Assignment 5 SPARQL queries

Authors: David García Valcarce and Jimena Martín Reina

I would like you to create the SPARQL query that will answer each of these questions. Please submit the answers as screenshots of the query plus the first 2-3 lines of response. Submit to Moodle. NO programming is required! Use whatever SPARQL client interface you want (Jupyter, Yasgui, etc.) Thanks!

For many of these you will need to look-up how to use the SPARQL functions 'COUNT' and 'DISTINCT' (we used 'distinct' in class), and probably a few others...

<u>UniProt SPARQL Endpoint: http://sparql.uniprot.org/sparql (note that you need to configure the endpoint to GET if you're using YASGUI)</u>

Q1: 1 POINT How many protein records are in UniProt?

```
Your SPARQL query

Add common prefixes

1  PREFIX X++
64  SELECT (COUNT(?protein) AS ?count)
65  WHERE {
?protein a up:Protein .
}

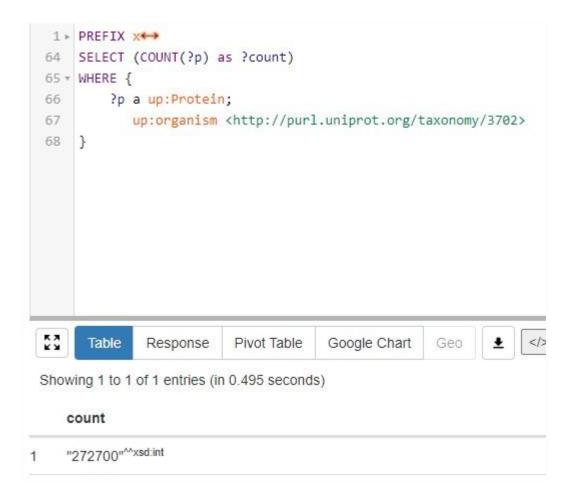
Submit Query
```

Results

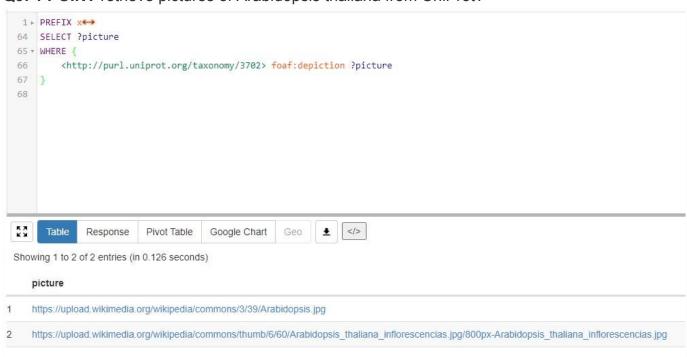
Sparql XML Sparql JSON CSV Show query Share

count

"412838422"xsd:int



Q3: 1 POINT retrieve pictures of Arabidopsis thaliana from UniProt?



Q4: 1 POINT: What is the description of the enzyme activity of UniProt Protein Q9SZZ8

Your SPARQL query Add common prefixes | PREFIX X++ | | 64 | | 65 | | 55 | | 66 | | WHERE { | WALUES | Protein | < | http://purl.uniprot.org/uniprot/095ZZ8> } | 7 | | 7 | | 7 | | 7 | | 7 | | 7 | | 7 | | 7 | | 7 | | 8 | | 7 | | 7 | | 7 | | 7 | | 8 | | 9 | | 9 | | 1 | | 7 | | 7 | | 7 | | 8 | | 9 | | 9 | | 1 | | 7 | | 7 | | 8 | | 9 | | 9 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | | 1 | |

Results

Sparql XML Sparql JSON CSV Show query Share

descriptior

"Nonheme diiron monooxygenase involved in the biosynthesis of xanthophylls. Specific for beta-ring hydroxylations of beta-carotene. Has also a low activity toward the beta- and epsilon-rings of alpha-carotene. No activity with acyclic carotenoids such as lycopene and neurosporene. Uses ferredoxin as an electron donor."xsd:string

Q5: 1 POINT: Retrieve the proteins ids, and date of submission, for 5 proteins that have been added to UniProt this year (HINT Google for "SPARQL FILTER by date")

Results

Sparql XML Sparql JSON CSV Show query Share

proteinID	submissionDate
http://purl.uniprot.org/citations/36876065❤	"2023"xsd:gYear
http://purl.uniprot.org/citations/35942639*	"2023" ^{xsd:g} Year
http://purl.uniprot.org/citations/SIPFCF38EF708ADD808♥	"2023-03"xsd:gYearMonth
http://purl.uniprot.org/citations/35640876	"2023"xsd:gYear
http://purl.uniprot.org/citations/36328097	"2023" ^{xsd:g} Year

Q6: 1 POINT How many species are in the UniProt taxonomy?

```
1 ▶ PREFIX x↔
64 SELECT (COUNT(DISTINCT ?species) as ?count)
65 ▼ WHERE {
66    ?species a up:Taxon
67 }
```

count

"2887005"xsd:int

Q7: 2 POINTS How many species have at least one protein record? (this might take a long time to execute, so do this one last!)

```
Your SPARQL query

Add common prefixes

1 PREFIX X SELECT (COUNT(DISTINCT ?0) AS ?count)
64
65 WHERE {
7 p a up:Protein;
up:organism ?0 .
}

Submit Query
```

Results

Sparql XML Sparql JSON CSV Show query Share

count "1305985"^{xsd:int} **Q8: 3 POINTS** find the AGI codes and gene names for all Arabidopsis thaliana proteins that have a protein function annotation description that mentions "pattern formation"

```
1 ► PREFIX X↔
64 SELECT DISTINCT ?geneName ?agi ?description
65 ▼ WHERE {
    ?protein a up:Protein ;
66
             up:organism <http://purl.uniprot.org/taxonomy/3702>;
              up:encodedBy ?gene ;
69
             up:annotation ?annotation .
    ?annotation a up:Function Annotation ;
70
                rdfs:comment ?description .
71
72 FILTER(CONTAINS(?description, "pattern formation"))
73
     ?gene skos:prefLabel ?geneName ;
          up:locusName ?agi .
74
75 }
```

geneName	agi	description
	ved-stri	"Transcription factor that is specifically required for the development of root hairs (PubMed:17556585). Acts with RSL1 to positively regulate root hair development (PubMed:17556585). Acts downstream of genes that
RHD6" ^{xsd:string}	"At1g664/0"*****	Isregulate epidermal pattern formation, such as GL2 (PubMed:17556585). Targets directly RSL4, another transcription factor involved in the regulation of root hair elongation (PubMed:20139979). Acts with RSL1 as transcription factor that integrates a Jasmonate (JA) signaling pathway that stimulates root hair growth (PubMed:31988260). "MSd.string"
		canscription (action in integraces a jointonic (pA) signaling partners (that signaling partners) that summates (not that grown (rubinets) 2700200). "Activates the ARF proteins by exchanging bound GDP for free GTP, Plays a role in vesicular protein sorting, Acts as the major regulator of endosomal vesicle trafficking but is also involved in the endocytosis process. Coul
		function redundantly with SNL1 in the retrograde Golgi to endoplasmic reticulum trafficking required for the coordinated polar localization of auxin efflux carriers which in turn determines the
N"XSd:string	"Δ+1σ13980"Xsd:strii	Redirection of auxin flow. Mediates the sorting of PIN1 from endosonal compartments to the basal plasma membrane and the polarization of PIN3 to the bottom side of hypocotyl endodermal cells. Involved in the specifica
	Aligi3700	of apical-basal pattern formation in the early embryo and during root formation. Required for correct cell wall organization leading to normal cell adhesion during seedling development. Also plays an essential role in
		hydrotropism of seedling roots.**xd-string
PK1"xsd:string	"At1g69270"xsd:strii	"Involved in the main abscisic acid-mediated (ABA) signaling pathway and in early ABA perception. Together with RPK2, required for pattern formation along the radial axis (e.g. the apical embryonic domain cell types that
		generate cotyledon primordia), and the apical-basal axis (e.g. differentiation of the basal pole during early embryogenesis).
		"Component of the cullin-RING ublquitin ligases (CRL), or CUL3-RBX1-BTB protein E3 ligase complexes which mediate the ublquitination and subsequent proteasomal degradation of target proteins. The functional specific
"CUL3B"xsd:string	"At1g69670"xsd:strii	of the CRL complex depends on the BTB domain-containing protein as the substrate recognition component. Involved in embryo potent formation and endosperm development. Required for the normal division appearance of the contract the college of the c
	-	organization of the foot stem cens and columnar foot cap cens, regulates primary foot growth by an unknown pathway, but in an ethylene-treperiterit manner, runctions in distart out patterning, by an ethylene-independent
		mechanism. Functionally redundant with CUL3A.**Sd.string
PK2"XSd:string	"At3g02130"xsd:strii	"Key regulator of anther development (e.g. lignification pattern), including tapetum degradation during pollen maturation (e.g. germination capacity). Together with RPK.1 required for pattern formation along the radial axis
		the apical embryonic domain cell types that generate cotyledon primordia), and the apical-basal axis (e.g. differentiation of the basal pole during early embryogenesis).
		"Component of the cullin-RING ubliquin ligases (CRL), or CUL3-RBX1-BTB protein E3 ligase complexes which mediate the ubliquitination and subsequent proteasomal degradation of target proteins. The functional speci
:UL3A"xsd:string	"At1g26830"xsd:strii	of the CRL complex depends on the BTB domain-containing protein as the substrate recognition component. Involved in embryo pattern formation and endosperm development between the account of the CRL complex depends on the BTB domain-containing protein as the substrate recognition component. Involved in embryo pattern formation and endosperm development development proteins and involved in embryo pattern formation and endosperm development development. Browned and account of the contraction of the contra
		organization of the root stem cells and columnia foot cap cells. Regulates primary root growth by an unknown pathway, but in an ethylene dependent manner, i directions in distant out patterning, by an ethylene independent
		mechanism. Functionally redundant with CUL3B. ^{excd string}
TL3"xsd:string	"At2g42580"xsd:strii	"Involved in osmotic and salt stress tolerance. May play a role in the control of meristematic cell size during osmotic stress. May function as an adapter protein for BRL2 and may be required for signaling affecting leaf vas
AMT1"XSd:string	"At5g55250"xsd:strii	"Catalyzes the methylation of the free carboxyl end of the plant hormone indole-3-acetic acid (IAA). Converts IAA to IAA methyl ester (MeIAA). Regulates IAA activities by IAA methylation. Methylation of IAA plays an
	0	important role in regulating plant development and auxin homeostasis. Required for correct leaf pattern formation. MeIAA seems to be an inactive form of IAA.
		"Probable transcription factor involved in cell specification and pattern formation during embryogenesis. Binds to the L1 box DNA sequence 5'-TAAATG[CT]A-3'. Plays a role in maintaining the identity of L1 cells, possibly
TML1"xsd:string	"At4g21750"XSd:Strii	Islinteracting with their L1 box or other target-gene promoters; binds to the LIP1 gene promoter and stimulates its expression upon imbibition (PubMed:24989044). Acts as a positive regulator of gibberellins (GAs)-regulated
		epidermal gene expression (e.g. LIP1, LIP2, LTP1, FDH and PDF1) (PubMed:24989044). Functionally redundant to PDF2 (PubMed:24989044). Seems to promote cell differentiation (PubMed:25564655). ****dd-string**
		"Functions in a MAP kinase cascade that acts as a molecular switch to regulate the first cell fate decisions in the zygote and the early embryo. Promotes elongation of the zygote and development of its basal daughter cell
'DA"XSd:string	"At1g63700"xsd:strii	the extra-emproyanic suspensor. In stomatal development, acts downstream of the LRR reception with the extra-emproyanic suspensor. In stomatal development, acts downstream of the LRR reception of the MKKK5-MPK3/MRK5-module to regulate operated legislate before cell
	8	is specified. Plays a certifiar for in both guard certification, which shapes the
		morphology of plant organs. Upon brassinosteroid signaling, is inhibited by phosphorylation of its auto-inhibitory N-terminal domain by the GSK3-like kinase ASK7."ssd.string
EX1"xsd:string	"At3g09090" ^{xsd:strii}	¹⁸⁻ Required for exine pattern formation during pollen development, especially for primexine deposition.**Sd ³ 5tring
		"Transcription factor required for quiescent center cells specification and maintenance of surrounding stem cells, and for the asymmetric cell division involved in radial pattern formation in roots. Essential for cell division in
CR"xsd:string	"At3g54220"xsd:strii	Isnot differentiation of the ground tissue. Also required for normal shoot gravitropism. Regulates the radial organization of the shoot axial organs. Binds to the promoter of MGP, NUC, RLK and SCL3. Restricts SHR movmer
		sequesters it into the nucleus of the endodermis." vsd.string
ODC A Dawxsd-stri	DENALO - 4 4 74 ONXSd-Strip	**Sequesces is time the nucleus of the endouerims.** Sequesces is time the nucleus of the
OPGAP3		
5L1"xsd:string	"A+E -27000"Xsd-strii	"Transcription factor that is specifically required for the development of root hairs (PubMed:17556585). Acts with RHD6 to positively regulate root hair development (PubMed:17556585). Acts downstream of genes that
DL1	AL3g3/800	eragulate epidermal pattern formation, such as GL2 (PubMed:17556585). Acts with RHD6 as transcription factor that integrates a jasmonate (JA) signaling pathway that stimulates root hair growth (PubMed:31988260). ***S
		"Guanine-nucleotide exchange factor (GEF) that acts as an activator of Rop (Rho of plants) GTPases by promoting the exchange of GDP for GTP. In postembryonic roots, modulates root stem cell maintenance by regulating
OPGEF7"xsd:stri	^{ng} "At5g02010" ^{xsd:strii}	Rexpression of PLT1 and PLT2, which are key transcription factors that mediate the patterning of the root stem cell niche. May connect RopGEF-regulated Rac/Rop signaling and auxin-dependent PLT-regulated root patter
	-	formation.**xsd-string
WEET8"XSd:string	"At5g40260"Xsd:strii	¹⁸ Mediates both low-affinity uptake and efflux of sugar across the plasma membrane. Required, in pollen, for microspore cell integrity and primexine pattern formation (PubMed:18434608, PubMed:25988582). ************************************
		"Transcription factor required for quiescent center cells specification and maintenance of surrounding stem cells, and for the asymmetric cell division involved in radial pattern formation in roots. Essential for both cell divi
SHR"XSd:string	na a na a convedetrio	Name cell specification. Regulates the radial organization of the shoot axial organs and is required for normal shoot gravitropism. Directly controls the transcription of SCR, and when associated with SCR, of MGP, RLK, TRI,

From the MetaNetX metabolic networks for metagenomics database SPARQL Endpoint: https://rdf.metanetx.org/sparql (this slide deck will make it much easier for you! https://www.metanetx.org/cgi-bin/mnxget/mnxref/MetaNetX RDF schema.pdf)

Q9: 4 POINTS: what is the MetaNetX Reaction identifier (starts with "mnxr") for the UniProt Protein uniprotkb:Q18A79

Virtuoso SPARQL Query Editor						
Default Data Set Nam	ne (Graph IRI)					
Query Text						
SELECT DISTINCT ?r WHERE{ ?peptide meta:pept ?catalyzes meta:pe ?gpr meta:cata ?ca meta:reac ?reactio ?reaction rdfs:lab	<pre>tttp://purl.un reaction_ident Xref uniprot: pt ?peptide . talyzes ; on . pel ?reaction_</pre>	iprot.org/uniprot/> ifier Q18A79 .				
Sponging:	Use only local data (including data retrieved before), but do not retrieve more					
Results Format:	HTML					
Execution timeout:	0	milliseconds (values less than 1000 are ignored)				
Options:	 ✓ Strict checking of void variables Log debug info at the end of output (has no effect on some queries and output formats) Generate SPARQL compilation report (instead of executing the query) 					
(The result can only be sen	nt back to browser, i	not saved on the server, see <u>details</u>)				
Run Query Reset						
4	4.0					
reaction_iden	tifier					
"mnxr165934"						
"mnxr145046c	3"					

FEDERATED QUERY - UniProt and MetaNetX

Q10: 5 POINTS: What is the official locus name, and the MetaNetX Reaction identifier (mnxr.....) for the protein that has "glycine reductase" catalytic activity in Clostridium difficile (taxon 272563). (this must be executed on the https://rdf.metanetx.org/sparql_endpoint)

Query Text

```
PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#>
PREFIX up: <http://purl.uniprot.org/core/>
PREFIX mnx: <https://rdf.metanetx.org/schema/>
PREFIX skos: <http://www.w3.org/2004/02/skos/core#>
PREFIX taxon: <http://purl.uniprot.org/taxonomy/>
SELECT DISTINCT ?locus ?reaction_id ?label
WHERE {
  SERVICE <a href="https://sparql.uniprot.org/sparql">https://sparql.uniprot.org/sparql>{</a>
    SELECT ?uniprot protein ?locus ?label
      ?uniprot protein a up:Protein ;
                up:organism taxon:272563 ;
                up:mnemonic ?locus ;
                up:enzyme ?enzyme .
      ?enzyme skos:prefLabel ?label .
      FILTER(CONTAINS(?label, "glycine reductase"))
    }
  }
  ?mnx protein mnx:peptXref ?uniprot protein .
  ?cata mnx:pept ?mnx protein .
?gpr mnx:cata ?cata ;
       mnx:reac ?reaction .
  ?reaction rdfs:label ?reaction id .
```

locus	reaction_id	label	
"Q185M4_CLOD6"	"mnxr157884c3"	"glycine reductase"	
"Q185M6_CLOD6"	"mnxr157884c3"	"glycine reductase"	
"Q185M3_CLOD6"	"mnxr157884c3"	"glycine reductase"	
"Q185M5_CLOD6"	"mnxr157884c3"	"glycine reductase"	
"Q185M1_CLOD6"	"mnxr157884c3"	"glycine reductase"	
"Q185M4_CLOD6"	"mnxr162774c3"	"glycine reductase"	
"Q185M6_CLOD6"	"mnxr162774c3"	"glycine reductase"	
"Q185M3_CLOD6"	"mnxr162774c3"	"glycine reductase"	
"Q185M1_CLOD6"	"mnxr162774c3"	"glycine reductase"	