



AgBioData SGV

Standards for Genetic Variation

Promoting the Use of variant Identifiers for Genetic Markers in Agricultural Research

Timothee Cezard EMBL - EBI

PAG 32



Challenges with (non)FAIR variation datasets



Only raw data is shared

Dataset is shared in paper's supplementary

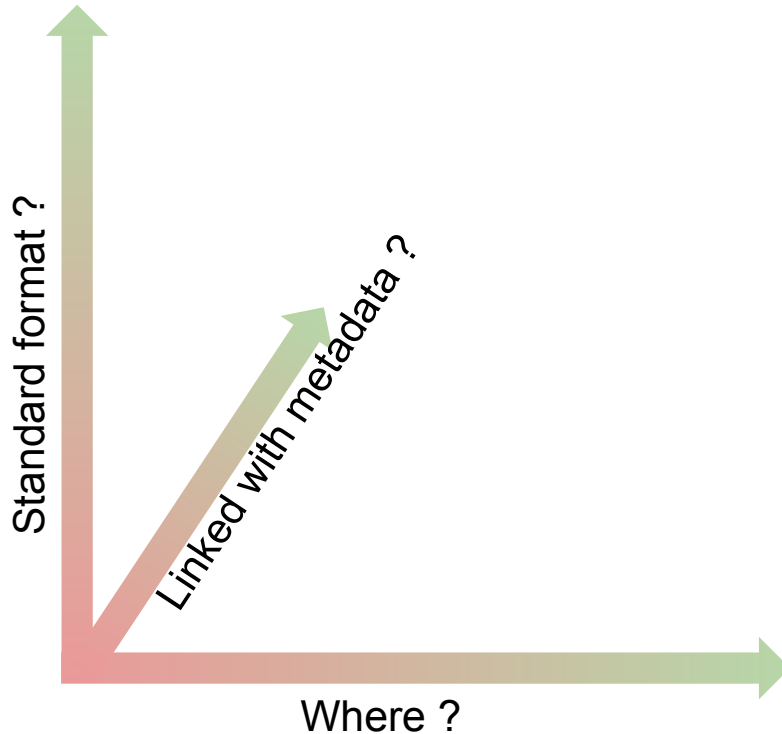
Dataset is share in generalist FAIR data repository

Dataset is share in specialist FAIR data repository





Challenges with (non)FAIR variation datasets



Increasingly complex context

- New tools and technologies
- New assemblies being generated
- Complex metadata

Standards for Genetic Variation WG



AgBioData SGV



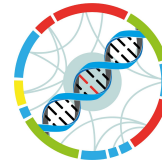
Co-Chairs:

- Marcela K. Tello-Ruiz
- Timothee Cezard

Most active members:

- Nahla Bassil
- Osman Gutierrez
- Rex Nelson
- Jodi Humann
- Sebastian Beier
- Moira Sheehan
- Sarah Dyer
- Melanie Harrison
- Irene Cobo
- Doreen Ware
- Sharon Wei

EMBL-EBI



TreeGenes

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

PAG32: AgBioData Workshop



GENOME DATABASE FOR VACCINIUM
Genomics, genetics, and breeding resources for blueberry,
cranberry, bilberry, and lingonberry research



Landscape analysis through curations



Curation Decision tree



Reuse



Cannot
Reuse

- GV data not in standardized format
- Missing reference genome assembly
- GV data not readily available (e.g., private FTP)
- Samples not using standard



AgBioData SGV

Improving Interoperability via standards

1. **Genomics variation dataset: VCF** Variant Call Format



[Specification](#) from GA4GH

1. Metadata: BioSamples

Variant detection: bcftools, GATK

Variant imputation: PLINK2

1. Markers: rs IDs

Validation: [vcf-validator](#)

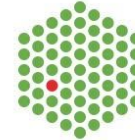


AgBioData SGV

Improving Interoperability via standards

1. Genomics variation dataset: VCF

Entity shared between NCBI, EBI, DDBJ



1. Metadata: BioSamples

Container for any metadata

1. Markers: rs IDs

Community defined checklist





AgBioData SGV

Improving Interoperability via standards

1. Genomics variation dataset: VCF

1. Metadata: BioSamples

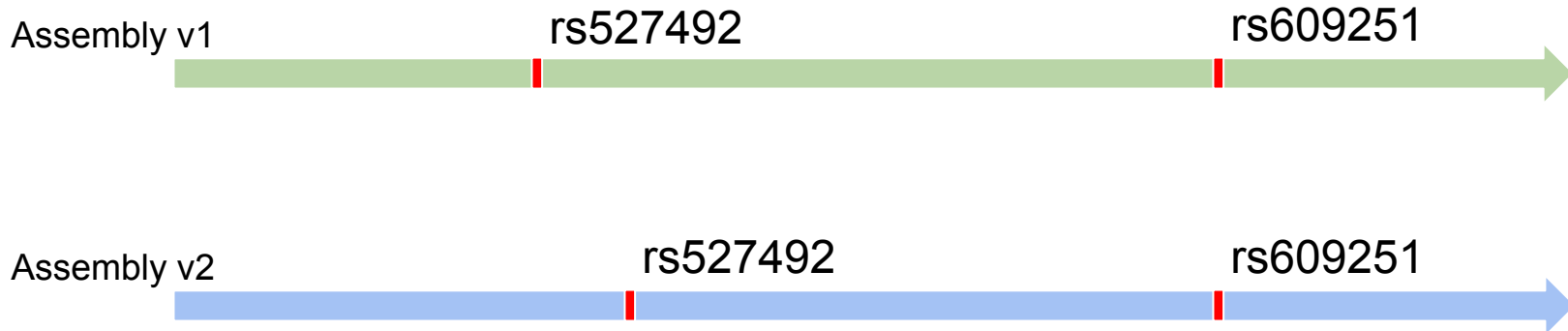
1. Markers: rs IDs



Ref SNP identifiers (rs IDs)

Location on a genome associated with a variation

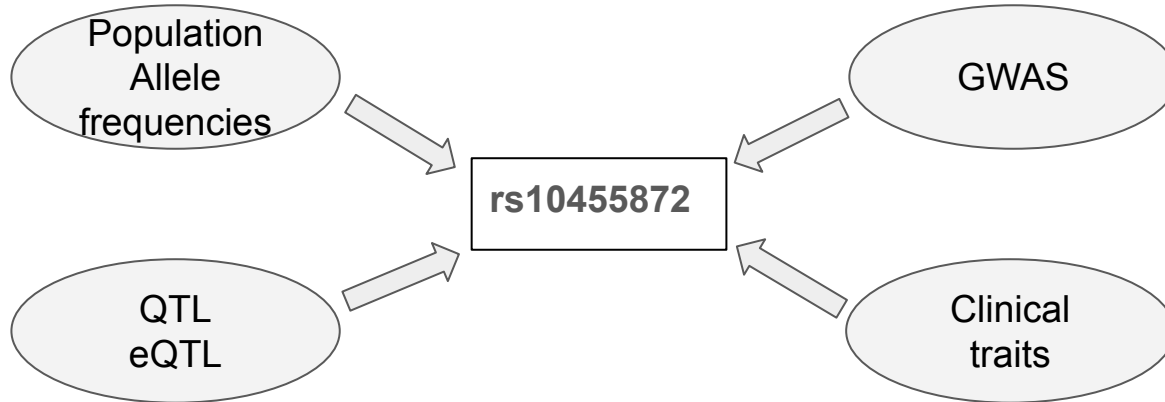
- Globally unique long term accessions
- Identify Variable loci on a genome
- Stable across genome assembly version





Power of using rs IDs

Natural candidate for aggregating several data around marker



Use in publication to describe makers



~ 1 M Publications Link to a RS ids



How to generate rs IDs ?

Submitters



Variation
dataset



Stable RS
release



Promoting use of RSids - Gramene / SorghumBase



SNP count (M)*	EVA release5	Gramene Pan-Genome Sites	Gramene Pan-Genome Sites rsID
Sorghum	50	59	40
Maize	78	50	47
Grape	0.36	0.46	0.32
Rice	32	28	27

M*: Million

The 4 pan-genome subsites of Gramene have been updated with the most recent rsIDs from EVA release version 5.





Promoting use of RSids - Soybase

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

Davoud Torkamaneh^{1,2}, Jérôme Laroche³, Aurélie Tardivel^{1,2,3}, Louise O'Donoghue³, Elroy Cober⁴, Istvan Rajcan⁵ and François Belzile^{1,2,*}

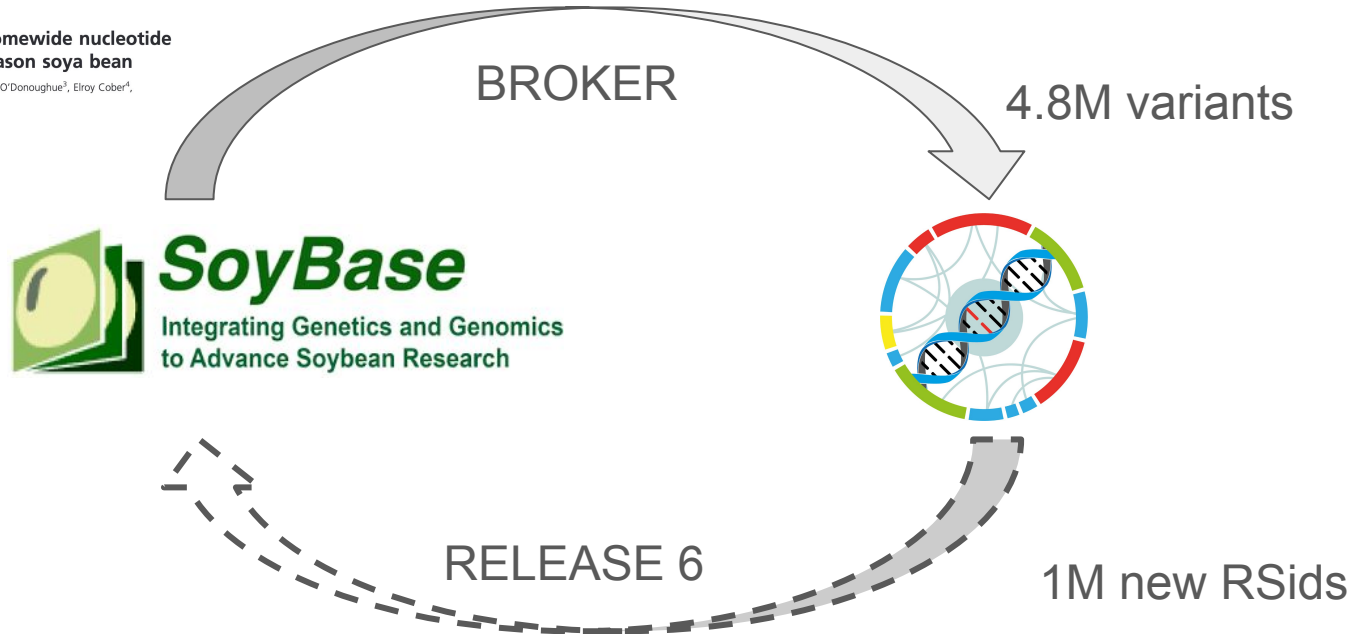
¹Département de Physiologie, Université Laval, Québec City, QC, Canada

²Institut de Biologie Intégrative et des Systèmes (IBIS), Université Laval, Québec City, QC, Canada

³CRISMA, Centre de Recherche Sur Les Grains Inc., Saint-Mathieu de Bellefleur, QC, Canada

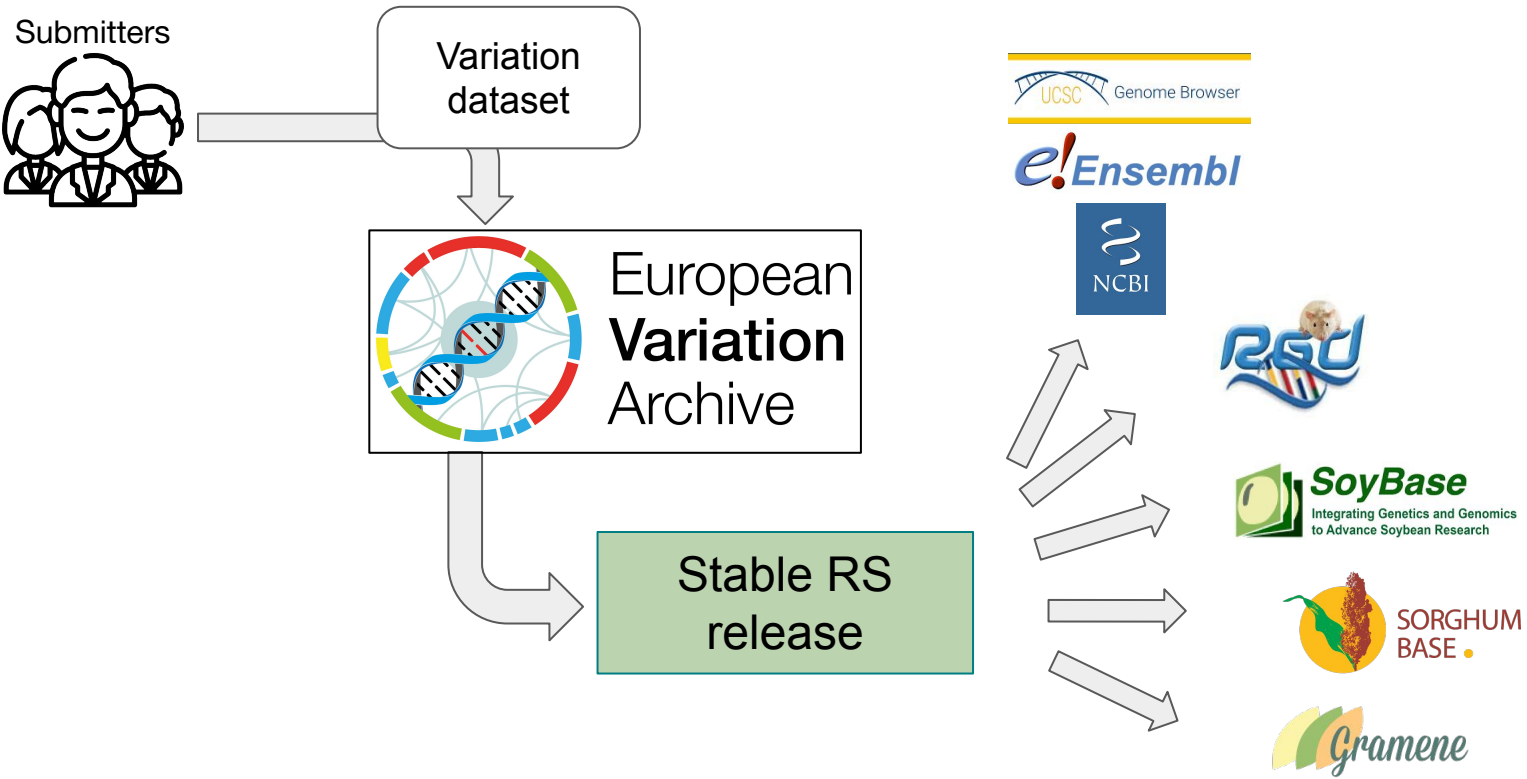
⁴Agriculture and Agri-Food Canada, Ottawa, ON, Canada

⁵Department of Plant Agriculture, Crop Science Bldg., University of Guelph, Guelph, ON, Canada





Integration of rs IDs in downstream resources





AgBioData SGV

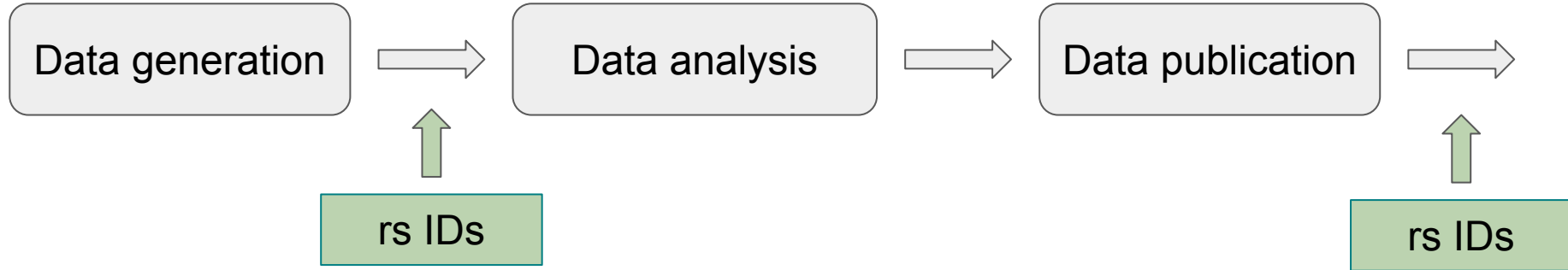
Promoting use of rs IDs

Workshop for AgBioData communities to highlight markers associated with rs ids.



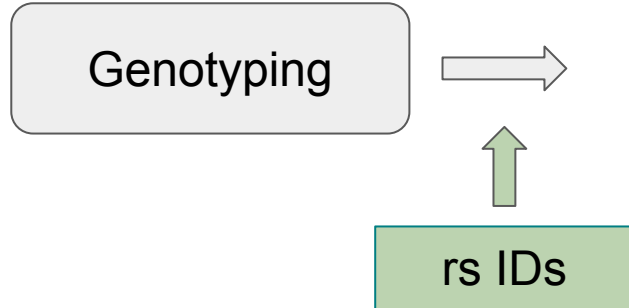


Promoting use of rs IDs





Promoting use of rs IDs



marker1	rs102743256	A	T
marker2	rs398477430	C	A
marker3	rs214779348	G	C
...			

Easier to convert to genomic coordinates

Improve interoperability with other DB downstream



AgBioData SGV

Promoting use of rs IDs - Industry collaboration

Develop a community marker panel with RS ids:

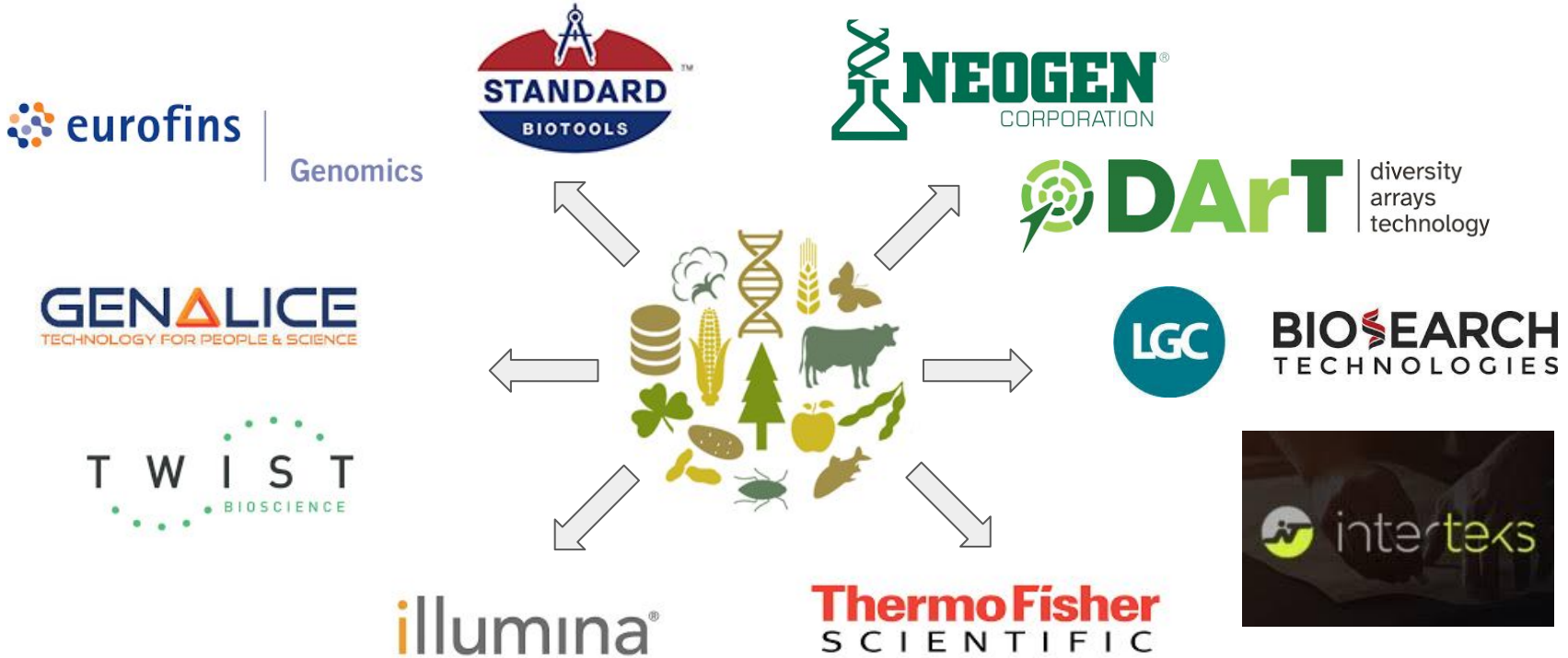
- Sorghum 2.4K SNPs (AgriPlex)
- 26 Markers without RS ids were assigned one





AgBioData SGV

Promoting use of rs IDs - Industry collaboration





Conclusions & Future work

- AgBioData Standard Genetic Variation meets monthly
- Recommendations for standard to be used
- Adoption of standards by community members
- Outreach internally, externally and to industry partners

⇒ Writing White paper





AgBioData SGV

Come and talk to us



Booth 207

EMBL-EBI



Workshop on Sunday at 8:00



European
Variation
Archive

Poster 345



Poster 108