

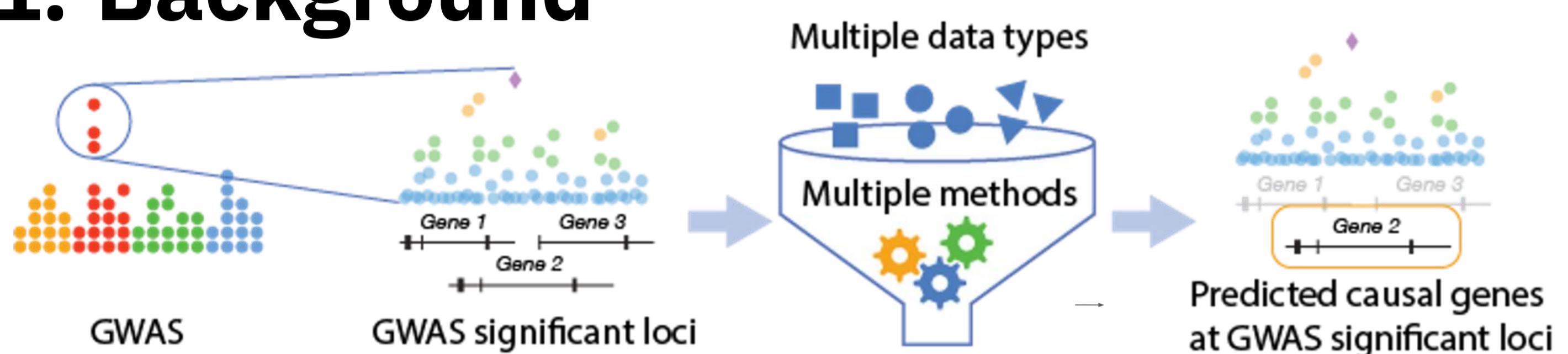
# The PEGASUS framework for Predicted Effector Gene (PEG) Reporting

## Predicted Effector Gene Aggregation, Standards and Unified Schema

Aoife McMahon<sup>1</sup>, Yue Ji<sup>1</sup>, Laura W. Harris<sup>1</sup>, Julie Jurgens<sup>2</sup>, Maria Costanzo<sup>2</sup>, Jason Flannick<sup>2,3,4</sup>,  
Helen Parkinson<sup>1</sup>, Noël P. Burtt<sup>2</sup>

1) European Molecular Biology Laboratory, European Bioinformatics Institute, Wellcome Genome Campus, Hinxton, Cambridge, UK. 2) Programs in Metabolism and Medical & Population Genetics, The Broad Institute of MIT and Harvard, Cambridge, MA, USA. 3) Department of Pediatrics, Boston Children's Hospital, Boston, MA, USA. 4) Department of Pediatrics, Harvard Medical School, Boston, MA, USA.

## 1. Background



Genome-wide association studies (GWAS) identify genomic regions (loci) where genetic variation is significantly associated with risk of a disease or magnitude of a trait

- To determine which gene is the most likely the effector gene (i.e. the gene mediating the effect of the trait-associated variant), researchers aggregate and integrate multiple types of evidence
- Effector gene prediction is a major output of post-GWAS analyses, aimed at identifying mechanistically relevant genes and potential drug targets**

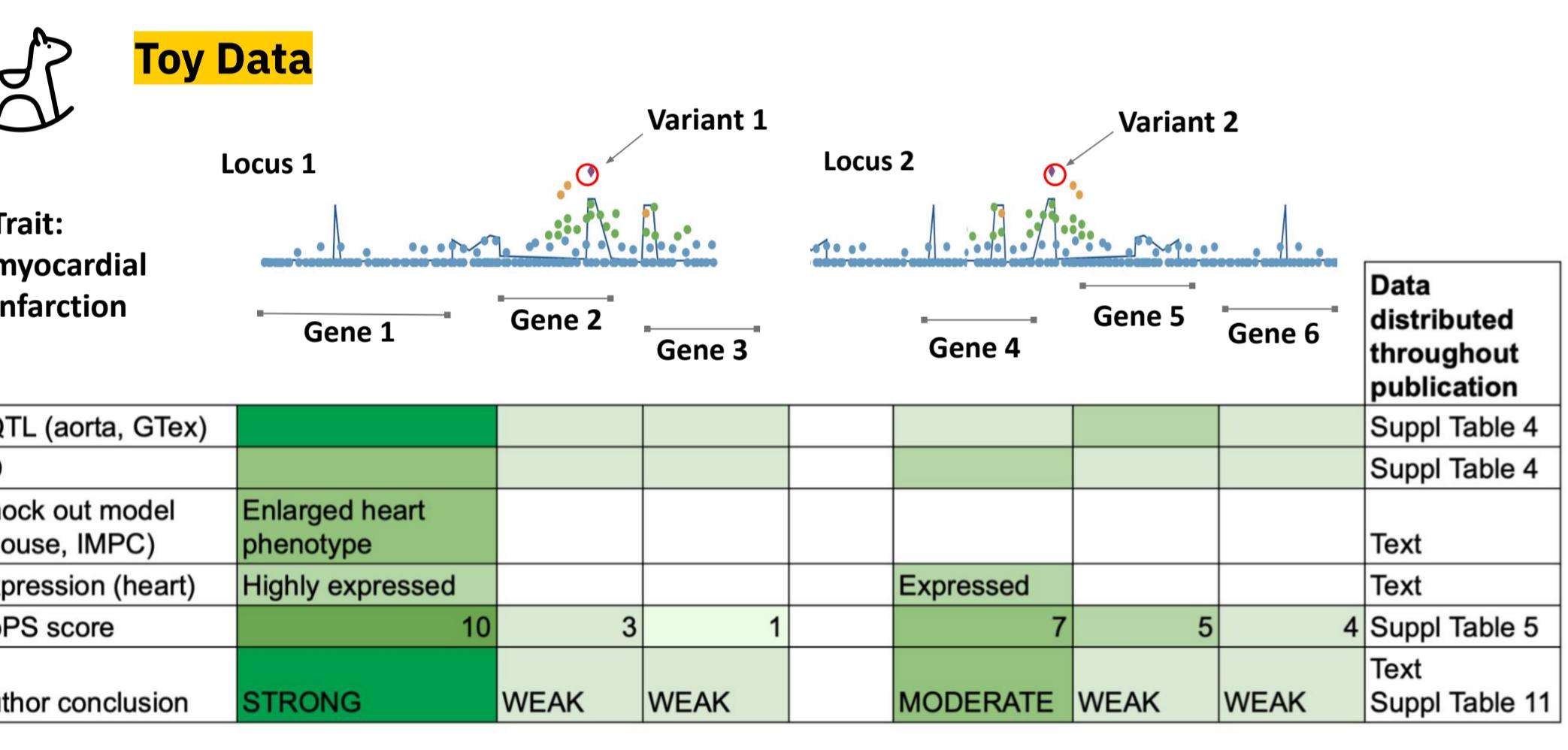


Illustration of the kind of data found in predicted effector gene publications, depth of green = strength of evidence

## 2. Problem

**Broad inconsistency in data content and format of PEG data impairs interpretation, integration and reusability**

- 10% presented as summary images without underlying data
- 29% present evidence for top gene per locus (not all genes considered)
- 19% don't identify the locus under investigation
- 29% use a scoring system to convey a conclusion
- Data distributed throughout publications

From landscape analysis of 169 publications, Costanzo et al, 2025

## 3. Goals

**Aim** - Findable Accessible Interoperable Reusable (FAIR) predicted effector gene data

**Ultimate long term aim** - enable meta analysis to define gold standard lists of genes involved in traits

### Use cases

**Data integration and reuse** – enable computational ingest, submission to knowledgebase, and support AI/ML/KG use.

**Research and hypothesis development** – generate hypotheses, confirm independent findings, prioritize genes for drug targets

**Consideration** balance \*data quality (FAIRness) \*burden on data generator \*need for data curation

## 4. Process

### Activities

#### Community Engagement

Workshop, Sept 2024, 80 attendees, at Broad and EBI

Working Group meetings monthly in 2025; Developed and iterated over a 'strawman' standard, benchmarked to assess suitability.

ASHG Ancillary, Oct 2025

#### Landscape analysis and general recommendations (Costanzo et al)

##### BOX 2

General recommendations emerging from a community workshop on standards and infrastructure for predicted effector gene lists

Provide genes, traits, loci definition, confidence measures. Share data in single machine-readable file. Establish a working group for standards and benchmarks. Create an open-source, FAIR-compliant catalog.



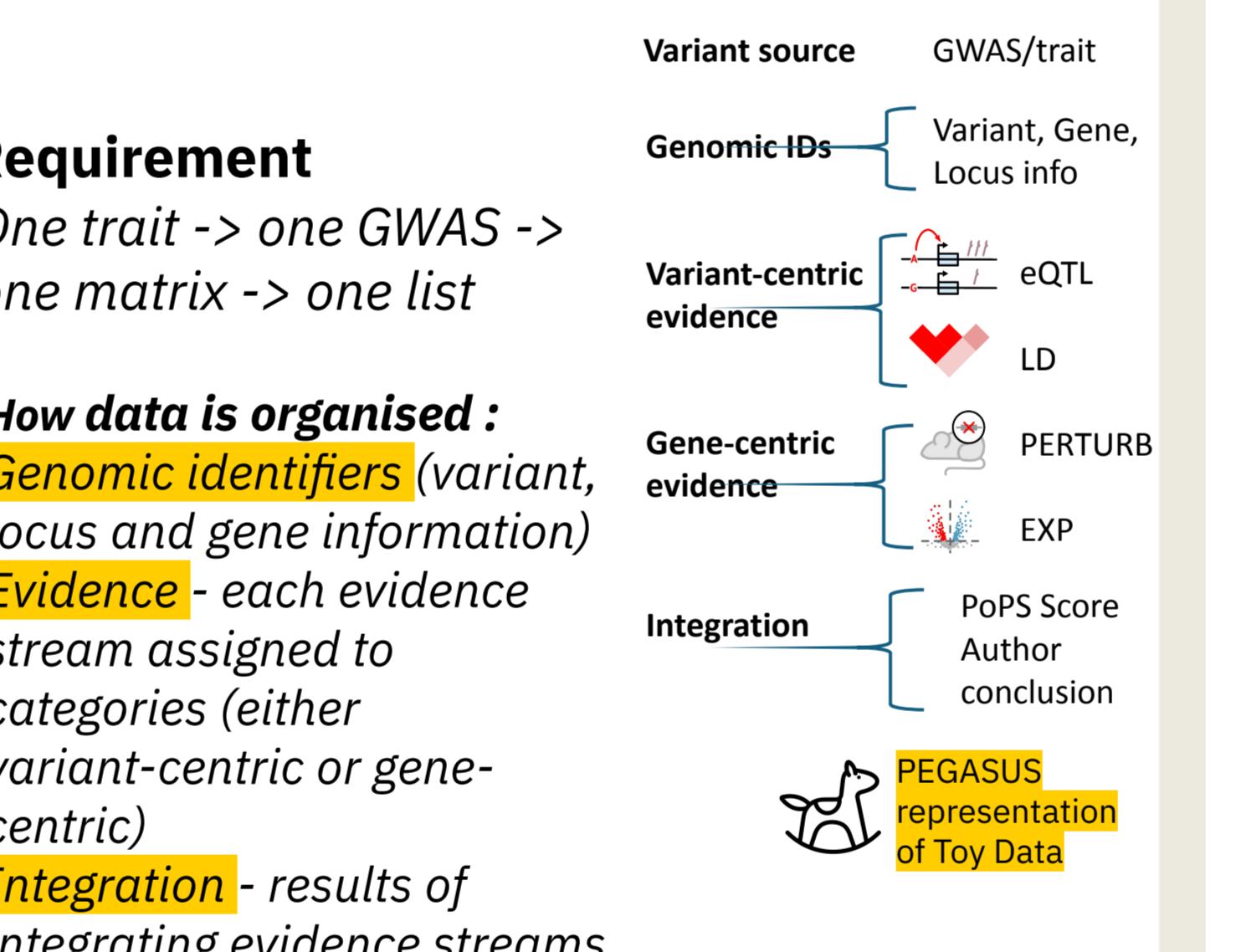
## 5. Proposed Framework

**PEG lists must be supported by an evidence matrix & described by metadata**

**Metadata** → **Evidence matrix** → **List**

1. Trait and provenance of the original GWAS data
2. Sources for each evidence type
3. Description of methods used
4. Column header definitions

1. Includes evidence for all genes considered (not just 'top' genes)
2. Lists the best predictions (top genes)
3. Indicates which evidence categories were included in analysis
4. Includes an author's conclusion (integration)



## 6. Conclusion

### Proposed framework

- facilitates submission to knowledge base (e.g. Predicted Effector Genes Knowledge Portal [pegkp.org](http://pegkp.org)) and linkage to GWAS Catalog
- enhances interpretability and enables interoperability of PEG data
- lays foundations for meta-analysis and generation of gold standard

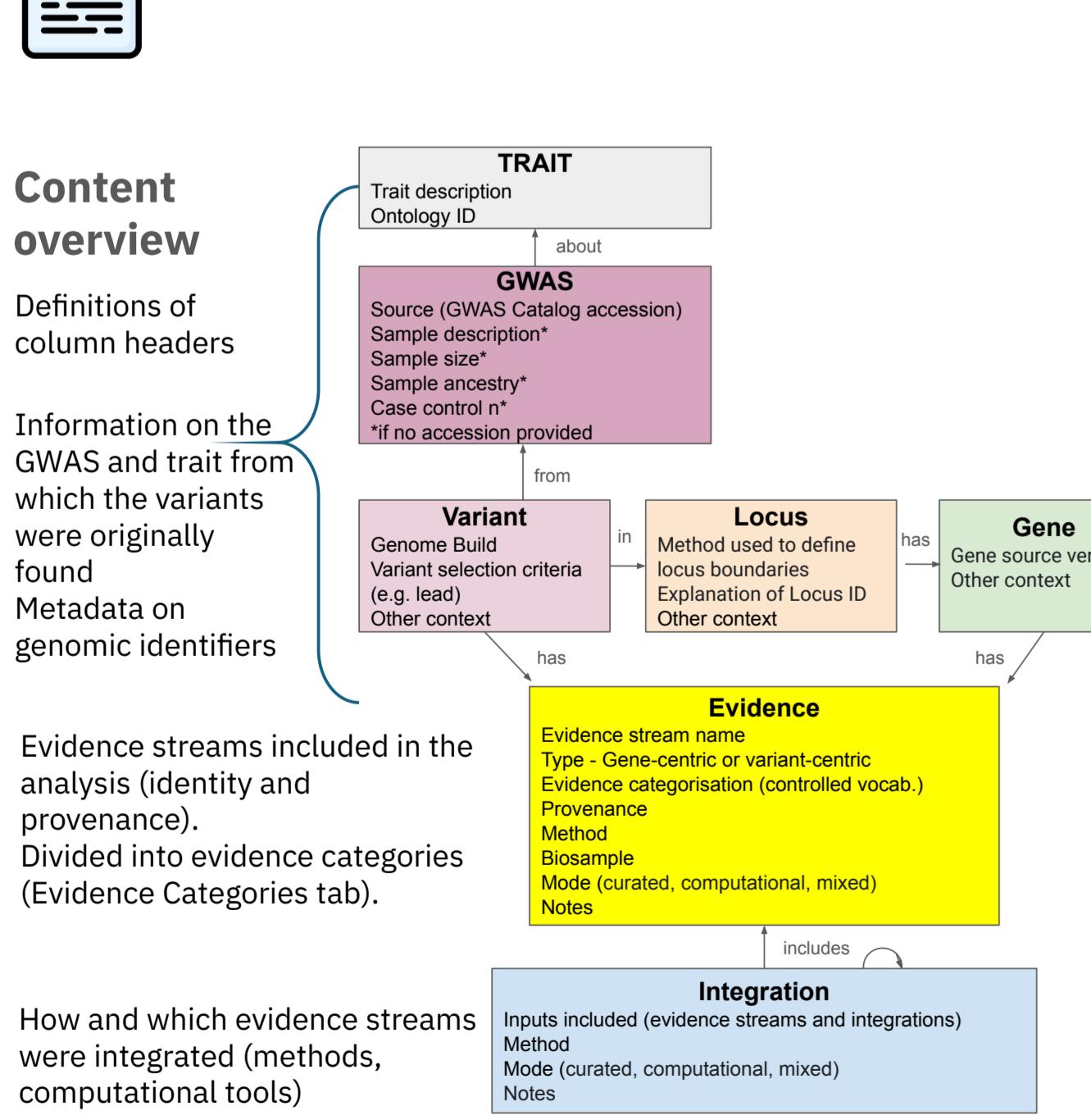
### Next steps

- community feedback on suitability
- uptake

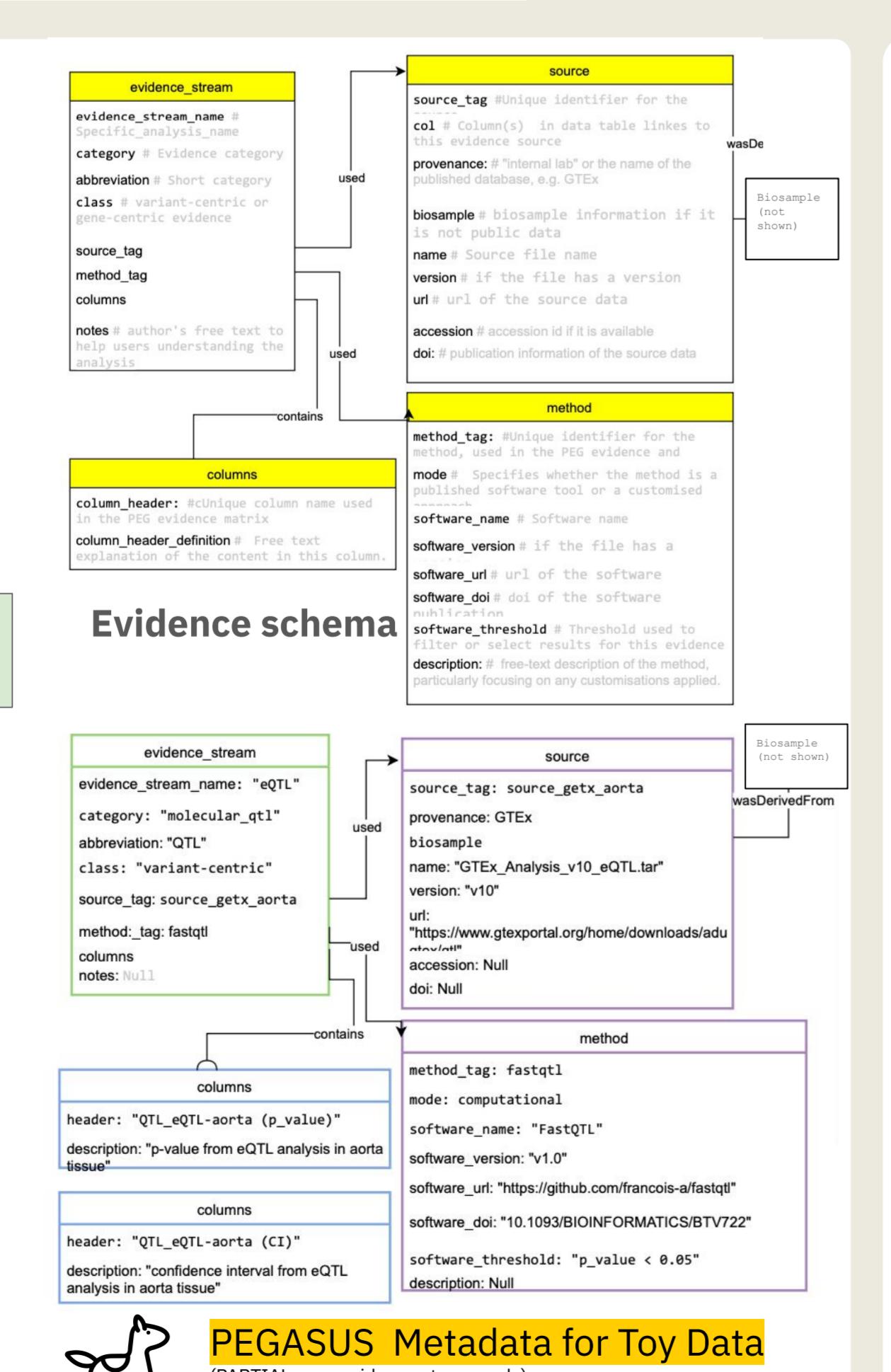
### Join us!

The PEG Working Group is an open community building standards, tools, and FAIR infrastructure for PEG lists. Email: [help@kp4cd.org](mailto:help@kp4cd.org)

## PEG Metadata



Full toy data example including yaml metadata representation (PEGASUS documentation materials)



## PEG Evidence Matrix

### Content

Column header	Variant information		Genomic identifiers		Gene information	Evidence	Integration	
	Primary Variant ID (chr1p)	Primary Variant ID (rsID)	Locus range	Locus ID			Category [xyz]	Category [xyz]
<b>Requirement</b>								
Description	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Method	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Source	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Category	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Method	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Source	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Category	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Method	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Source	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Category	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Method	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Source	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Category	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Method	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Source	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Category	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Method	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Source	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Category	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Method	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Source	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Category	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Method	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Source	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Category	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Method	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Source	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Category	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Method	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Source	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Category	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Method	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Source	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Category	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Method	Variant to which variant-centric evidence relates.		Range around the primary variant considered in this analysis.	Internal or curated to the region considered.	Gene to which gene-centric evidence relates.	Column relating to the evidence, defined in the metadata file.	Column relating to the integration, defined in the metadata file.	
Source	Variant to which variant-centric evidence relates.		Range around the primary variant considered					