# Workshop Series Strategies to overcome your challenges in multi-omics data integration

Monday 22<sup>nd</sup> June 2020 Linked data in practice: An RDF-based approach with SPARQLing Genomics



# Data standards and multi-omics data integration

- **12.10** General intro
- 12.30 RDF and SPARQL intro and demo
- **13.10** Interactive Quiz
- **13.20** Ontologies and FAIR data
- **13.40** Q&A + Open Discussion









# Monday 22<sup>nd</sup> June 2020 Linked data in practice

#### **General Introduction**





# **General introduction**

- 1. Reasoning
- 2. Data and metadata
- 3. RDF and SPARQL
- 4. Current omics data landscape
- 5. SPARQLing Genomics project







- Interoperability across data, platforms, research lines, disciplines
- Harmonizing data and metadata
- FAIR data requirements
- Reusable and trustworthy data and information





	SPREADSHEET.xlsx Info
	SPREADSHEET.xlsx 8 KB Modified: Today, 15:10
	Add Tags ▼ General: Kind: Microsoft Excel Workbook (.xlsx) Size: 8.209 bytes (12 KB on disk) Where: Macintosh HD + Users + jbohmer + Desktop
	Created: Tuesday, 9 June 2020 at 15:10 Modified: Tuesday, 9 June 2020 at 15:10 Stationery pad Locked
XLS	<ul> <li>▼ More Info:</li> <li></li> <li>▶ Name &amp; Extension:</li> </ul>
	<ul> <li>Comments:</li> <li>Open with:</li> </ul>
	Preview:
	Sharing & Permissions:











	А	В	С	D	E	F
1	Sample	Population	Has Omni Genotypes	Has Axiom Genotypes	Has Affy 6.0 Genotypes	Has Exome/LOF Genotypes
2	HG00096	GBR	1			1
3	HG00097	GBR	1			1
4	HG00098	GBR	1			1
5	HG00099	GBR	1			1



https://www.flaticon.com/authors/smashicons ftp.1000genomes.ebi.ac.uk/vol1/ftp/technical/working/20130606\_sample\_info/20130606\_sample\_info.xlsx









### **METADATA-Standard**



Standardised METADATA:		
Creator:		
Creation Date:		
Gender:		
Sample-ID:		
Data-Type:		
Access Condition:		
:		





### **METADATA-Standard**



Standardised <b>METADATA:</b>		
Creator:	Free Text	
Creation Date:	DD.MM.YYY	
Gender:	CONTROLLED	
Sample-ID:	SYNTAX	
Data-Type:	CONTROLLED	
Access Condition:	CONTROLLED	





# **Semantic ambiguity**

#### How many ways can you say 'female'?

18-day pregnant females	female (lactating)	individual female	worker caste (female)
2 yr old female	female (pregnant)	lgb*cc females	sex: female
400 yr. old female	female (outbred)	mare	female, other
adult female	female parent	female (worker)	female child
asexual female	female plant	monosex female	femal
castrate female	female with eggs	ovigerous female	3 female
cf.female	female worker	oviparous sexual females	female (phenotype)
cystocarpic female	female, 6-8 weeks old	worker bee	female mice
dikaryon	female, virgin	female enriched	female, spayed
dioecious female	female, worker	pseudohermaprhoditic female	femlale
diploid female	female(gynoecious)	remale	metafemale
f	femele	semi-engorged female	sterile female
famale	female, pooled	sexual oviparous female	normal female
femail	femalen	sterile female worker	sf
female	females	strictly female	vitellogenic replete female
female - worker	females only	tetraploid female	worker
female (alate sexual)	gynoecious	thelytoky	hexaploid female
female (calf)	healthy female	female (gynoecious)	female (f-o)
hen	probably female (based on m	orphology)	

female (note: this sample was originally provided as a \"male\" sample to us and therefore labeled this way in the brawand et al. paper and original geo submission; however, detailed data analyses carried out in the meantime clearly show that this sample stems from a female individual)",

EMBL-EBI





https://www.slideshare.net/ConnectedDataLondon/ontology-services-for-the-biomedical-sciences; Courtesy of N. Silvester, European Nucleotide Archive, EMBL- EBI

# Semantic ambiguity and ontology terms

Meta-data element	Preferred ontology term for meta-data	Value types	lssue Number	lssue (open / closed)
Individual ID	NCIT:C164337	ID [string]	#3	closed
Gender	SIO:010029	Male SIO:010048 Female SIO:010052 Unknown / Undetermined	#4	closed
Genotypic sex	PATO:0020000	UNKNOWN_KARYOTYPE, XX, XY, XO, XXY, XXX, XXYY, XXXY, XXXX, XYY, OTHER_KARYOTYPE	#69	closed





# Linked Data with RDF and SPARQL



























# **RDF triplets in a knowledge graph**



A Scale-Out RDF Molecule Store for Improved Co-Identification, Querying and Inferencing - Scientific Figure on ResearchGate. Available from: https://www.researchgate.net/figure/RDF-triples-about-a-yeast-protein\_fig1\_43516365 [accessed 17 Dec, 2019]



### **Current Omics data landscape**

	LC-MS	GO			EFO	
Python	Bash		.mzTAB		.VCF	
Java	R	NCIT		.NMR		DUO
MIABIS	C++		.BAM		.GC-MS	
.mzXML		.FAST	Q	Perl	DUBLIN CORE	
	LISP		.dta		Scheme	
ISO					RI	<b>)F</b>
	SI	NOMED				





# **SPARQLing Genomics – Pilot 2019**





https://www.internationalgenome.org/wiki/Analysis/vcf4.0/ "Graph of klick.jÖrg - Die Homepage" by jÖrg is licensed under CC BY-SA 2.0 https://www.sparqling-genomics.org/



# **SPARQLing Genomics – Pilot 2019**

## Todividuals
##
:vcf2rdf-0.99.9
rdf:type :SoftwareProgram :
rdfs:label "vcf2rdf" ·
referencement "Tool to convert VCE files to DDE "
Turs: comment foot to convert ver files to KDF
:table2rdf-0.99.9
rdf:type :SoftwareProgram ;
rdfs:label "table2rdf" ;
rdfs:comment "Tool to convert tabular data to RDF." .
:folder2rdf-0.99.9
rdf.tyne ·SoftwareProgram ·
rdfs.jobol "foldoroddf" -
ruts tabet Totterzitut ;
rdfs:comment "lool to extract RDF from files in a directory." .
:sg-web-0.99.9
<pre>rdf:type :SoftwareProgram ;</pre>
rdfs:label "SPARQLing genomics web interface";
rdfs:comment "Tool for guerving multiple SPAROL endpoints using a web interface".

VCF2RDF TABLE2RDF XML2RDF JSON2RDF

....



https://www.sparqling-genomics.org/ https://dublincore.org/ http://geneontology.org/ https://github.com/EBISPOT/DUO

#### Dublin Core Metadata Initiative







# **SPARQLing Genomics – Portal 2020**

sparqling-genomics	
Dashboard     Manual     +     Log out       Overview     →     Collect     →     Structure     →     Query     →     Automate	spargling-genomics
	version 0.99.11, June 19, 2020
Name #Queries Actions	
roel 1	
Assigned graphs +	
Members of this project have access to the following graphs.	
Graph Actions	
https://manual.sparqling-genomics.org (moonstone)	
Queries	
Query         Connection         Executed by         Duration (in seconds)         Actions ( <u>Remove unselected</u> )	
SELECT 7s 7p 7o FROM <https: manual.sparqling-genomics.org=""> { ?s ?p ?o }     moonstone     roel     0</https:>	
v0.99.11   <u>Source code</u>	









# Monday 22<sup>nd</sup> June 2020 Linked data in practice

#### **RDF intro and SPARQLing Genomics demo**





# Insights from the first XO workshop

Go to www.menti.com and use the code 45 38 35

#### Mentimeter Which skills are needed for Big data and multiomics data integration? Knowledge on data 1st structure 2nd **Programming Skills** Knowledge on common 3rd experimental design Following multi-omics 4th analytical developments Identification burdens and 5th alaorithmic details 6th Statistical Knowledge 25







# Monday 22<sup>nd</sup> June 2020 Linked data in practice

#### **QUIZ TIME!**









# Monday 22<sup>nd</sup> June 2020 Linked data in practice

# We need to talk about -data standards -vocabularies ontologies





# **Back to the semantic ambiguity**

#### How many ways can you say 'female'?

18-day pregnant females	female (lactating)	individual female	worker caste (female)
2 yr old female	female (pregnant)	lgb*cc females	sex: female
400 yr. old female	female (outbred)	mare	female, other
adult female	female parent	female (worker)	female child
asexual female	female plant	monosex female	femal
castrate female	female with eggs	ovigerous female	3 female
cf.female	female worker	oviparous sexual females	female (phenotype)
cystocarpic female	female, 6-8 weeks old	worker bee	female mice
dikaryon	female, virgin	female enriched	female, spayed
dioecious female	female, worker	pseudohermaprhoditic female	femlale
diploid female	female(gynoecious)	remale	metafemale
f	femele	semi-engorged female	sterile female
famale	female, pooled	sexual oviparous female	normal female
femail	femalen	sterile female worker	sf
female	females	strictly female	vitellogenic replete female
female - worker	females only	tetraploid female	worker
female (alate sexual)	gynoecious	thelytoky	hexaploid female
female (calf)	healthy female	female (gynoecious)	female (f-o)
hen	probably female (based on mor	phology)	

female (note: this sample was originally provided as a \"male\" sample to us and therefore labeled this way in the brawand et al. paper and original geo submission; however, detailed data analyses carried out in the meantime clearly show that this sample stems from a female individual)",







### **RDF version**







### **RDF version**







# SPARQL query to select "female"







# FAIR genome definition of gender

Meta-data element	Preferred ontology term for meta-data	Value types	lssue Number	lssue (open / closed)
Individual ID	NCIT:C164337	ID [string]	#3	closed
Gender	SIO:010029	Male SIO:010048 Female SIO:010052 Unknown / Undetermined	#4	closed
Genotypic sex	PATO:0020000	UNKNOWN_KARYOTYPE, XX, XY, XO, XXY, XXX, XXYY, XXXY, XXXX, XYY, OTHER_KARYOTYPE	#69	closed





# **Using SPARQL to ontologise**







### Unifying omics data one term at a time







# FAIR data with linked data







# FAIR data work around



Subjects		optional 💙
Specify subjects from a taxonomy or controll	led vocabulary. Each term must be uniquely identified	l (e.g. a URL). For free form text, use the keywords field in basic information section.
Subjects	female germ-line cyst formation	G0:0048135 🗢 🗙
	Term	GENEONTOLOGY ×
	+ Add another subject	Unifying Biology







# Find standards on FAIRsharing.org







# SUMMARY

- We can create a knowledge graph of diverse omics data using RDF
- Ontologies help unifying terminology across research domains
- We can transform a spreadsheet into ontologically defined data
- Flexible data management requires easy data transformations, which can be done with SPARQL
- SPARQLing Genomics is a starting point for interoperable omics data





# Monday 22<sup>nd</sup> June 2020 Linked data in practice

### Q&A Open Discussion









This work is licensed under a Creative Commons Attribution 4.0 International License.