



University of  
Central Lancashire  
UCLan

# CO3409

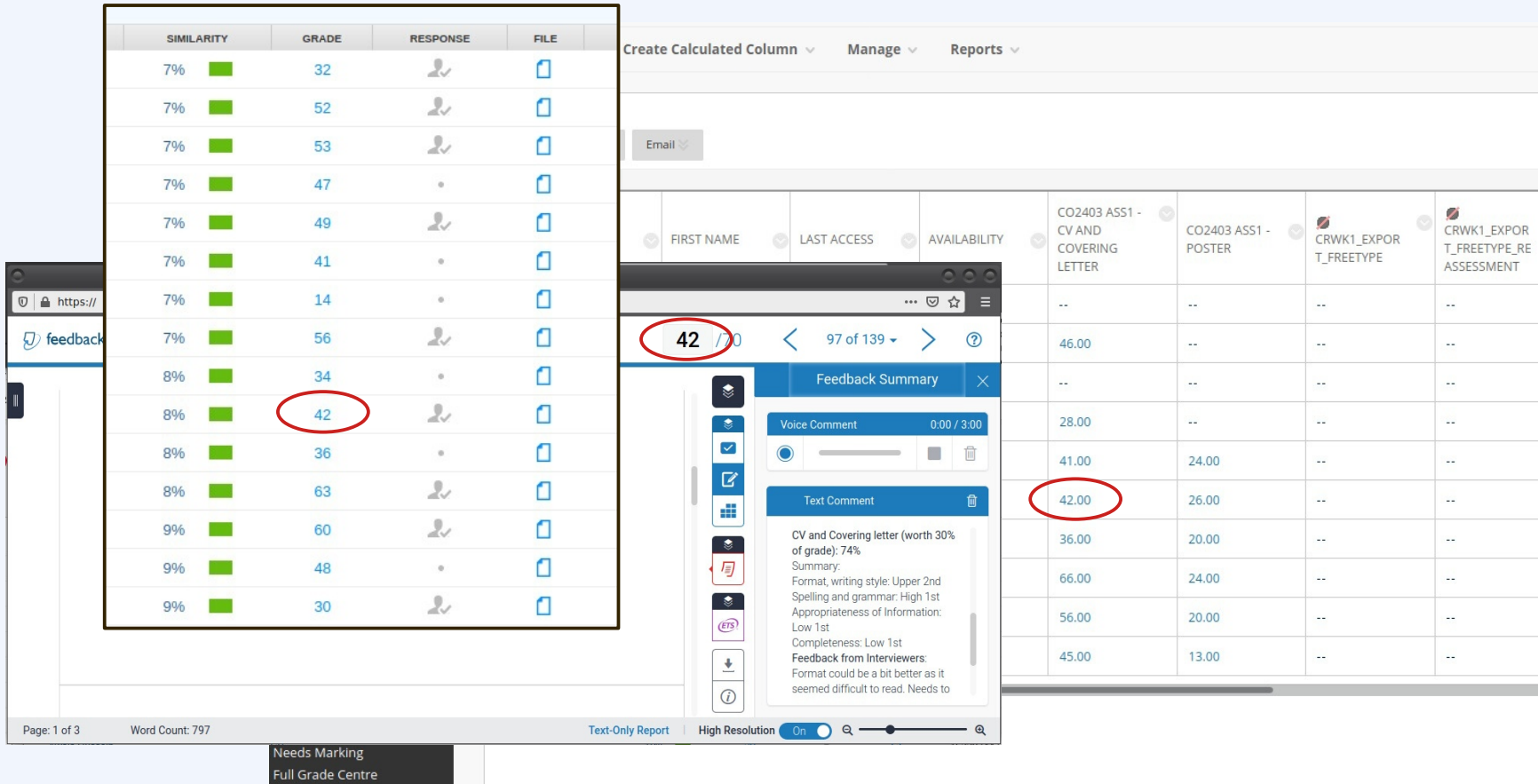
# Distributed Enterprise Systems

FAIR principles  
Metadata schemas

Where opportunity creates success

# FAIR principles

# Data management & distributed enterprise systems



The screenshot displays a web-based assessment system interface. On the left, a table lists student results with columns for Similarity, Grade, Response, and File. The grade 42 is circled in red. On the right, a detailed feedback summary for a student is shown, with the score 42.00 also circled in red. The feedback includes a voice comment section and a text comment section detailing the student's performance on a CV and covering letter.

SIMILARITY	GRADE	RESPONSE	FILE
7%	32	[User Icon]	[File Icon]
7%	52	[User Icon]	[File Icon]
7%	53	[User Icon]	[File Icon]
7%	47	[User Icon]	[File Icon]
7%	49	[User Icon]	[File Icon]
7%	41	[User Icon]	[File Icon]
7%	14	[User Icon]	[File Icon]
7%	56	[User Icon]	[File Icon]
8%	34	[User Icon]	[File Icon]
8%	42	[User Icon]	[File Icon]
8%	36	[User Icon]	[File Icon]
8%	63	[User Icon]	[File Icon]
9%	60	[User Icon]	[File Icon]
9%	48	[User Icon]	[File Icon]
9%	30	[User Icon]	[File Icon]

FIRST NAME	LAST ACCESS	AVAILABILITY	CO2403 ASS1 - CV AND COVERING LETTER	CO2403 ASS1 - POSTER	CRWK1_EXPOR T_FREETYPE	CRWK1_EXPOR T_FREETYPE_RE ASSESSMENT
..	..	..	..	..	..	..
46.00	..	..	..	..	..	..
..	..	..	..	..	..	..
28.00	..	..	..	..	..	..
41.00	24.00	..	..	..	..	..
42.00	26.00	..	..	..	..	..
36.00	20.00	..	..	..	..	..
66.00	24.00	..	..	..	..	..
56.00	20.00	..	..	..	..	..
45.00	13.00	..	..	..	..	..

**Feedback Summary**

Voice Comment: 0:00 / 3:00

Text Comment:

CV and Covering letter (worth 30% of grade): 74%

Summary:

- Format, writing style: Upper 2nd
- Spelling and grammar: High 1st
- Appropriateness of Information: Low 1st
- Completeness: Low 1st

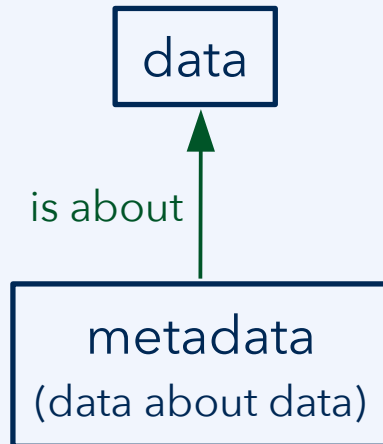
Feedback from Interviewers: Format could be a bit better as it seemed difficult to read. Needs to

What do we (and the system) need to know to **use and reuse the data** correctly?

# FAIR data and metadata

Metadata are “descriptive data about an object” (ISO 11179), usually a digital object

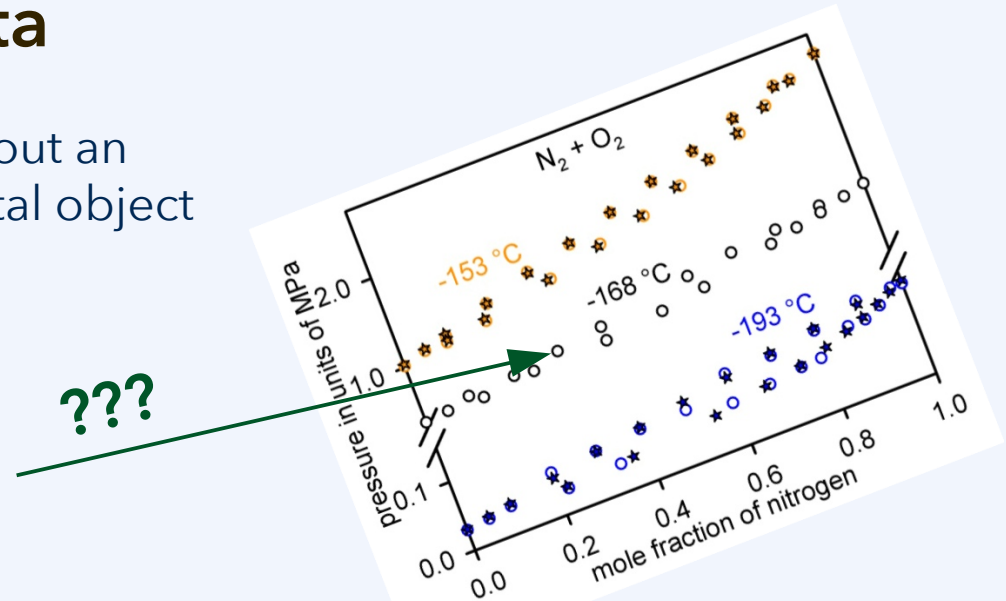
Digitization



Digitalization

Industry 5.0<sup>2</sup>

“How was the data point obtained?”  
“How precise is it?” “Who did it?” ...



## Competency questions:<sup>1</sup>

Representative queries about data (e.g., for metadata), to be competently answered by a knowledge base.

Good practice in data management (**FAIR principles**):

Make all data **findable**, **accessible**, **interoperable**, and **reusable**.

<sup>1</sup>M. Grüninger, M. S. Fox, in *Benchmarking: Theory and Practice*, doi:10.1007/978-0-387-34847-6\_3, **1995**.

<sup>2</sup>M. Breque, L. De Nul, A. Petridis, *Industry 5.0*, EC policy brief, doi:10.2777/308407, **2021**.

# FAIR principles of data management<sup>1</sup>

## Findability

- F1. Globally unique **persistent identifiers (PID)**
- F2. **Enriched with metadata**
- F3. Data identifier included in metadata
- F4. **Registered/indexed in searchable platform**

## Interoperability

- I1. **Formal language** used for **knowledge representation**
- I2. Metadata use **vocabularies** that are themselves FAIR
- I3. Semantic web principles, **data can refer to other data**

## Accessibility

- A1. **Retrievable from PID** via a standard protocol
  - A1.1. Open and freely implementable protocol
  - A1.2. ... **authentication/authorization** if necessary
- A2. Metadata remain accessible (beyond data)

## Reusability

- R1. Metadata include a plurality of accurate and relevant attributes
  - R1.1. Release data and metadata with an accessible **data usage licence**
  - R1.2. Data are annotated with a detailed **provenance description**
  - R1.3. Relevant **disciplinary and community standards** are fulfilled

<sup>1</sup>M. D. Wilkinson *et al.*, "The FAIR Guiding Principles ...," doi:10.1038/sdata.2016.18, **2016**.

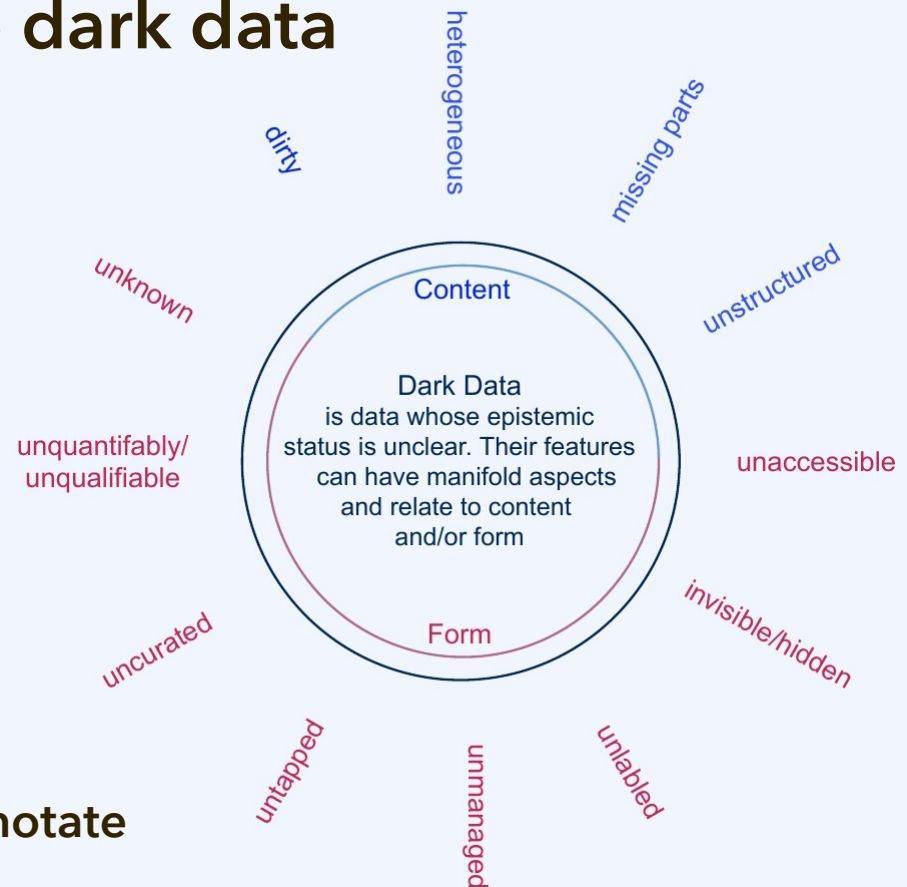
# FAIR data as opposed to dark data

To be FAIR, and therefore also **reusable**, the **epistemic status** (= knowledge status) of data needs to be characterized: Beyond interoperating via some I/O mechanisms, **we must know in what way the data constitute knowledge.**

For data reusability, it is crucial to **annotate** data with all the required metadata.

The opposite of reusable data are **dark data**: Data with an uncharacterized epistemic status. Today there is a “**deluge of dark data**” – most data are dark.<sup>1</sup>

<sup>1</sup>B. Schembera, *J. Supercomput.* 77: 8946 – 8966, doi:10.1007/s11227-020-03602-6, **2021**.



(Figure: Courtesy B. Schembera.)

# Semantic interoperability

Three branches of the theory of formal languages:

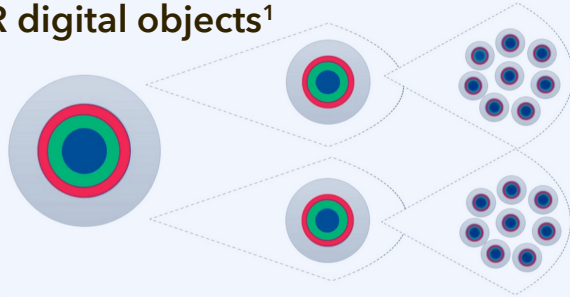
- **Syntax**, theory of the **structure** of language → data formats (e.g., JSON)
- **Semantics**, theory of the **meaning** of language → knowledge graph
- **Pragmatics**, theory of the **use** of language → processes and protocols

Generally speaking, **semantics** refers to “**meaning**,” as opposed to syntax, which refers to “proper grammar and notation.” Normally, there can only be a semantic content if there is a correct syntax, but the same content (e.g., knowledge graph) can be represented in arbitrarily many different formats.

**Semantic interoperability** allows us to coherently use many different formats, web service and API specifications, DB schemas and architectures. This is necessary whenever a distributed enterprise system is **internally heterogeneous**, or when it needs to **exchange information** with external systems.

# Semantic interoperability

## FAIR digital objects<sup>1</sup>



### Problems

Lack of (or overabundance of)

- P1: explicit definitions
- P2: common semantics (general ontologies)
- P3: reference repository
- P4: common metadata scheme across communities
- P5: metadata models



### Needs

- N1: principle approaches/tools for ontology and metadata schemes
- N2: harmonisation across disciplines
- N3: harmonisation of data of the same type
- N4: federated access to existing research data repositories



### Recommendations

- R1: definitions of concepts, metadata and data schemes
- R2: creating semantic artefacts with open licenses
- R3: associated documentation for semantic artifacts
- R4: repositories of semantic artefacts
- R5: minimum metadata model and cross walks discovery
- R6: extensible options for disciplinary metadata
- R7: apply a broad definition of data (datasets, workflows, lab protocols, software, methods, hardware design, etc.)
- R8: clear protocols and building blocks for catalogues



**EUROPEAN OPEN  
SCIENCE CLOUD**

EOSC Interoperability Framework<sup>1</sup>

<sup>1</sup>EOSC Executive Board, *EOSC Interoperability Framework*, doi:10.2777/620649, 2021.



# Semantic interoperability

Now: How does this look in the context of concrete data?

Let us imagine we receive a file "sigma.dat" (on the right).

Three modes of interoperability (*i.e.*, agreements) are needed:

```
# Model 1
# A      sigma      sigma_err
40.0    1.17745    0.167
60.0    3.03579    0.3592
80.0    3.62384    0.3797
100.0   4.30474    0.3719

# Model 2
# A      sigma      sigma_err
40.0    1.25022    0.1238
60.0    2.75247    0.2723
80.0    4.05209    0.2691
100.0   4.05401    0.2726
```

## 1) Syntactic:

formal relations between signs



The file format

(ASCII text file, tab separated columns, etc.).

## 2) Semantic:

meaning, relations between signs and what they refer to



Info about the content: e.g., what each column and block means, the data provenance, etc.

## 3) Pragmatic:<sup>1,2</sup>

relation between signs and their use, environment, users, and practices



*E.g.:* We type "rm sigma.dat" in a terminal.

Depending on our rights on the file, it will be removed or not.

<sup>1</sup>The **EOSC Interoperability Framework** calls this **technical**, **organizational** and **legal** interoperability.

<sup>2</sup>On pragmatic interoperability for enterprise systems, see doi:10.1007/978-3-030-81200-3\_4, **2021**.

# Metadata schemas

# What do you see?



Use only simple sentences consisting of:

- A **subject**
- A **predicate**
- An **object**

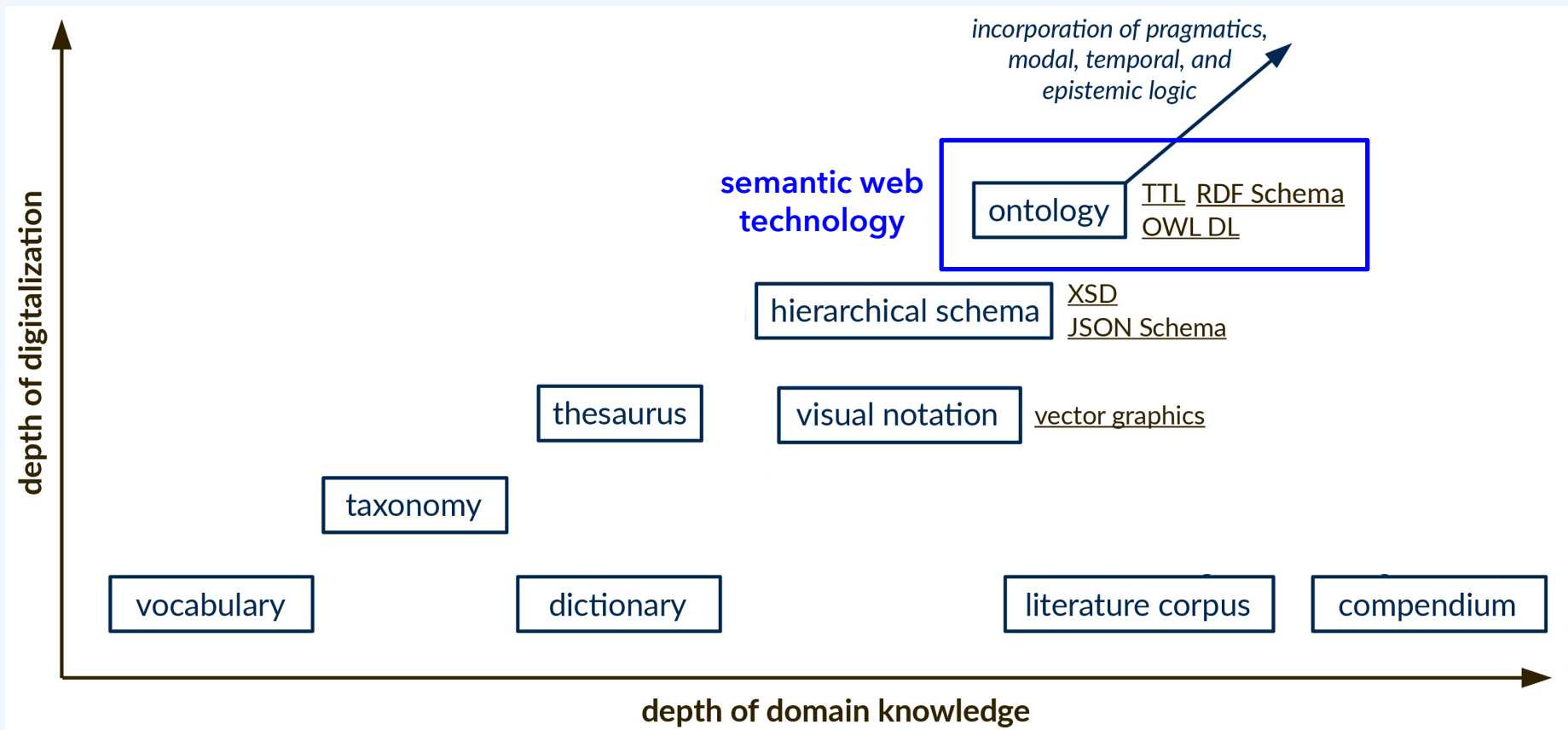
Such as:

"The-elephant  
is-dancing-in  
the-room."

"The-wheel  
is-part-of  
the-car."

# Metadata standardization

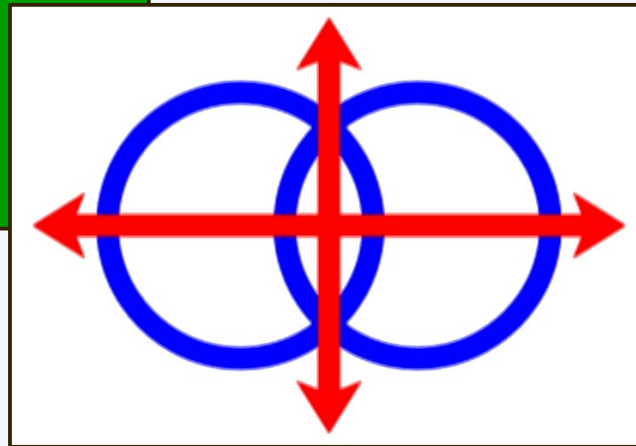
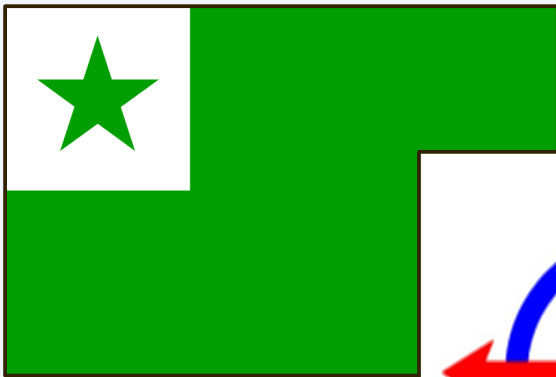
Hierarchy of **semantic artefacts** (*i.e.*, metadata standards)



# Metadata standardization

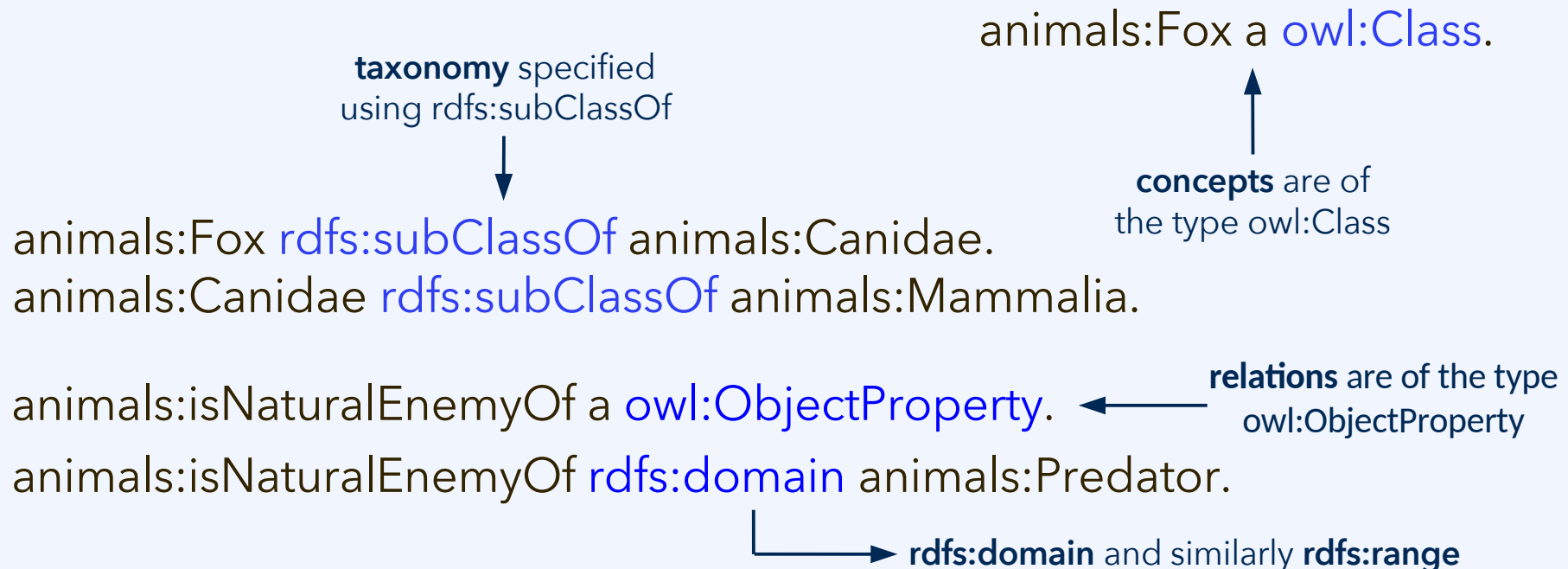
## „One World Language“

*or now, “Web Ontology Language (OWL),”  
based on RDF Schema (RDFS)*



# RDF Schema + Web Ontology Language

**Ontologies** are (data and) metadata **schemas for linked data**. They define what kinds of **knowledge graphs** are permitted. They specify what **concepts** can be instantiated by **individuals**, and what **relations** there can be between them; languages: **RDF Schema (RDFS)** and **Web Ontology Language (OWL)**.



# Schema/ontology design based on scenarios



What did you see?

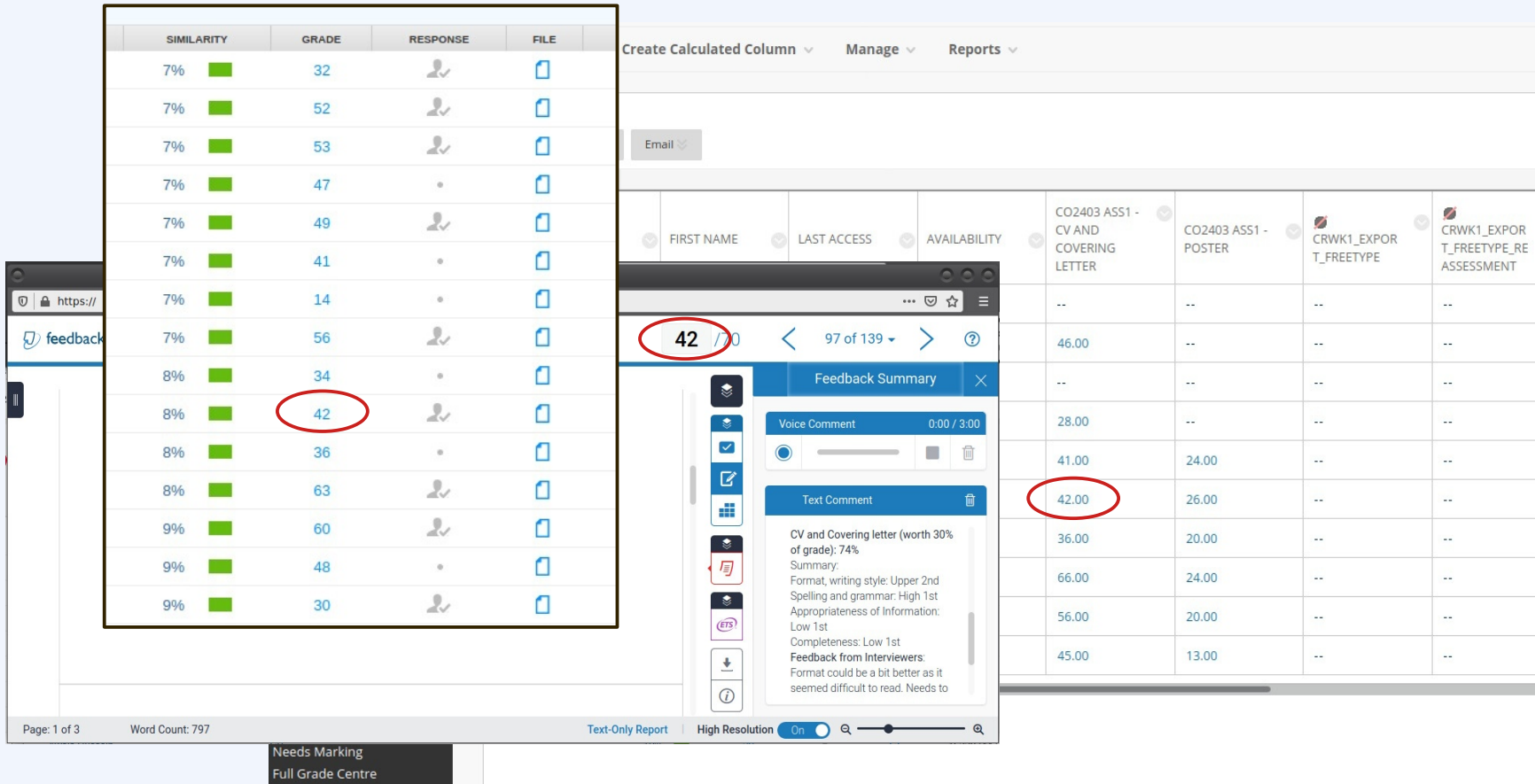
One approach to designing ontologies/schemas consists in **describing example scenarios.**

Usually, different people describe the same scenario in different terms, causing **semantic heterogeneity.**

**Concepts instantiated by individuals** should be in the ontology/RDF schema. **Relations** occurring as edges in the knowledge graph should be included; the **domains and ranges** of these relations should be included as concepts.

Different platforms may use different metadata schemas. To facilitate interoperability, an alignment is needed (e.g., an **ontology alignment**).

# Schema/ontology design via competency questions



The screenshot displays a learning management system interface. On the left, a table lists student performance metrics. On the right, a 'Feedback Summary' window is open, showing a text comment for a student with a score of 42.00.

SIMILARITY	GRADE	RESPONSE	FILE
7%	32	[User Icon]	[File Icon]
7%	52	[User Icon]	[File Icon]
7%	53	[User Icon]	[File Icon]
7%	47	[User Icon]	[File Icon]
7%	49	[User Icon]	[File Icon]
7%	41	[User Icon]	[File Icon]
7%	14	[User Icon]	[File Icon]
7%	56	[User Icon]	[File Icon]
8%	34	[User Icon]	[File Icon]
8%	42	[User Icon]	[File Icon]
8%	36	[User Icon]	[File Icon]
8%	63	[User Icon]	[File Icon]
9%	60	[User Icon]	[File Icon]
9%	48	[User Icon]	[File Icon]
9%	30	[User Icon]	[File Icon]

FIRST NAME	LAST ACCESS	AVAILABILITY	CO2403 ASS1 - CV AND COVERING LETTER	CO2403 ASS1 - POSTER	CRWK1_EXPOR T_FREETYPE	CRWK1_EXPOR T_FREETYPE_RE ASSESSMENT
--	--	--	46.00	--	--	--
--	--	--	28.00	--	--	--
--	--	--	41.00	24.00	--	--
--	--	--	42.00	26.00	--	--
--	--	--	36.00	20.00	--	--
--	--	--	66.00	24.00	--	--
--	--	--	56.00	20.00	--	--
--	--	--	45.00	13.00	--	--

**Feedback Summary**

Voice Comment 0:00 / 3:00

Text Comment

CV and Covering letter (worth 30% of grade): 74%  
 Summary:  
 Format, writing style: Upper 2nd  
 Spelling and grammar: High 1st  
 Appropriateness of Information:  
 Low 1st  
 Completeness: Low 1st  
**Feedback from Interviewers:**  
 Format could be a bit better as it seemed difficult to read. Needs to

Another strategy for building an ontology consists in gathering **competency questions** and including the employed concepts and relations in the ontology.



# Knowledge graph validation using SHACL

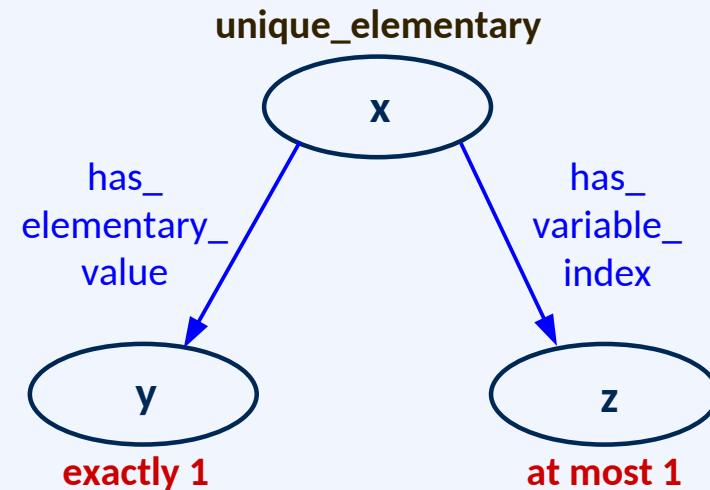
RDF schema, combined with the open world assumption, is very liberal in what knowledge graphs it permits. However, an API will usually need to specify a concrete kind of information content to be exchanged for a particular action.

**Shapes Constraint Language (SHACL)** can be used for such specifications.<sup>1</sup>

```

:unique_elementary_shape a sh:Shape;
  sh:targetClass :unique_elementary;
  sh:property [
    sh:path :has_elementary_value;
    sh:minCount 1;
    sh:maxCount 1
  ], [
    sh:path :has_variable_index;
    sh:maxCount 1
  ].

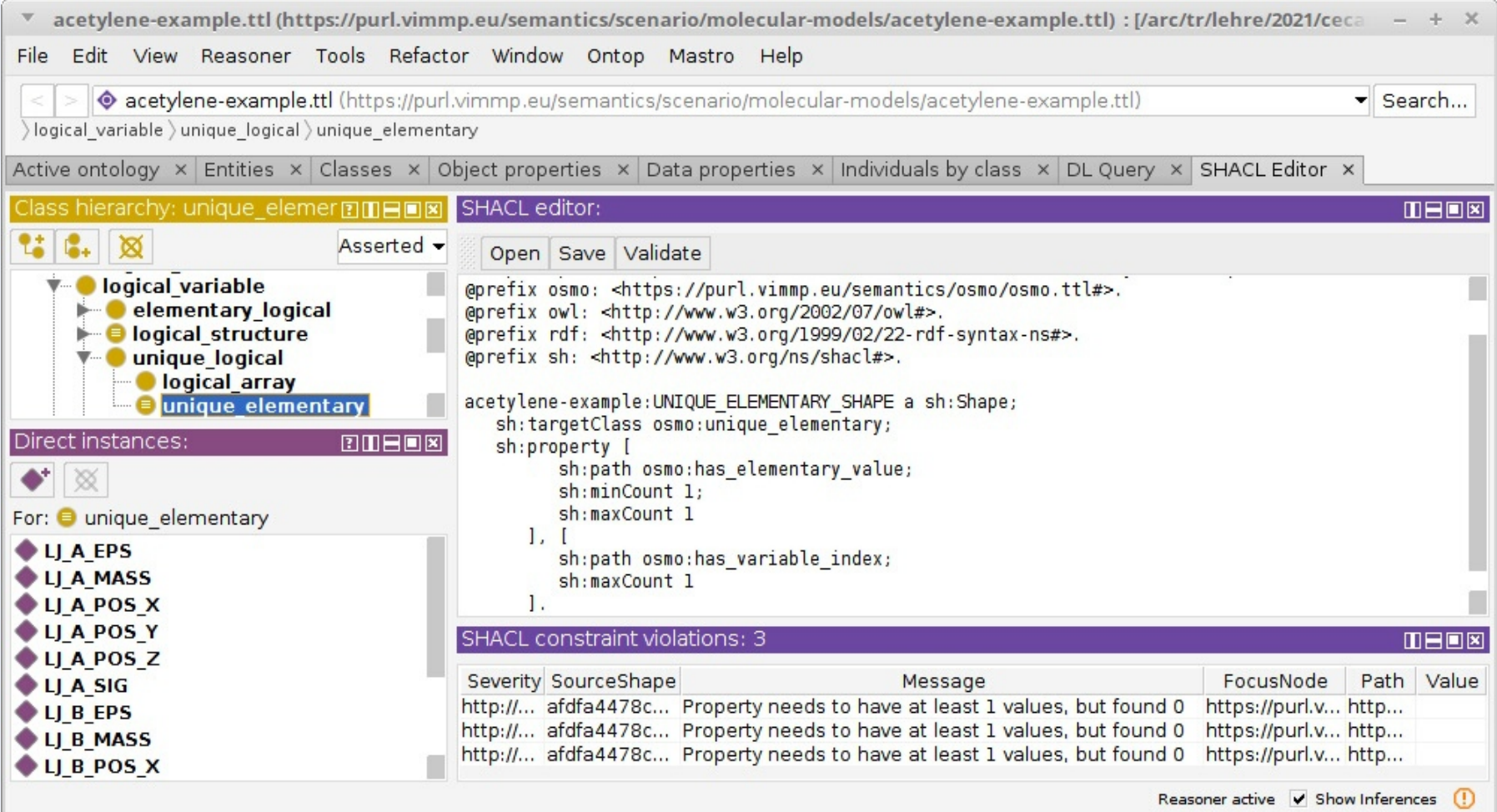
```



The open world assumption is **not** applied when evaluating SHACL constraints!

<sup>1</sup>W3C recommendation, <https://www.w3.org/TR/shacl/>, 2017.

# Knowledge graph validation using SHACL



The screenshot shows the Protégé SHACL Editor interface. The main editor displays the following SHACL constraint:

```

@prefix osmo: <https://purl.vimmp.eu/semantics/osmo/osmo.ttl#>.
@prefix owl: <http://www.w3.org/2002/07/owl#>.
@prefix rdf: <http://www.w3.org/1999/02/22-rdf-syntax-ns#>.
@prefix sh: <http://www.w3.org/ns/shacl#>.

acetylene-example:UNIQUE_ELEMENTARY_SHAPE a sh:Shape;
  sh:targetClass osmo:unique_elementary;
  sh:property [
    sh:path osmo:has_elementary_value;
    sh:minCount 1;
    sh:maxCount 1
  ], [
    sh:path osmo:has_variable_index;
    sh:maxCount 1
  ].
  
```

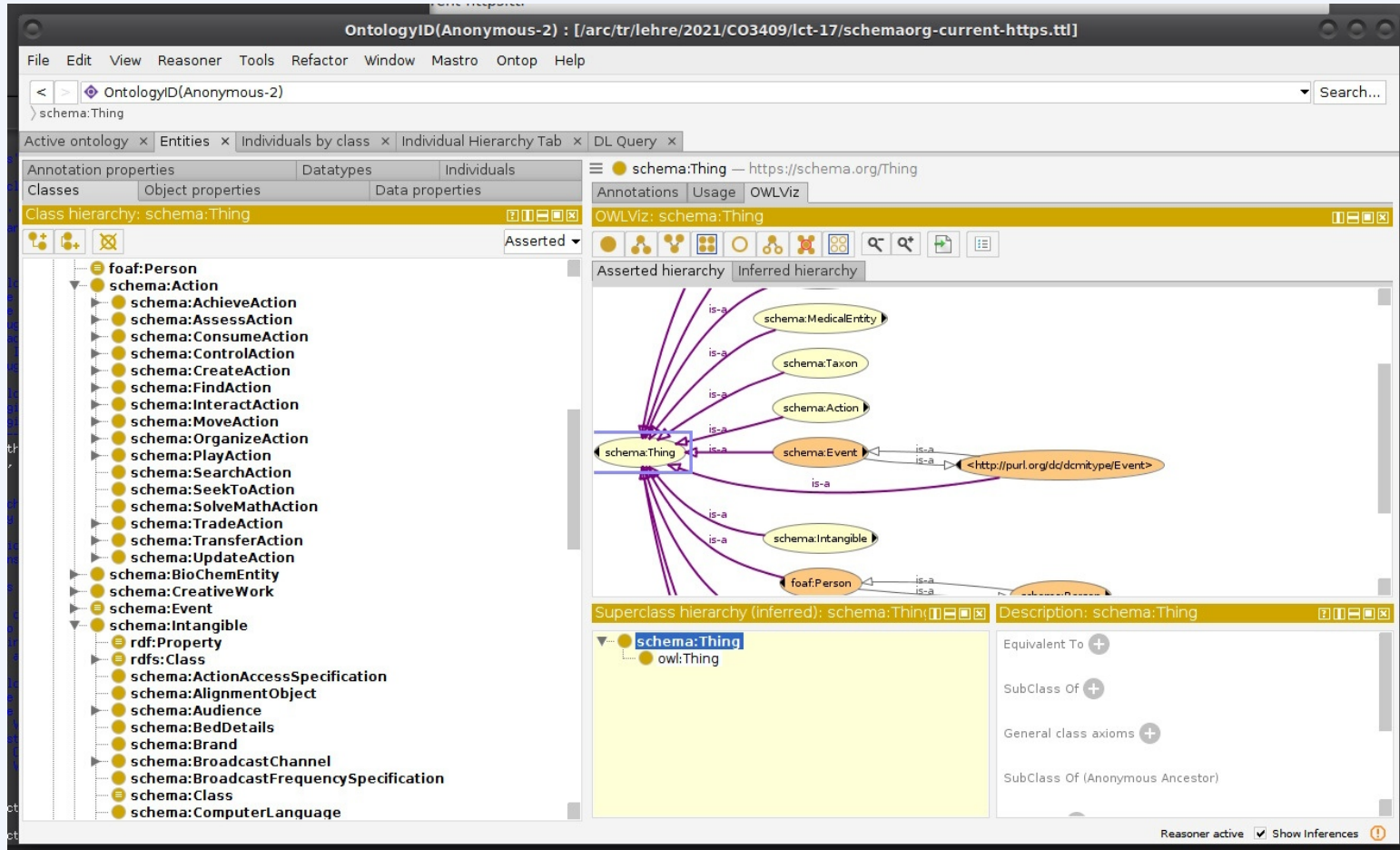
The SHACL constraint violations table shows 3 violations:

Severity	SourceShape	Message	FocusNode	Path	Value
http://...	afdfa4478c...	Property needs to have at least 1 values, but found 0	https://purl.v...	http...	
http://...	afdfa4478c...	Property needs to have at least 1 values, but found 0	https://purl.v...	http...	
http://...	afdfa4478c...	Property needs to have at least 1 values, but found 0	https://purl.v...	http...	

(example: SHACL in Protégé)

# Discussion

# Schema.org: A metadata schema used by Google<sup>1, 2</sup>



<sup>1</sup>Schema.org definitions and documentation: <https://schema.org/docs/full.html>.

<sup>2</sup>Ontology in TTL format at <https://schema.org/version/latest/schemaorg-current-https.ttl>.



University of  
Central Lancashire  
UCLan

**CO3409**

# **Distributed Enterprise Systems**

**FAIR principles**  
**Metadata schemas**

Where opportunity creates success