

Making Data Deeply FAIR through Lightweight Standards, the KnetMiner and ELIXIR Cases

Marco Brandizi, Rothamsted Research

Nov 20th, 2023, ELIXIR Plant Biology Community Meeting



Let's start from a Use Case...



KnetMiner®

 jupyter KnetMiner_SPARQL_EA Last Checkpoint: 37 minutes ago

File Edit View Run Kernel Settings Help

         Code 

Choose from the list of studies related to the chosen Tax ID:

Study_Ti... 

Run

Study Accession is: E-MTAB-8520

Total Number of Genes in study = 1717

Let's start from a Use Case...

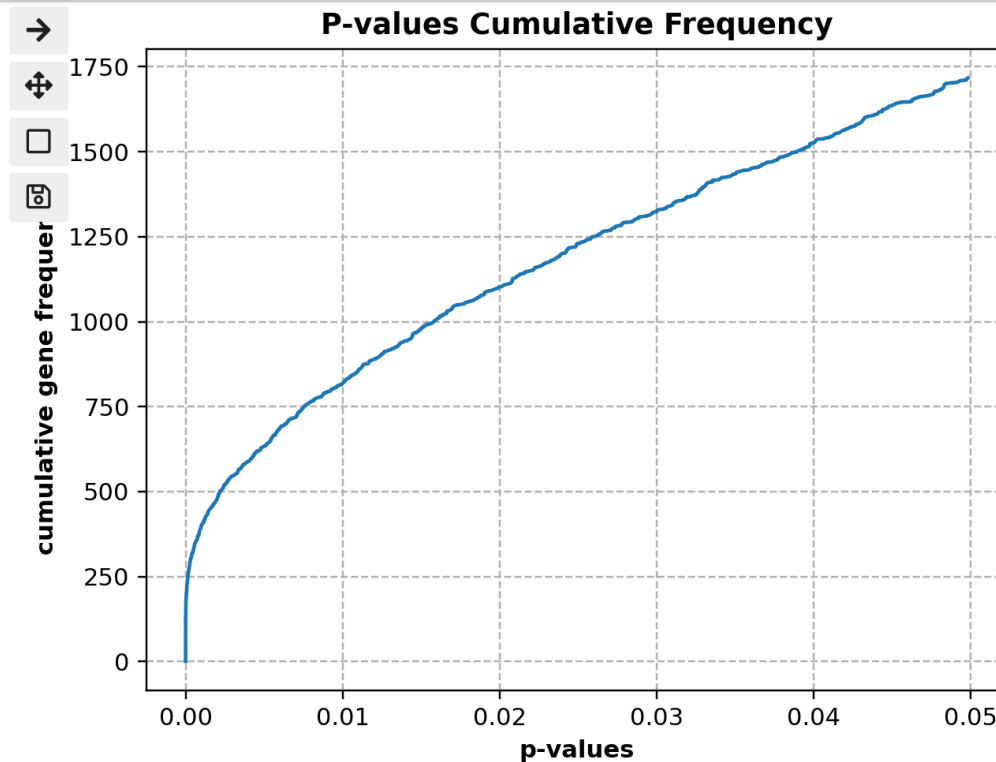


jupyter KnetMiner_SPARQL_EA Last Checkpoint: 39 minutes ago

File Edit View Run Kernel Settings Help

Code

JupyterLa



After choosing the p-value and number of genes with the slider, click the 'Run Analysis' s results.

Pvalues 0.0001

p-value = 0.0001, Number of genes = 228

Run Analysis

Let's start from a Use Case...



The enrichment table below has 153 rows.

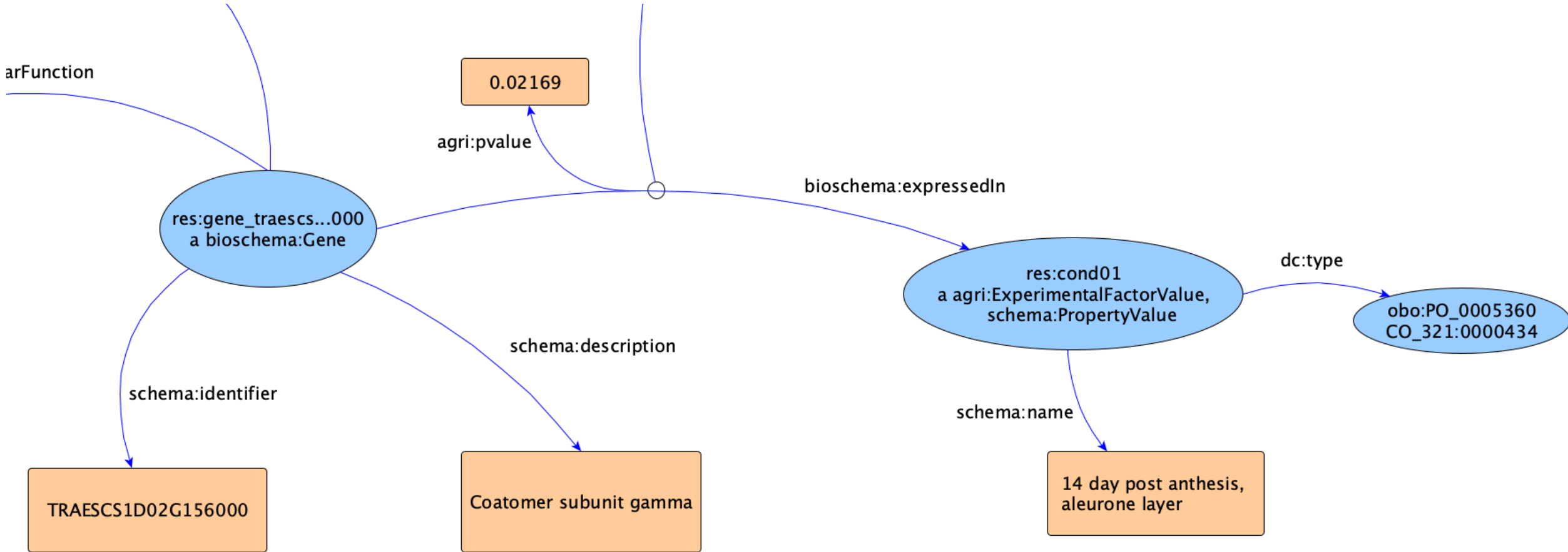
To view the whole table, see the 'View whole tables section' or click on the download link below:

[Download enrichment table CSV file](#)

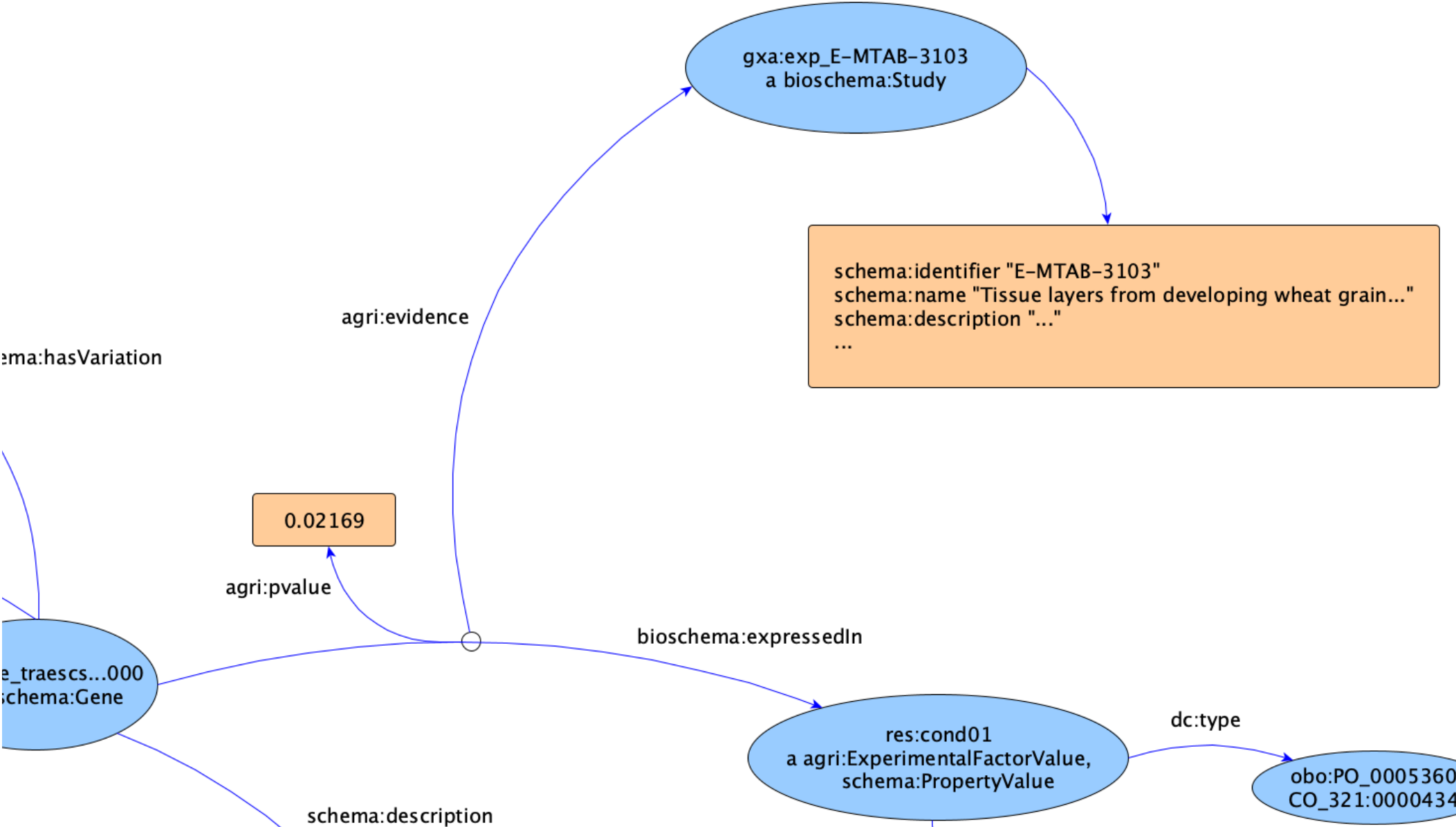
	Ontology Term	Preferred Name	odds ratio	exact p-value	adj p-value	Reference Genes	User/Study Genes
0	TO_0000430	germination rate	17.746196	2.489665e-77	3.809187e-75	5626	107
1	TO_0006002	proline content	10.730036	3.963930e-57	3.032406e-55	8960	107
2	TO_0000276	drought tolerance	5.943980	1.637385e-36	8.350665e-35	16360	112
3	TO_0000190	seed coat color	14.368467	4.736225e-22	1.811606e-20	1150	28
4	TO_0002661	seed maturation	4.699963	4.384546e-16	1.341671e-14	6291	48
5	TO_0000253	seed dormancy	4.632023	5.420875e-14	1.382323e-12	5296	41
6	TO_0000112	disease resistance	2.867433	1.414227e-11	3.091096e-10	15631	70
7	TO_0000043	root morphology trait	2.772096	4.127774e-10	7.894368e-09	13868	62
8	TO_0000495	chlorophyll content	3.268355	6.032035e-10	1.025446e-08	7762	22
9	TO_0000259	heat tolerance	4.252257	1.052225e-10	1.025446e-08	11111	22

Try it! <https://github.com/Rothamsted/knetgraphs-gene-traits>

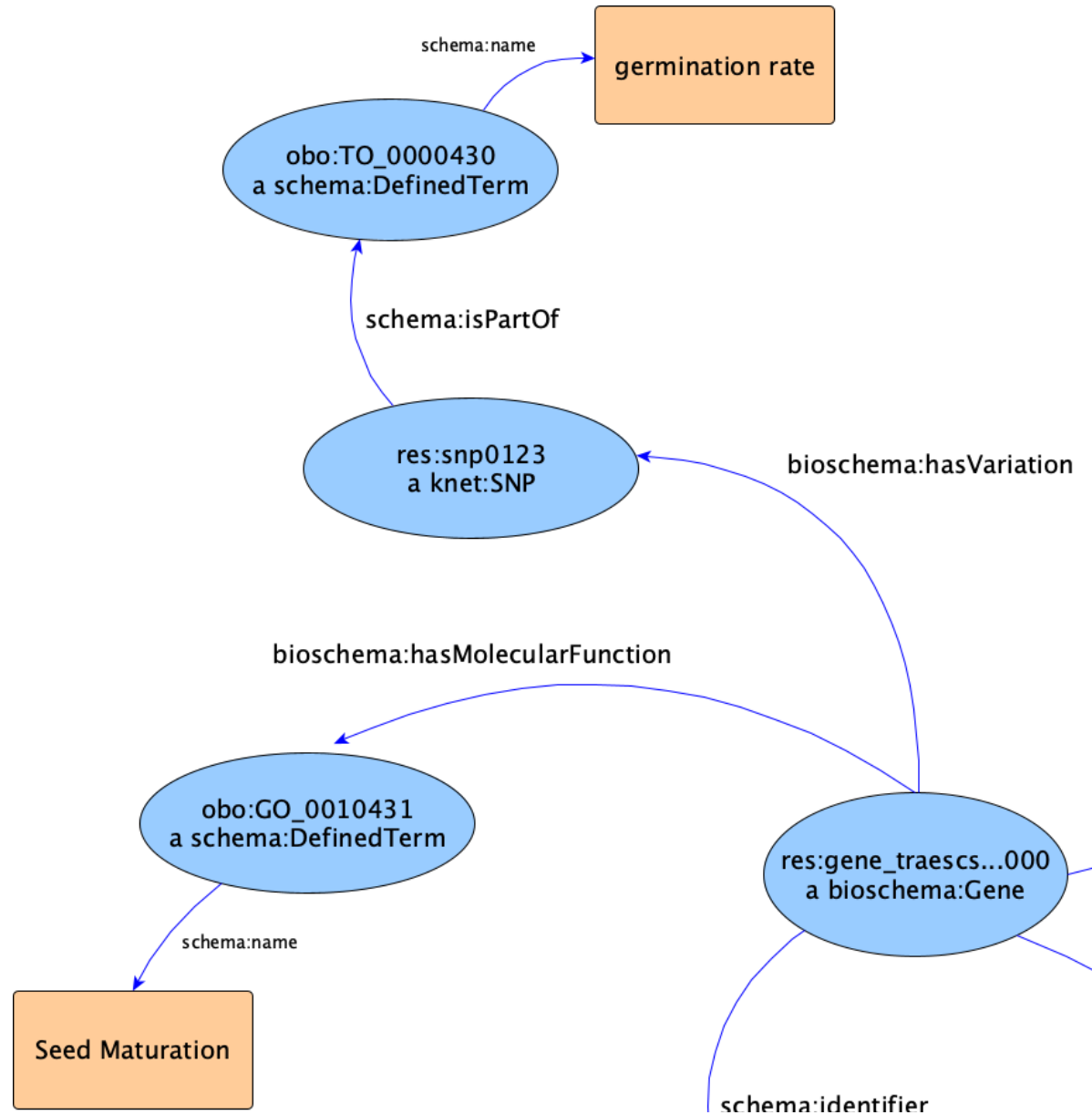
Behind the Scenes



Behind the Scenes



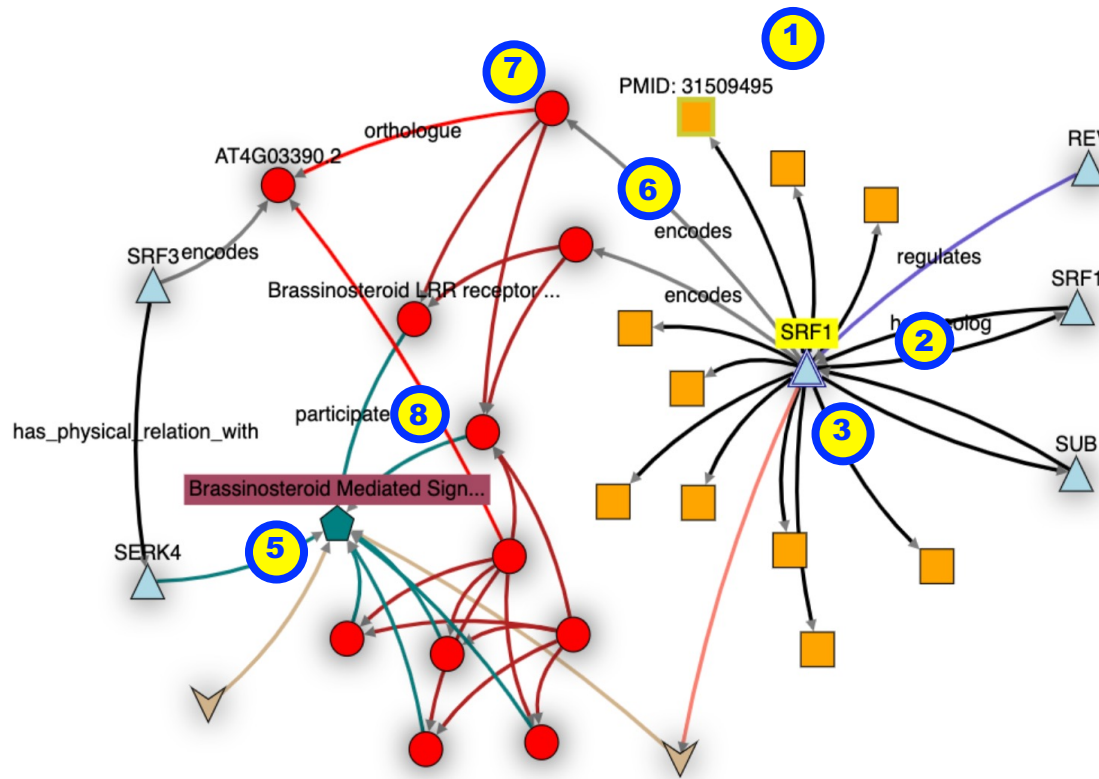
Behind the Scenes



Where it started from...



Based on **publications**, which **genes** are *related* to the **yellow rust disease**? In which **biological processes** are their *encoded* proteins *involved*?



Info box: ✕

1	Concept Type:	Publication
	Source:	NLM
	Evidence:	Imported from database
	Synonyms:	PMID: 31509495
	Attributes:	
	JOURNAL_REF	Plant disease
	YEAR	2019
	AbstractHeader	Genome-Wide Linkage Mapping Reveals Stripe Rust Resistance in Common Wheat (<i>Triticum aestivum</i>) Xinong1376.
	Abstract	Stripe rust, also known as yellow rust , is a significant threat to wheat yield worldwide. Adult plant resistance (APR) is the preferred way to obtain durable protection in these winter wheat cultivars. Xinong1376 has maintained acceptable APR to stripe rust in field environments. To characterize APR in this cultivar, 190 F ₁₀ recombinant inbred lines (RILs) developed from Xiaoyan81 × Xinong1376 were evaluated for infection type and disease severity in fields either artificially or naturally inoculated. The population along with parents were genotyped using the Illumina 90K single-nucleotide polymorphism arrays. Six quantitative trait loci (QTL) were detected using the inclusive composite interval mapping method. <i>Qyr.nwafu-4AL</i> and <i>QYr.nwafu-6BL.3</i> conferred stable resistance in all environments, and likely corresponded to a gene-rich region on the long arm of chromosomes 4A and 6B. <i>QYr.nwafu-5AL</i> , <i>QYr.nwafu-6BL.1</i> , and <i>QYr.nwafu-6BL.2</i> were also detected.

Interactive Legend:

- Gene: 6/131
- Publication: 11/16
- Protein: 11/31
- Trait: 0/299
- SNP: 0/265
- Domain: 0/9
- CoExpCluster: 2/6
- PO: 0/117
- BioProc: 1/656
- MolFunc: 0/37
- CellComp: 0/35
- EC: 0/2
- Phenotype: 0/56
- SNPEffect: 0/1
- CoExpStudy: 0/2

Concepts: 31 (1663); Relations: 45 (4024)

Have a try! knetminer.com

Integrating with EBI GXA



KnetMiner®

Based on **publications**, which **genes** are *related* to the **yellow rust disease**?

Knetminer data (include PubMed, ENSEMBL, Text Mining on mentions)

In which **conditions** *are expressed*?

EBI Gene Expression Atlas

Try it with SPARQL: <https://tinyurl.com/2qerv5wn>



Adopt the Knowledge **graph data model**

- Network of nodes and relationships, each with properties (Ehrlinger, Wöß, 2016)
- Typical features coming from Knowledge representation, formal logics
- Made with a variety of techniques, eg, manual curation, imports, machine learning (Gabrilovich, Usunier, 2016)

Focus on **exploratory research**

- Yes: find gene candidates, interesting articles
- No: perform a precise ANOVA analysis of gene significance, compare gene expression across experiments in a precise way

Coherently, focus on **lightweight schematisation**, integrate with more formal models, ie, OBO ontologies, when useful

- Simple and complementary approach
- Suitable for integrating a high number of heterogeneous datasets, web sources, “noisy” data from the world **wild** web
- Data can be presented directly to the user

Simple, informal



“transmembrane receptor protein serine/threonine kinase activity (GO:0004675)”

subClassOf: “protein serine/threonine kinase activity (GO:0004674)”

AND “transmembrane receptor protein kinase activity (GO:0019199)”

AND Restriction:

onProperty: part of (BFO_0000050)

someValueFrom: “transmembrane receptor protein serine/threonine kinase signalling pathway (GO:0007178)”

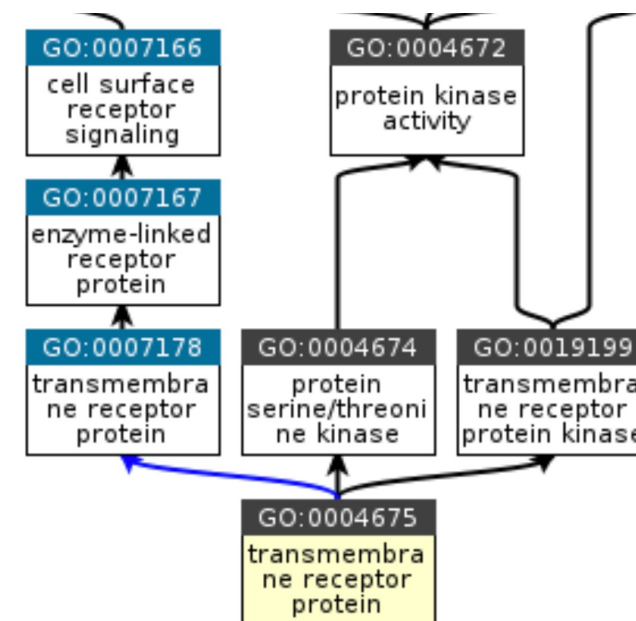
versus:

“transmembrane receptor protein serine/threonine kinase activity (GO:0004675)”

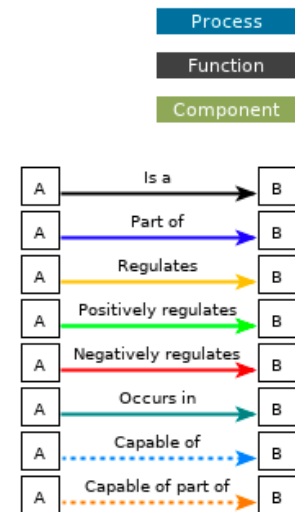
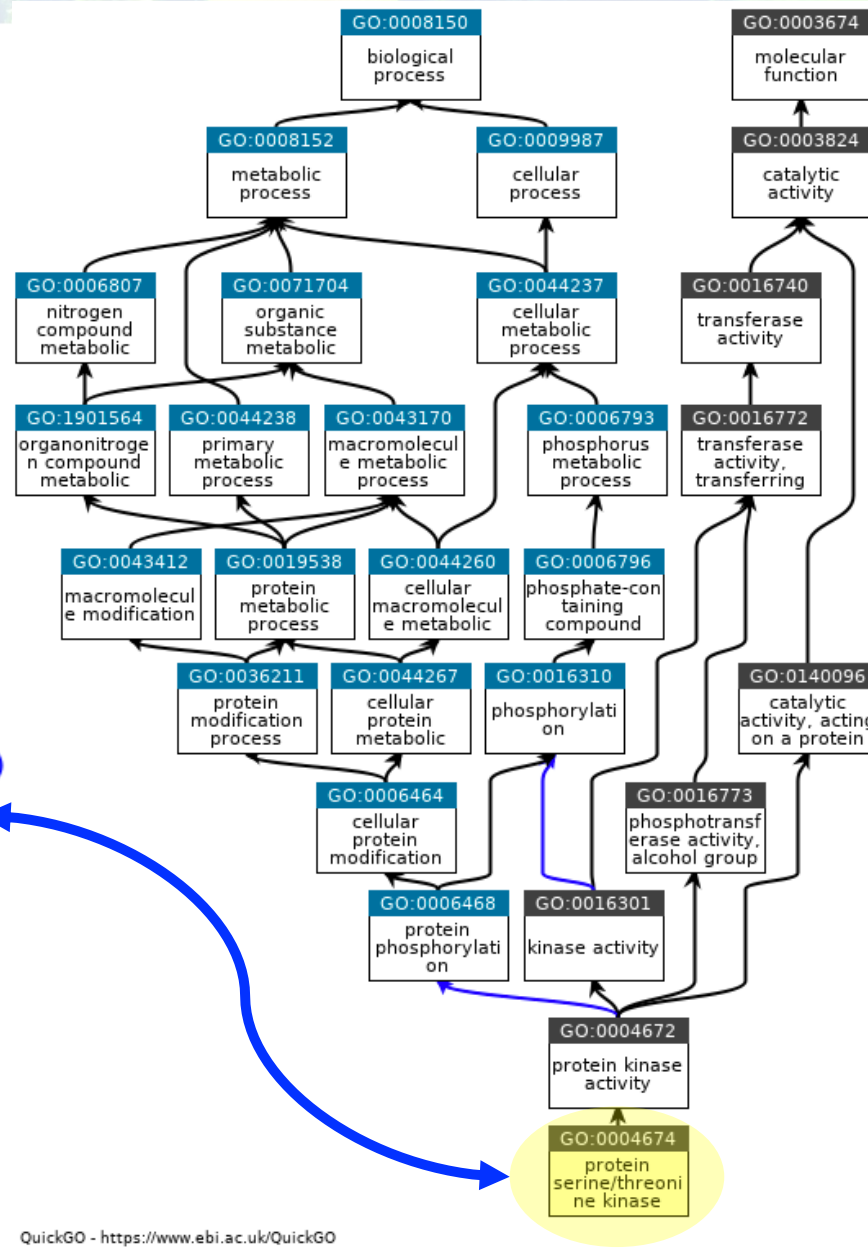
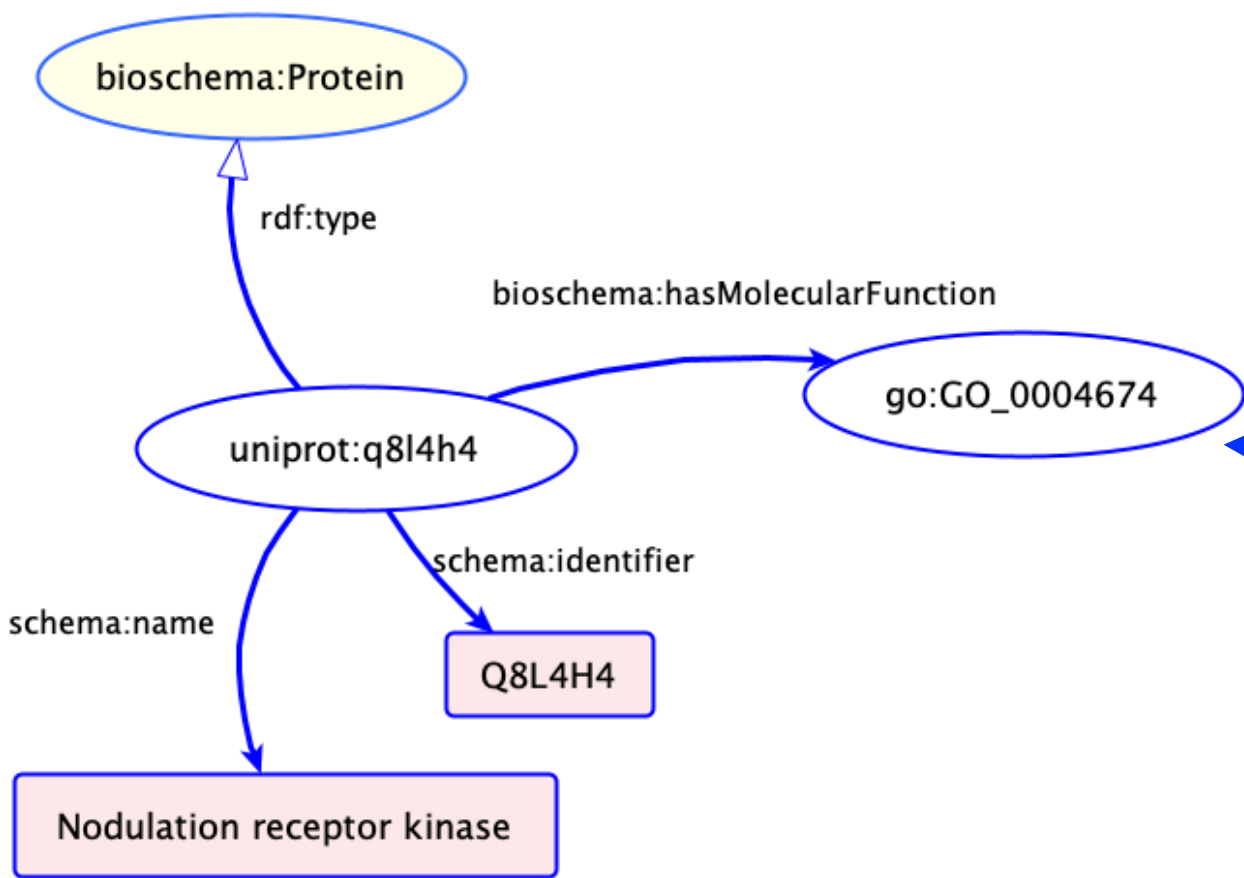
schema:partOf: “transmembrane receptor protein serine/threonine kinase signalling pathway (GO:0007178)”

subClassOf: “protein serine/threonine kinase activity (GO:0004674)”,

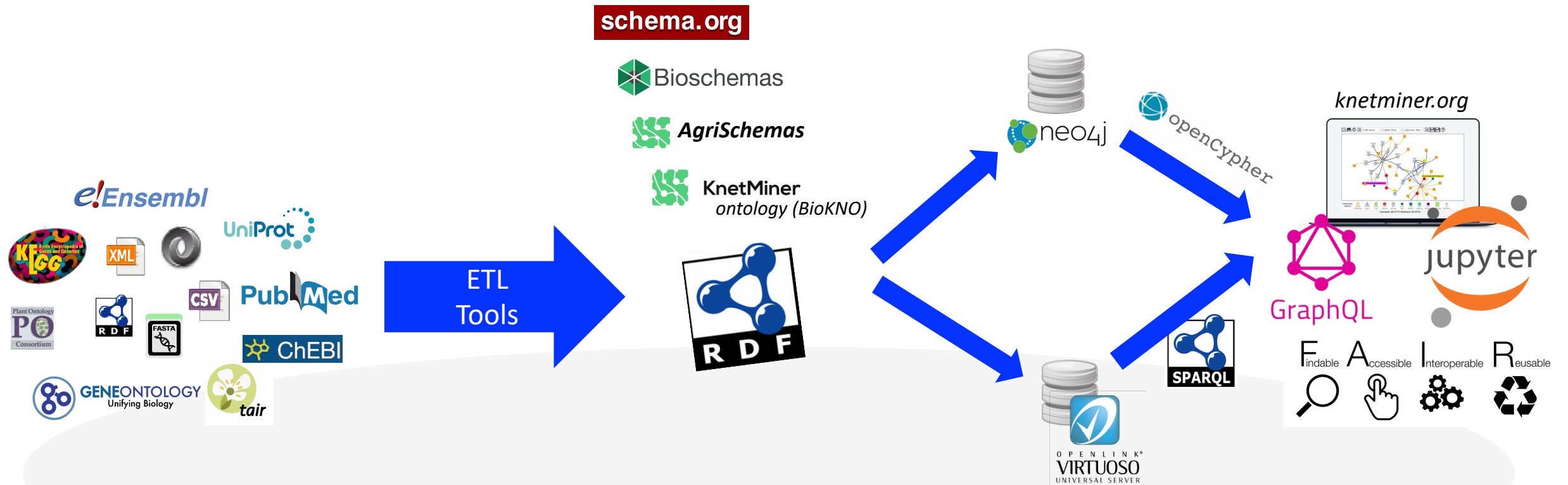
“transmembrane receptor protein kinase activity (GO:0019199)”



Complementary to "real" ontologies



And it can be FAIRer



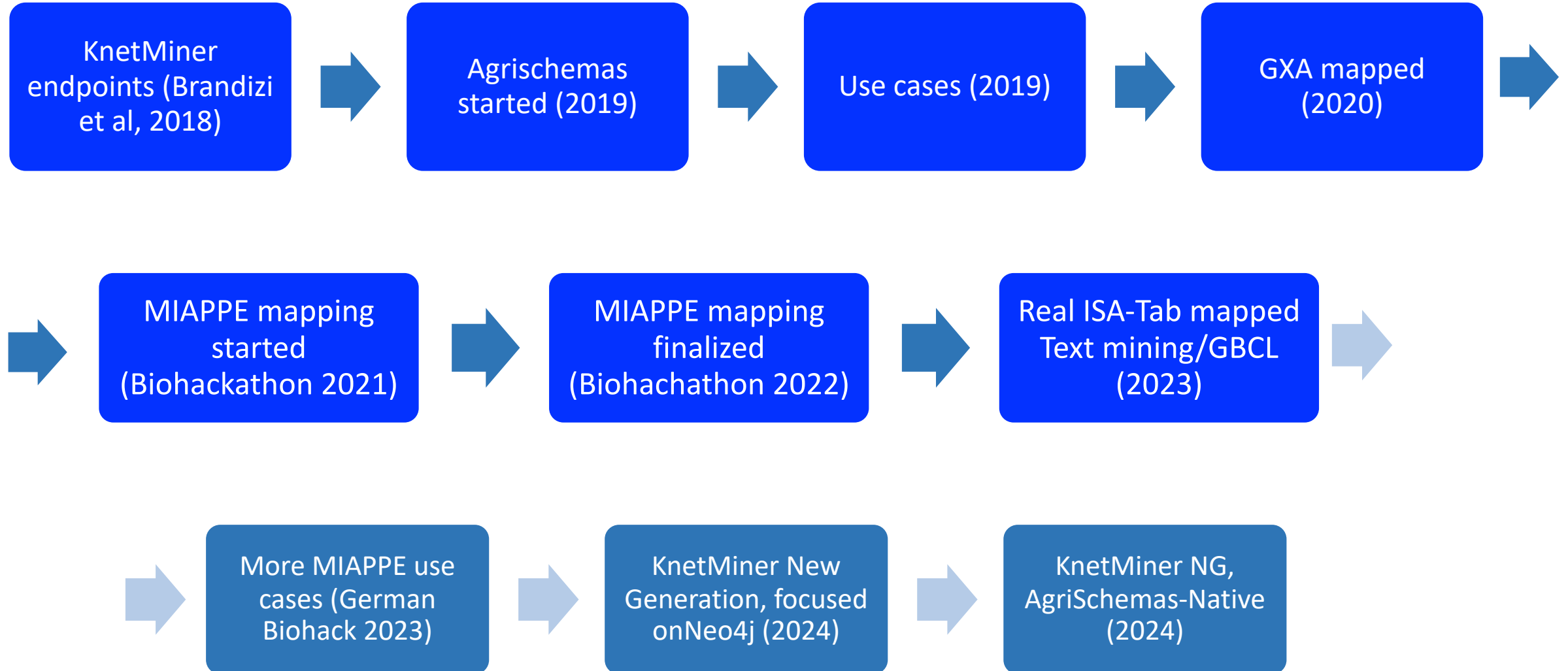
Based on **publications**, which **genes** are *related* to the **yellow rust disease**? In which **biological processes** are their *encoded proteins* *involved*?

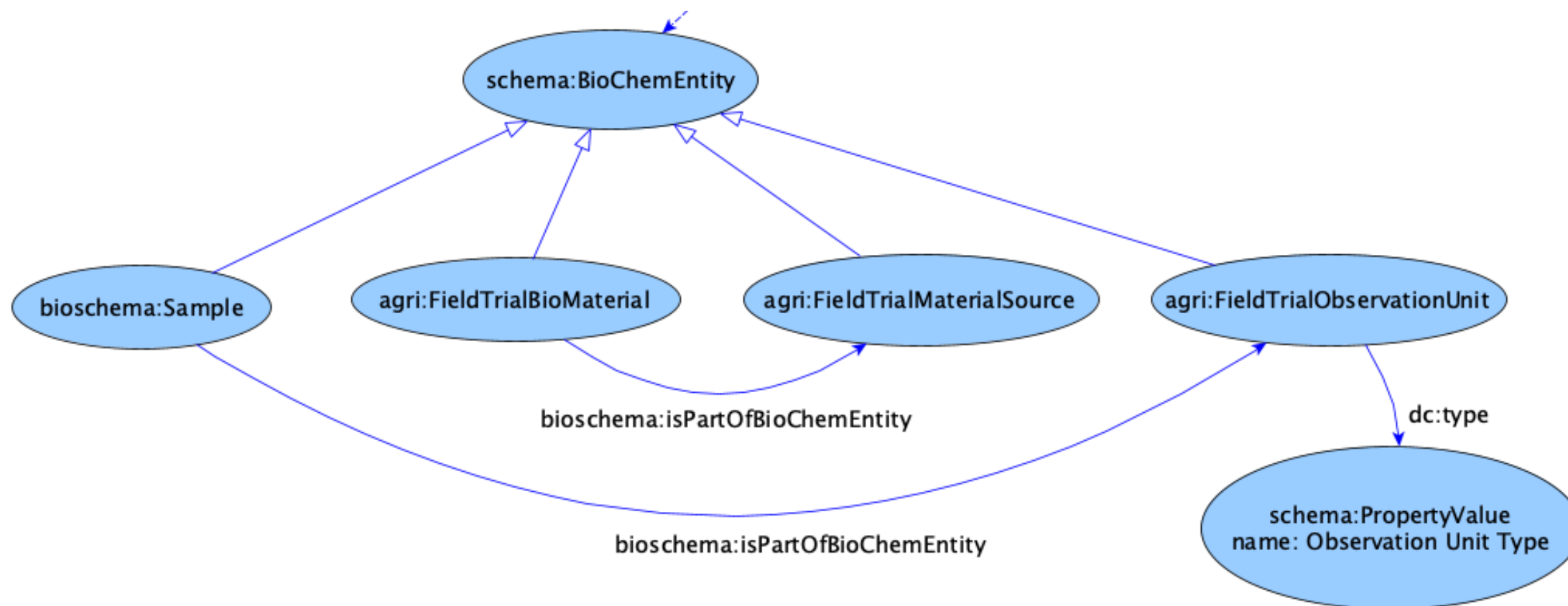
- 🌿 Knowledge Graph Patterns and use cases about plant biology, agronomy, food, forestry, weather
- 🌿 Based on standards, mainly schema.org, bioschemas
- 🌿 Used to integrate knetminer data and other relevant data (more later)
- 🌿 Allows for publishing datasets and data endpoints (knetminer.com/data)
- 🌿 Together with ETL and consuming tools

The story so far



KnetMiner®





Based on field trials, which germplasms have least yield loss under drought? Which genes (variants, markers) are associated to them? Which bioproc, traits, mol functions are most significant?

Acknowledgements

KnetMiner and Rothamsted

- Mennatullah Shehata, Jupyter project lead
- Jeremy Parsons, bioinformatics engineer, ETL author
- Arne De Klerk, product owner
- Keywan Hassani-Pak, team leader and CEO
- Lawal Olaotan, UI developer
- Alumni and past collaborators
- Chris Baker, IDE department director
- Chris Rawling, dept director, consultant
- Brett Drury, GCBL text mining project

Empats, contractor developers

ELIXIR

- Cyril Pommier
- Sebastian Beier
- Daniel Arend

Bioschemas

- Alasdair Gray

Germany Hackathon

- Gabriel Schneider (and FAIRAgro)

