

European Bank for induced pluripotent Stem Cells

Semantic Queries in EBiSC

1.4 31.05.2023

The EBISC – European Bank for induced pluripotent Stem Cells project has received support from the Innovative Medicines Initiative Joint Undertaking under grant agreements n° 115582 and 821362, resources of which are composed of financial contribution from the European Union and EFPIA companies' in kind contribution. www.imi.europa.eu

Contents

Contents	. 2
Introduction	. 3
EBiSC Ontology	. 3
Semantic Linkage to Diseases	. 4
Cell Line Identifiers	. 4
IRIs in the EBiSC Ontology	. 5
SPARQL Queries	. 5
SPARQL Interface EBiSC Platform	. 5
SPARQL Examples	. 6
Appendix 1 – Introduction to Ontologies and SPARQL	. 9
Ontologies	. 9
SPARQL	. 9
Appendix 2 – Protégé	11
Protégé GUI	11



Introduction

Each cell line in the European Bank for induced pluripotent Stem Cells (EBiSC) catalogue is described by a detailed dataset, including data and metadata provided by the cell line depositor. We focused out, that a simple dataset might not be sufficient to display all associated data to characterise a cell line. Some items, like diseases or gene mutation, requires a more complex and comprehensive data description method.

We decided to use a semantic data description by an ontology to archive this goal. For more information about ontologies, please consult Appendix 1.

EBiSC Ontology

The aim of the EBiSC Ontology (available at <u>https://ebisc.org/ontologies/ebisc.owl</u>) is to provide fully semantic descriptions of the data and metadata of pluripotent Stem Cells registered in the EBiSC platform and to make the cell lines more discoverable for EBiSC users.

As this ontology describes cell lines, it is based on the Cell Line Ontology¹. Several commonly available ontologies have been imported to enable the most comprehensive possible descriptions of all important metadata. They include information about cell types, cell lines, diseases, employed experimental methods, anatomical entities, genes and proteins.

The following picture shows a short excerpt of the global description of a cell line including some associated metadata.



Figure 1: Semantic description of a specific cell line

¹ https://www.ebi.ac.uk/ols/ontologies/clo



Semantic Linkage to Diseases

An important information is the connection of a cell line to a certain disease.

This feature is particularly important to provide users who search for cell lines with the most appropriate matches, e.g. matches that relate to a specific disease context or genetic mutation/variant.

This connection can exist in two different ways. On the one hand, we have information about the donor of the line and his/her diseases (affected or unaffected). Thus, cell lines can be linked to diseases, which have been diagnosed in the donor, or cell lines can possess disease-related mutations, which have been typed in the donor, who carries the disease mutation.

On the other hand, a line itself can be genetically modified and in this way serve as a role model (or "experimental tool") for investigating disease mechanisms (see Figure 2).



Figure 2: Linkage and detailed sematic information of a disease

Cell Line Identifiers

Every cell line in the EBiSC Ontology is described by a CLO_ID, because of its relation to the Cell Line Ontology. This CLO_ID is also part of the cell line's metadata in the EBiSC user interface (see Figure 3).



	External Databases
hPSCreg	BIONi015-A
BioSamples	SAMEA4342649
Cellosaurus	CVCL_LE13
ECACC	66540266
CLO	CLO_0100561
Wikidata	Q54796794

Figure 3: CLO_ID of cell line BIONi015-A in the EBiSC catalogue

IRIs in the EBiSC Ontology

Every class in the ontology has a unique identifier called IRI².

The following example explains the IRI behaviour in EBiSC:

- Cell line name: STBCi098-A
- IRI: http://purl.obolibrary.org/obo/CLO_0101911
 - General part: http://purl.obolibary.com/obo/ (will not change for ontologies, that a part of the OBO-Foundry³)
 - Variable part: CLO_0101911 (specific ID)

An easy way to analyse the content of the EBiSC ontology to use the software programme Protégé (see Appendix 2)

SPARQL Queries

Instead of using Protégé, the related information of a cell line can also be accessed by SPARQL (see Appendix 1 for a short introduction in SPARQL).

SPARQL Interface EBiSC Platform

The SPARQL interface of the EBiSC platform can be reached via <u>https://ebisc.org/sparql</u> (see Figure 4).

In the field "Query Text", you can enter your SPARQL Query

"Run Query" will show you the result of Query in the Result view (see Figure 4).

³ https://obofoundry.org/



² Internationalized Resource Identifiers (IRIs) (w3.org)

$\leftarrow \ \rightarrow \ G$	https://ebisc.org/sparql	URL SPARQL interface	Ē] 120% 公	⊻ 🧿 ≡
Virtuoso SPARC	L Query Editor]		
Default Data Set Name https://hpscreg.eu/ontolo	(Graph IRI) gies/hpscreg.owl Graph IRI (do not	t change!)	About <u>N</u>	amespace Prefixes Inferenc	<u>æ rules</u> <u>RDF views</u>
select distinct ?Con	cept where {[] a ?Concept} LIMIT 100	SPARQL qu	iery		
(Security restrictions of this a Results Format: Execution timeout:	erver do not allow you to retrieve remote RDF data, see de HTML 0 milliseconds (values less than	<u>italis.)</u> n 1000 are ignored)			<i>II</i> ;
Options:	 Strict checking of void variables Log debug info at the end of output (has no effective of the second second	ct on some queries and output formats) of executing the query)			
(The result can only be sent	back to browser, not saved on the server, see details)				
Run Query Re Run	SPARQL query				
	Virtuoso version	Copyright © 2023 OpenLink Software 07.20.3229 on Linux (x86_64-pc-linux-gnu), Single Se	erver Edition		

Figure 4: SPARQL interface for EBiSC

SPARQL Examples

The following lines will show some SPARQL examples for querying relevant cell line information available in the EBiSC Ontology. The queries can be easily adopted by changing the highlighted part.

• Get all cell lines with a related donor disease

IRIs of donor diseases in the ontology (can be replaced in owl:onProperty part):

- Donor has disease: http://purl.obolibrary.org/obo/CLO_0000015
- Patient is carrier of disease: http://purl.obolibrary.org/obo/CLO_0000003

PREFIX rdfs: <http://www.w3.org/2000/01/rdf-schema#> PREFIX owl: <http://www.w3.org/2002/07/owl#> SELECT DISTINCT (STR(?clname) AS ?line) (STR(?dislab) AS ?disease) WHERE { ?dis rdfs:label ?label. ?label bif:contains "**neurodegenerative disease**'". ?sub rdfs:subClassOf* ?dis. ?cell rdfs:subClassOf ?rest. ?rest owl:onProperty <http://purl.obolibrary.org/obo/CLO_0000015>. ?rest owl:someValuesFrom ?sub. ?sub rdfs:label ?dislab. ?cell rdfs:label ?clname. ?cell rdfs:seeAlso ?so. filter contains(STR(?so),"ebisc").

} GROUP by ?cell ORDER by ?line



Result (excerpt):

line	disease
BIONi010-C-41	myotonic dystrophy type 1
BIONi010-C-42	myotonic dystrophy
BIONi010-C-43	myotonic dystrophy
CBRCULi002-A	myotonic dystrophy type 1
CENSOi008-A	myotonic dystrophy

• Get all cell lines with a genetically modified gene related to a specific disease

GROUP by ?cell ORDER by ?line

Result (excerpt):

line	disease	
BIONi010-C-17	Alzheimer disease	
BIONi010-C-2	Alzheimer's disease	
BIONi010-C-25	Alzheimer's disease	
BIONi010-C-3	Alzheimer disease	
BIONi010-C-4	Alzheimer's disease	



- <u>Get all cell lines with a modified gene that plays a role in a specific biological</u>
 <u>process</u>
 - Get ID of modifying gene from Gene Ontology The ID can retrieve this URL: <u>https://www.ebi.ac.uk/ols/ontologies/go</u>. Type in the name and select "search".
 - 2. Copy the ID from the result page and paste it in the Query below.

Inal transduction Search		
Jump to	signal transduction	Search GO
lignal transduction GO GO:0007165	년 http://purl.obolibrary.org/obo/GO_0007165 (집 Copy	
ignal transduction involved in filamentous growth GO GO.0001402 ignal transduction by p53 class mediator GO GO.007/2331	covers signaling from receptors located on the surface of the cell and signaling via covers signaling from receptors located on the surface of the cell and signaling via to events at and within the receiving cell. [GCC: mtg_signaling_feb11 GOC: go_	e.g. regularison on unicomprovi or regularison of the fitted doll prodess. Signal transduction a molecules located within the cell. For signaling between cells, signal transduction is res curators]
Search (1) S for signal transduction	-1: Tree view (3) Term mappings	Term information
	biological_process	Graph view Graph view Graph view Graph view Witingdia-Simpal transduction
	biological regulation de regulation of biological process	Roset tree Subsets
	eregulation of cellular process	Show all siblings goslim_chembl, goslim_metagenomics,
	e cellular process	gosim_ban, gosim_candida
	cell communication	Note that signal transduction is defined broad
	cellular response to stimulus	include a ligand interacting with a receptor,
	P signal transduction	being triggered. A change in form of the sign
	eresponse to stimulus	every step is not necessary. Note that in man
	ellular response to stimulus	initiation of transcription. Note that specific
	E ingran transduction	transcription factors may be annotated to this

Figure 5: Get Id from Gene Ontology

Result (excerpt):

line BIONi010-C-9 BIONi010-C-5



Appendix 1 – Introduction to Ontologies and SPARQL

The following lines will give a short overview about ontologies explained by a simplified example (see Figure 6).

Ontologies

An ontology consist of elements (*classes*) that exist in a specific domain and *properties* to describe them. Properties are relationships to link two classes or *attributes* to describe a class.

The easiest way to link a class to another is the *subClassOf* property. A subClass is a more precise description to a superclass, like creature -> animal -> dog -> poodle.



Figure 6: ontology example

It is also possible to link more than two classes. Dogs and cats are both animals.

Properties can also be a bit more complex. As you can see in the example, dogs and cats can have fur. However, fur is a subClass neither of dogs nor of cats. This connection can be realised by a specific property, which is called "hasSome" here.

So, the elements in an ontology are represented by a *graph* structure. Each element of ontology can be described by *triples* (class – property – class – property....).

Ontologies provide more features, which are going beyond this example. Detailed information can be found here <u>https://www.w3.org/standards/semanticweb/ontology</u>.

SPARQL

SPARQL is a query language to receive information from ontologies in RDF format (<u>https://www.w3.org/RDF/</u>). As these datasets are described in triples, its queries have to be constructed in that manner.

The next lines shows some simplified examples.



• all triples of a dataset:

```
SELECT * WHERE {
    graph ?g {
        ?class ?property ?linkedClass .
    }
}
```

Result:

?class	?property	?linkedClass
animal	subClassOf	creature
dog	hasSome	fur
Cat	subClassOf	animal

classes with linked by a specific property

```
SELECT ?class
WHERE {
?class hasSome fur.
}
```

Result:

?class
cat
dog

• All subclasses

```
SELECT ?class
WHERE {
?class subClassOf dog.
}
```

Result:

?class	
poodle	

SPARQL is a very comprehensive query language. More information can be found here: <u>https://www.w3.org/TR/rdf-sparql-query/</u>.



Appendix 2 – Protégé

Protégé GUI

An easy way to analyse the content of the EBiSC ontology is the tool Protégé (<u>https://protege.stanford.edu/software.php</u>), as displayed in the following screen. This tool can be installed on every computer and run as a standalone application.

ebic, cle (https://bpicetg.eu/ontologies/bbic, cls.ow) ; [ChUken/sohnik/Download/ebic.cwi]	- 0 >
File Edit View Reasoner Tools Refactor Window Help	
ebisc clo (https://hpscreg.eu/ontologies/ebisc clo.owl)	- Search.
apprinter for the factor (watering entry) containing encryoners) only experimentally modified cells in etcs) (subured cell (solid) and the cell () hadced pulpations starm cells the cell () hanse induced pulpations starm cells (secold) #2004.	
Active ontology × Entities × Individuals by class × DL Query ×	
Tree – "nure" class 04-A – CLO:0102650 – http://purl.obolibrary.org/obo/CLO_0102650	
ons Usage	
hierarchy (linkage by moth	208)
subClassOf property)	
type: xsd:string]	@80
RCi004-A	
RCI002-B Attributes of cell line	080
RCI003-A Second Second B Se	inec
RC1004-A"	ipsc O O O
RC1004-A-1 SeeAlso	@ × C
RCi004-B https://www.cellosaurus.org/CVCL_9S47	
-• RCI005-A seeAlso	@×0
RCI006-A https://www.wikidata.org/wiki/Q54949504	
-• RC1007-A Description: RC6004-A	208
-• RCI007-C Equivalent To 😳	
- RC1009-A	
- Sigilou -A-1 Subclass Of C	
- Signo1 A-11	7000
- SIGI001-A-12	7000
 SIGi001-A-14 'derived from patient carrier of disease' some 'Huntington's disease' 	7080
-• SIGi001-A-15	0080
-• SIGi001-A-16 "derives from' some 'female organism' LITIKAGE to Classes with	0080
SIGi001-A-17 'derives from' some 'human tissue donor' specific properties	7@×0
SiGl001-A-18 "has potential to develop into' some ectoderm	- 7080
• SIGIDUT-A-19 • SIGIDUT-A-2	?@×0
* the notential to develop into' some mesoderm	2 @ X O
• SIGIO01-A-3	

Figure 7: Details of cell line "RCI004-A" in Protégé.

