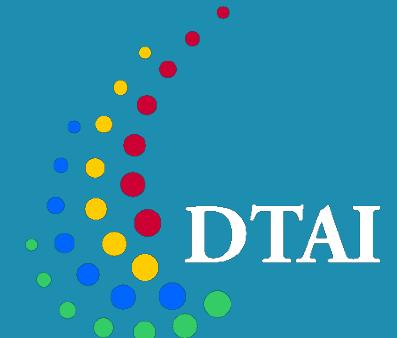




Aligning heterogeneous semi-structured data and knowledge graphs



Anastasia Dimou

 anastasia.dimou@kuleuven.be

 [@natadimou](https://twitter.com/natadimou)

my story so far...

Tenure-track assistant professor
Data & Cost-efficient ML & AI

Declarative Languages & Artificial Intelligence
(DTAI) dtai.cs.kuleuven.be/
Computer Science Department
wms.cs.kuleuven.be
KULeuven kuleuven.be/

Leuven.AI ai.kuleuven.be/
KULeuven Institute for Artificial Intelligence

Flanders Make
flandersmake.be/



Knowledge Graphs construction

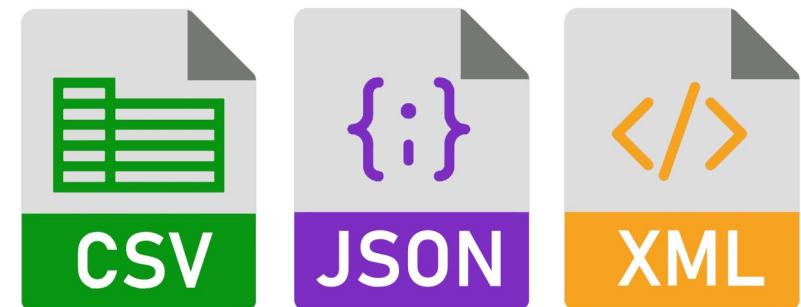
Knowledge Graphs are
crowdsourced, e.g., Wikidata.org



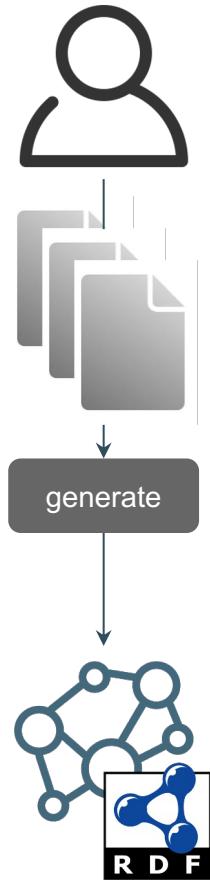
extracted from
unstructured data, e.g., plain text



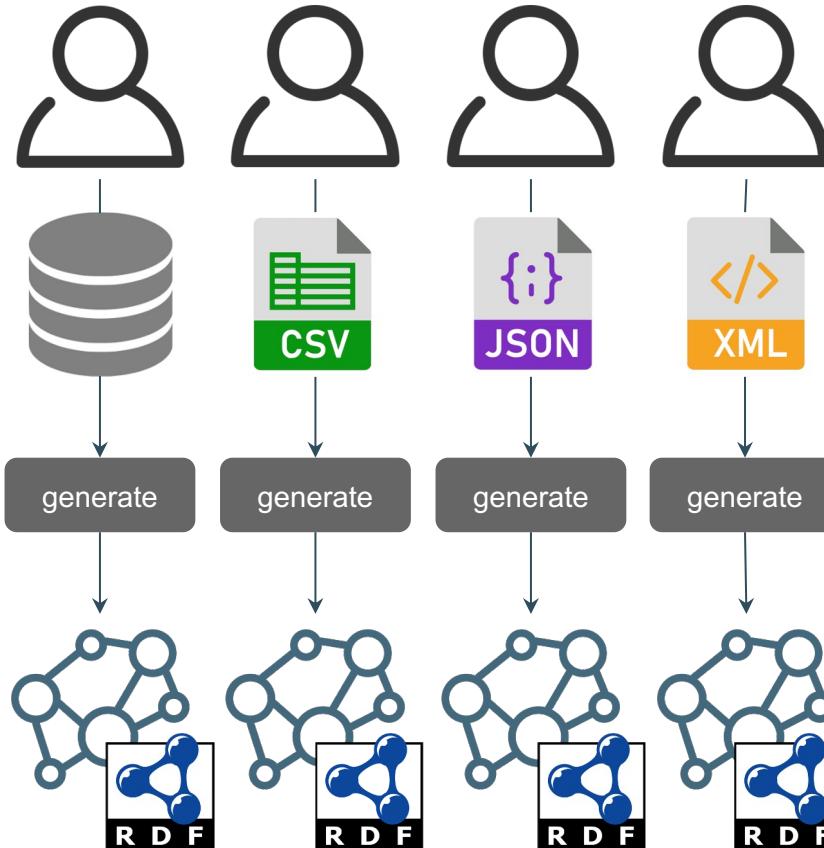
(semi-)structured data
tabular data, e.g., DBs, CSV
hierarchical data, e.g., XML, JSON



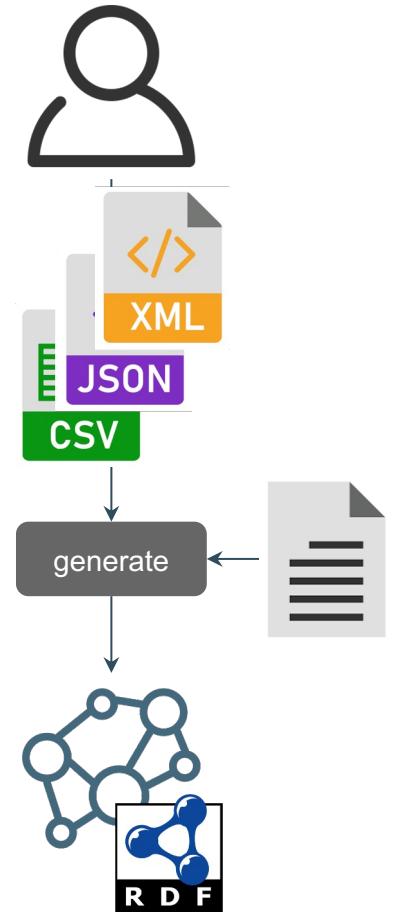
hard-coded



format-specific



declarative



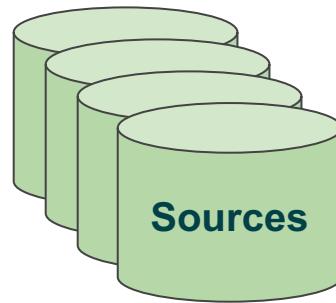
(-) new development cycle
for every modification

(-) learn & maintain multiple solutions
if data in different formats

(-) learn & maintain single solutions
(+) declarative describe rules

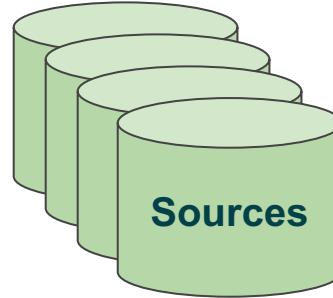
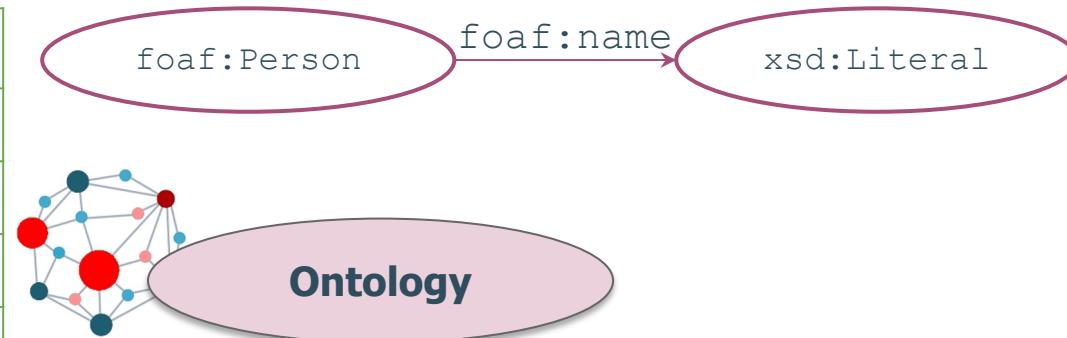
Declarative Knowledge Graph Construction

rank	name	mark
1	Anzhelika Sidorova	4.95
2	Sandi Morris	4.90
3	Katerina Stefanidi	4.85
4	Holly Bradshaw	4.80
5	Alysha Newman	4.80
6	Angelica Bengtsson	4.80



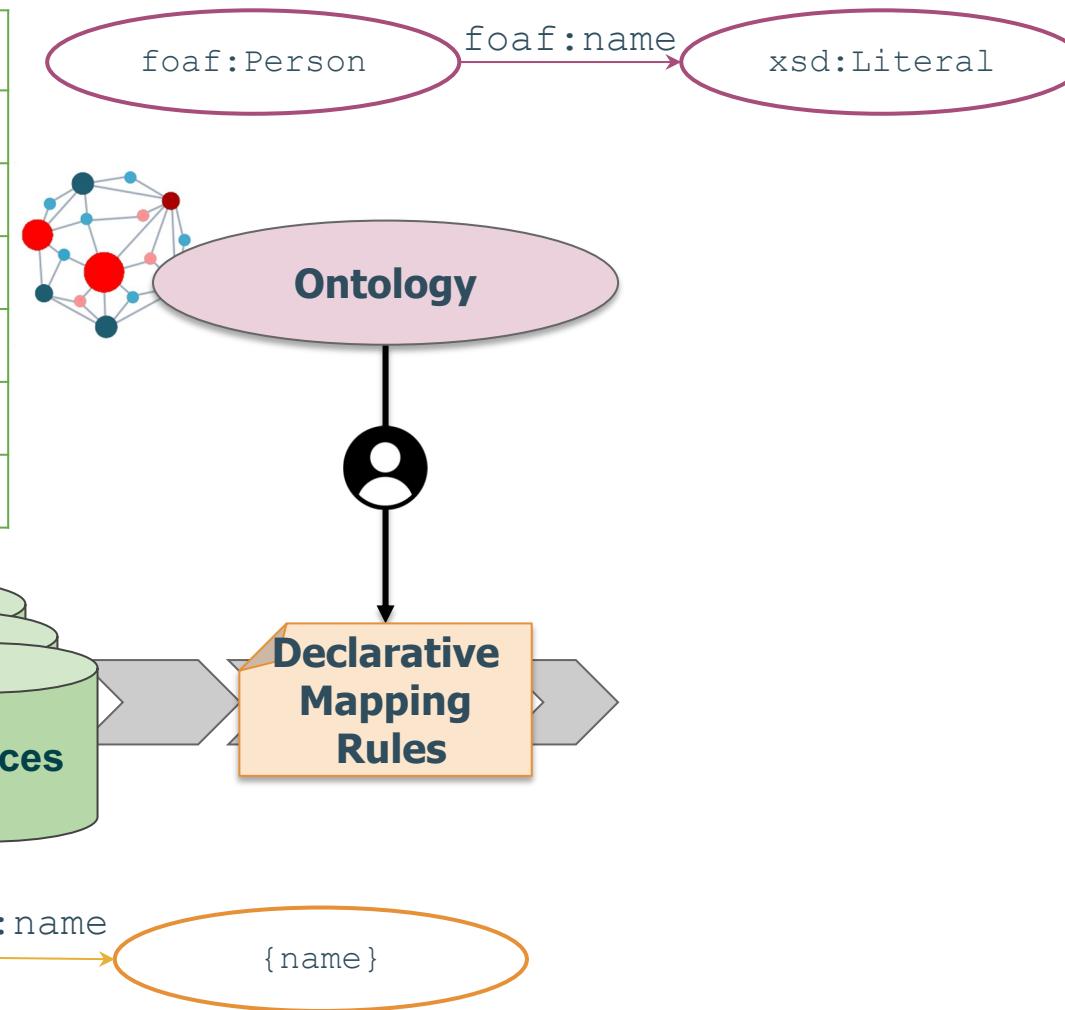
Declarative Knowledge Graph Construction

rank	name	mark
1	Anzhelika Sidorova	4.95
2	Sandi Morris	4.90
3	Katerina Stefanidi	4.85
4	Holly Bradshaw	4.80
5	Alysha Newman	4.80
6	Angelica Bengtsson	4.80



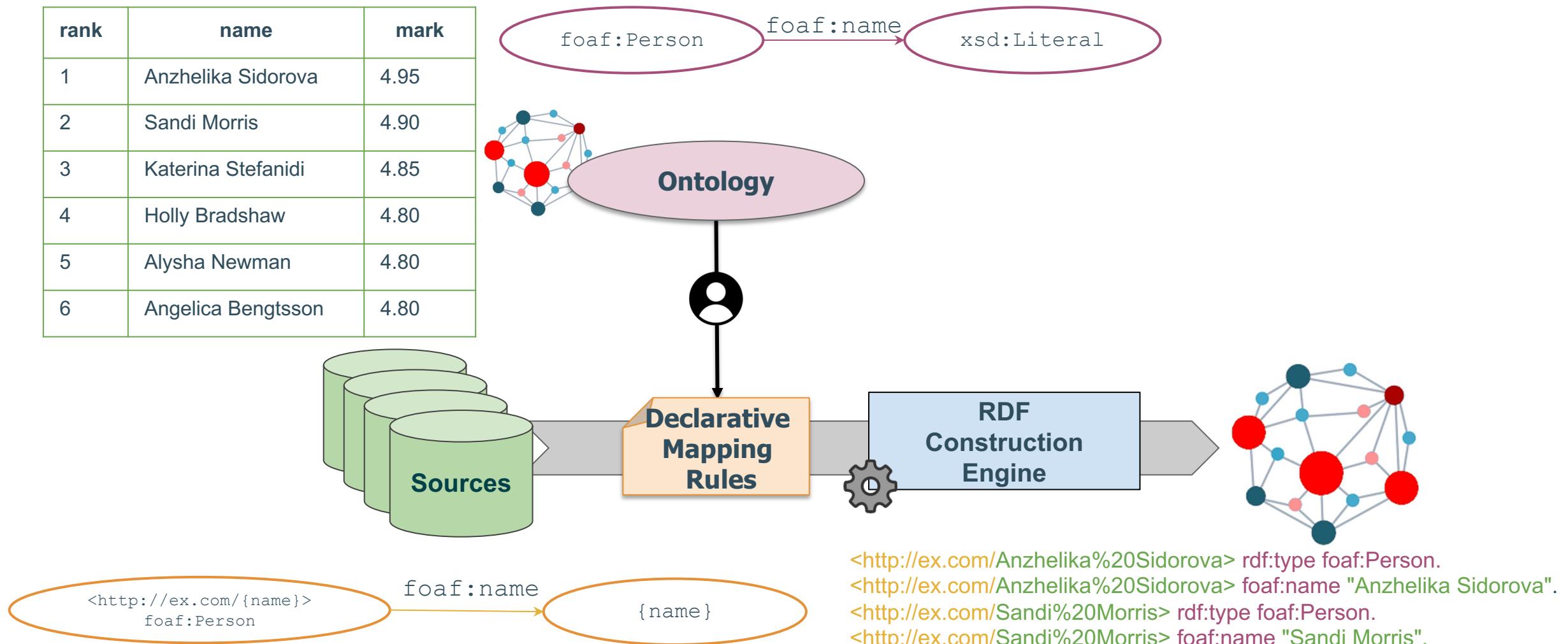
Declarative Knowledge Graph Construction

rank	name	mark
1	Anzhelika Sidorova	4.95
2	Sandi Morris	4.90
3	Katerina Stefanidi	4.85
4	Holly Bradshaw	4.80
5	Alysha Newman	4.80
6	Angelica Bengtsson	4.80



Declarative Knowledge Graph Construction

rank	name	mark
1	Anzhelika Sidorova	4.95
2	Sandi Morris	4.90
3	Katerina Stefanidi	4.85
4	Holly Bradshaw	4.80
5	Alysha Newman	4.80
6	Angelica Bengtsson	4.80



What did I do?

Mapping Languages & RML

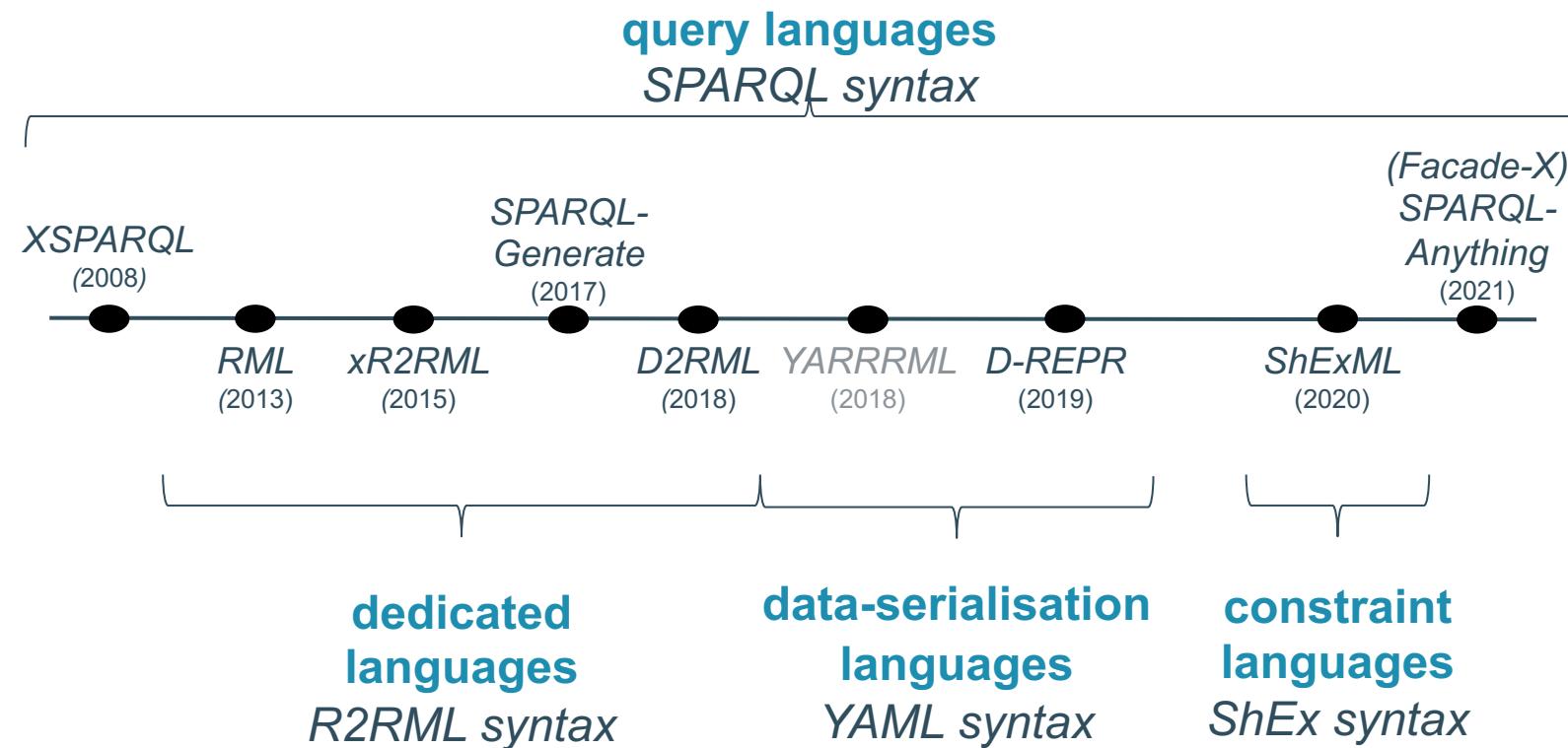
Mapping rules & User support

Mapping rules & RDF graphs validation

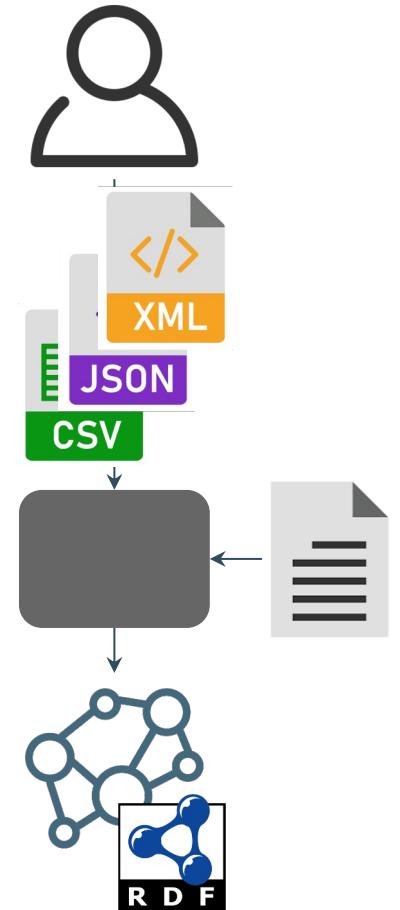
Mapping rules & SHACL shapes

Implementations & benchmarks

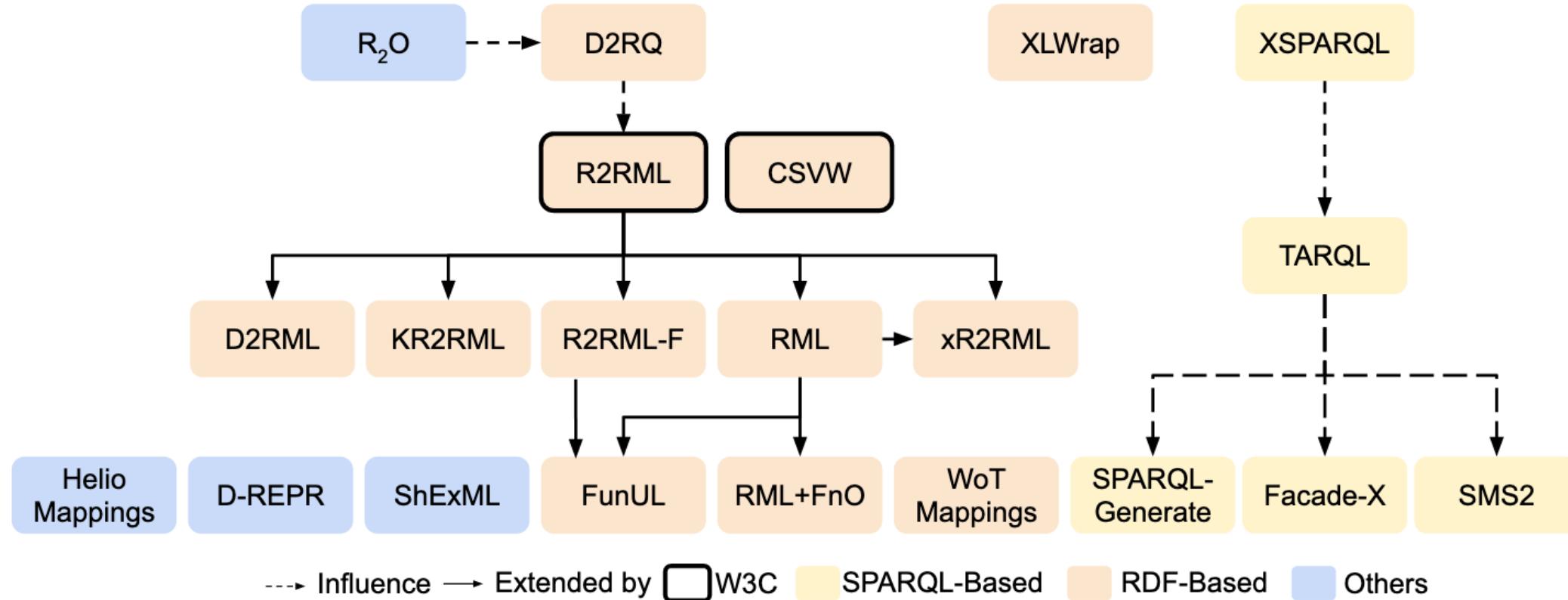
Declarative Mapping Languages



declarative



Relationships of Declarative Mapping Languages



cross-mapping-language alignments

YARRRML-to-RML

<https://github.com/RMLio/yarrmml-parser>

<https://github.com/oeg-upm/yarrmml-translator>

RML-to-SPARQL-Generate

<https://github.com/sparql-generate/rml-to-sparql-generate>

ShExML-to-RML

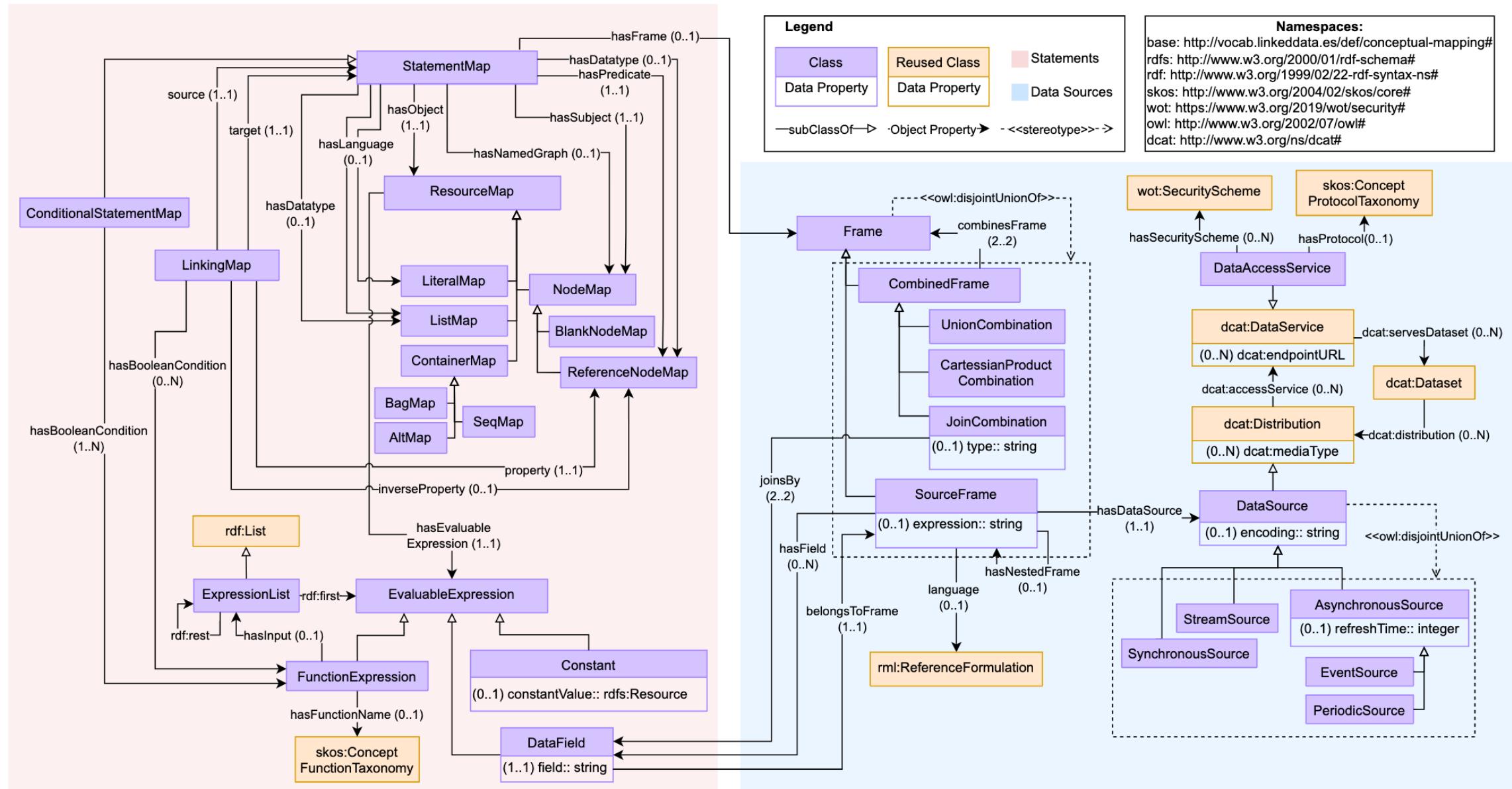
<https://github.com/herminiogg/ShExML>

RML-to-SPARQL-Anything

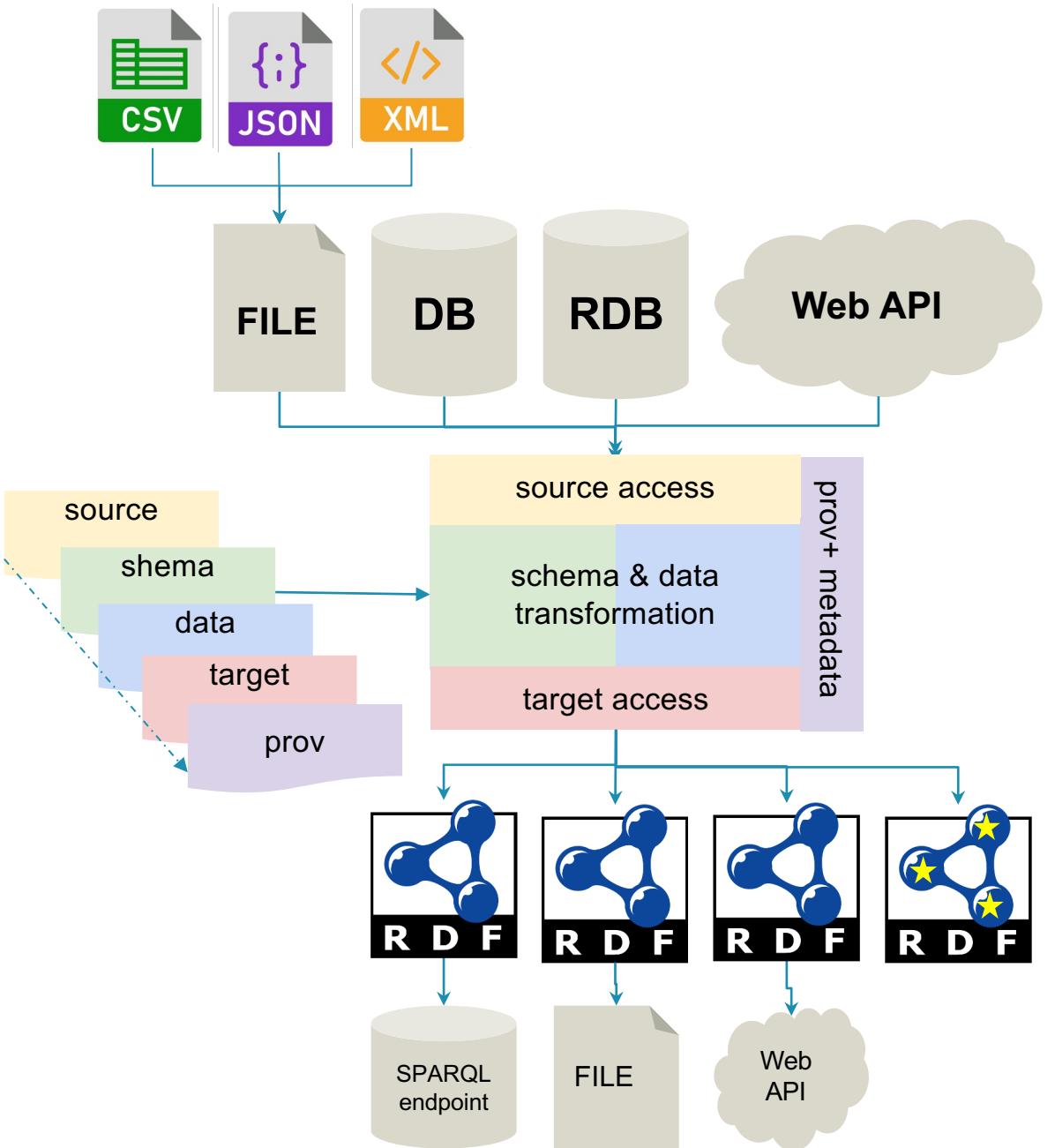
<https://github.com/DylanVanAssche/rml2sparqlanything>

D2RQL-to-RML

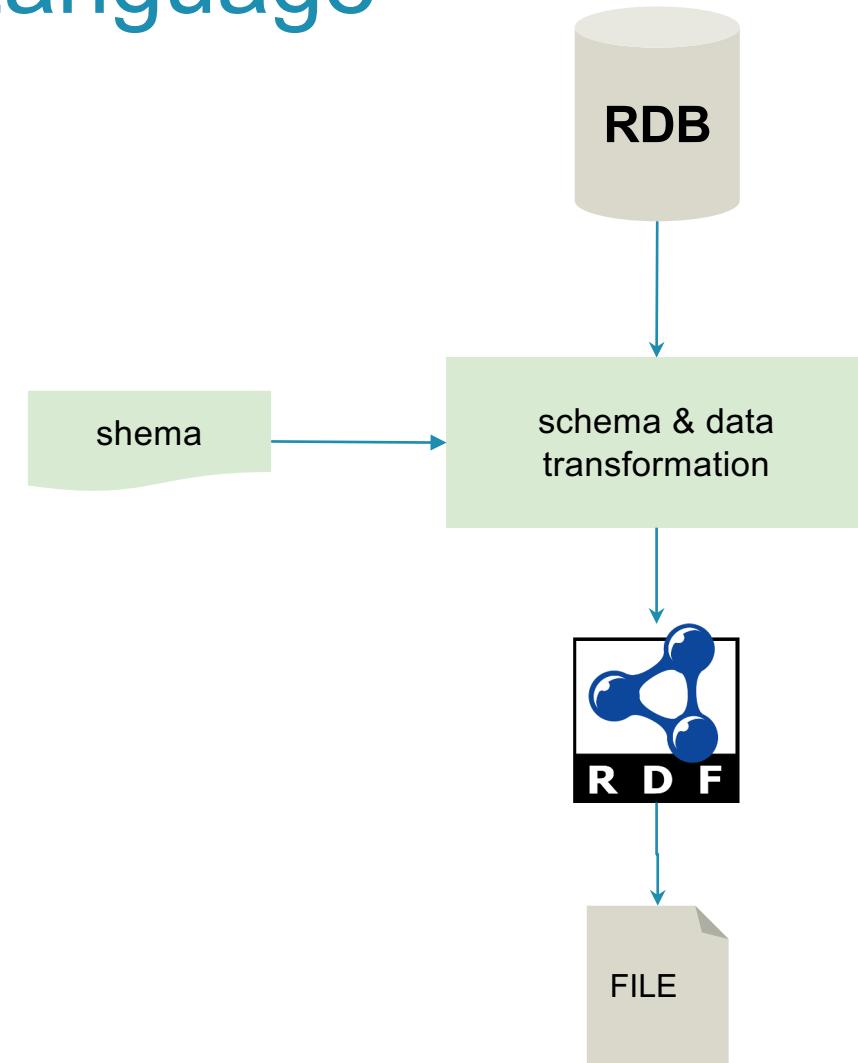
https://github.com/RMLio/D2RQ_to_R2RML



RDF Mapping Language (RML)

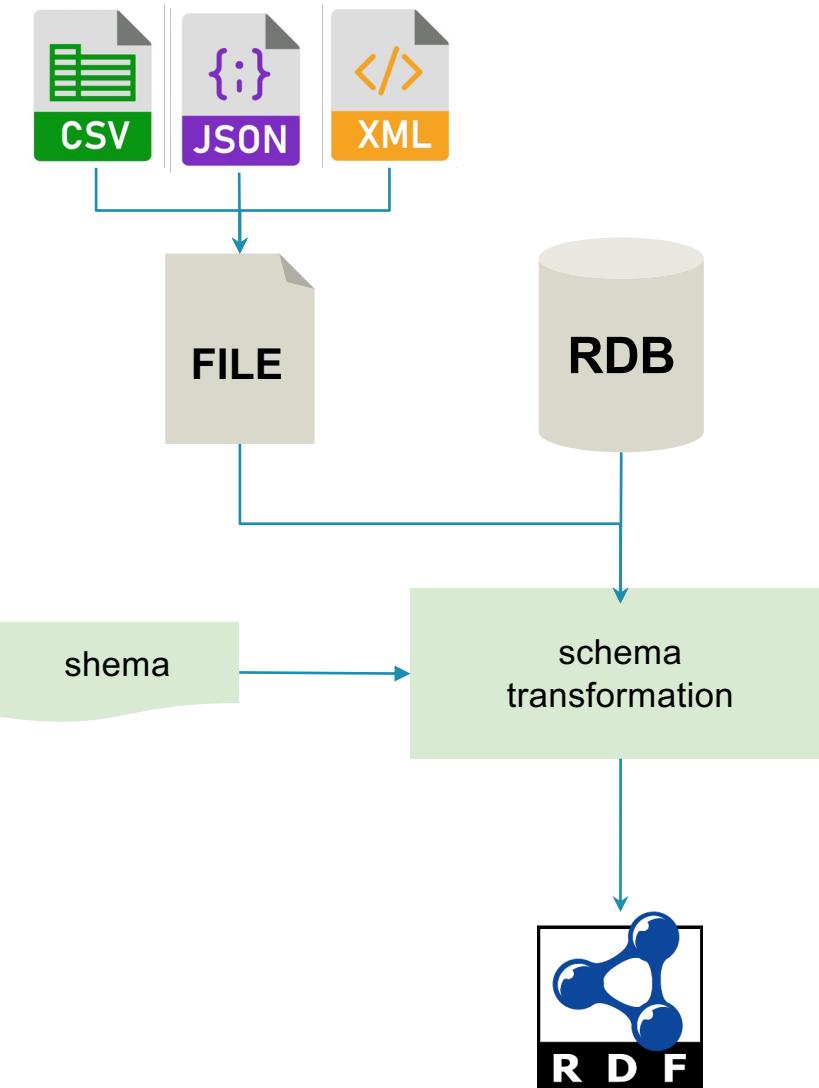


Relational to RDF Mapping Language (R2RML)



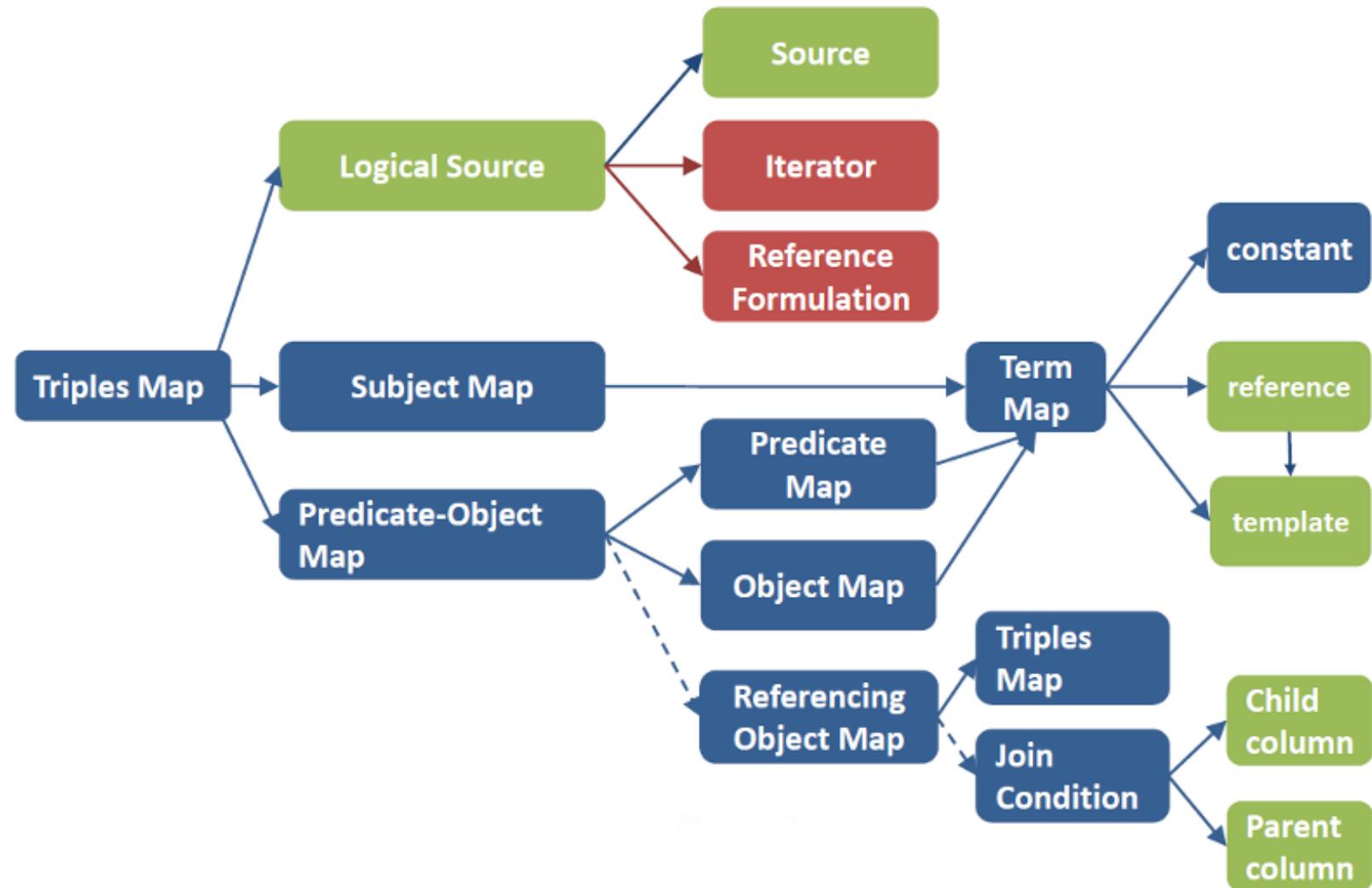
RDF Mapping Language (RML)

rank	name	mark
1	Anzhelika Sidorova	4.95
2	Sandi Morris	4.90
3	Katerina Stefanidi	4.85
4	Holly Bradshaw	4.80
5	Alysha Newman	4.80
6	Angelica Bengtsson	4.80



ex:Anzhelika%20Sidorova rdf:type foaf:Person.
ex:Anzhelika%20Sidorova foaf:name "Anzhelika Sidorova".
ex:Sandi%20Morris rdf:type foaf:Person.
ex:Sandi%20Morris foaf:name "Sandi Morris".

R2RML Vs RML



RML: source description

Local File(s)

Database connectivity

D2RQ www.d2rq.org/d2rq-language

Web source(s) (Web API/service)

DCAT www.w3.org/TR/vocab-dcat-2/

CSVW www.w3.org/ns/csvw

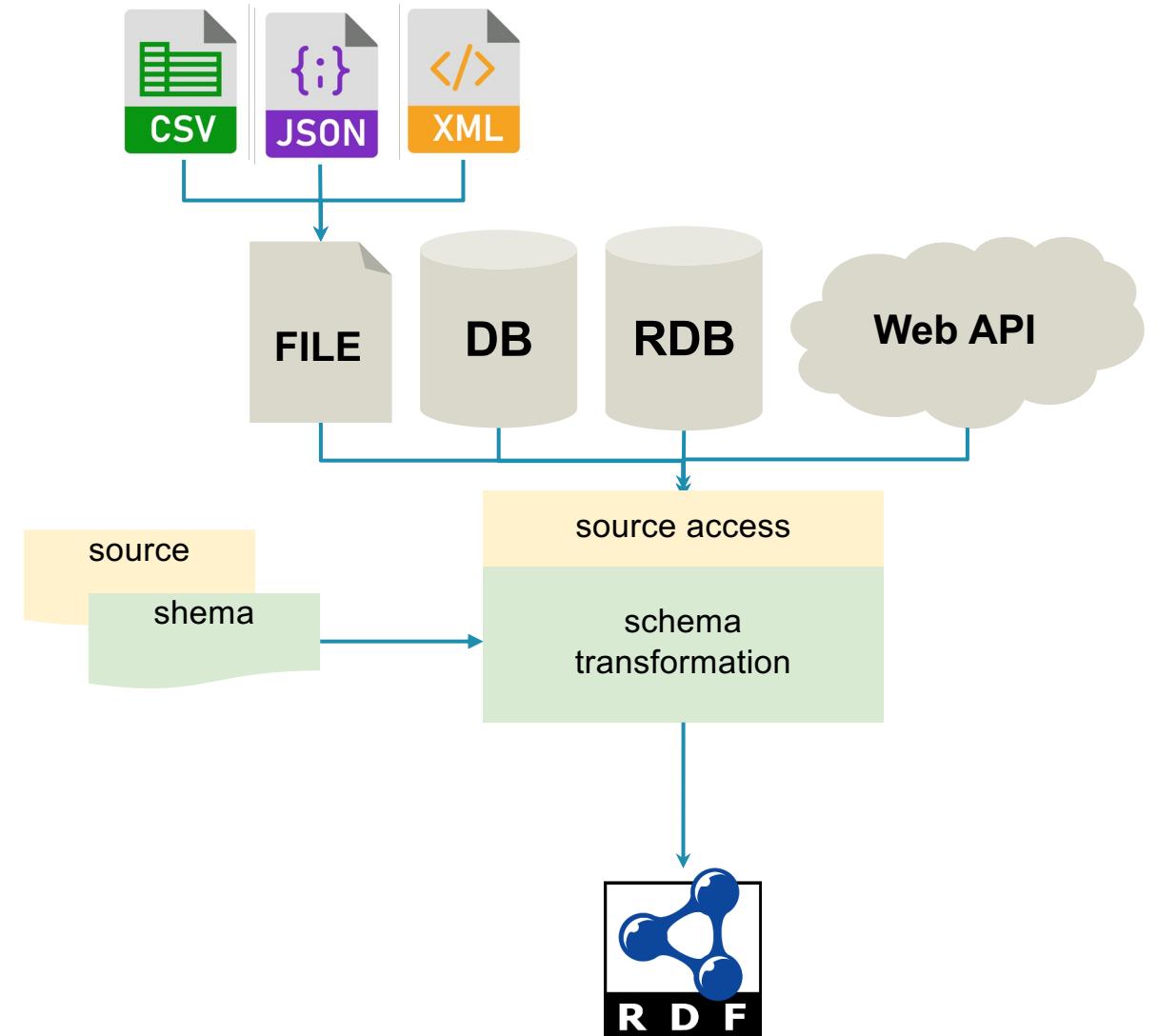
Hydra www.hydra-cg.com/spec/latest/core/

VOiD www.w3.org/TR/void/

RDF source(s)

VOiD (endpoint) www.w3.org/TR/void/

SPARQL-SD www.w3.org/TR/sparql11-service-description/



RML+FnO: data transformations

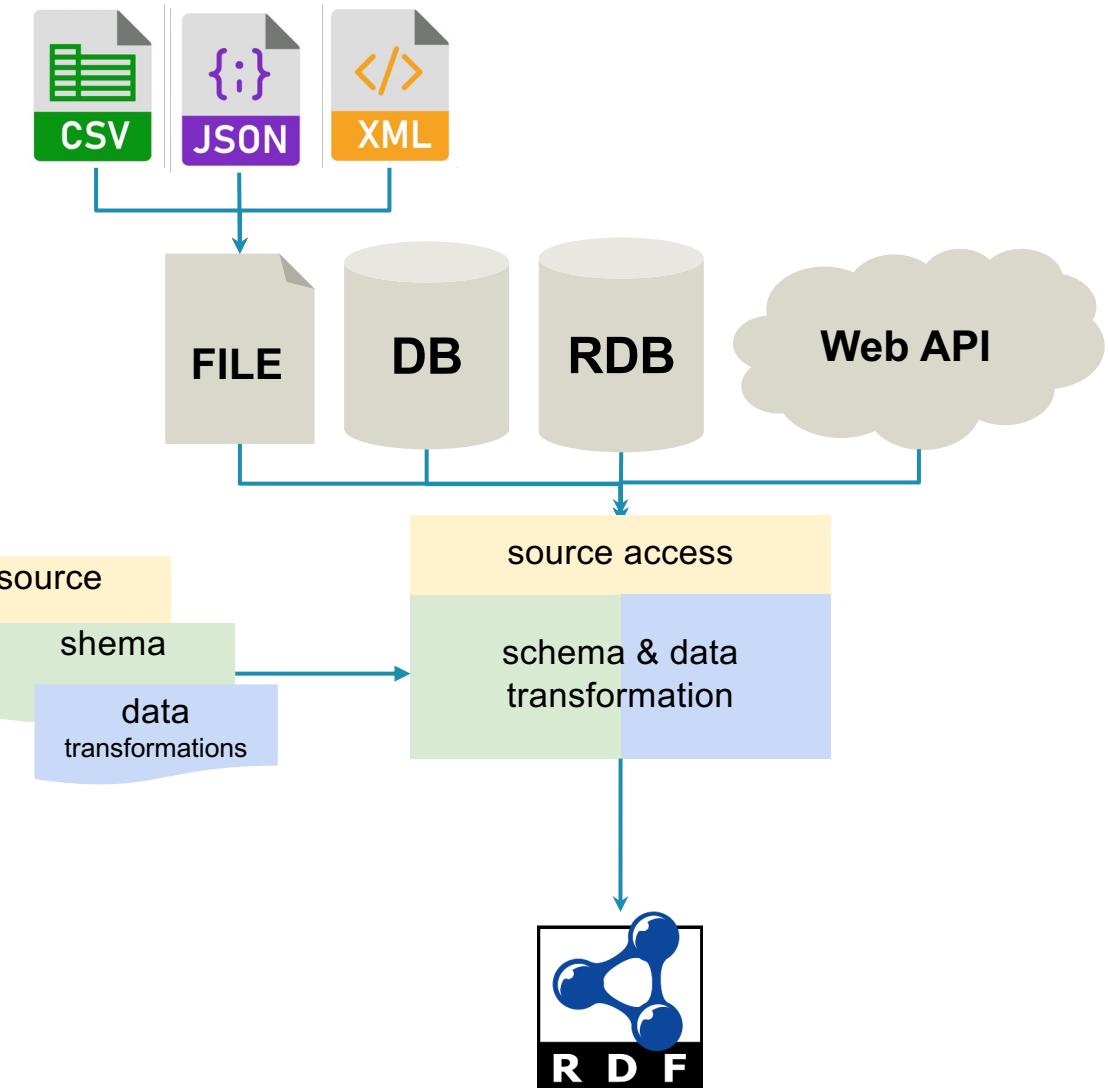
<https://fno.io/spec/>

alternatives:

R2RML-F

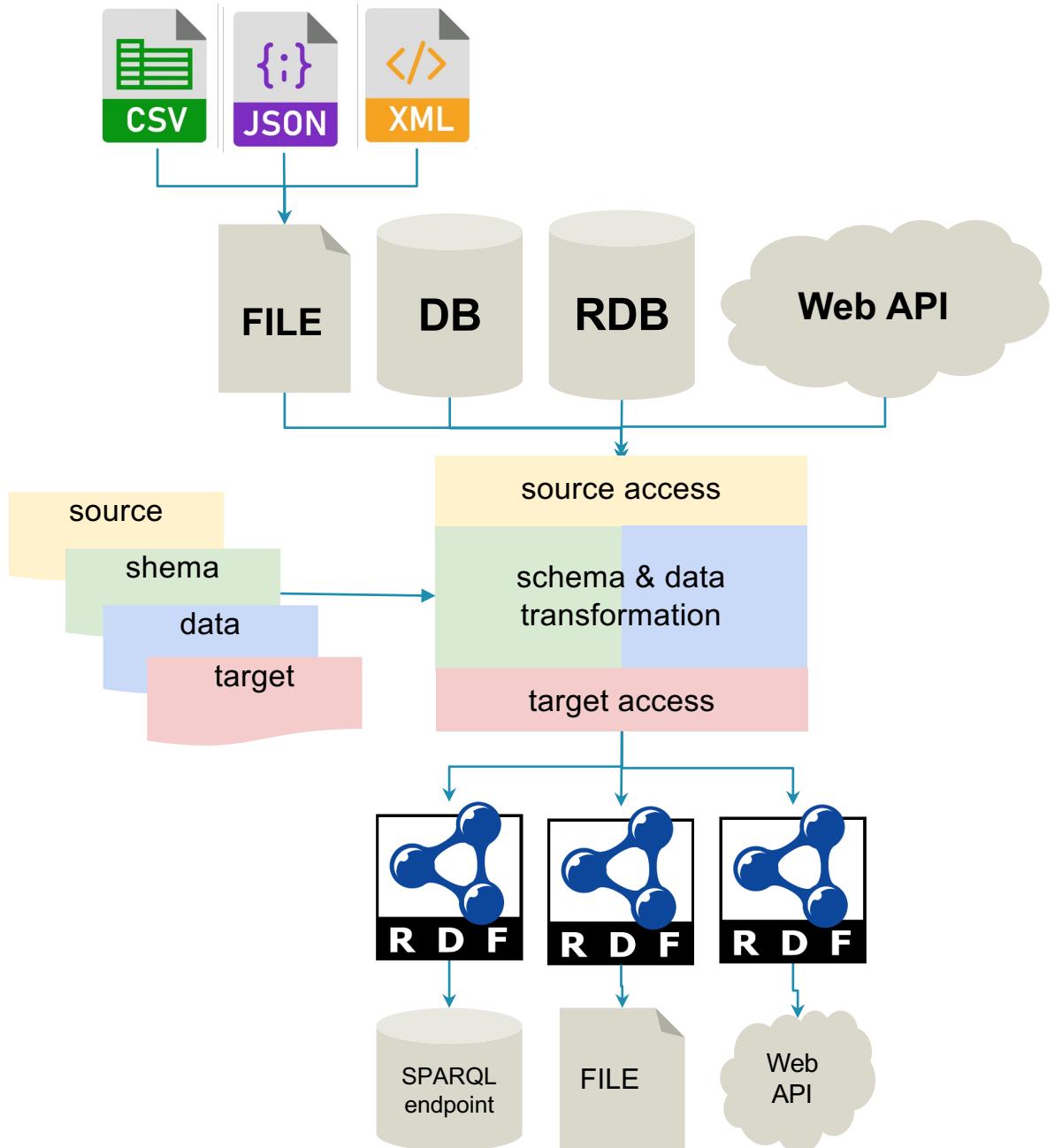
FunUL

rank	name	mark
1	Anzhelika Sidorova	4.95
2	Sandi Morris	4,90
3	Katerina Stefanidi	4.85
4	Holly Bradshaw	4.80
5	Alysha Newman	4.80
6	Angelica Bengtsson	4.80



`ex:Anzhelika%20Sidorova rdf:type foaf:Person.`
`ex:Anzhelika%20Sidorova foaf:name "ANZHELIKA SIDOROVA".`
`ex:Sandi%20Morris rdf:type foaf:Person.`
`ex:Sandi%20Morris foaf:name "SANDI MORRIS".`

RML: target description

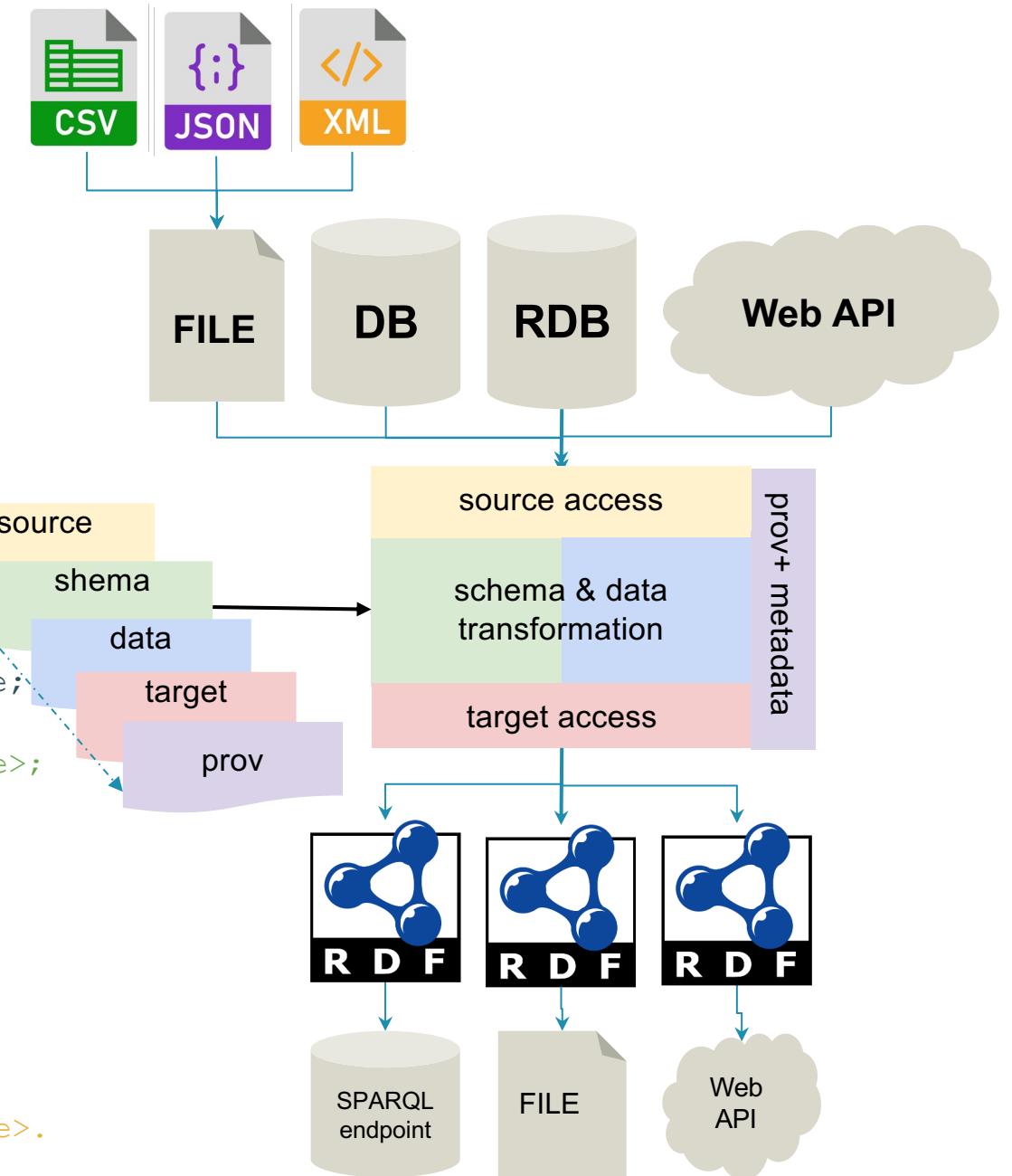


RML: provenance & metadata

```
ex:Anzhelika%20Sidorova rdf:type foaf:Person.  
ex:Anzhelika%20Sidorova foaf:name "Anzhelika Sidorova".  
ex:Sandi%20Morris rdf:type foaf:Person.  
ex:Sandi%20Morris foaf:name "Sandi Morris".
```

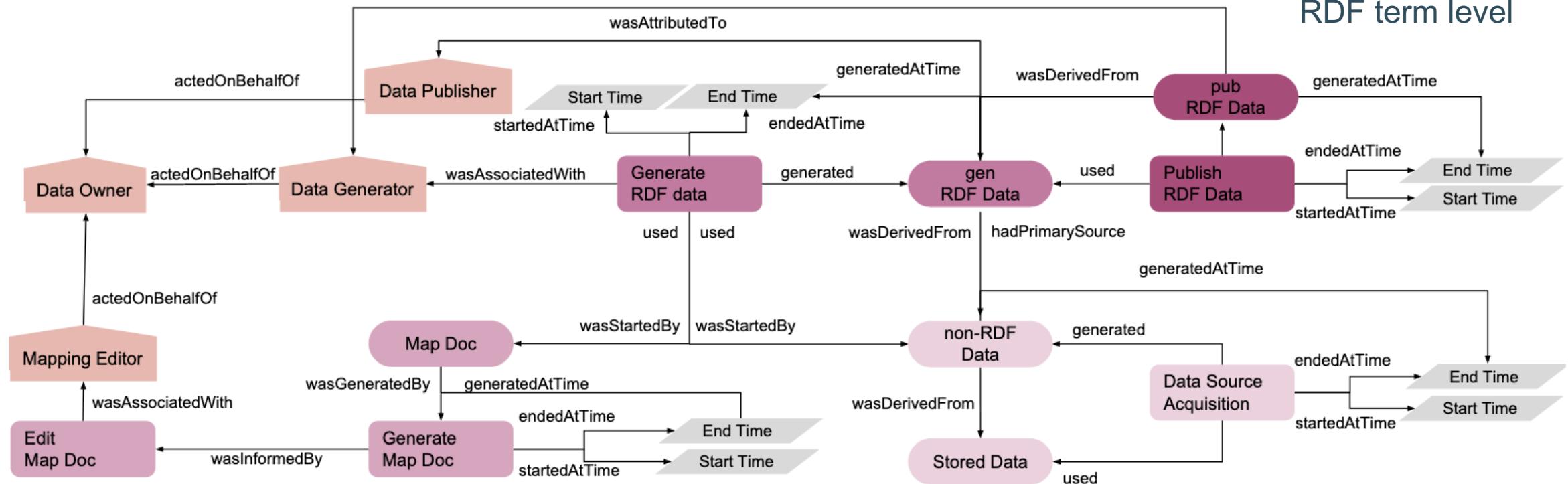
```
<#RDF_Dataset> a prov:Entity, void:Dataset;  
prov:generatedAtTime "2022-01-05T17:10:00Z"^^xsd:dateTime;  
prov:wasGeneratedBy <#RDFdataset_Generation>;  
prov:wasDerivedFrom <#DB_LogicalSource>, <#DCAT_LogicalSource>;  
prov:wasAssociatedWith <#RMLProcessor>;  
prov:wasAttributedTo <http://natadimou.com>.
```

```
<#RDFdataset_Generation> a prov:Activity;  
prov:generated <#RDF_Dataset>;  
prov:startedAtTime "2022-01-05T17:00:00Z"^^xsd:dateTime;  
prov:endedAtTime "2022-01-05T17:10:00Z"^^xsd:dateTime;  
prov:wasInformedBy <#MapDoc_Generation>;  
prov:used <#MapDoc>, <#DB_LogicalSource>, <#DCAT_LogicalSource>.
```

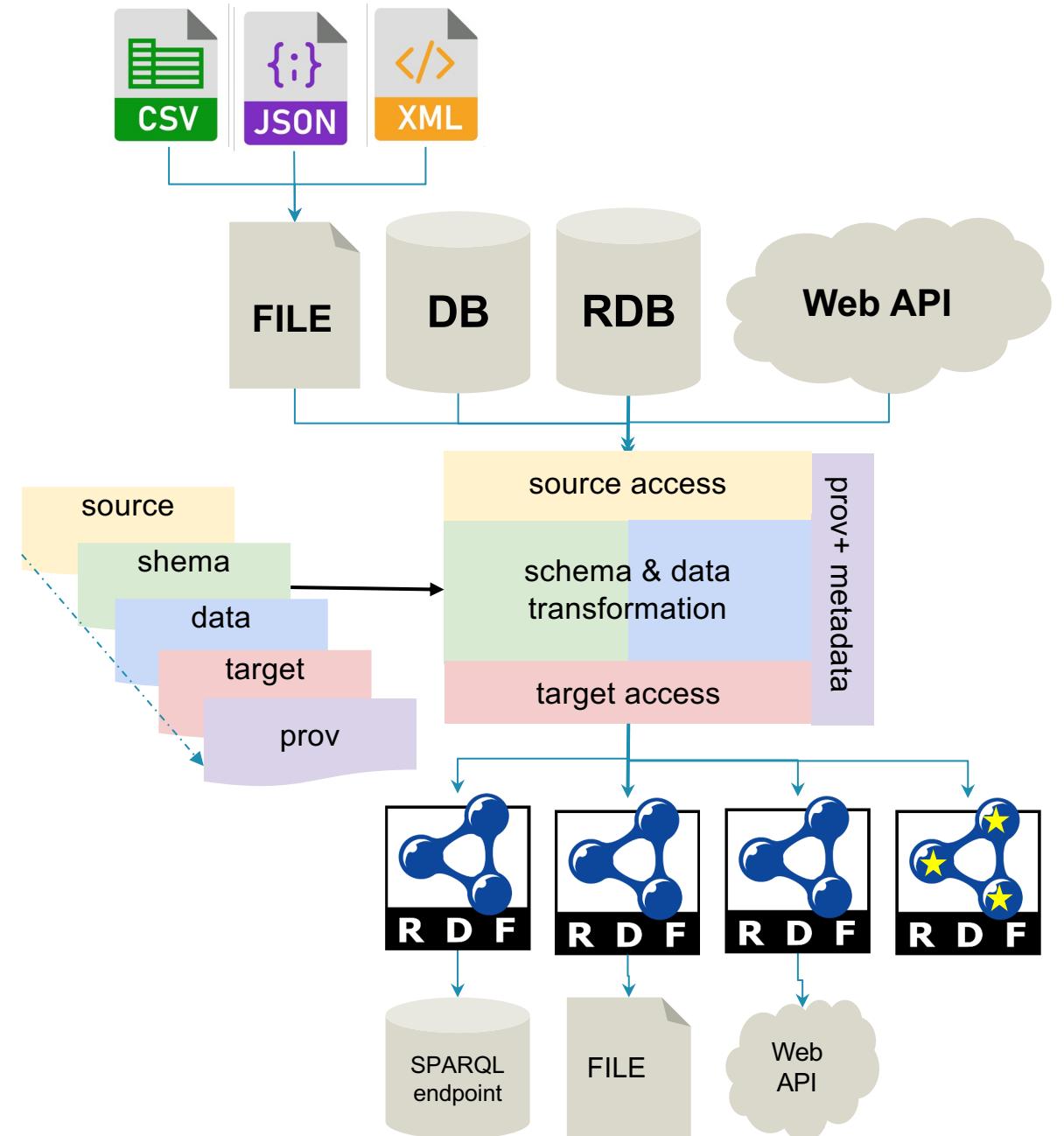
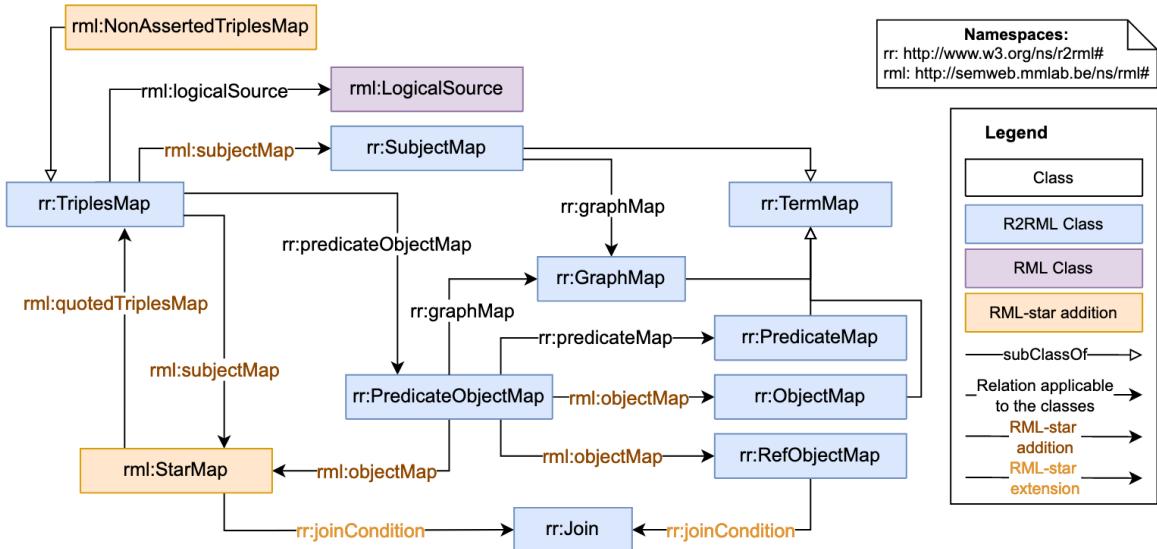


RML: provenance

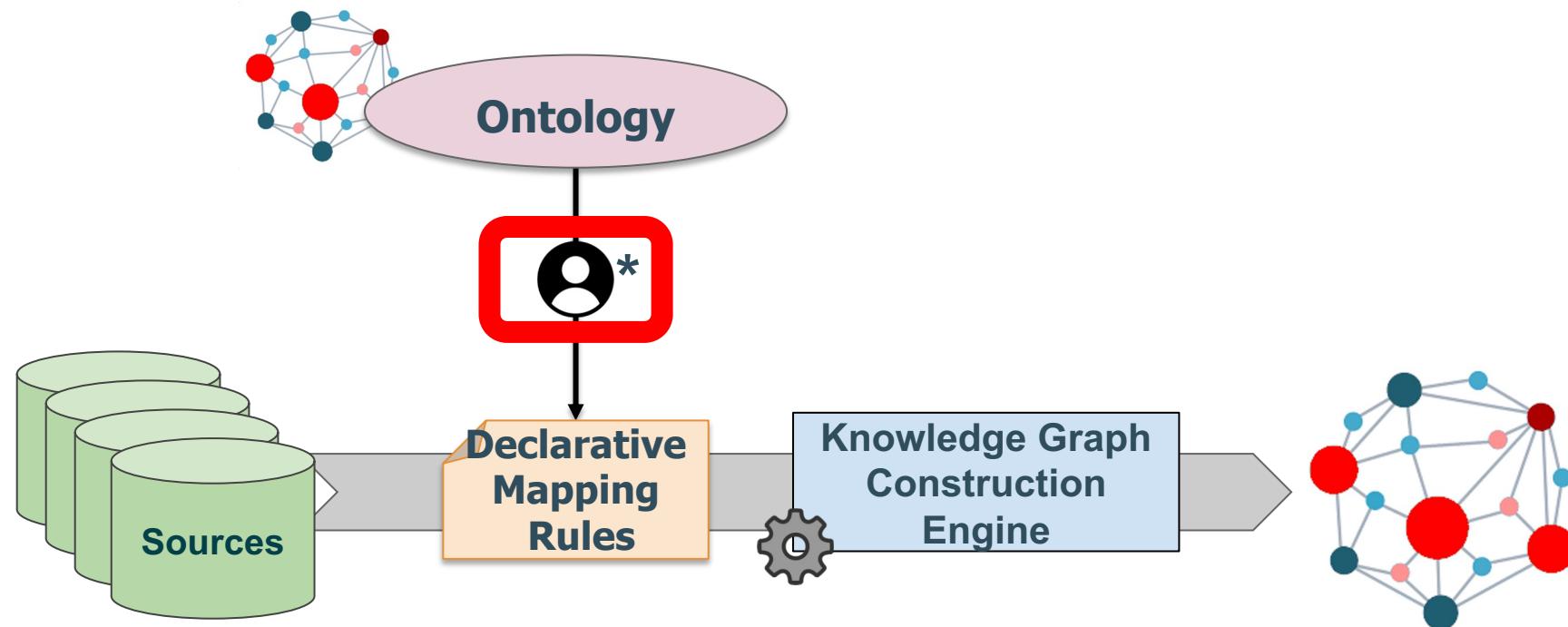
Metadata details levels:
dataset level
named-graph level
subgraph level
RDF triple level
RDF term level



RML-star



Declarative Knowledge Graph Construction



* On average it takes 6 Person-Month to create a knowledge graph

What did I do?

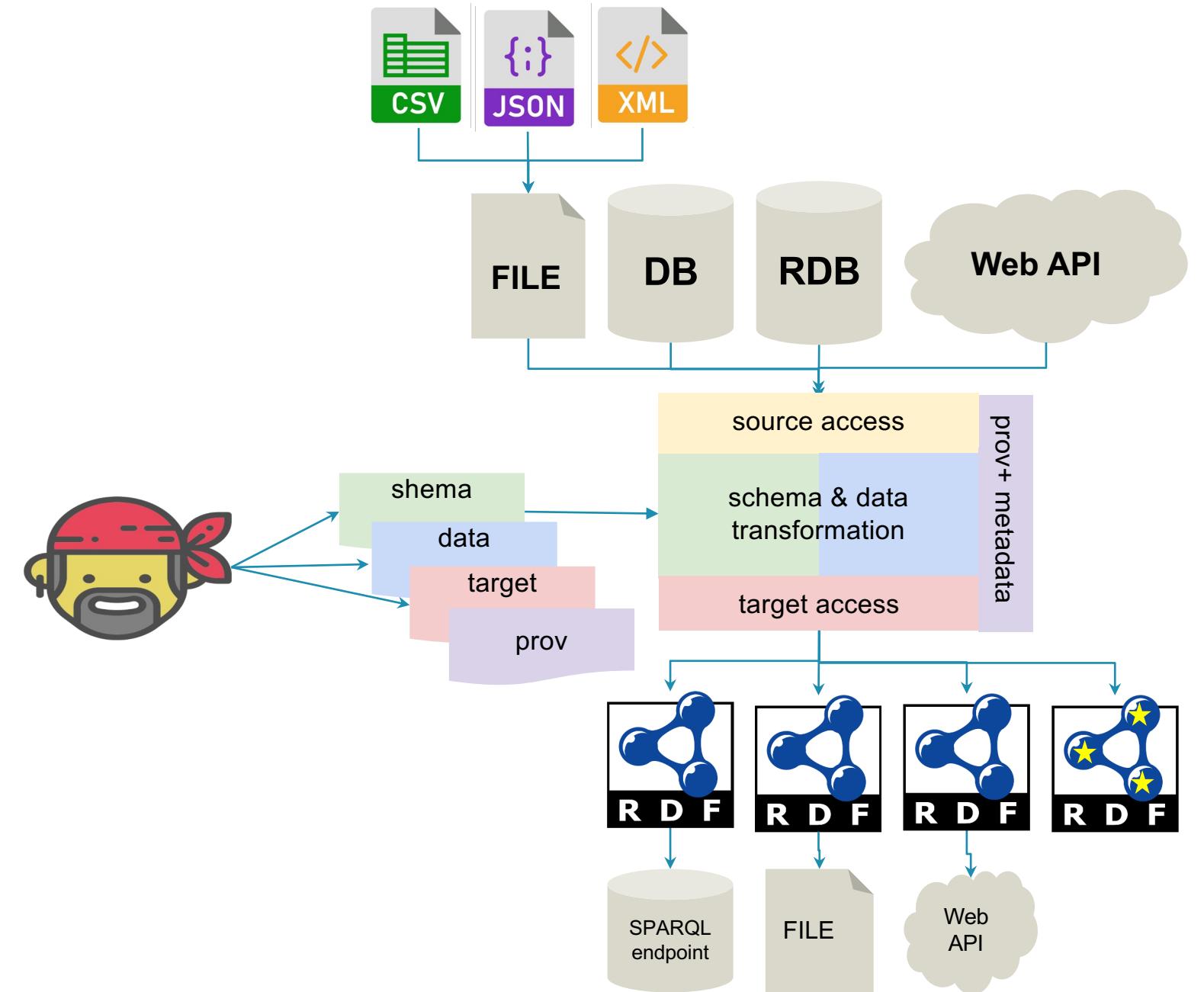
Mapping Languages & RML

Mapping rules & User support

Mapping rules & RDF graphs validation

Mapping rules & SHACL shapes

Implementations & benchmarks





Matey

Everyone need's a matey, this is **YARRRML's Matey!**

See [below](#) to start editing YARRRML-documents!

Or, check the [screencast!](#)

Reload example: [People \(JSON\)](#)

[Advanced](#)

[Facebook](#)

Actions: [Generate RML](#)

[Generate LD](#)

Layout:



Input: data.json ▾



```
1 - {  
2 -   "persons": [  
3 -     {  
4 -       "firstname": "John",  
5 -       "lastname": "Doe"  
6 -     },  
7 -     {  
8 -       "firstname": "Jane",  
9 -       "lastname": "Smith"  
10 -    },  
11 -    {  
12 -      "firstname": "Sarah",  
13 -      "lastname": "Bladinck"  
14 -    }  
15 -  ]  
16 }
```

Input: YARRRML ▾

```
1 - prefixes:  
2 -   ex: "http://example.com/"  
3 -  
4 - mappings:  
5 -   person:  
6 -     sources:  
7 -       - ['data.json~jsonpath', '$.persons[*]']  
8 -       s: http://example.com/${firstname}  
9 -     po:  
10 -      - [a, foaf:Person]  
11 -      - [ex:name, ${firstname}]
```

Output: Turtle/TriG ▾

```
1 |  
2 |  
3 |
```

File | Edit | Mapping | View | Help

RML EDITOR

Detail Lowest Low Moderate High Highest

directors.csv

directors.csv **[+]**

ID

first_name

last_name

date_of_birth

country

Detail **[X]** **[D]**

{title}

rdfs:label

...om/{ID} schema:Movie

...t_name{}

...om/{ID} ...vieDirector

Level Subject ^ Predicate Object

	http://example.com/0	rdf:type	schema:Movie
	http://example.com/0	rdfs:label	The Shawshank Redem
	http://example.com/1	rdf:type	schema:Movie
	http://example.com/1	rdfs:label	The Godfather
	http://example.com/2	rdf:type	schema:Movie
	http://example.com/2	rdfs:label	The Godfather: Part II
	http://example.com/3	rdf:type	schema:Movie
	http://example.com/3	rdfs:label	The Dark Knight
	http://example.com/4	rdf:type	schema:Movie
	http://example.com/4	rdfs:label	12 Angry Men
	http://example.com/5	rdf:type	schema:Movie
	http://example.com/5	rdfs:label	Schindlers List
	http://example.com/6	rdf:type	schema:Movie
	http://example.com/6	rdfs:label	Pulp Fiction
	http://example.com/7	rdf:type	schema:Movie
	http://example.com/7	rdfs:label	The Good the Bad and the Ugly
	http://example.com/8	rdf:type	schema:Movie
	http://example.com/8	rdfs:label	The Lord of the Rings
	http://example.com/9	rdf:type	schema:Movie

ID first_name last_name date_of_birth

dir_AAA	Frank	Darabont	28-01-1955
dir_AAB	Francis Ford	Coppolla	07-04-1939
dir_AAC	Christopher	Nolan	30-07-1970

Ready

Editors

Matey

<https://github.com/rmlio/matey>

RMLEditor

<https://app.rml.io/rmleditor/>, <https://rml.io/tools/rmleditor/>, <https://github.com/RMLio/rmleditor-ce>

Mapeauthor

<https://morph.oeg.fi.upm.es/tool/mapeauthor>, <https://github.com/oeg-upm/morph-website>

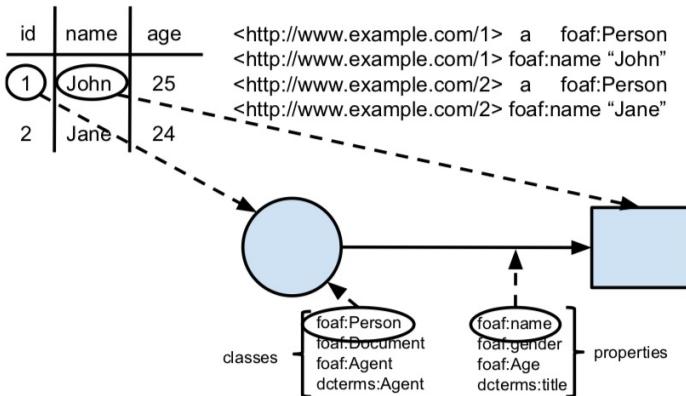
Map-On

<http://semanco-tools.eu/map-on>, <https://github.com/arc-lasalle/Map-On>

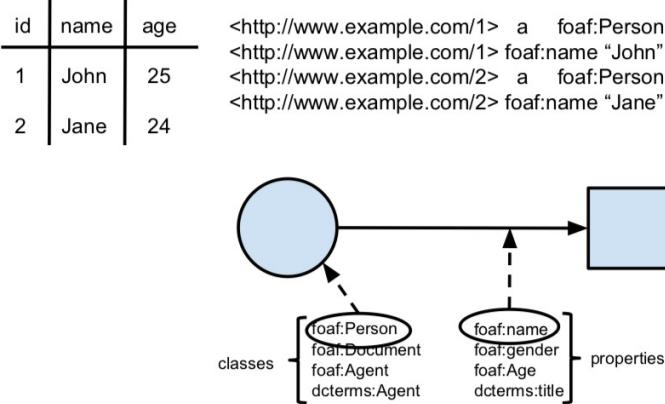
RMLx Visual Editor

<http://pebbie.org/mashup/rml>

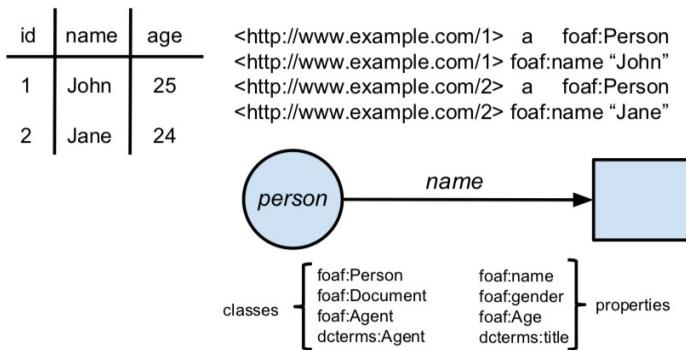
Approaches for generating mapping rules



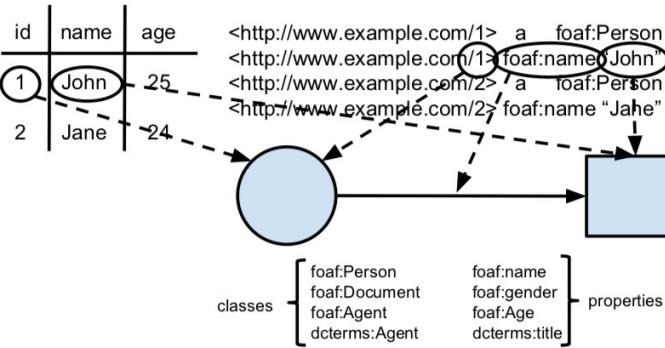
(a) data-driven



(b) schema-driven

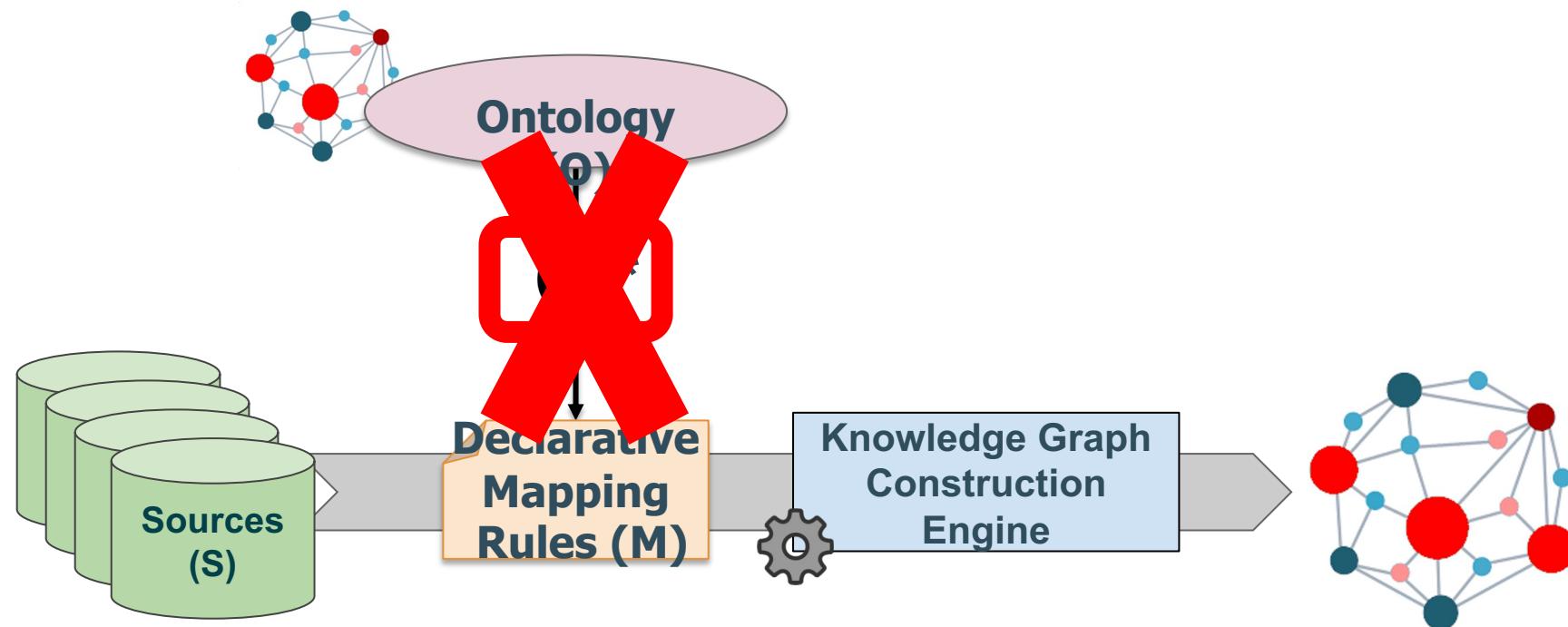


(c) model-driven



(d) result-driven

Declarative Knowledge Graph Construction



* On average it takes 6 Person-Month to create a knowledge graph

The million dollar question!

How can we automate the knowledge graph construction and reduce its time and effort?



Direct Knowledge Graph Construction

wrapper systems

FDR2 (2004)

D2RQ (2004)

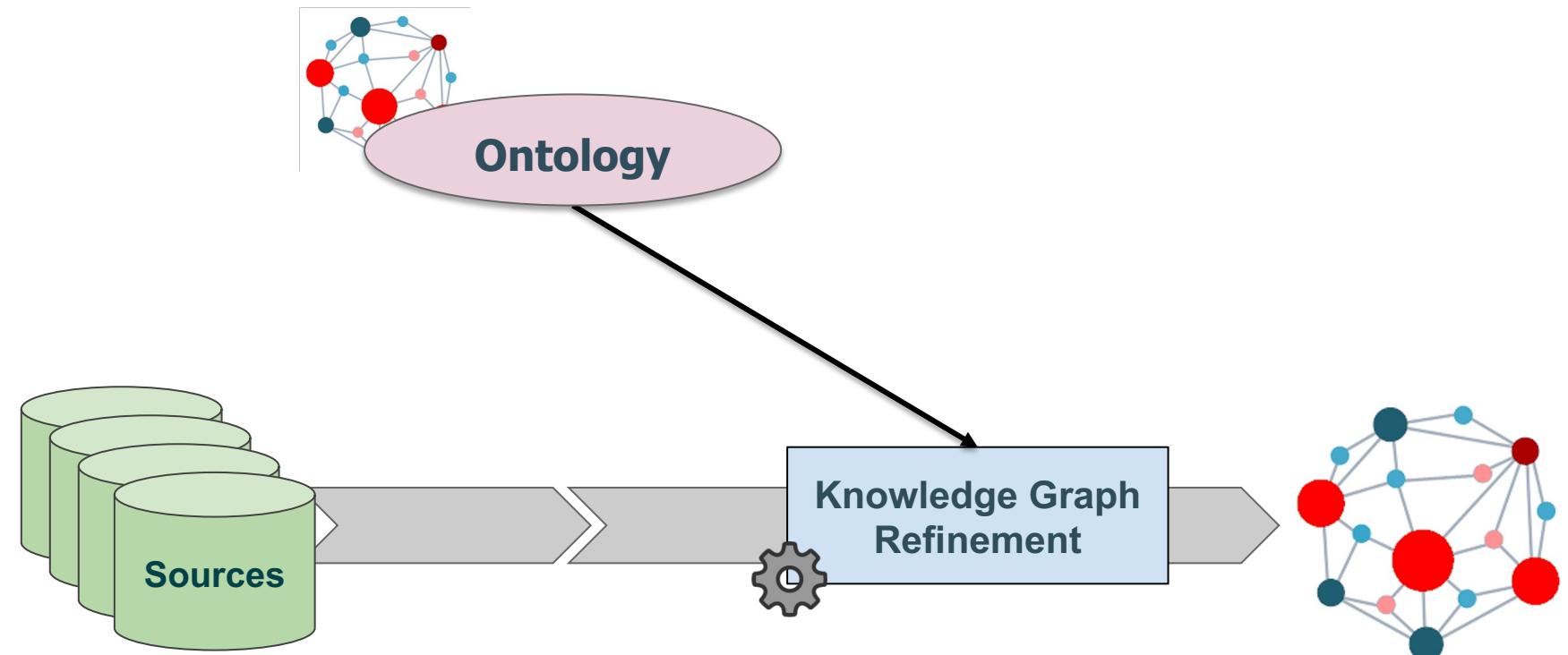
METAmorphoses (2004)

Relational.OWL (2005)

R₂O (2006)

RDB2Onto (2006)

D2OMapper (2006)

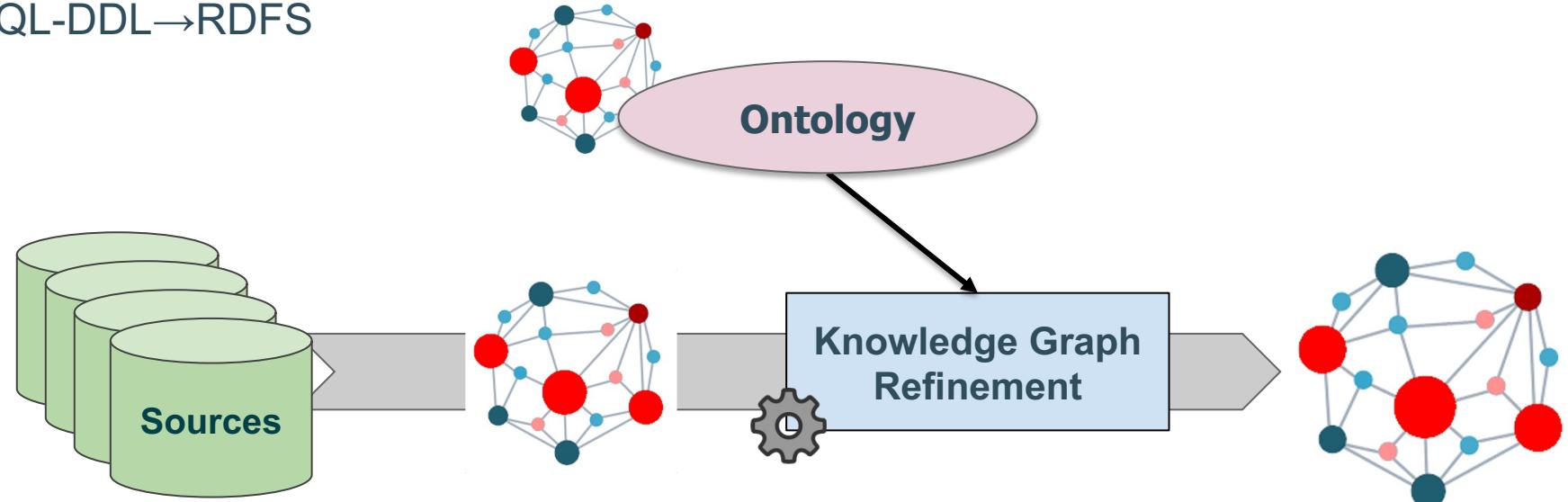


Direct Knowledge Graph Construction

direct mapping

Buccella et al. (2004)
Li et al. (2005)
Shen et al. (2006)
Stojanovic et al. (2002)
Astrova et al. (2004)

} SQL-DDL→OWL
} SQL-DDL→RDFS



Direct Knowledge Graph Construction

declarative automation

Mirror (2015)

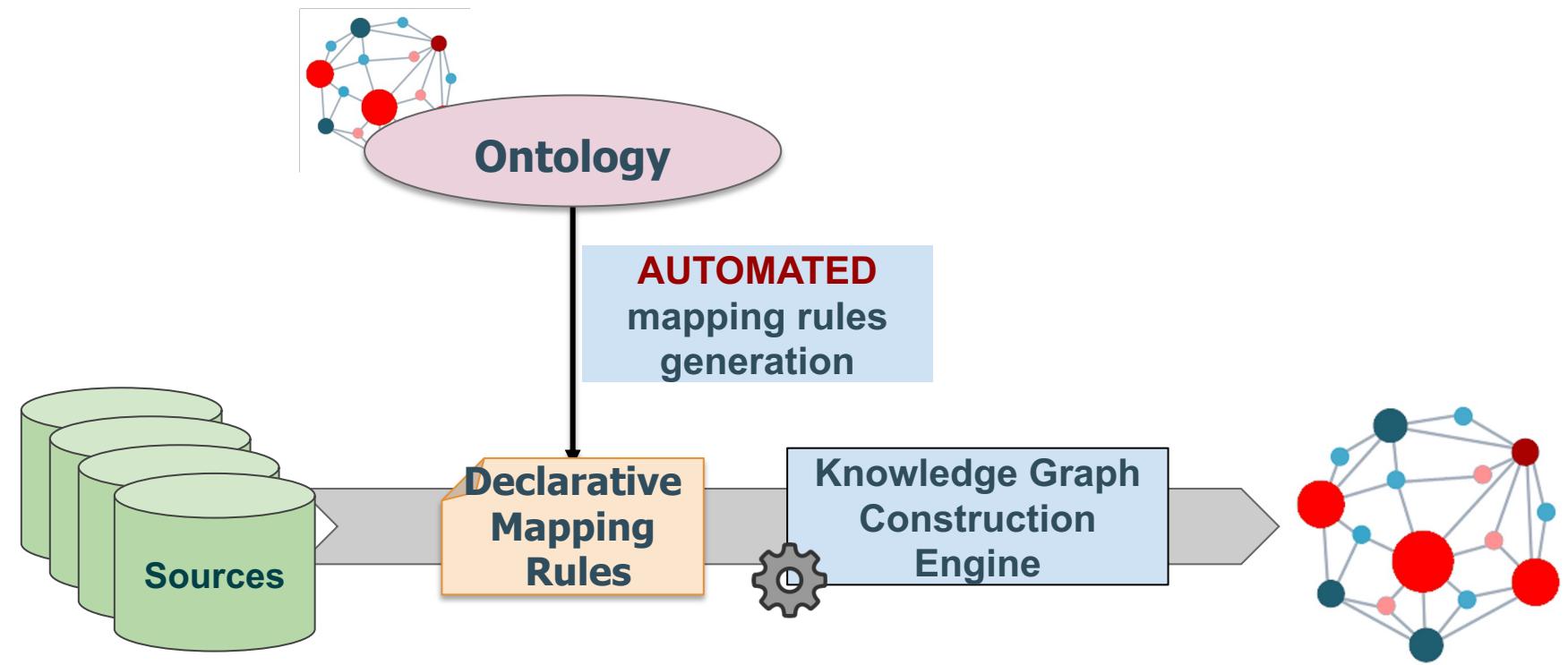
AutoMap4OBDA (2016)

BootOX (2015)

IncMap (2013)

COMA++ (2005)

Ontop (... - nowadays)



RODI benchmark → disappointing results

Direct Knowledge Graph Construction

SemTab challenge

MTab (2019,2020,2021)

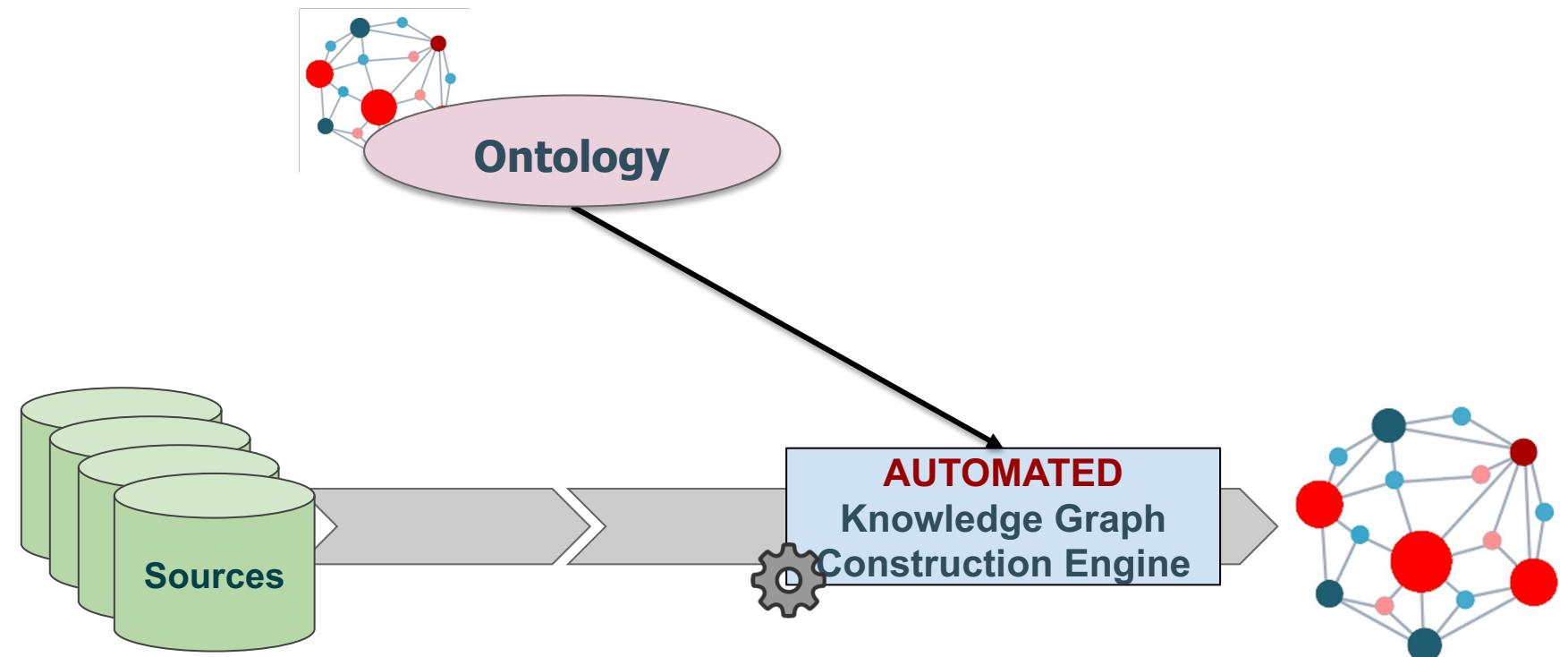
MAGIC (2020)

MantisTable (2021)

JenTab (2021)

DAGOBAH (2019,2020,2021)

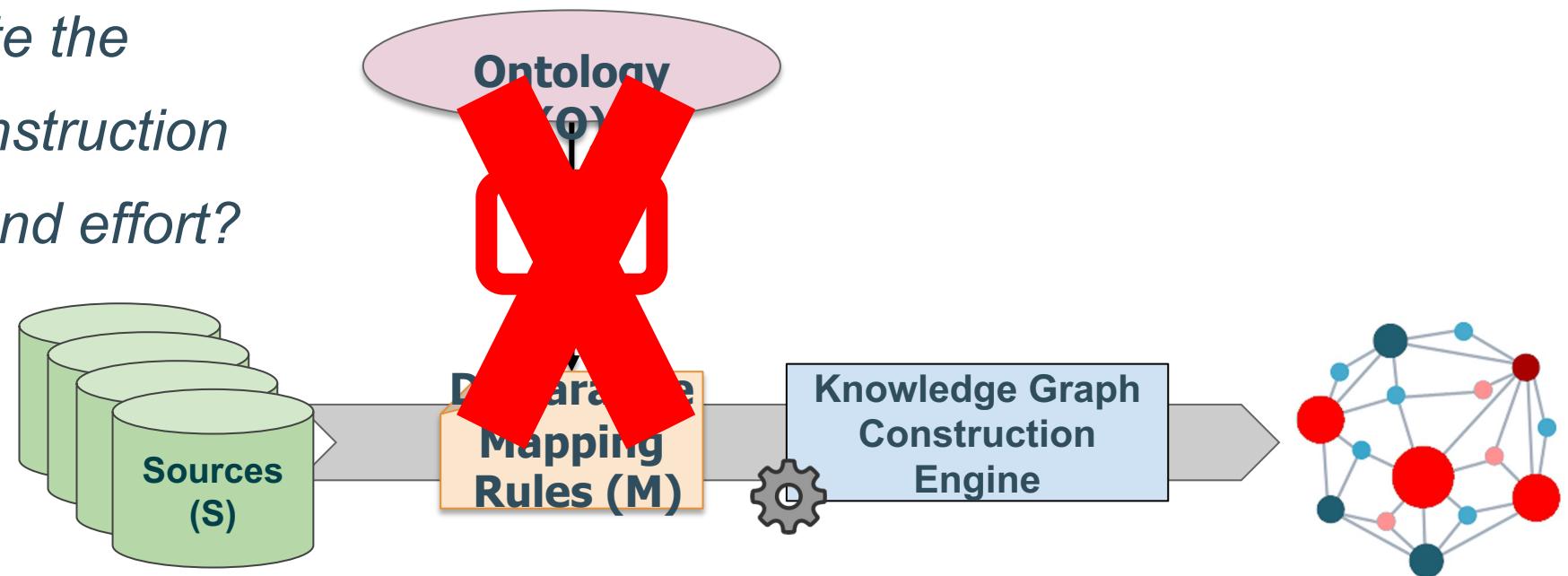
....



SemTab Challenge → overfit ⇒ lack of generalisation

The million dollar question!

How can we automate the knowledge graph construction and reduce its time and effort?



What did I do?

Mapping Languages & RML

Mapping rules & User support

Mapping rules validation

Mapping rules & SHACL shapes

Implementations & benchmarks

rank	name	mark
1	Anzhelika Sidorova	4.95
2	Sandi Morris	4.90
3	Katerina Stefanidi	4.85
4	Holly Bradshaw	4.80
5	Alysha Newman	4.80
6	Angelica Bengtsson	4.80

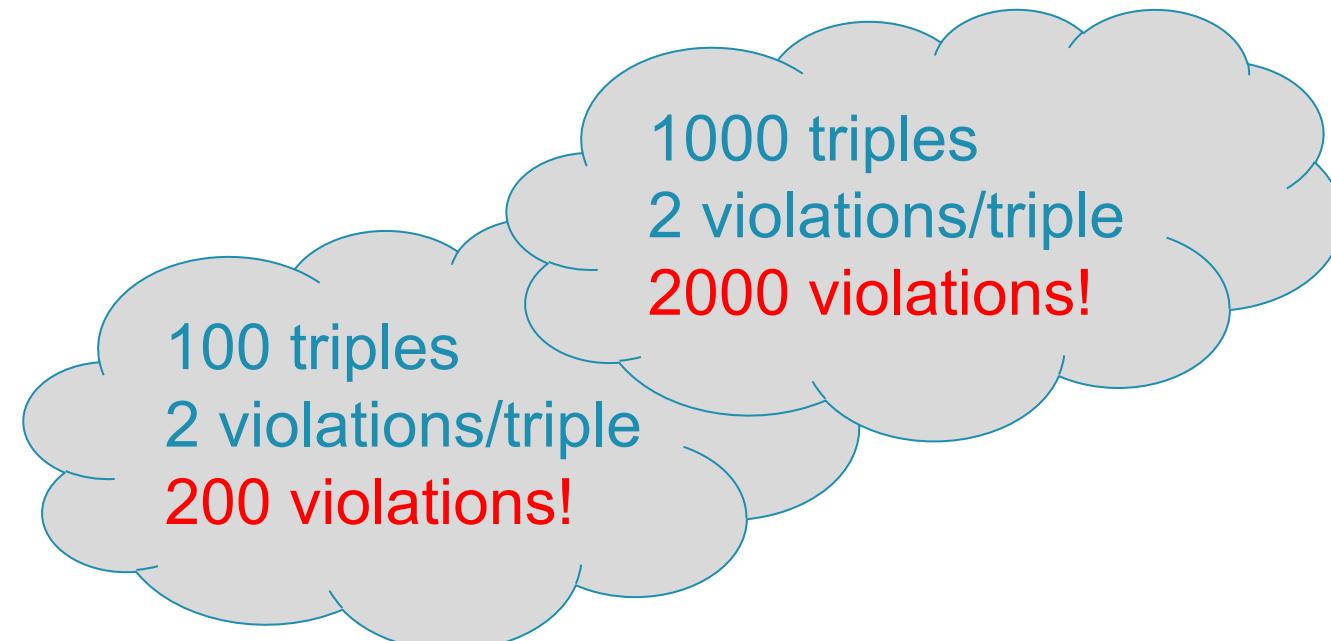


```

ex:Anzhelika%20Sidorova rdf:type foaf:Person bibo: Document .
ex:Anzhelika%20Sidorova foaf:name "Anzhelika Sidorova"^^xsd:integer.
ex:Sandi%20Morris rdf:type foaf:Person bibo: Document .
ex:Sandi%20Morris foaf:name "Sandi Morris"^^xsd:integer.

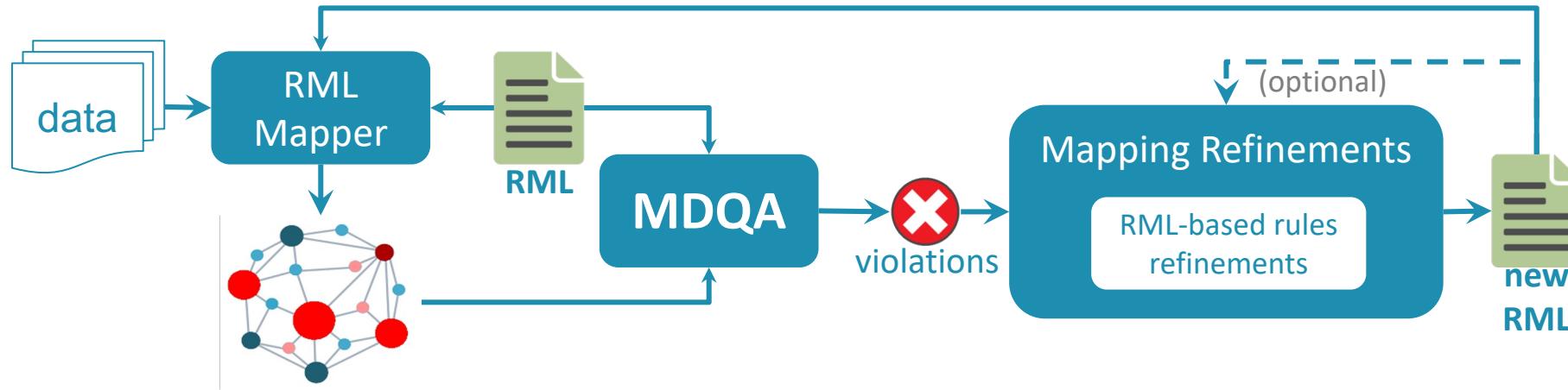
```

rank	name	mark
1	Anzhelika Sidorova	4.95
2	Sandi Morris	4.90
3	Katerina Stefanidi	4.85
4	Holly Bradshaw	4.80
5	Alysha Newman	4.80
6	Angelica Bengtsson	4.80



ex:Anzhelika%20Sidorova rdf:type foaf:Person bibo: Document .
 ex:Anzhelika%20Sidorova foaf:name "Anzhelika Sidorova"^^xsd:integer.
 ex:Sandi%20Morris rdf:type foaf:Person bibo: Document .
 ex:Sandi%20Morris foaf:name "Sandi Morris"^^xsd:integer.

Mapping rules quality assessment

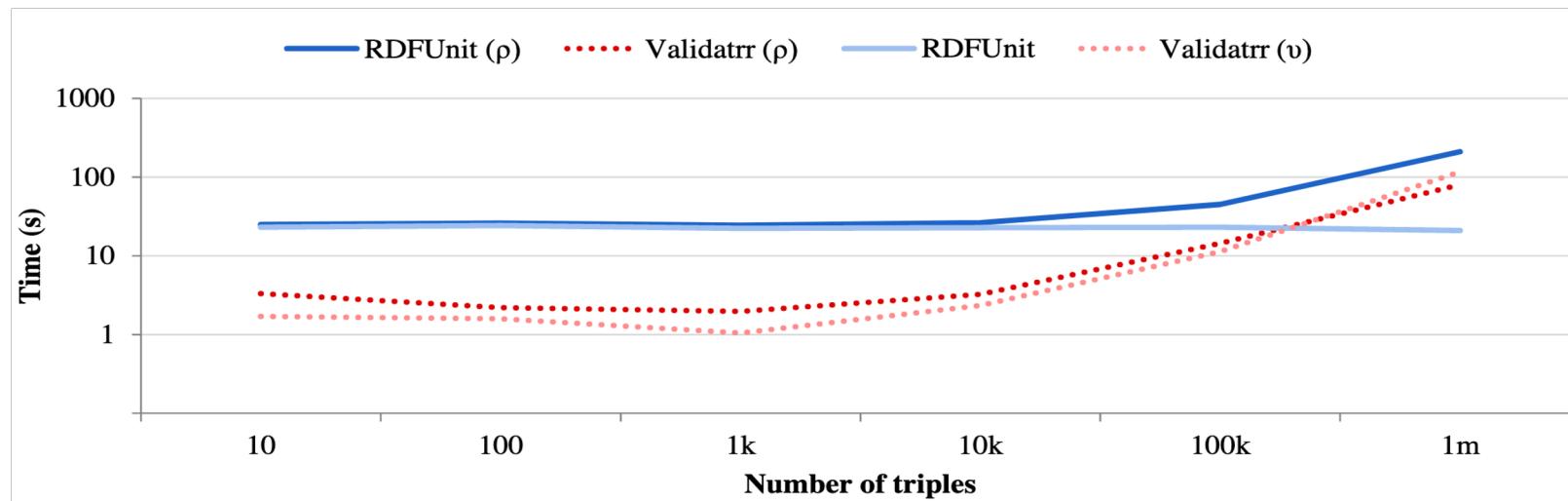
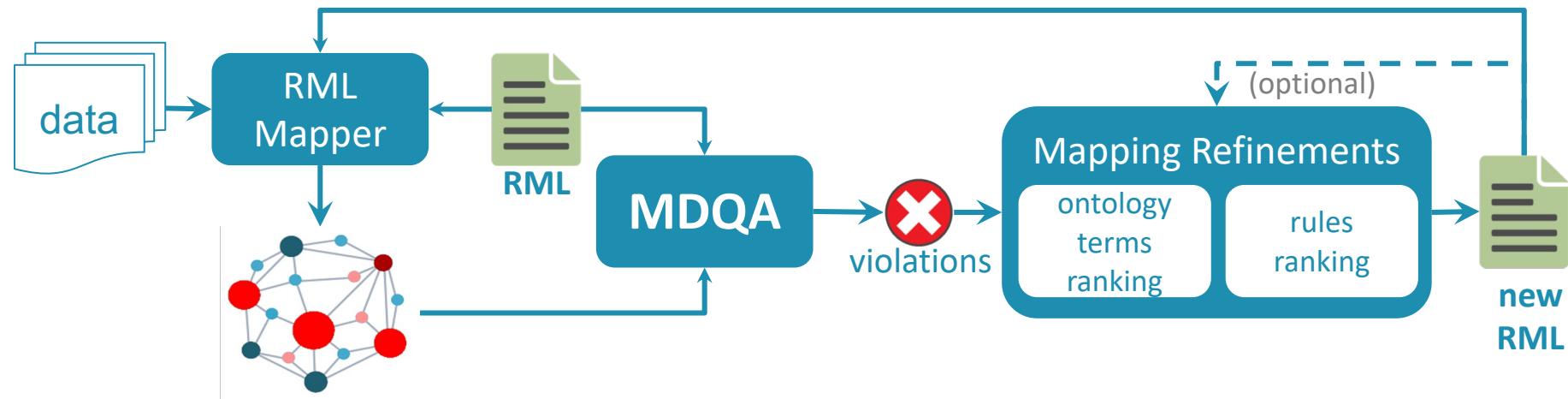


✓ DBpedia Quality Assessment
knowledge graph: 16h
RML rules: 32s

✗ Certain test cases require a complete Linked Data set
e.g., (qualified) cardinality, (inverse) functionality, (a)symmetry, irreflexivity

Assumption
ontology terms are perfect!

Mapping rules quality assessment



datasets have more than 100k triples
lodstats: 94%
lodlaundromat: 60%

What did I do?

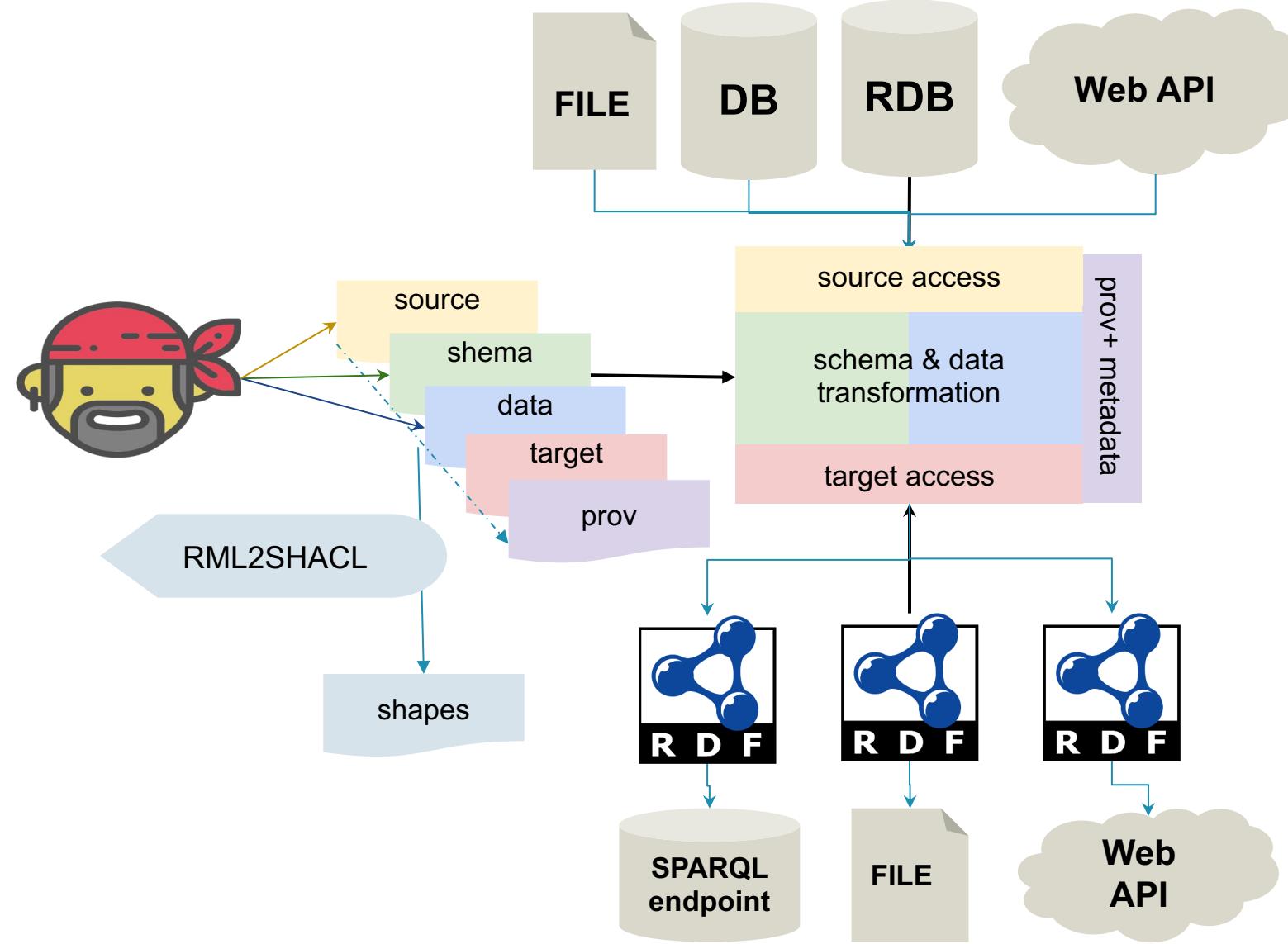
Mapping Languages & RML

Mapping rules & User support

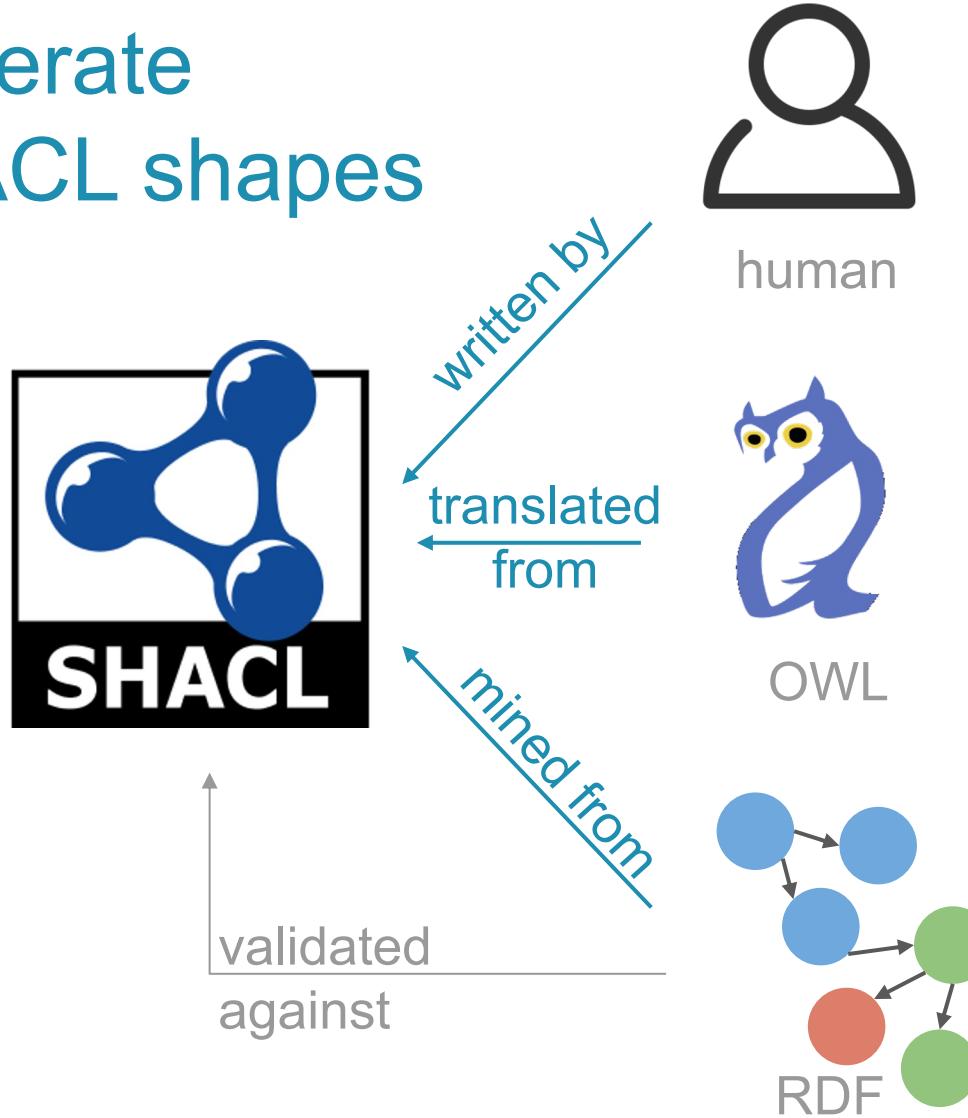
Mapping rules & RDF graphs validation

Mapping rules & SHACL shapes

Implementations & benchmarks



Generate SHACL shapes

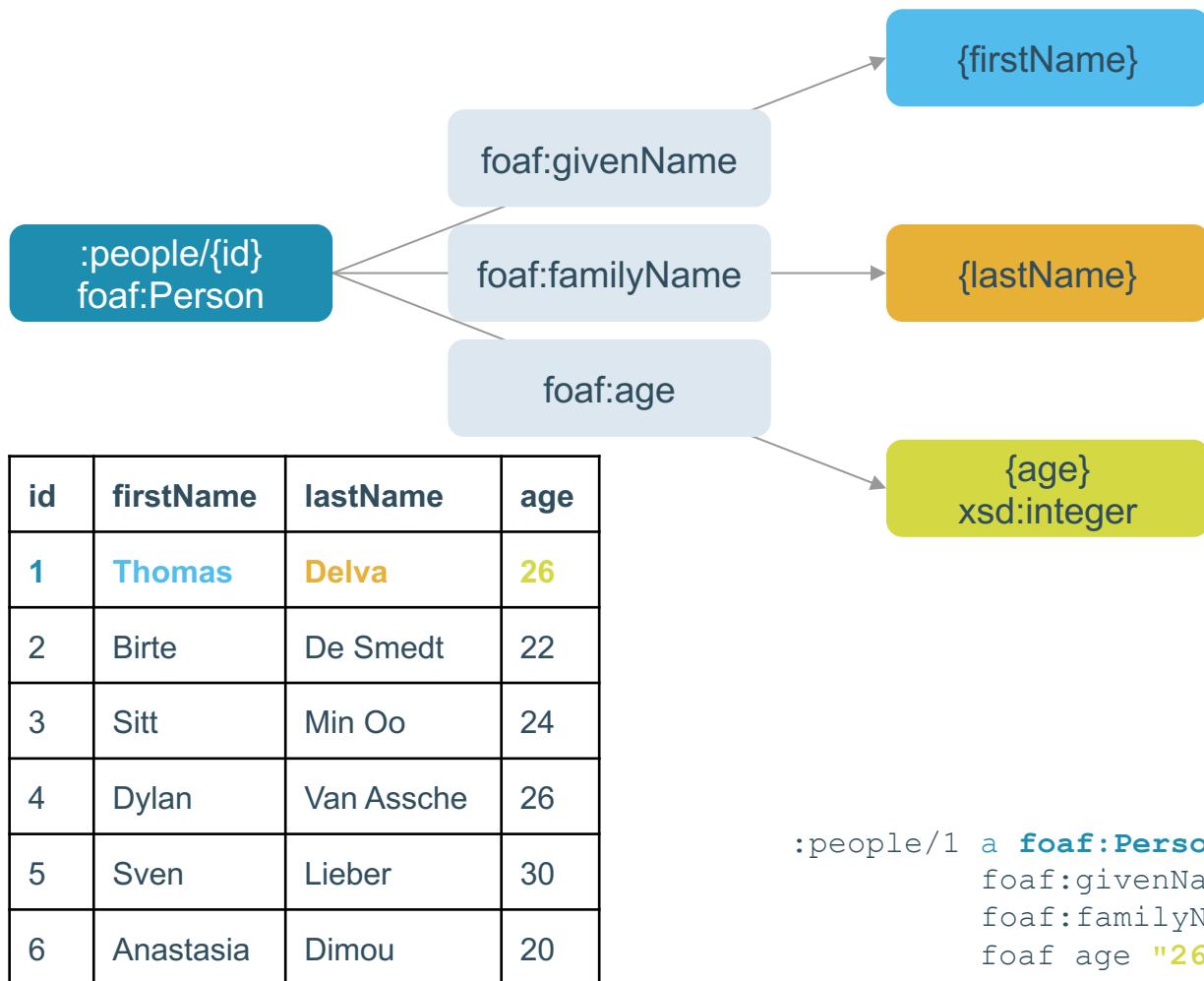


- + most control
- human effort needed

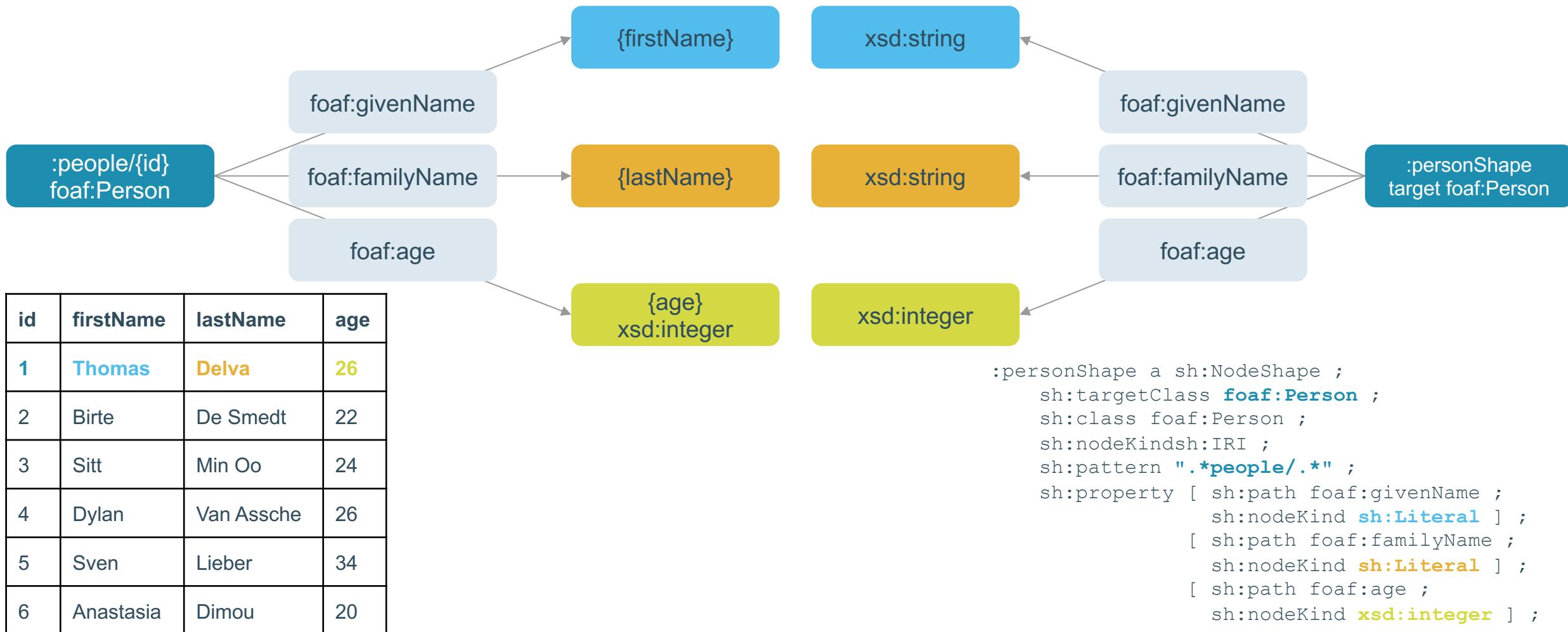
- + constant in data-size
- limited to ontological constraints,
ontology required

- + always possible
- expensive for large data,
uncertain in nature

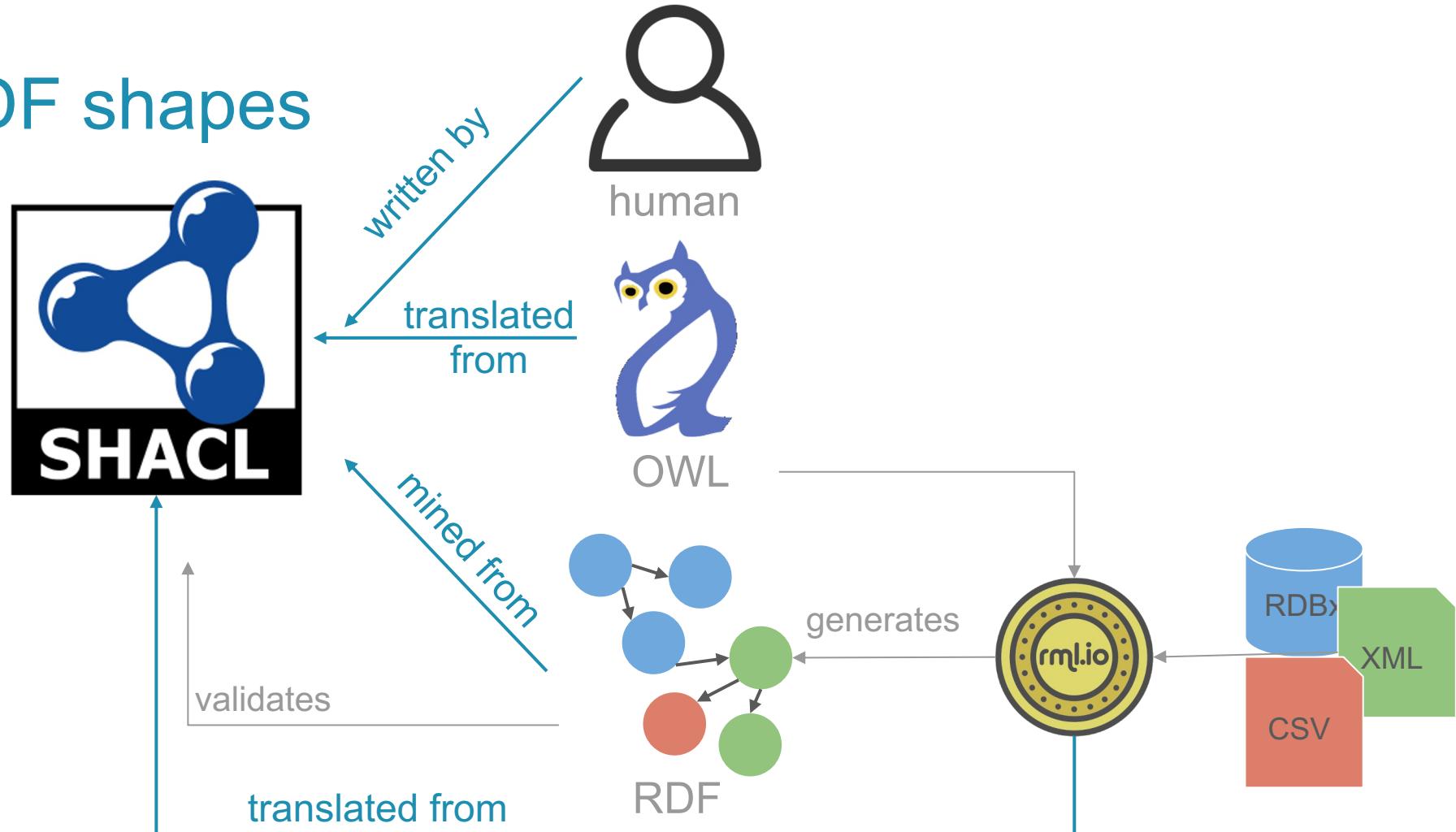
Generate RDF with RML



RML mapping rules & SHACL shapes



Generate RDF shapes from RML

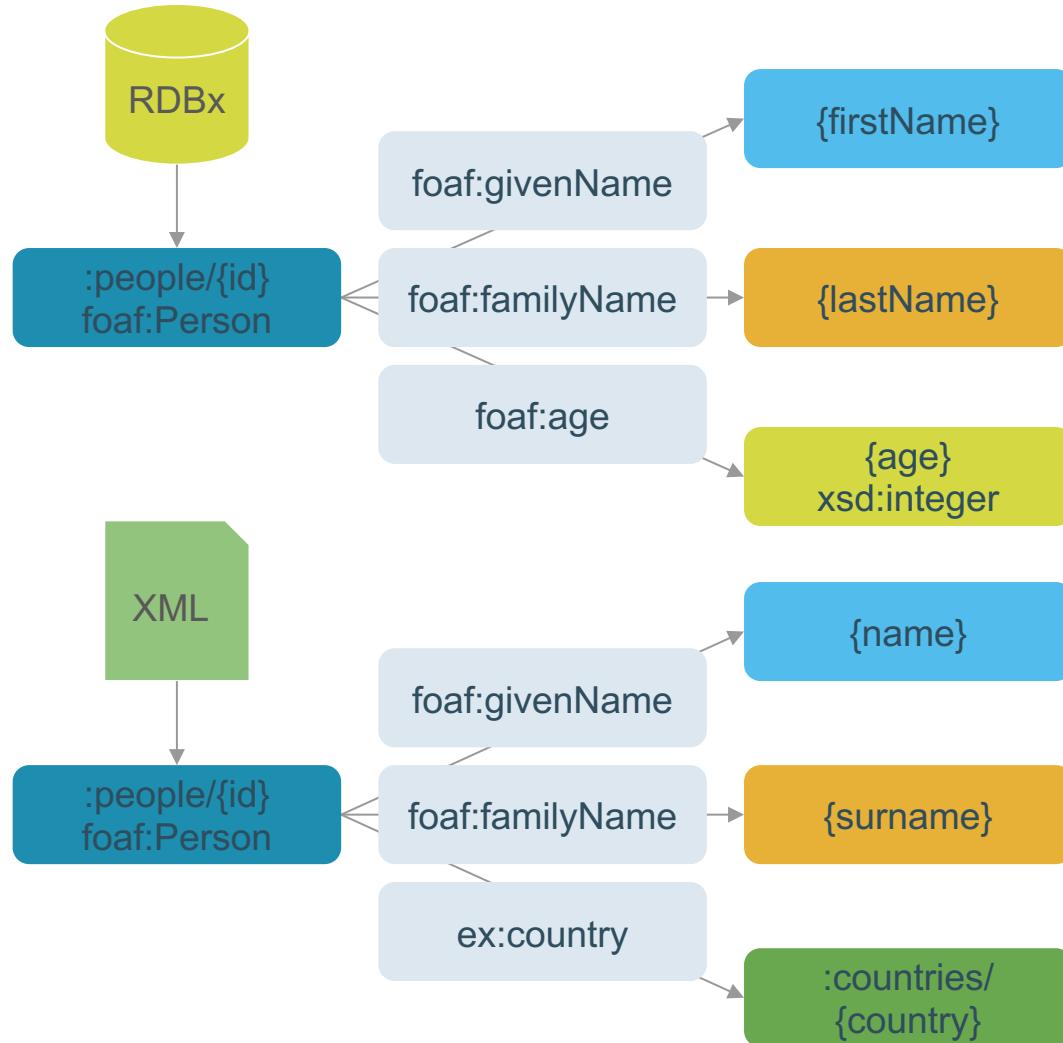


RML & SHACL alignment

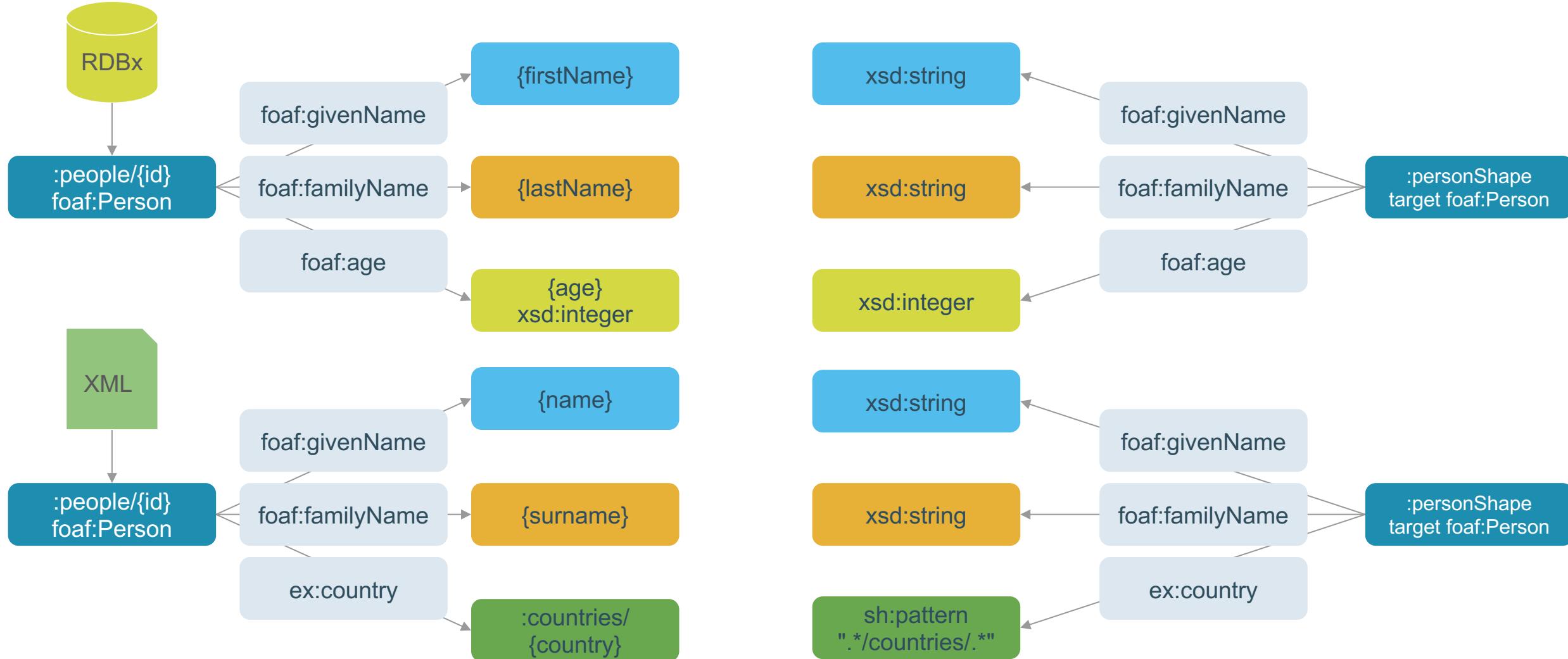
RML	SHACL
rr:subjectMap & rr:SubjectMap	sh:NodeShape
rr:predicateObjectMap	sh:property
rr:PredicateObjectMap	sh:PropertyShape
rr:class	sh:class, sh:targetClass
rr:predicate	sh:path
rr:referencingObjectMap	sh:node
rr:termType	sh:nodeKind
rr:datatype	sh:datatype
rr:language	sh:languageIn
rr:constant	:in
rr:template	sh:pattern

note: the alignment
is not formally
described yet

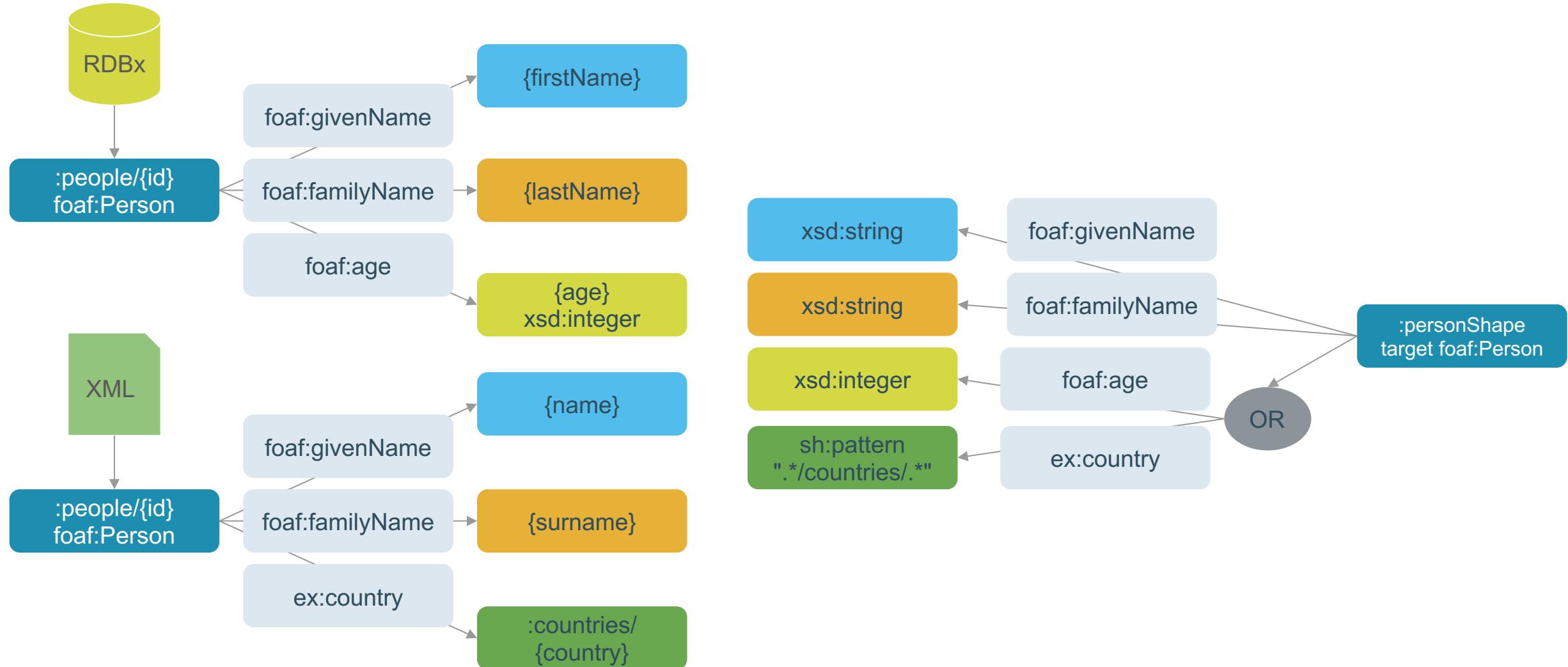
RML mapping rules

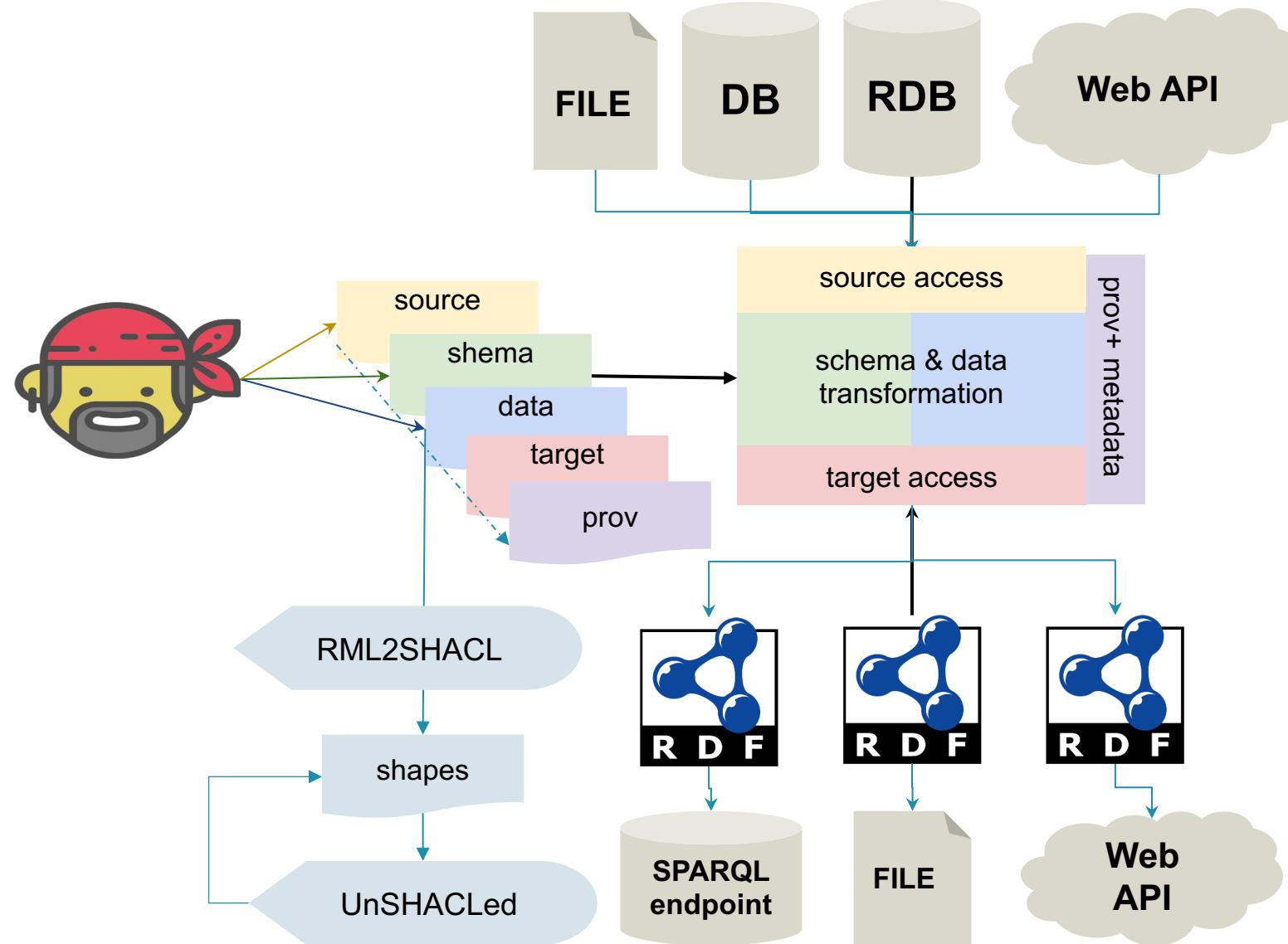


RML mapping rules & SHACL shapes

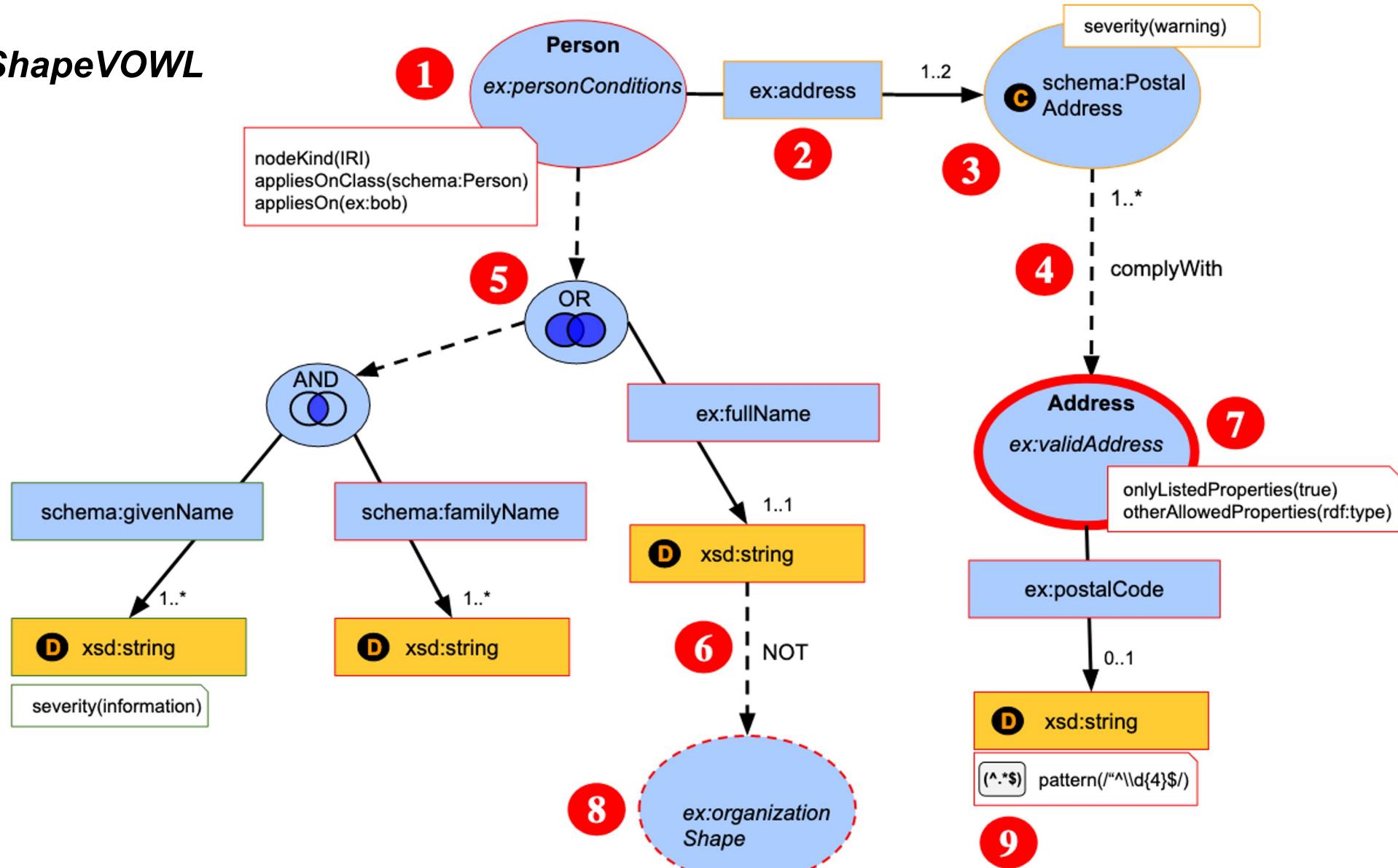


RML mapping rules → SHACL shapes

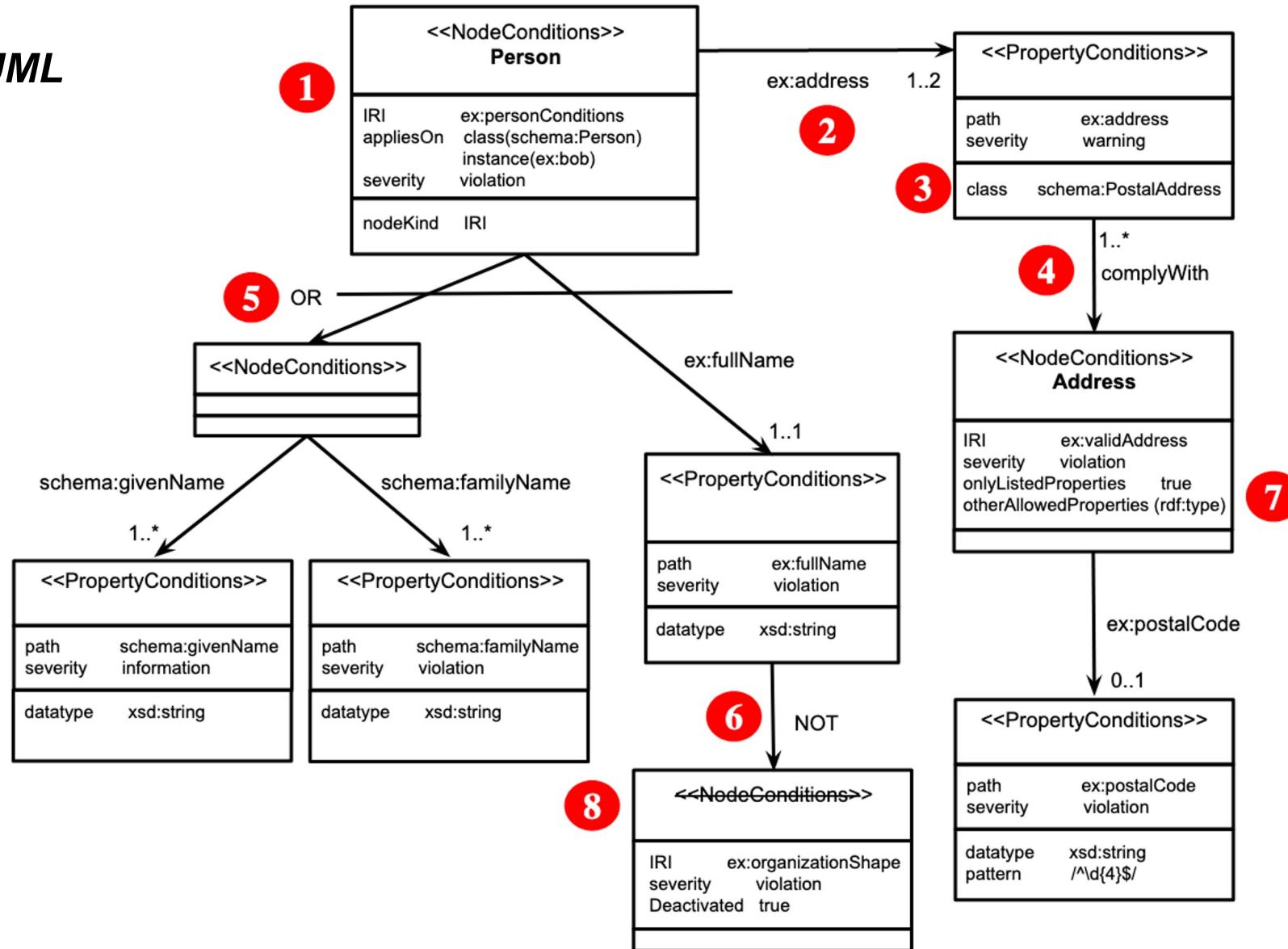




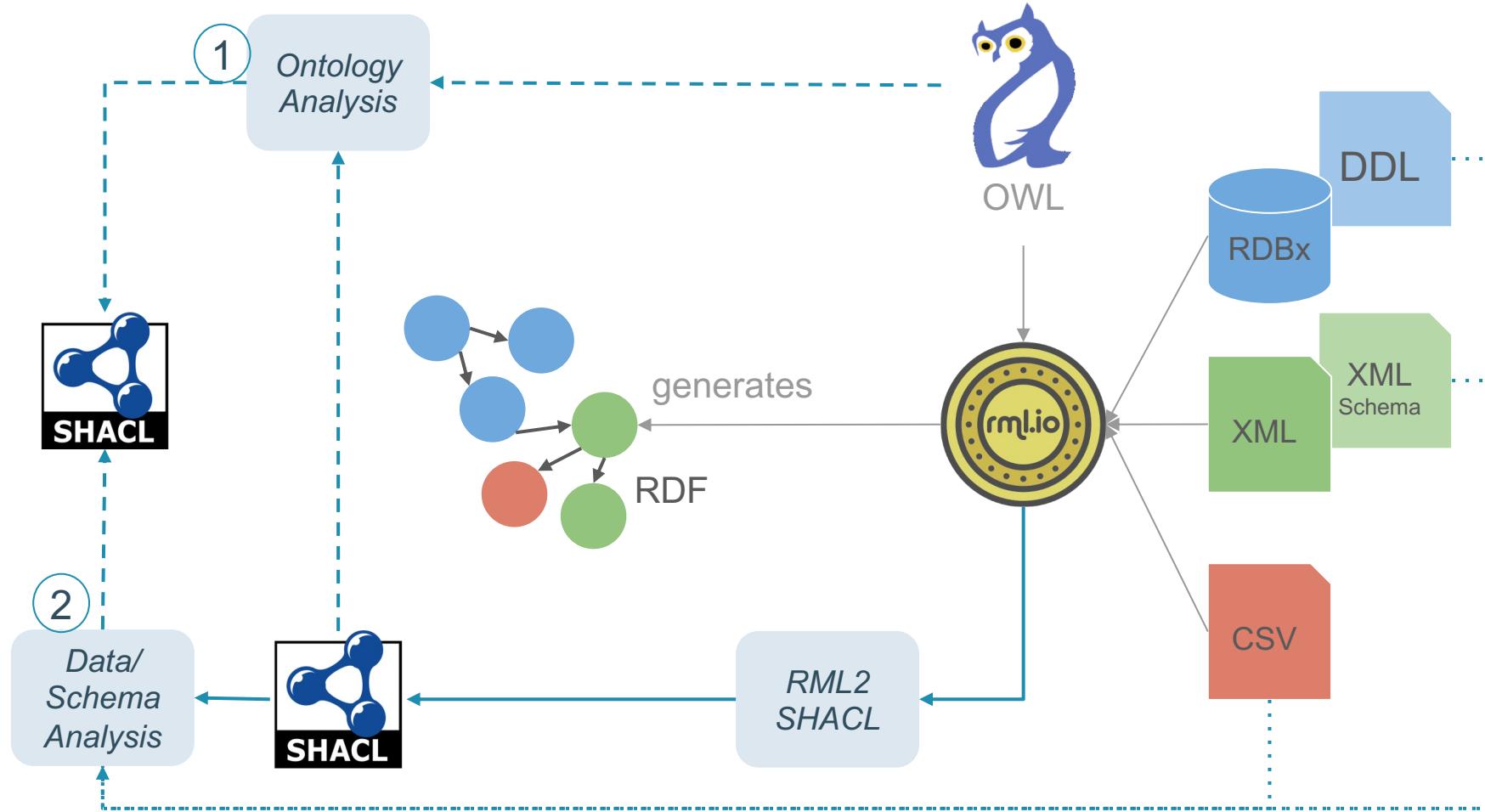
ShapeVOWL



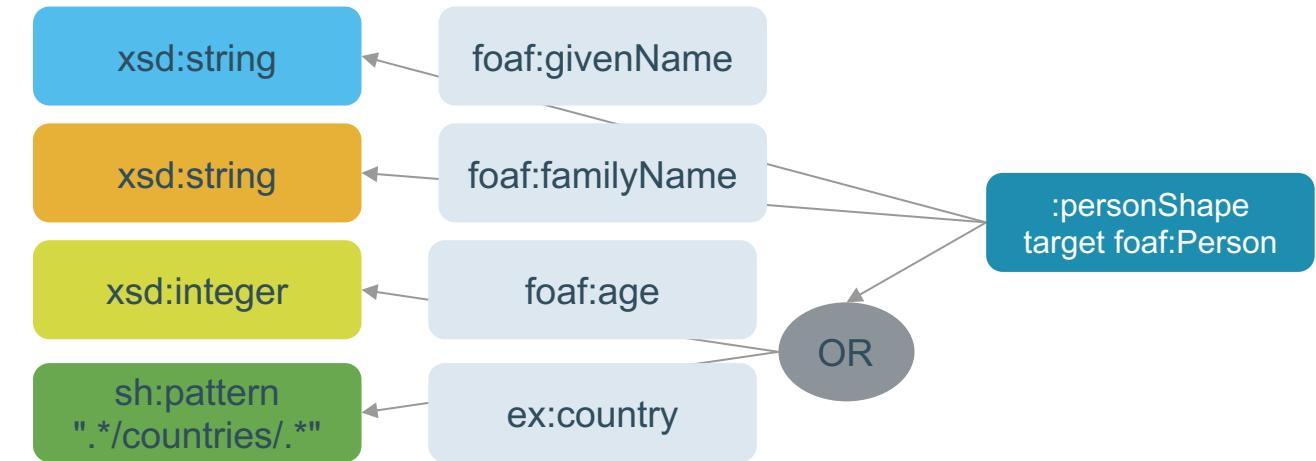
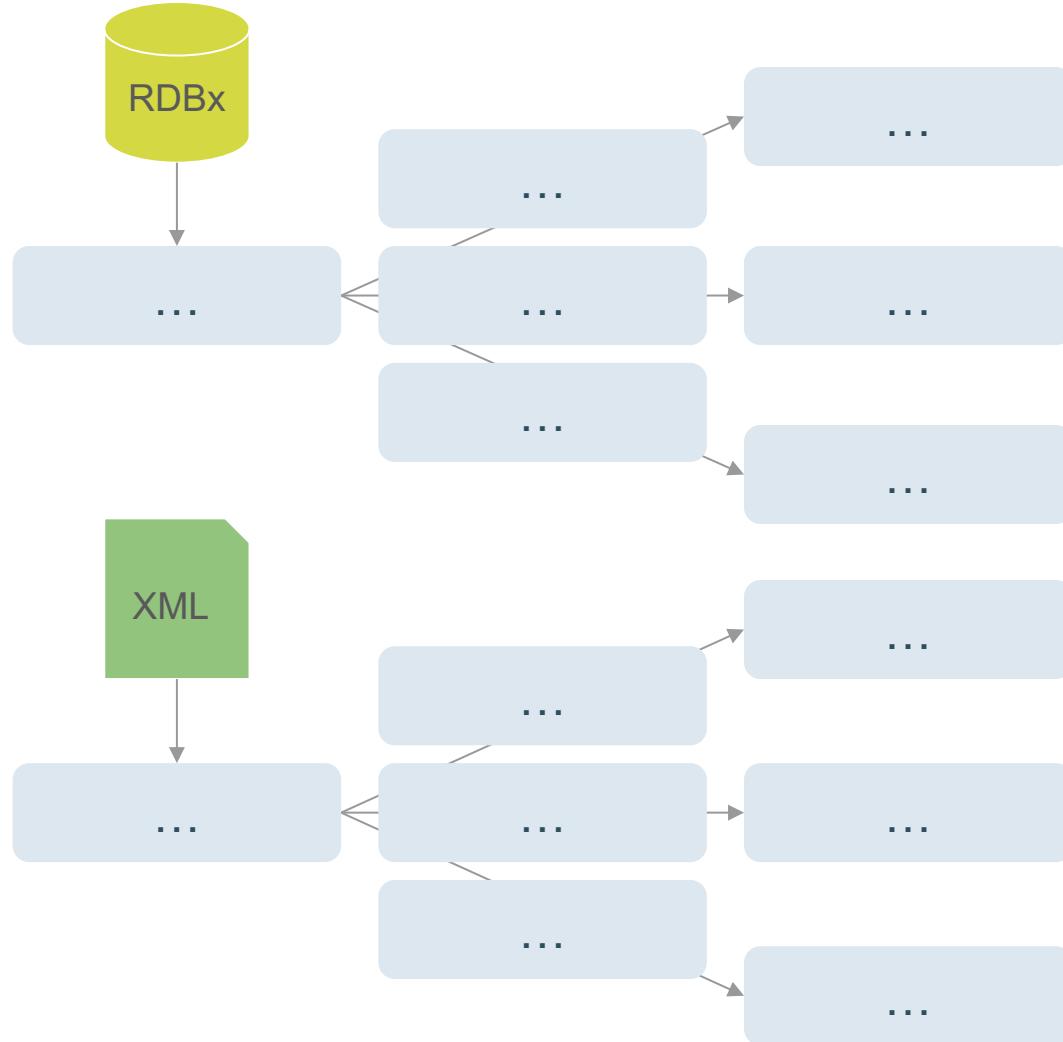
ShapeUML



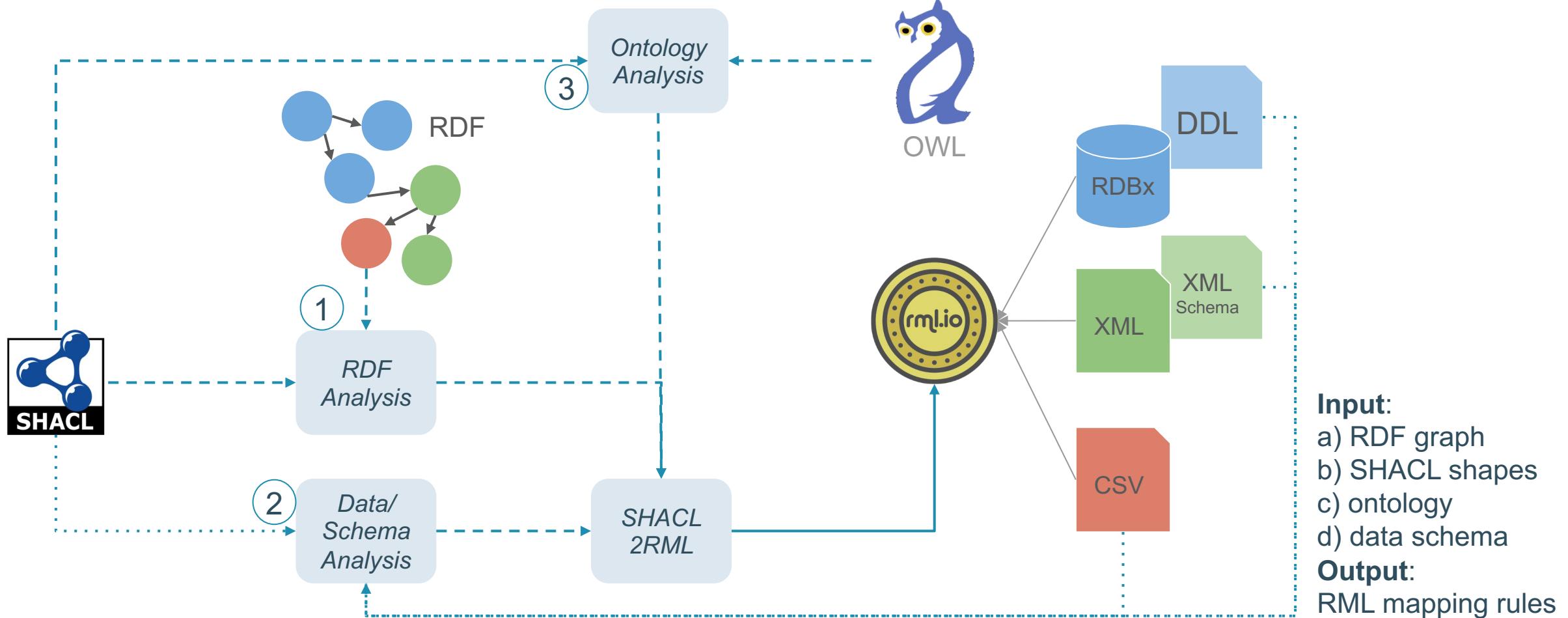
Enriching SHACL shapes



RML mapping rules ← SHACL shapes



Generating RML mapping rules from SHACL shapes



What did I do so far?

Mapping Languages & RML

Mapping rules & User support

Mapping rules & RDF graphs validation

Mapping rules & SHACL shapes

Implementations & benchmarks

RML processors

materialisation

DB2triples	(github.com/antidot/db2triples)
R2RML Parser	(github.com/nkons/r2rml-parser)
XSPARQL	(xsparql.sourceforge.net/)

*homogeneous
data sources*

Morph-RDB	(github.com/oeg-upm/morph-rdb)
Ontop	(github.com/ontop/ontop)

RMLMapper
CARMEL
RocketRML
Morph-xR2RML
RMLStreamer
Chimera
SDM-RDFizer
FunMap
MapSDI
Morph-KGC
GeoTriples

Java (github.com/RMLio/rmlmapper-java)
Java (github.com/carml/carml)
JavaScript (github.com/semanlifyit/RocketRML)
Scala (github.com/frmichel/morph-xr2rml)
Flink (github.com/RMLio/RMLStreamer)
Camel (github.com/cefriel/chimera)
heuristic-based planning (github.com/SDM-TIB/SDM-RDFizer)
function-free planning (github.com/SDM-TIB/FunMap)
deduplication-based optimizations (github.com/SDM-TIB/MapSDI)
mapping partitions (github.com/oeg-upm/morph-kgc)
geospatial data (github.com/LinkedEOData/GeoTriples)

*heterogeneous
data sources*

Morph-xR2RML	(github.com/frmichel/morph-xr2rml)
Squerall	(github.com/EIS-Bonn/Squerall)
Ontario	(github.com/SDM-TIB/Ontario/)

virtualisation

Choose your implementation

Choose yourself the best tool for your needs!

<http://rml.io/test-cases/>

<http://rml.io/implementation-report/>

Benchmarks:

GTFS: evaluate tools generating RDF graphs with RML mapping rules

(<https://github.com/oeg-upm/gtfs-bench>)

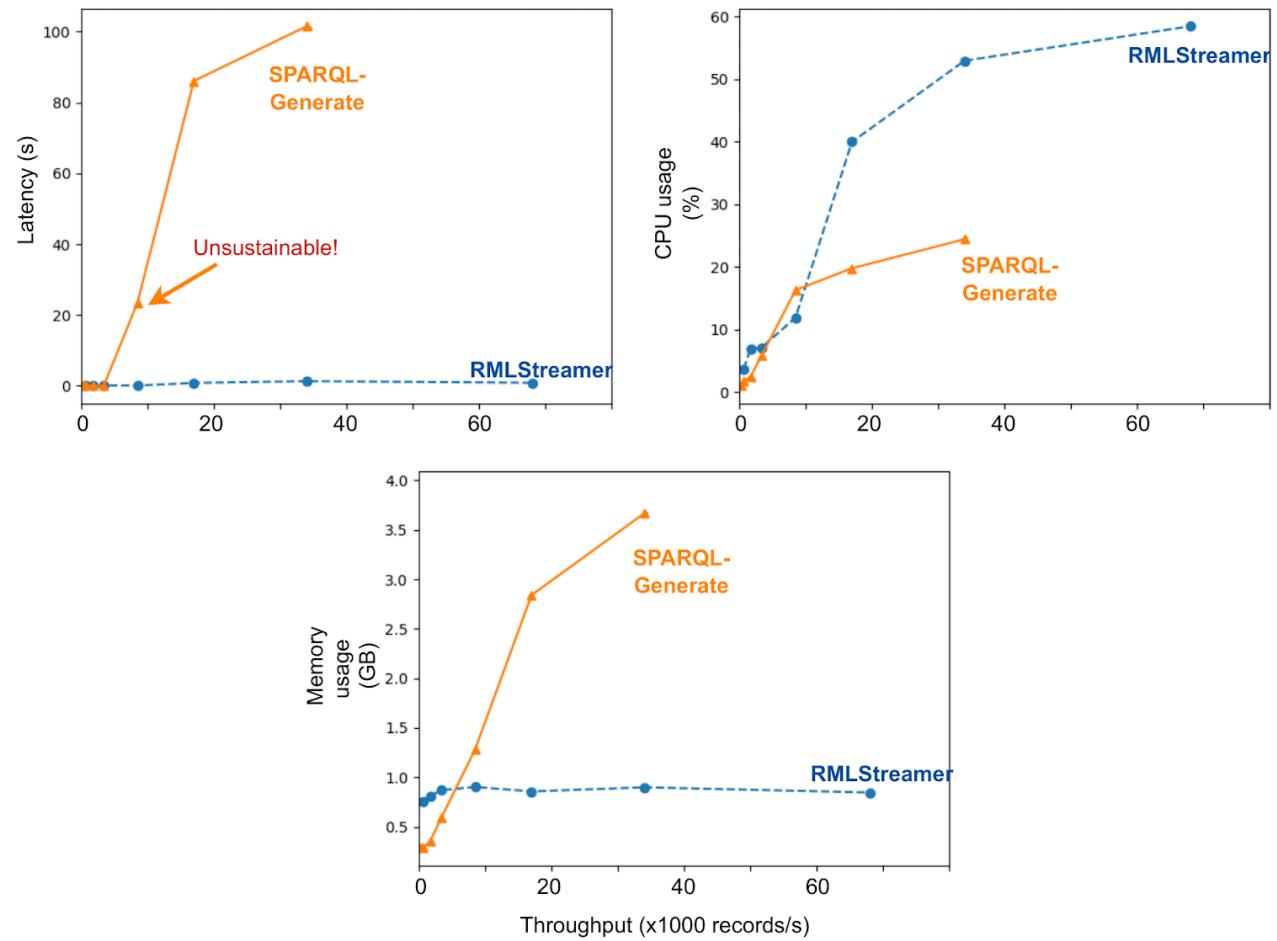
RODI:

test the quality of (semi-)automatically generated mapping rules

(<https://github.com/chrpin/rodi>)

Test Case	RMLMapper	CARML	RocketRML	SDM-RDFizer	RMLStreamer	Chimera	Morph-KGC
RMLTC0000-CSV	passed	passed	passed	passed	passed	passed	passed
RMLTC0000-JSON	passed	passed	passed	passed	passed	passed	failed
RMLTC0000-MYSQL	passed	inapplicable	inapplicable	passed	inapplicable	inapplicable	passed
RMLTC0000-POSTGRESQL	passed	inapplicable	inapplicable	passed	inapplicable	inapplicable	passed
RMLTC0000-SPARQL	passed	inapplicable	inapplicable	inapplicable	inapplicable	inapplicable	inapplicable
RMLTC0000-SQLSERVER	passed	inapplicable	inapplicable	passed	inapplicable	inapplicable	inapplicable

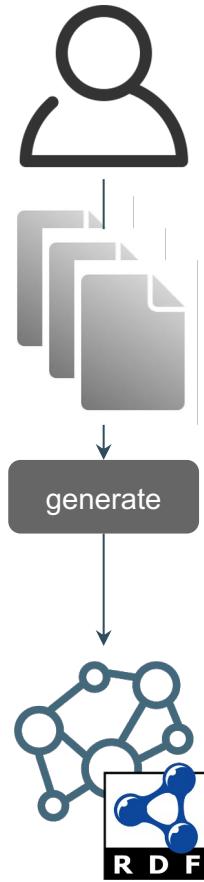
RMLStreamer-SISO



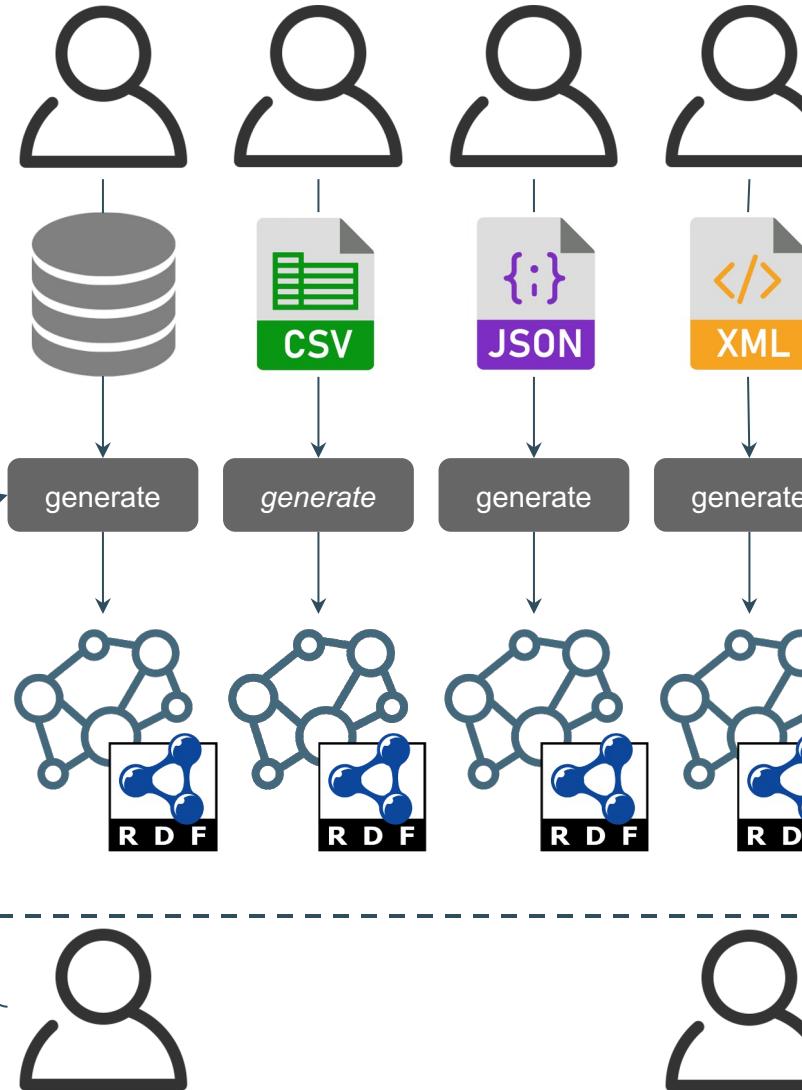
shameless advertisement:

Tuesday, 25 October 15:10-16:10

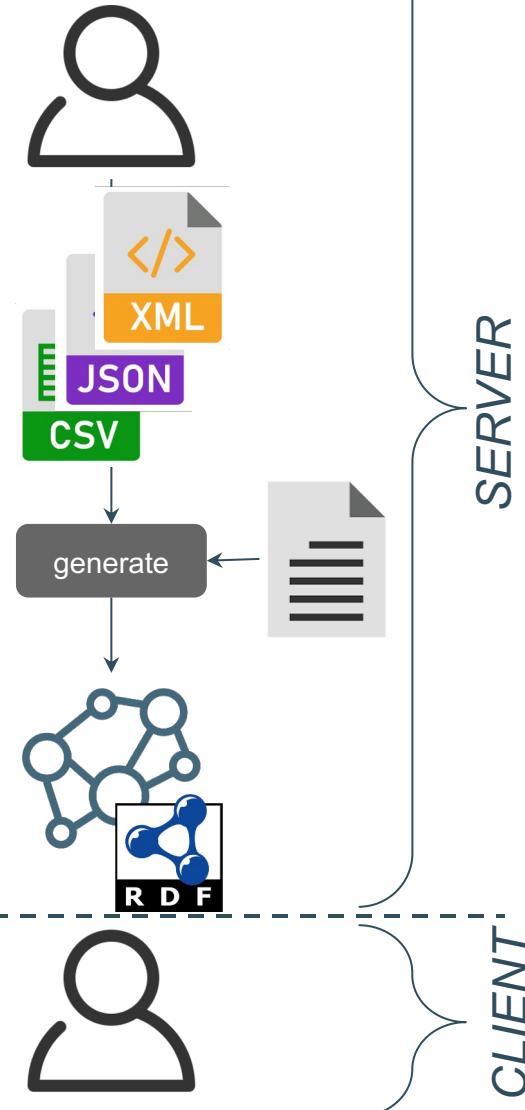
hard-coded



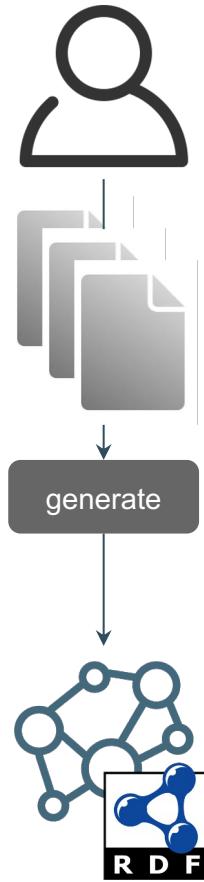
format-specific



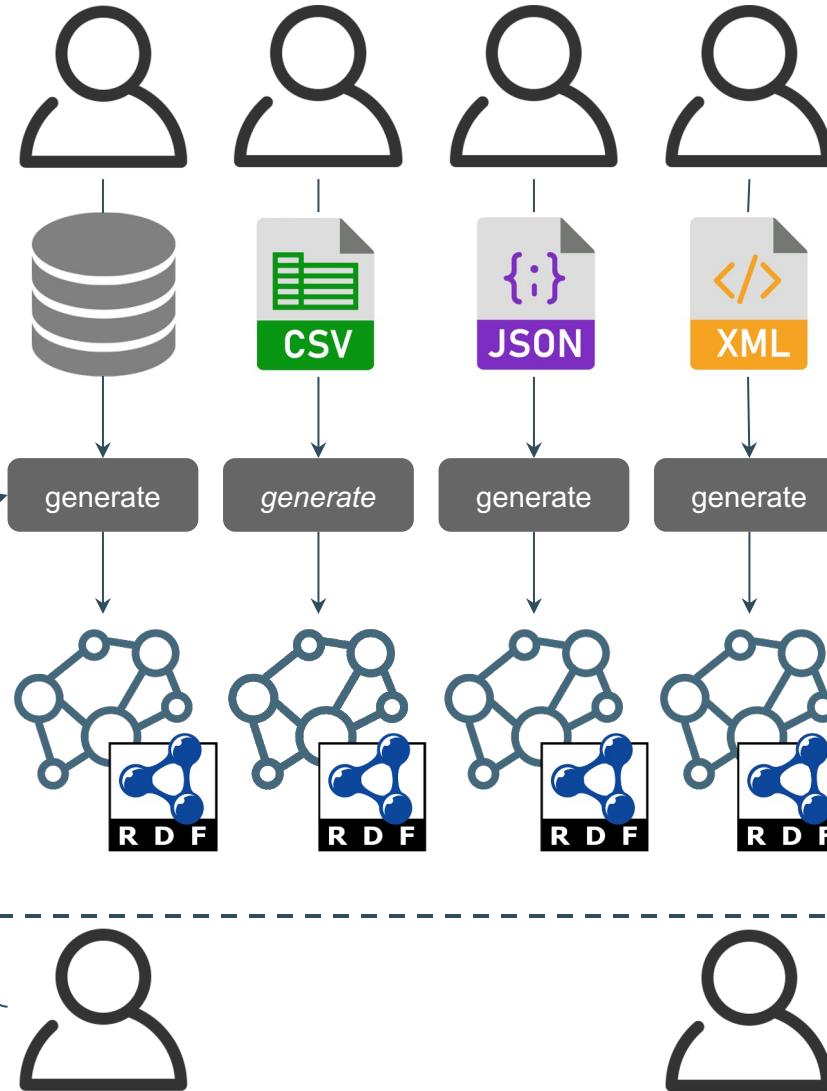
declarative



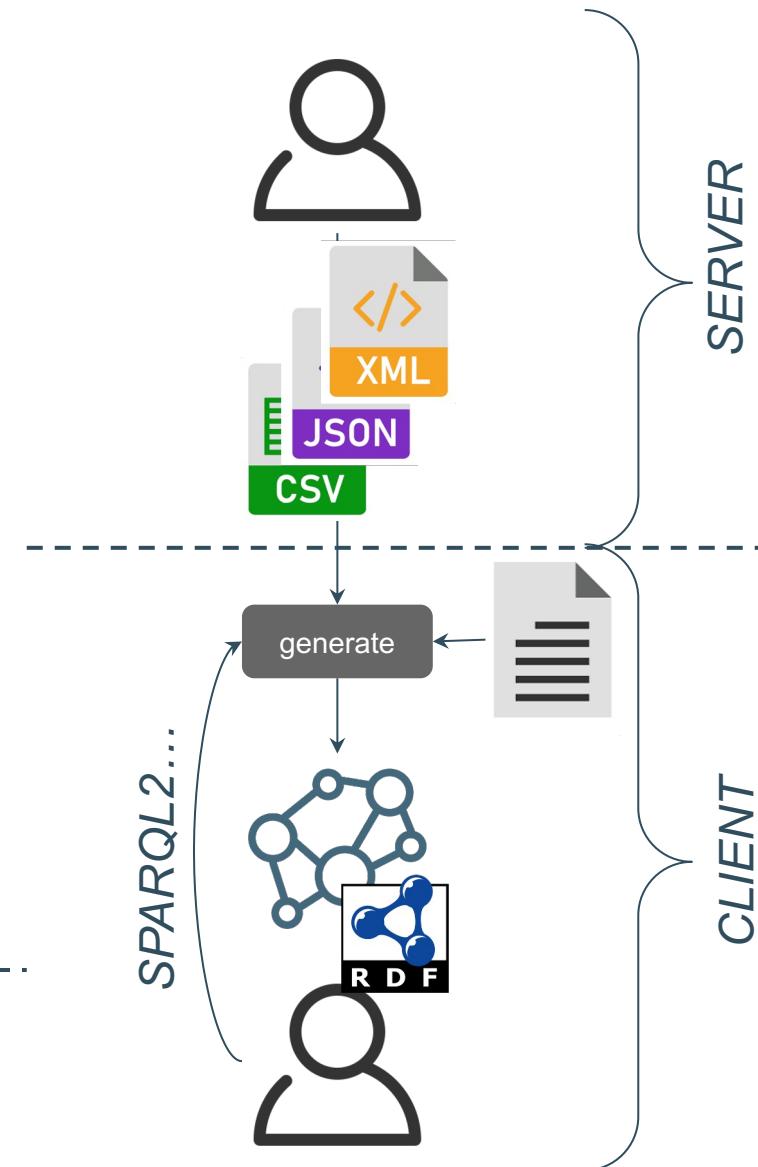
hard-coded

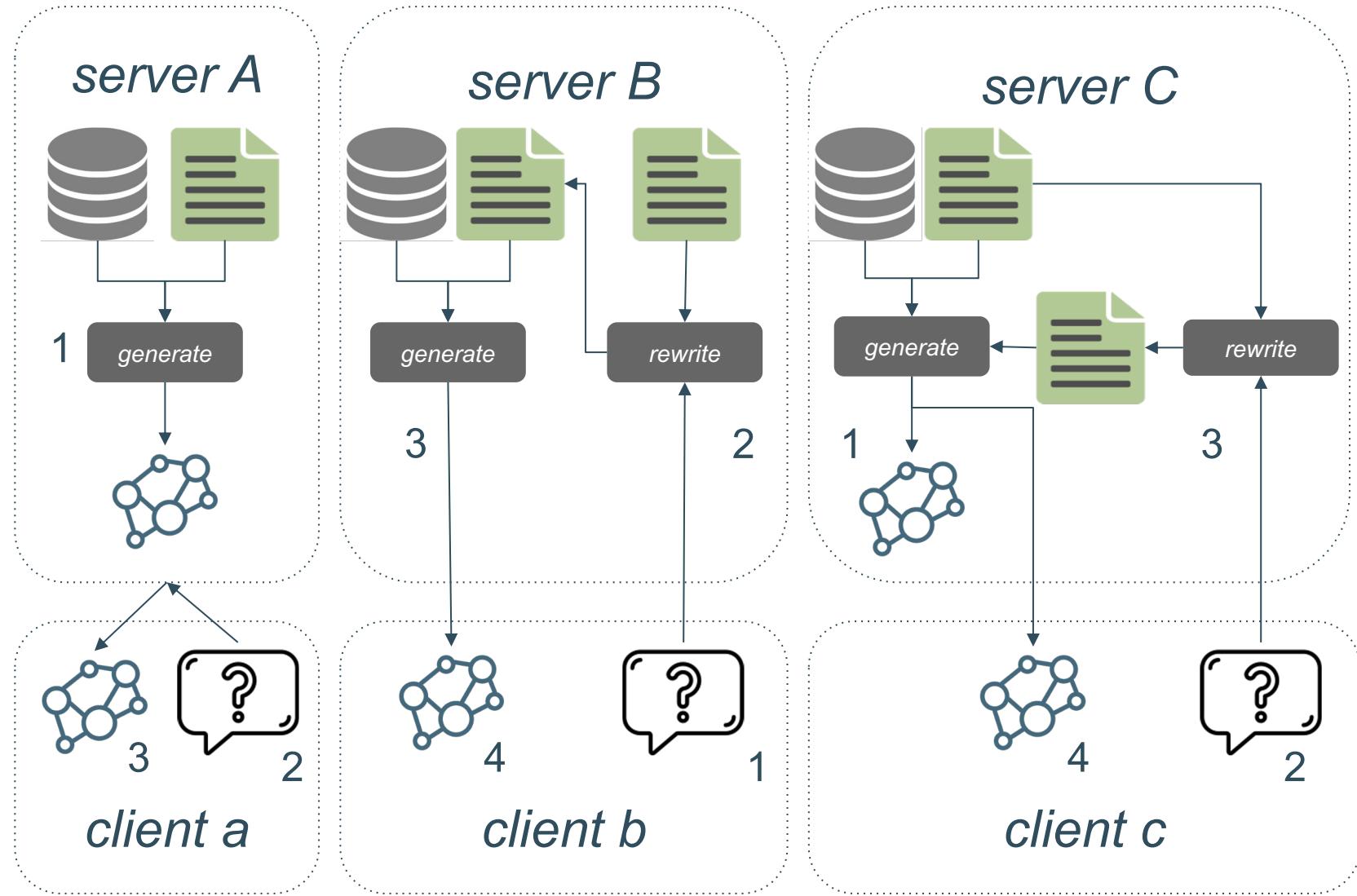


format-specific



declarative

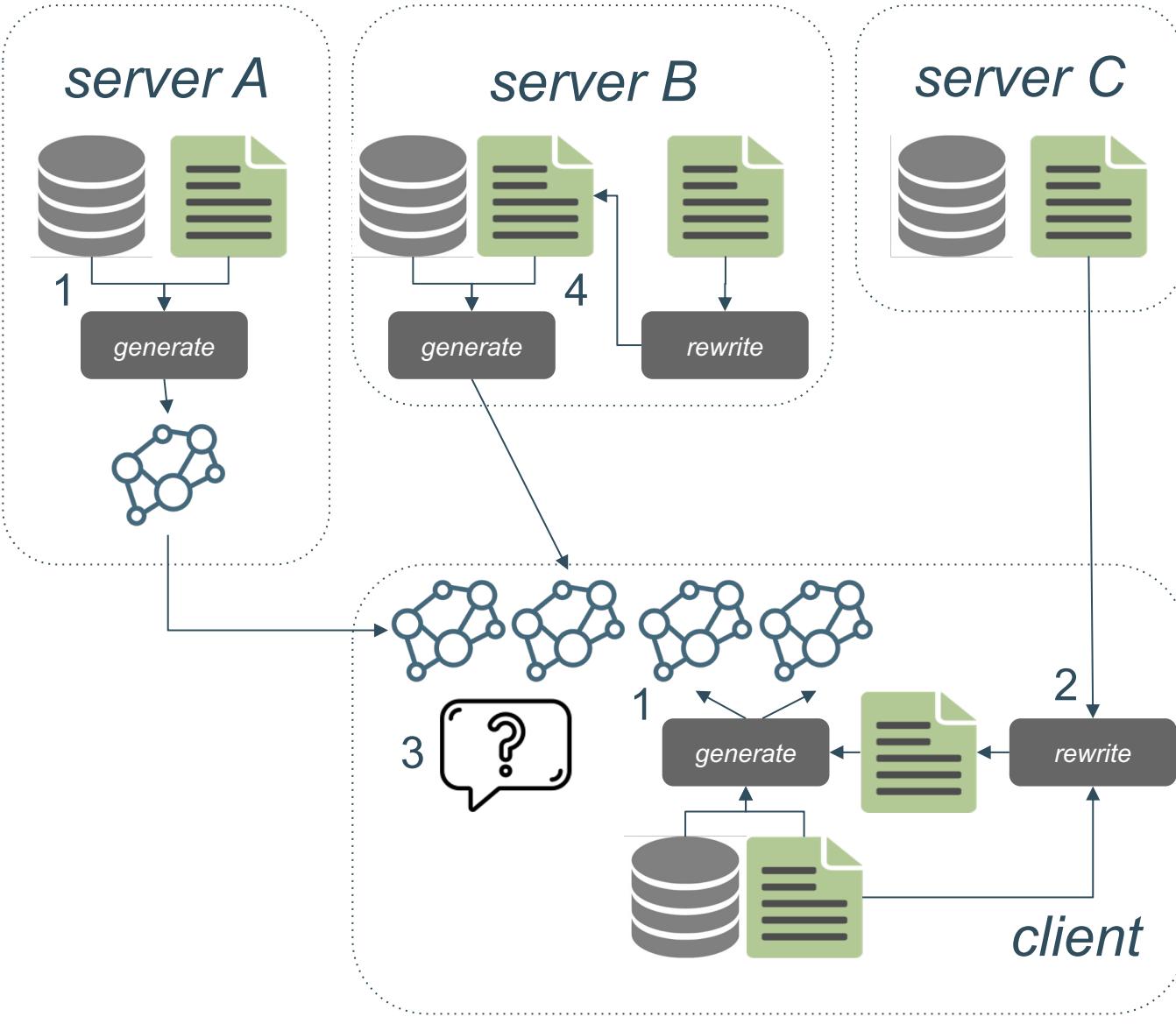




*knowledge graph
materialisation*

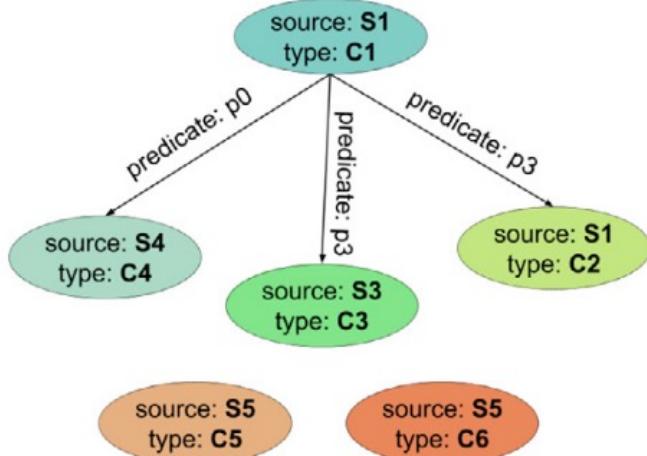
*virtual knowledge graphs generation
over heterogeneous data*

*hybrid materialisation & virtualisation
for knowledge graph generation*

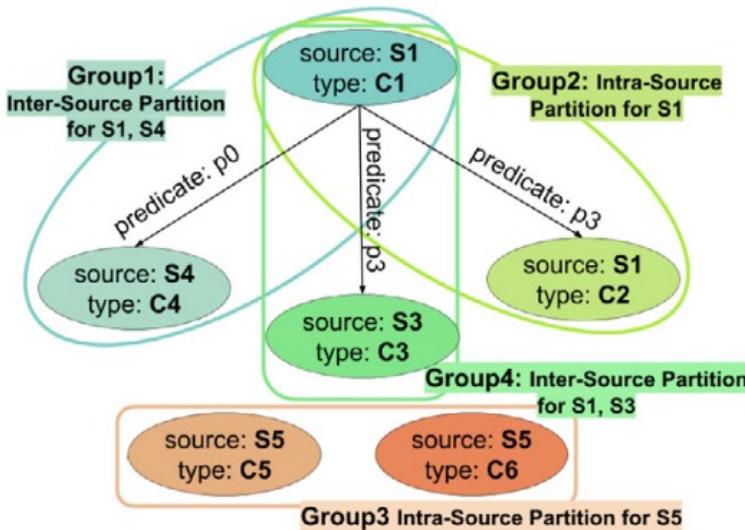


server A materializes its knowledge graph
 server B generates a virtual knowledge graph
 server C does not generate a knowledge graph
 client materialises its own knowledge graph
 & the knowledge graph for server C's data
 receives a virtual knowledge graph from server B
 receives a subgraph from server A
 combines all subgraphs to perform its tasks

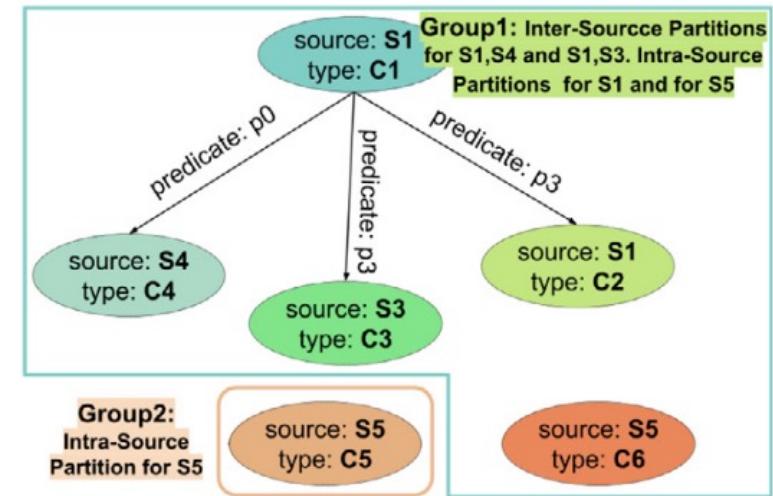
Mapping rules partition



No Partitioning
RMLMapper: Five-hour timeout (No results)
RocketRML: Out of memory (No results)
SDM-RDFizer: 441.18 sec (100% results)
Morph-KGC: 42.32 sec (100% results)



Optimized Partitioning
RMLMapper: Five-hour timeout (92.65% results)
RocketRML: 133.82 sec (100% results)
SDM-RDFizer: 211.02 sec (100% results)
Morph-KGC: 36.81 sec (100% results)



Random Partitioning
RMLMapper: Five-hour timeout (5.41% results)
RocketRML: Out of memory (5.41% results)
SDM-RDFizer: 369.27 sec (100% results)
Morph-KGC: 40.21 sec (100% results)

What did I do so far?

Mapping Languages & RML

Mapping rules & User support

Mapping rules & RDF graphs validation

Mapping rules & SHACL shapes

Implementations & benchmarks

What do we still miss?

Many RML implementations by now – focus on **materialisation implementations**
still **missing virtualisation & hybrid** (materialisation+virtualisation) solutions

Many RML **implementations for static data**
still **missing solutions for dynamic/streaming & versioned data**

Community **converges on mapping language** for heterogeneous semi-structured data
still **missing automated solutions** for mapping rules definition

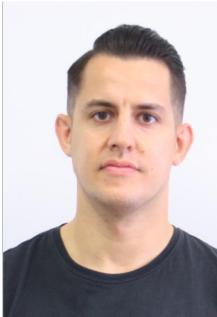
Only very **few GUIs** to support the humans
still **missing thorough research on methodologies for knowledge graphs construction & human-in-the-loop**

Preliminary benchmarks
still missing more fine-grained **benchmarking for all aspects of knowledge graph construction**

Preliminary quality assessment solutions
still **missing systematic solutions for quality assessment** & feedback loop

Thank you! former team @UGhent/imec

presentation based on works with
my former team at Ghent University
and my current team at KULeuven
and
many references to different members
of the Knowledge Graph Construction
Community Group



Sitt Min Oo

Gerald
Haesendonck

Gertjan
De Mulder

Thomas
Delva



Ioannis
Chrysakis

Dylan
Van Assche



Thank you! team @KULeuven



prof Anastasia Dimou Dr David Chaves Fraga
anastasia.dimou@kuleuven.be postdoctoral fellow
@natadimou (KU Leuven & UPM)
cleansing & benchmarking KGs



Ioannis Chrysakis joint PhD student (UGhent & KU Leuven)
KGs for dataspaces interoperability & privacy



Dylan Van Assche PhD student (UGhent)
trade offs of KGs materialisation & virtualisation



Ioannis Dasoulas PhD student (KU Leuven)
(June 2022 - now)
fusion of KGs & ML



Xuemin Duan PhD student (KU Leuven)
(Sep 2022 - now)
automated KG generation



(to start Dec 2022)
KGs for manufacturing



(to start Jan 2023)
KG & Decision Making



Knowledge Graphs Construction

W3C Community Group

w3.org/community/kg-construct/
github.com/kg-construct/
>150 members world-wide

new RML specifications
(work in progress)

<https://kg-construct.github.io/rml-core/>
<https://kg-construct.github.io/rml-target-source-spec/>
<https://w3id.org/kg-construct/collections-containers/>
<https://kg-construct.github.io/rml-star-spec/>



Knowledge Graph Construction Workshop

<https://kg-construct.github.io/workshop/>

1st edition theme: mapping languages

2nd edition theme: human-in-the-loop

3rd edition: automated KG construction

4th edition (2023): benchmarking implementations?

networks on Knowledge Graphs

graph theory
data management
knowledge representation
programming languages



Knowledge Graphs for Data Integration (KG4DI)
w3id.org/kg4di/
Scientific Research Network
kick off: November 30th

Distributed Knowledge Graph COST Action (DKG)

DKG – cost-dkg.eu

>100 members from >30 countries in Europe

Are Knowledge Graphs Ready for the Real World? Challenges and Perspective

dagstuhl.de/24061

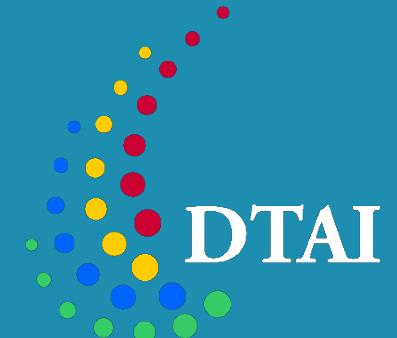
February 2024

*The adoption of Semantic Web technologies
depends on the availability of knowledge graphs*

*It's on our hands to investigate how to
efficiently generate high quality knowledge graphs*



Aligning heterogeneous semi-structured data and knowledge graphs



Anastasia Dimou

 anastasia.dimou@kuleuven.be

 [@natadimou](https://twitter.com/natadimou)