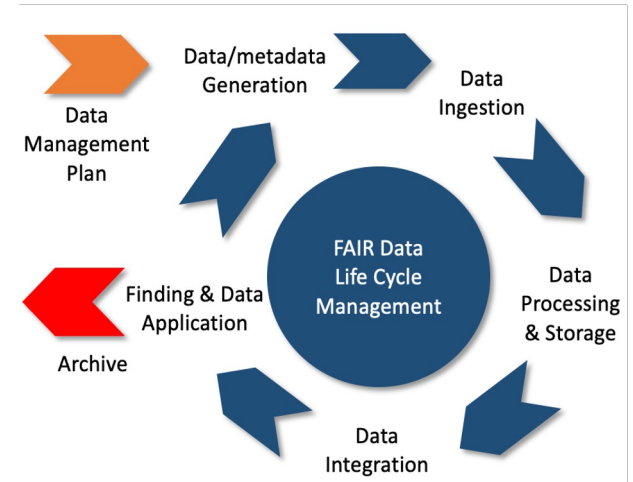- Integration between studies (new studies, meta-analyses, etc.)



Image credit: FAIRToolkit

=> Solution: Submit GV to EVA to get rsIDs, unique genetic variant identifiers

# Challenges associated with biosamples



| | Study 1 | Study 2 Study 3 | |
|---|---|---|---|
| | PI534138 | SC62C | SAP_PI534138.B064FABXX.2.**F7** |
| | PI534138 | SC62C | SAP_PI534138.B064FABXX.2.**G1** |

**Other names**

Matchikah       SAP-416

**AgBioData SGV**

| DB sample with multiple ids because of naming conventions | | |
|---|---|---|
| | --> *Join words* | --> *Use underscores* |
| **Study 1 at NCBI** Study 3 at DB2 | **Study 1 at DB1** | |
| EarlyHegari | EarlyHegari | Early_Hegari |
| IBC/E-38432 | 38432 | IBC_E38432 |
| Karper 669 | Karper669 | Karper_669 |
| Malisor 84-7 | Malisor84-7 | Malisor_84-7 |
| RTx7000 | RTx7000 | RT_7000 |
| *S. bicolor (PI226096)* | S.bicolor.subsp.Verticilliflorum(PI226096) | PI226096 |
| *S. bicolor subsp. drummondii (PI330272)* | S.bicolor.subsp.drummondii | PI330272 |
| *S. bicolor subsp. verticilliflorum (AusTRCF 317961)* | S.arundinaceum | AusTRCF_317961 |
| *S. bicolor subsp. verticilliflorum (PI300119)* | S.bicolor.subsp.Verticilliflorum(PI300119) | PI300119 |
| Cherekit (IBC/E-460) | Cherekit(S) | Cherekit_IBC_E460 |
| Kilo (IBC/E-382) | Kilo | Kilo_IBC_E382 |
| Yik.solate (IBC/E-339) | Yik.solate | Yik_IBC_E339 |
| Zengada (IBC/E-308) | Zangeda | Zengada_IBC_E308 |

=> Solution: Use standard germplasm identifiers (BioSample / Genebank IDs)

PAG31

# Recommendations for data standards for plants

- FAIRification of Plant Genotyping Data (& linking it to Phenotyping)
- First guidelines on FAIR handling of GV data published in 2022
- Provide a checklist to classify and validate the data to support iits submission to EVA (and BioSamples)

# Recommendations from Biocurator Meetings

Additional Suggestions for Plant Samples Metadata associated with VCFs

| Metadata field | Field Name | Definition | Format | Example | Cardinality |
|---|---|---|---|---|---|
| ##SAMPLE | | Metadata about a single sample genotype that is part of the genotyping experiment in the VCF file | Composite (see below) | ##SAMPLE=<ID=SAMN04168247, DOI=doi:10.18730/NBYG*, ext_ID=grin-global.org:USA126:PI 276837> | 1:N |
| | BioSample ID | Refers to a biological sample used as a 'reference' (e.g. to sequence its genome) or used in an assay database such as ENA, EVA, ArrayExpress. Always begin with SAM. The next letter is either E or N or D depending if the sample information was originally submitted to EMBL-EBI or NCBI or DDBJ, respectively. After that, there may be an A or a G to denote an Assay sample or a Group of samples. Finally, there is a numeric component that may or may not be zero-padded. | /(SAM)(E2N\|D)(A\|G)(\d+)/ | ID=SAMN04168247 | 1 |
| | External identifiers | - Primary accession - **One mandatory external ID for plants**. Impractical to enter metadata for each biosample; easier to add as a metadata line in VCF. Impractical for huge data sets as this would significantly increase the size of the VCF file. --- Source of accession [Genebank Name, Original Collection (not in genebank), etc.] Examples: GRIN, ICRISAT, WEIWS code:Species code (IPK), CNGB, GBIS, ORIGINAL COLLECTION --- Accession prefix. Examples: PI, IS, NSSL, GRIF, SOR, Collector ID --- Accession unique identifier or number. Example: six-digit PI number, five-digit IS number, four-digit following WEIWS:species number, collector number - Secondary accession - Sample inventory **if applicable**. Example: CR02, CR03, 07PL. Note: USDA germplasm repositories provide inventory accessions. - Other - Not necessary. Example: Population panel identifiers such as SAP-391, a member of the Sorghum Association Panel are not necessary and are well captured in germplasm registries like GRIN. Identifiers under which this genotyping sample is registered in other databases (either 'FAO-WIEWS_instcode:genus:accession_number' or 'DNS:database_identifier:identifier_scheme:identifier') | ext_ID=registry:identifier | ext_ID=grin-global.org:USA126:PI 276837 | 1:N |
| | Study sample identifier | Identifies specific plant/genotype used, **when available** . This will usually be specific to an individual research project and not publicly available. However, the plant or DNA sample may be shared between researchers. Different plant numbers from the same lot. Example: SC103 and SC103-14E share the same PI533752 accession. | | | 0-1 |
| | DOI, URL | DOI for the passport information of the genotyping sample. | URL, DOI | DOI=doi.org/10.18730/NBYG* | 0-1 |

=> BioSamples entries:

- ○ Require primary external identifier from major germplasm repository (e.g., GRIN, CGIAR, IPK, CNGB) with doi/url

- ○ Recommend including inventory or local number & identifier for the specific plant/genotype used in the study
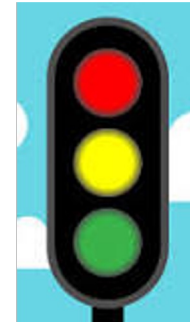
AgBioData SGV

# Technical challenges revealed through biocuration

- Missing reference genome assembly

- Reference genome not registered at INSDC

- GV data not readily available (e.g., private FTP)

- GV data not in standardized format (e.g., VCF)

  - Non-standard format at community DB (e.g., tabular output .xls)

  - No format conversion method provided

  - Only precursor sequencing reads provided

14

# FAIRifying public plant GV data sets

**AgBioData SGV**

| Species | Reference assembly in INSDC | VCF available | Sample IDs with DOI/URL from major germplasm repo | VCF in EVA & BioSamples | Samples qualified for cross-linking to other DBs | Recommended action |
|---|---|---|---|---|---|---|
| cranberry, raspberry, blackberry | ☐ (red) | ☐ | ☐ | ☐ | ☐ | Authors will need to submit assembly to INSDC |
| pear | ☐ (red) | ☐ (red) | ☐ | ☐ | ☐ | Authors will need to submit assembly to INSDC |
| strawberry | ☐ (red) | ☑ (green) | ☑ (yellow) | ☐ | ☐ | Authors will need to submit assembly to INSDC |
| grape | ☑ (yellow) | ☐ (red) | ☐ | ☐ | ☐ | Contacted authors to submit reference assembly to INSDC & provide VCF. Next contact Journal |
| poplar | ☑ (yellow) | ☑ (green) | ☐ | ☐ | ☐ | INSDC updated assembly. Next EVA to coordinate with CartograPlant /TreeGenes |
| apple, peach, cherry, hazelnut, kiwi | ☑ (green) | ☐ (red) | ☐ | ☐ | ☐ | Unknown whether VCFs are available. NCGR might follow up |
| maize | ☑ (green) | ☑ (green) | ☐ | ☐ | ☐ | Gramene Maize looking to coordinate with MaizeGDB |
| sorghum | ☑ (green) | ☐ (red) | ☐ | ☐ | ☐ | Contacted multiple authors/studies unsuccessfully |
| sorghum | ☑ (green) | ☑ (green) | ☑ (green) | ☑ (green) | ☑ (green) | SorghumBase coodination with EVA & GRIN |

STOP

…

GO

# Working towards solutions

- ❖ Assembly submissions to INSDC
  - ➢ Education & training
  - ➢ [Elixir cookbook recipe](Elixir cookbook recipe)

- ❖ Standard file format
  - ➢ Converter tools (e.g., excel => VCF)

- ❖ Data sharing
  - ➢ Minimum standards
  - ➢ File validation (community DBs effort)
  - ➢ Journals
  - ➢ Funding agencies

- ❖ BioSamples with germplasm IDs + sample doi/url
  - ➢ FAANG project extension
    - ■ Experimental, metadata & bioinformatics standards
    - ■ Reuse tools



FAANG
Functional Annotation of Animal Genomes

# Summary of Outcomes

- FAIRifying pilot studies (replaced tmp SNP IDs with rsIDs):
  - SorghumBase & Gramene: 41M sorghum rsIDs
  - Gramene Vitis: 0.3M grape rsIDs
- Standardized germplasm identifiers
  - Gramene, SorghumBase
- Recruited 14 new members
- Discussed synergy with Education & Sci Lit WGs
- Merged with Public Genetic Resources WG

Gramene Workshop
Tuesday, Jan. 16
Palm 8, 4 pm

# Future work

- Ensure relevant reference assemblies registered at INSDC by active participation of WG members to:
  - Promote data submission to EVA
  - Lower barrier for biocuration through training, SOPs, etc.
  - Convert historical data into current reference assembly
- Biosamples metadata biocuration hackathon
  - Cross-link accessions to germplasm repositories
  - Cross-link passport data (germplasm synonyms)
  - Index widely used population panels

# THANKS - Join our working group, chat with us…!

AgBioData SGV

**AgBioData Booth #422**
Sunday Opening, Monday & Tuesday lunch



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