



Virtual Knowledge Graph Generation from Heterogeneous Data Sources

**David Chaves-Fraga, Ontology Engineering Group
Universidad Politécnica de Madrid, Spain**

Freddy Priyatna, OEG-UPM

Oscar Corcho, OEG-UPM

`schema:email = lower(substr({name},1,1) || {surname} || '@fi.upm.es')`

✉ dchaves@fi.upm.es

🐦 [@dchavesf](https://twitter.com/dchavesf)

📅 24-25/10/2018

📍 [@imec/Ghent University](https://imec.ugent.be/)

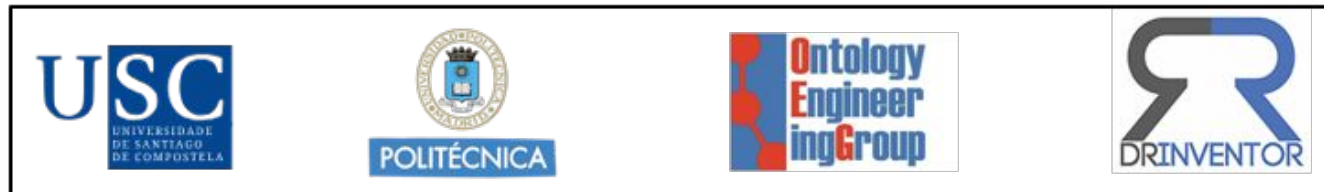
PhD Student and Researcher at OEG-UPM since 2016 (Data Integration team):

- MSc Thesis (2016): Methods and Techniques for the Evaluation of Ontology Learning
- PhD Thesis (2016-2020): Virtual Knowledge Graph Generation from heterogeneous resources

Interests:

- **OBDA**
- **Heterogeneous data**
- SPARQL
- Federated queries
- **Data Integration**
- Public Transport
- Linked Connections
- **R2RML- RML**
- **Virtualization - Access**

2011-2016



2017



2018



@dchavesf



dchaves@fi.upm.es

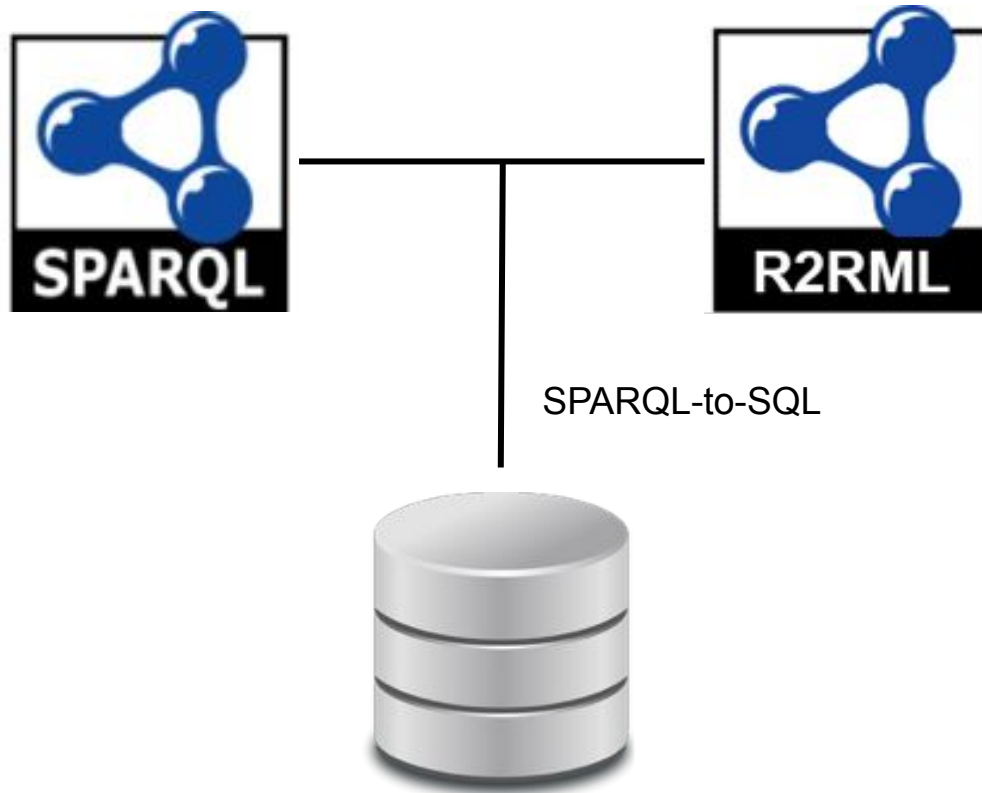


dchaves.oeg-upm.net



dachafra

OBDA... **O**ntology **B**ased **D**ata **A**ccess



Focused on optimizing the generated SQL query to improve the performance

But we are working on... Semantic **WEB**

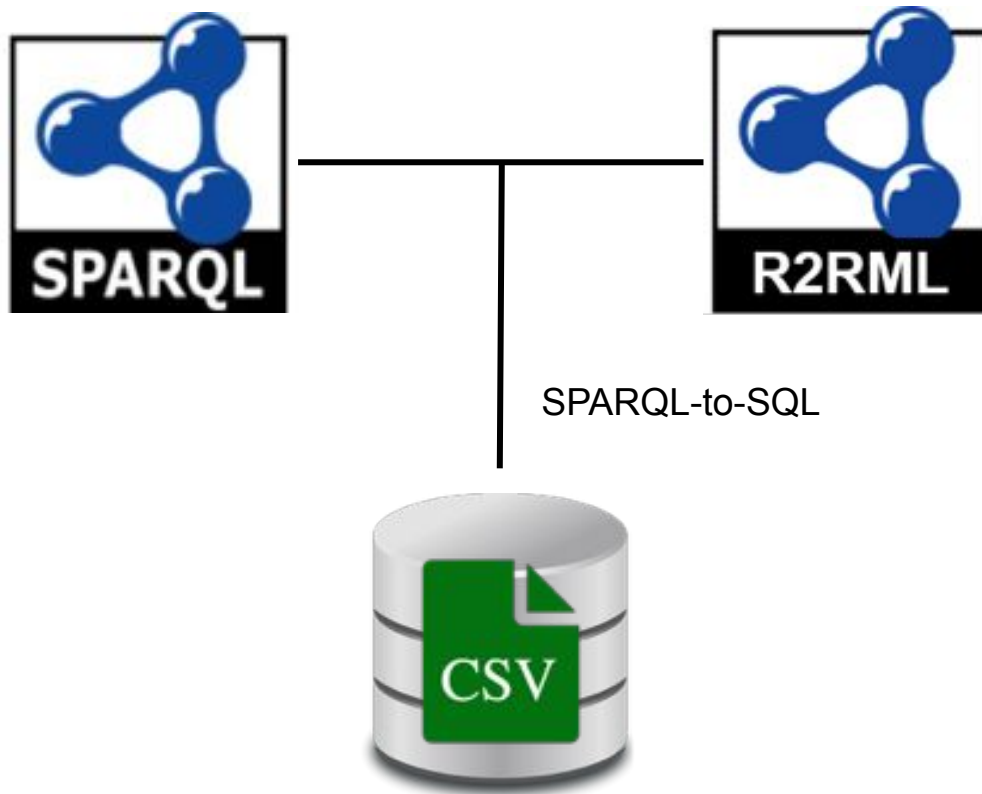
How is the data exposed on the Web?

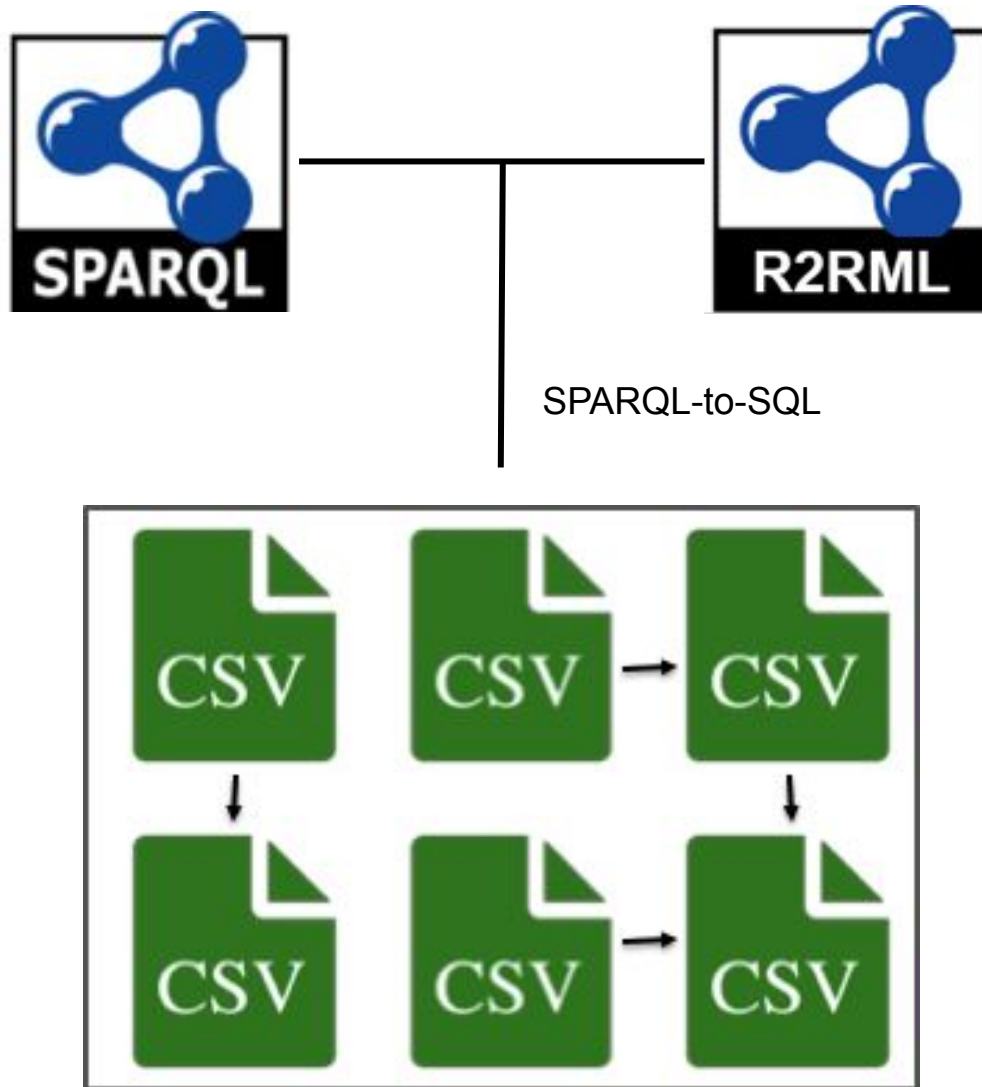
Formats
CSV (114629)
TXT (80014)
JSON (50676)
ZIP (50070)
HTML (45706)
GMZ (44712)
PDF (34770)
XLS (26356)
SHP (19778)
XML (19311)

Formato
CSV (10581)
XLS (7474)
JSON (7234)
HTML (6245)
PDF (3909)
XML-APP (2721)
XLSX (2649)
PC-Axis (2490)
XML (1951)
ASCII (1909)
JPG (1774)
KMZ (1504)
ZIP (1309)

Formatos	—
CSV (367)	
XML (130)	
XLS (128)	
XLSX (88)	
WMS (29)	
RDF (21)	
GeoJSON (7)	
JSON (7)	
prj (7)	
SHP (7)	
SHX (7)	

OBDA... Ontology **B**ased **D**ata **A**ccess



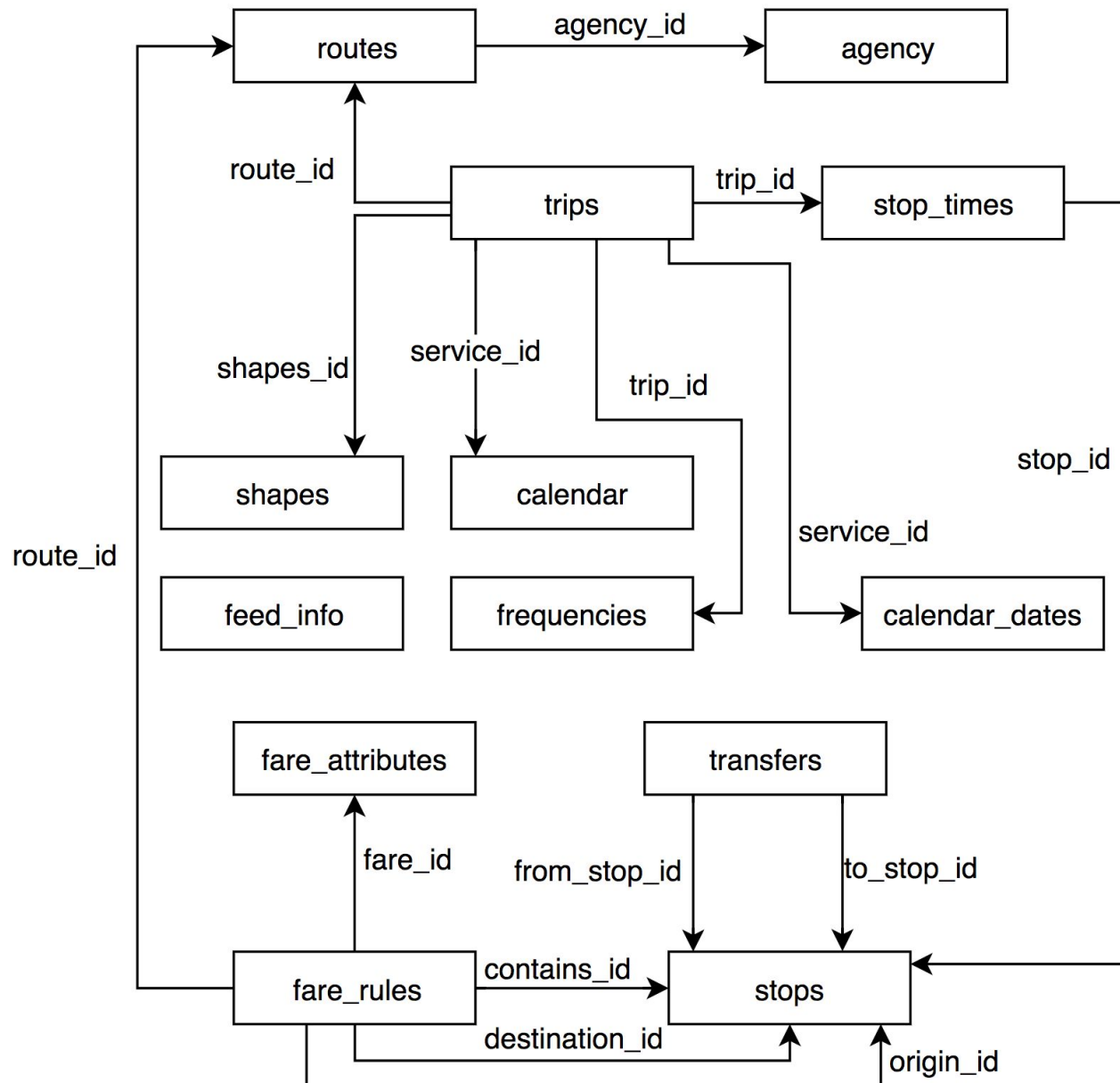
OBDA... **O**ntology **B**ased **D**ata **A**ccess

Multiple CSV files with relations among them:

1. Joins are not explicit
2. Constraints are not defined explicitly in the CSVs (PK, FKs)
3. The data may not be in the desirable format (e.g. dates)
4. CSVs are not in 3NF:
 - a. PK may be repeated
 - b. FKs may not be explicated
 - c. FKs could not have a 1:1 cardinality
 - d. Lists in column

R2RML is not enough for dealing with CSV(s) in an
OBDA approach

Let's give an example...



LD Generation from GTFS to LinkedGTFS (in hours)

Dataset (size mg)	Morph-R2RML	RML-Mapper
D1 (2.3)	0.004	3.739
D2 (2.6)	0.026	2.587
D3 (2.9)	0.068	0.778
D4 (3.4)	0.118	7.026
D5 (4.2)	0.115	7.026
D6 (4.7)	0.217	12.218
D7 (31)	1.153	151.541
D8 (96)	12.496	>160

Our mission as researchers it to provide solutions for:

- Generate Linked Data when:
 - o Quality is important
 - o The data is static

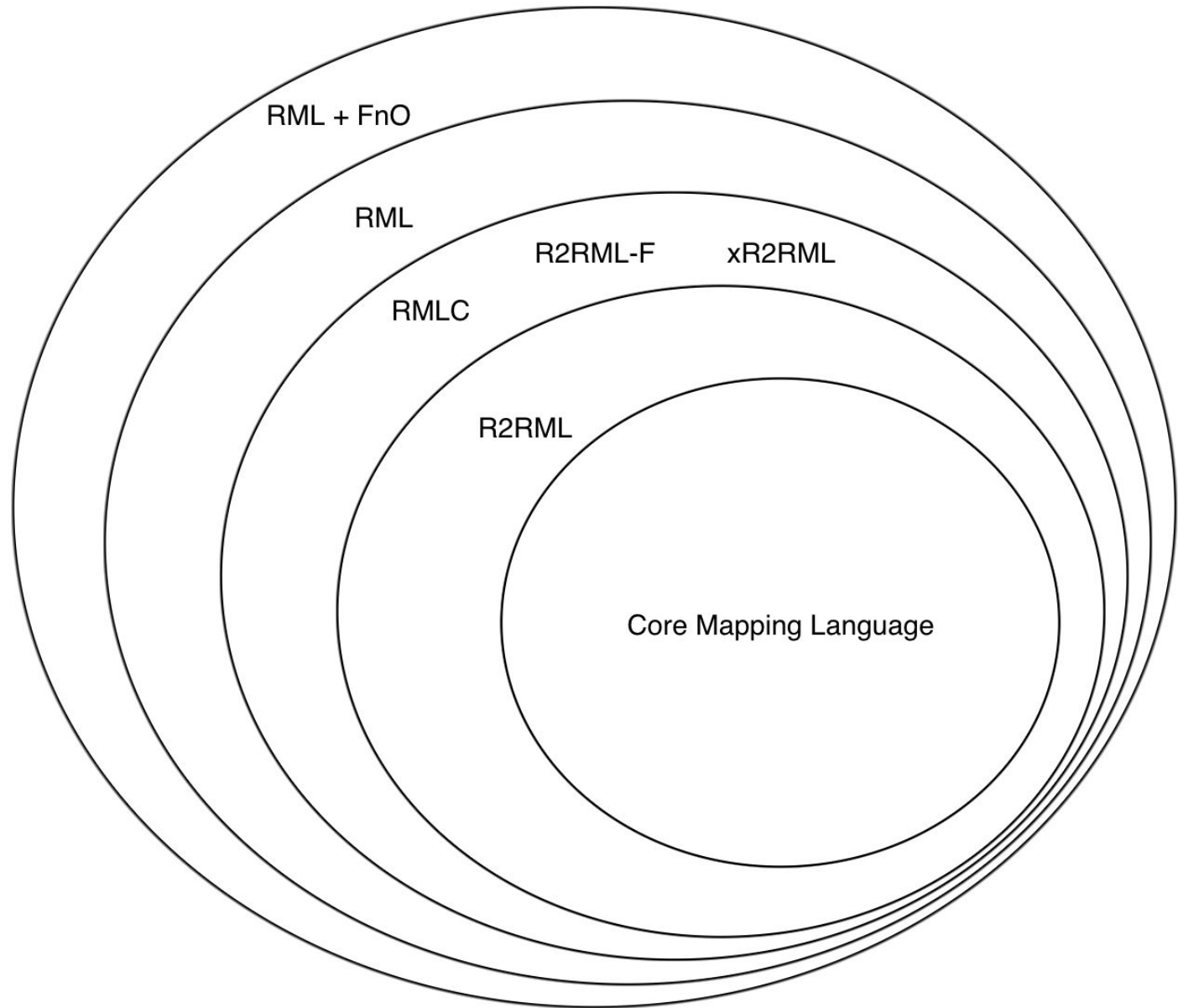
- Access to data using a graph query language when:
 - o Data is volatile
 - o Performance is relevant
 - o Underlying query engine for translation exists

Common point: The mapping language!

How to provide **access**/generation to heterogeneous data exposed on the web with relations among them using semantic technologies?



How can we extend standard Mapping Languages maintaining their semantics for using **OBDA engines** or LD generators?



Feature	R2RML	RML	RML-C
Data format	RBD	JSON,CSV, XML	CSV
Materialization	Yes	Yes	Yes
Virtualization	Yes	No	Yes
Functions	No	Yes (FnO)	Yes (SQL Functions)
Specification	Yes	Partially? (FnO+RML?)	Partially

“Virtual Statistics Knowledge Graph Generation from CSV files” D. Chaves-Fraga, F. Priyatna, I. Santana-Perez and O.Corcho at *SemStats Workshop co-located with ISWC18* (Best Paper Award)

“SATET: Providing access to multiple CSV on the Web using OBDA” D. Chaves-Fraga and O.Corcho (on-going work)



Virtual Statistics Knowledge Graph Generation from CSV files

**David Chaves-Fraga, Ontology Engineering Group
Universidad Politécnica de Madrid, Spain**

Freddy Priyatna, OEG-UPM

Idafen Perez-Santana, OEG-UPM

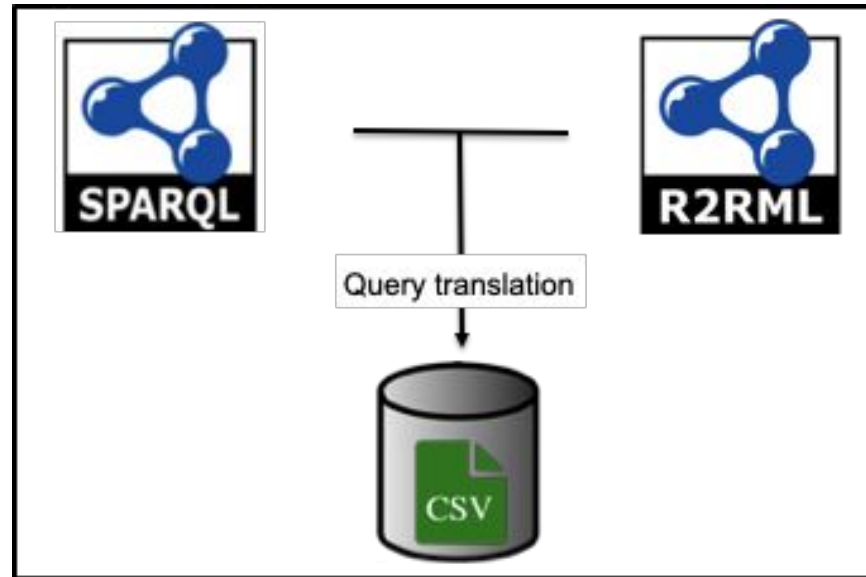
Oscar Corcho, OEG-UPM

✉ dchaves@fi.upm.es

🐦 [@dchavesf](https://twitter.com/dchavesf)

📅 08/10/2018

📍 ISWC18-SemStats

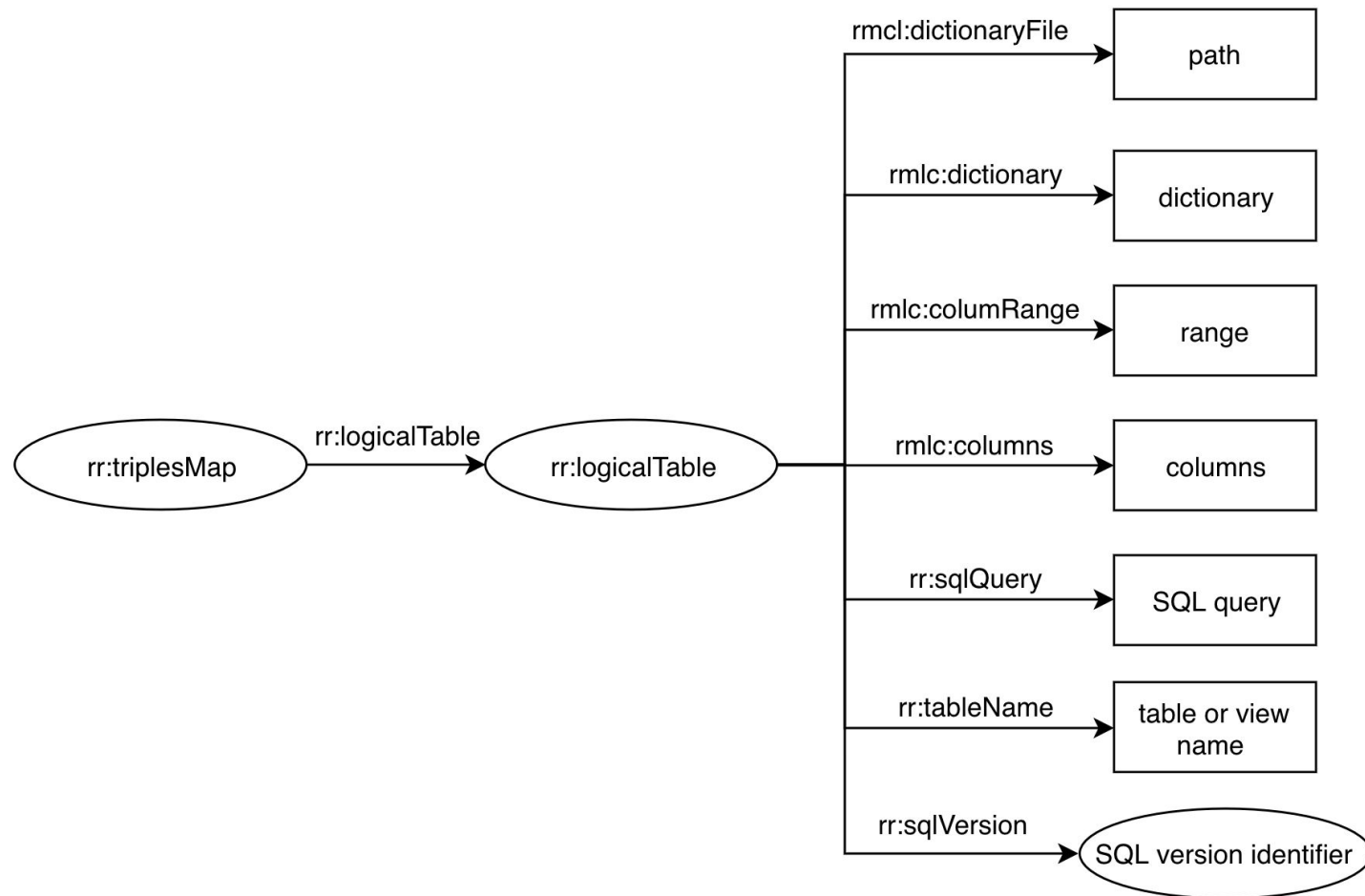


The **size** of the R2RML mapping depends on the **number of columns** in the CSV

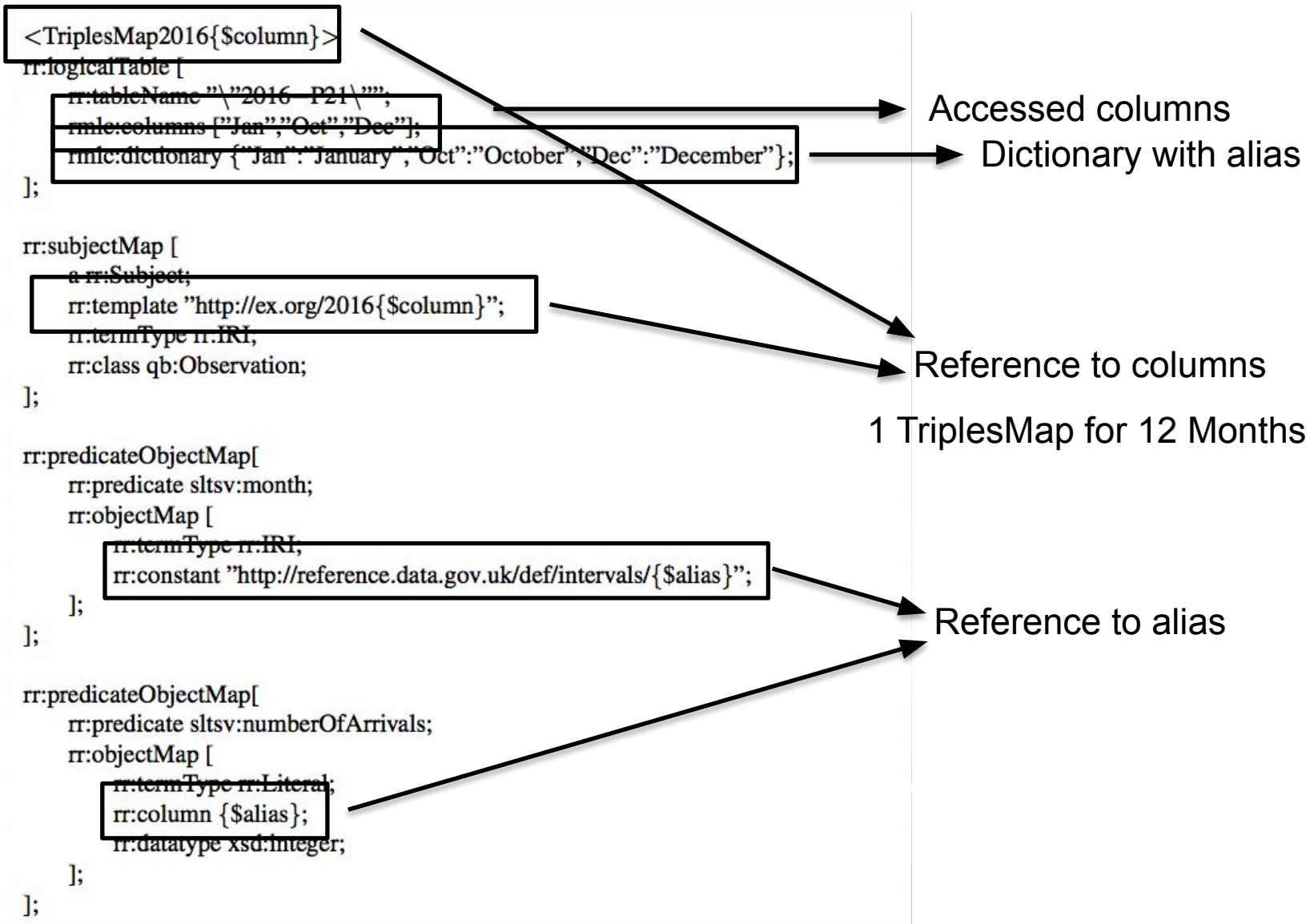


Difficulty of maintenance and creation

PROBLEM



Two variables for identifying independently each TriplesMap and provide access to the CSV data: `{ $column }, { $alias }`



Outputs:

- RMLC-Iterator for transforming the mappings to R2RML
- Morhp-RDB as OBDA engine for the query translation

Results:

D1

Features	R2RML	RMLC
Total Lines	~700	74
#TriplesMaps / #SubjectMaps	12	1
#PredicateObjectMaps	60	5

D2

Features	R2RML	RMLC
Total Lines	>2800	<70
#TriplesMaps / #SubjectMaps	>40	1
#PredicateObjectMaps	>170	4



SATET: Providing access to multiple CSV on the Web using OBDA

**David Chaves-Fraga, Ontology Engineering Group
Universidad Politécnica de Madrid, Spain**

Freddy Priyatna, OEG-UPM

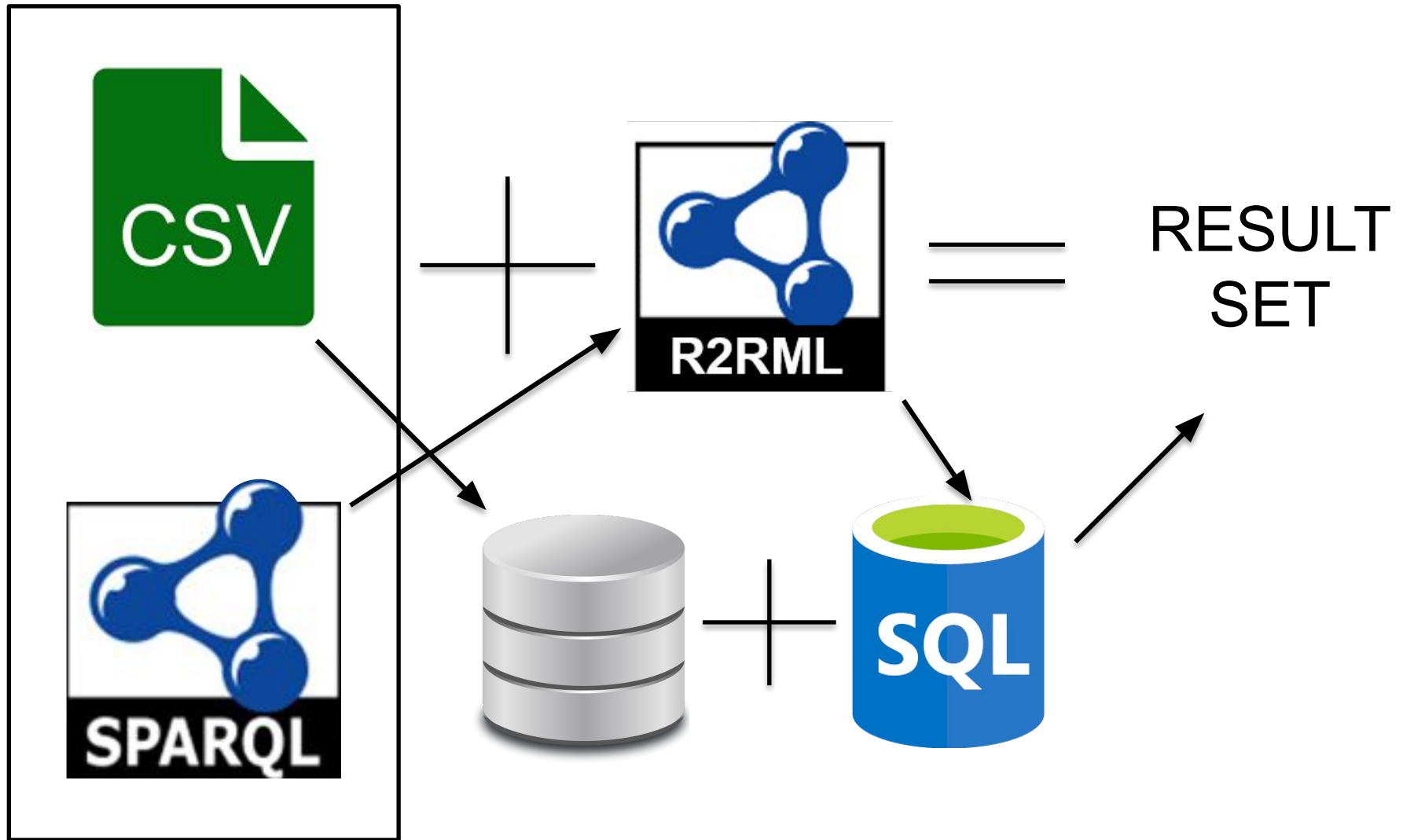
Oscar Corcho, OEG-UPM

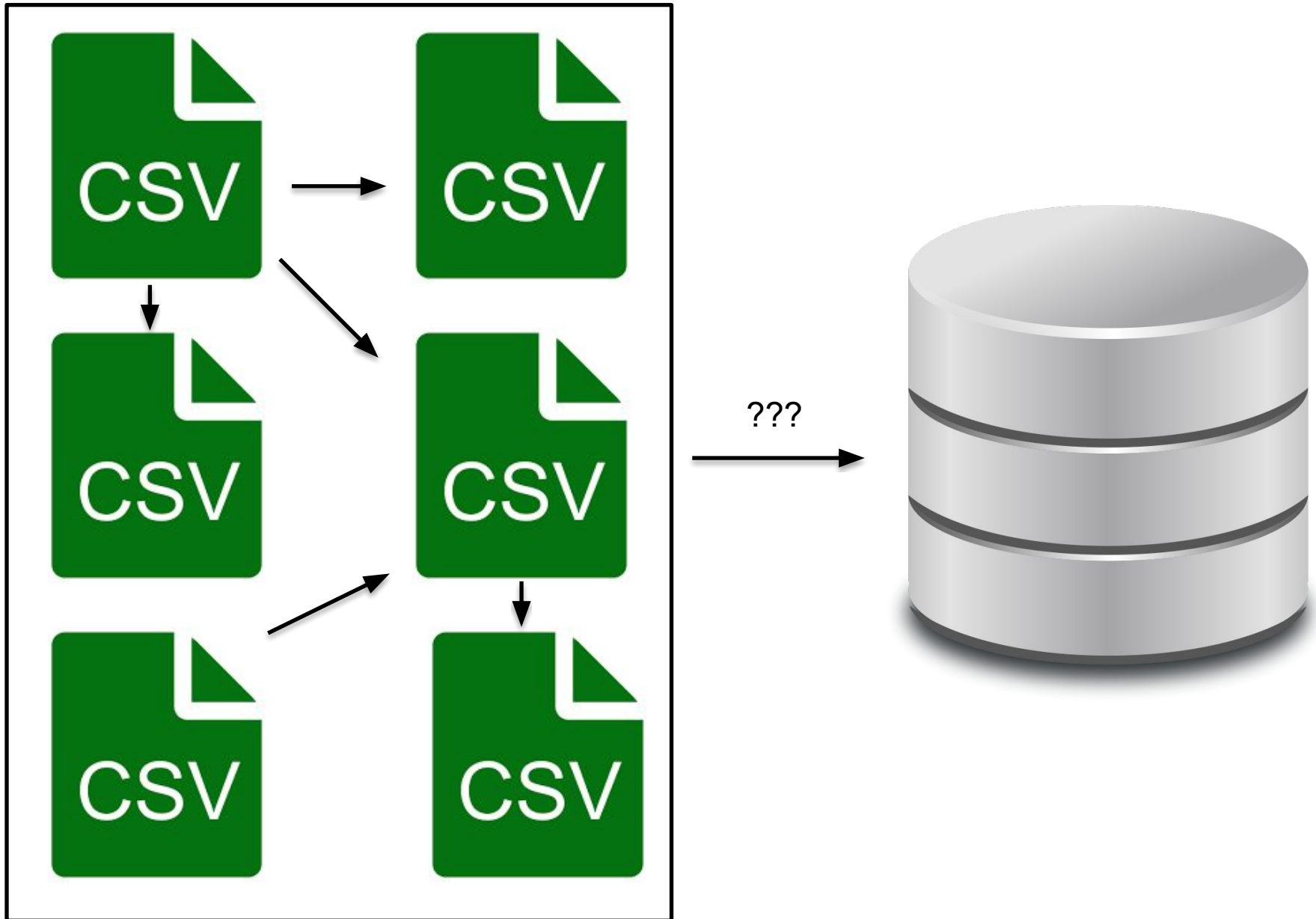
✉ dchaves@fi.upm.es

🐦 [@dchavesf](https://twitter.com/dchavesf)

📅 ???

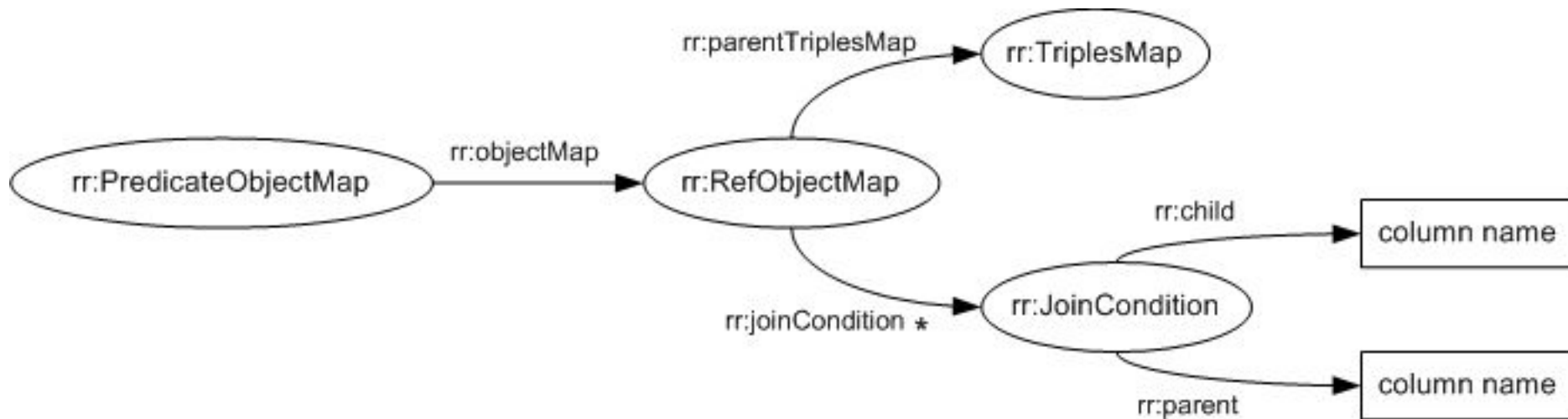
📍 ???





SATET: Semantic Access for heTEroogeneous Tabular data

- RMLC: Extension of R2RML for including SQL functions
 - Discover implicit joins among CSV files
 - Transforming CSV columns to RDF objects
- Generation of an enriched database schema using the mapping info (optimization)
- Semantic preservation of R2RML



Discovering implicit joins between CSV files

Relational Database

```
id,name,surname,birthdate,location  
1,david,chaves-fraga,27-11-1993,SDC
```

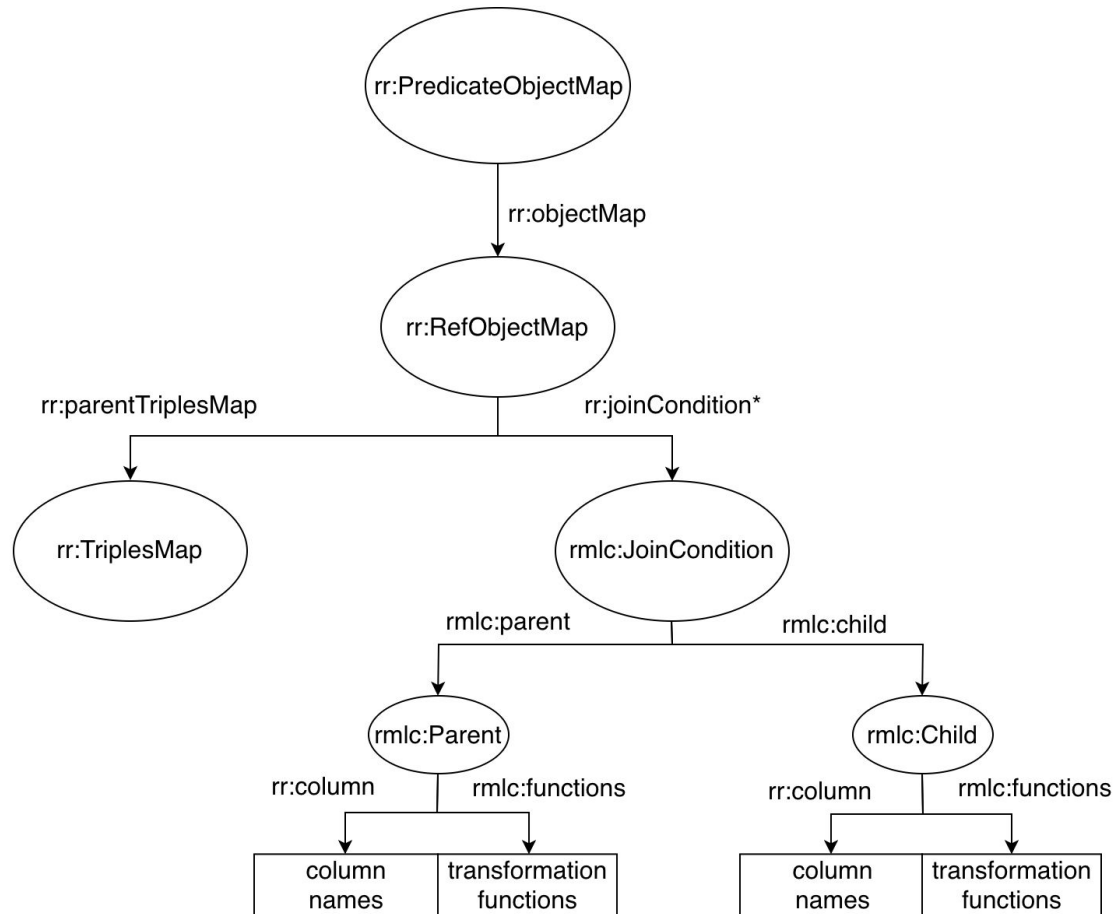
```
id,job  
1,phd_student
```

CSV files

```
name,surname,birthdate,location  
david,chaves_fraga,27111993,SDC
```

```
full_name,job  
"David Chaves Fraga","phd_student"
```

RMLC: RDF Mapping Language extension for heterogeneous CSV files



The functions are SQL basic transformation functions

Table 1

```
name,surname,birthdate,location  
david,chaves_fraga,27111993,SDC
```

Table 2

```
full_name,job  
"David Chaves Fraga","phd_student"
```

```
SELECT ?name ?birthday ?job WHERE {  
  ?name ?p1 ?birthday.  
  ?name ?p2 ?job .  
}
```

```
<#TriplesMap1>
```

```
....
```

```
rr:predicateObjectMap[
  rr:predicate foaf:name;
```

```
rr:objectMap [
```

```
  rr:parentTriplesMap <#TriplesMap2>;
```

```
  rr:joinCondition [
```

```
    rmlc:child [
```

```
      rmlc:functions "LOWER({FULL_NAME})";
```

```
    ];
```

```
    rmlc:parent [
```

```
      rmlc:functions "CONCAT({NAME},' ',REPLACE({SURNAME},' ',' '))";
```

```
    ];
```

```
  ];
```

```
];
```

```
];
```

```
.
```



```
SELECT ?name ?birthday ?job
WHERE {
  ?name ex:birthday ?birthday.
  ?name ex:job ?job .
}
```

```
SELECT name, birthday, table2.job FROM table1
INNER JOIN table2 ON
CONCAT(table1.name,' ',REPLACE(table1.surname,' ',' ')) = LOWER(table2.full_name)
```

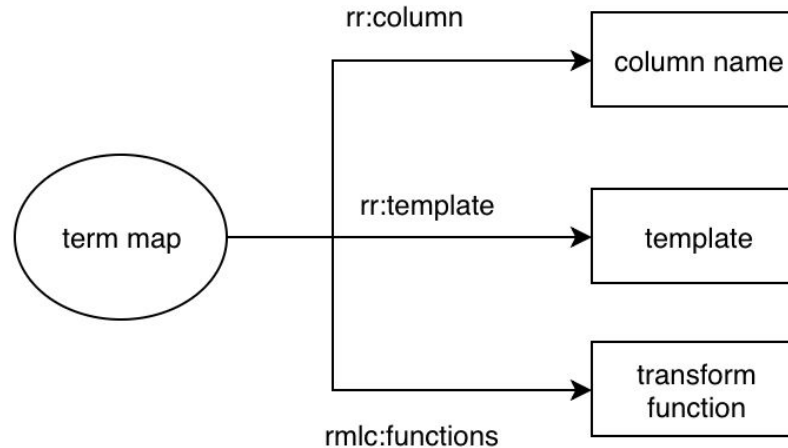
```
{
  listSocialMediaPosting {
    identifier
    comment
    author {
      identifier
      email
      familyName
      givenName
      name
      telephone
    }
  }
}
```

```
{
  "data": {
    "listSocialMediaPosting": [
      {
        "identifier": "http://ex.org/1",
        "comment": "Hallo Dunia@id",
        "author": {
          "identifier": "http://ex.org/Person/1",
          "email": "fpriyatna@fi.upm.es",
          "familyName": "Priyatna",
          "givenName": "Freddy",
          "name": "Freddy Priyatna",
          "telephone": "8141"
        }
      },
      {
        "identifier": "http://ex.org/2",
        "comment": "Hola Mundo@es",
        "author": {
          "identifier": "http://ex.org/Person/2",
          "email": "dchaves@fi.upm.es",
          "familyName": "Chaves",
          "givenName": "David",
          "name": "David Chaves",
          "telephone": "9063"
        }
      },
      {
        "identifier": "http://ex.org/3",

```

```
SELECT
  "listSocial"."id" AS "id",
  'http://ex.org/' || "listSocial".id || ' ' AS "identifier",
  "listSocial"."mensaje" AS "comment",
  "author"."id" AS "author__id",
  'http://ex.org/Person/' || "author".id || ' ' AS "author__identifier",
  lower(substr("author"."nombre",1,1) || "author"."apellido || '@fi.upm.es') AS "author__email",
  "author"."apellido" AS "author__familyName",
  "author"."nombre" AS "author__givenName",
  ' ' || "author"."nombre || ' ' || "author"."apellido || ' ' AS "author__name",
  "author"."telephone" AS "author__telephone"
FROM comentarios "listSocial"
LEFT JOIN personas "author" ON "listSocial".usuario = lower(substr("author"."nombre",1,1) || "author"."apellido)
```

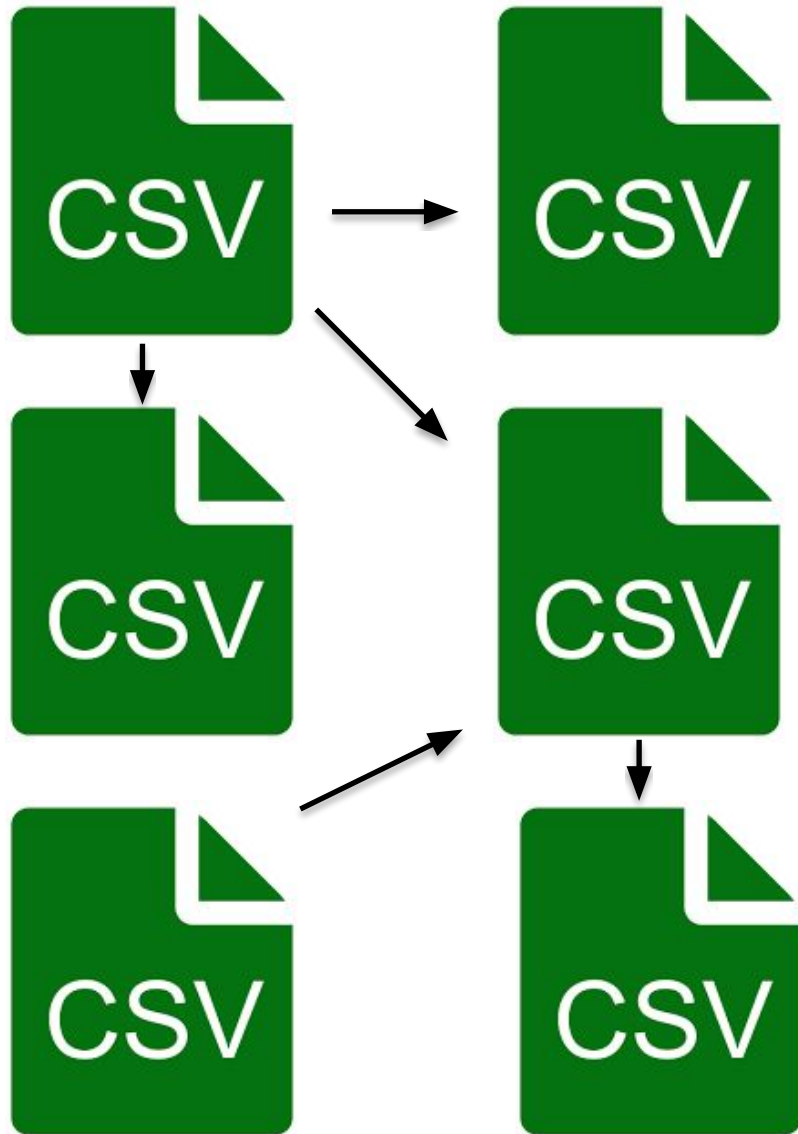
Transforming CSV columns to RDF objects



<TriplesMap1>

...

```
rr:predicateObjectMap[
  rr:predicate ex:shortName;
  rr:objectMap [
    rr:datatype xsd:string;
    rmlc:functions "REPLACE(SUBSTRING(LOWER{FULL_NAME},1,5),' ','-')";
  ];
]
rr:predicateObjectMap[
  rr:predicate ex:yearofBirthday;
  rr:objectMap [
    rmlc:functions "YEAR({birthday})"; ];
];
```



+ RMLC =

SATET

++



- Primary Keys
- Foreign Keys
- Datatypes

- RMLC maintains the semantics of R2RML
- It's aligned with R2RML:
 - ObjectMaps with Functions → new column in the table with the name of the predicate
 - Joins with Functions → new columns in the tables
 - SATET transforms RMLC to R2RML
- SATET can be introduced on the top of state-of-art OBDA engines (morph/ontop) for using their optimizations to efficiently access to CSV files

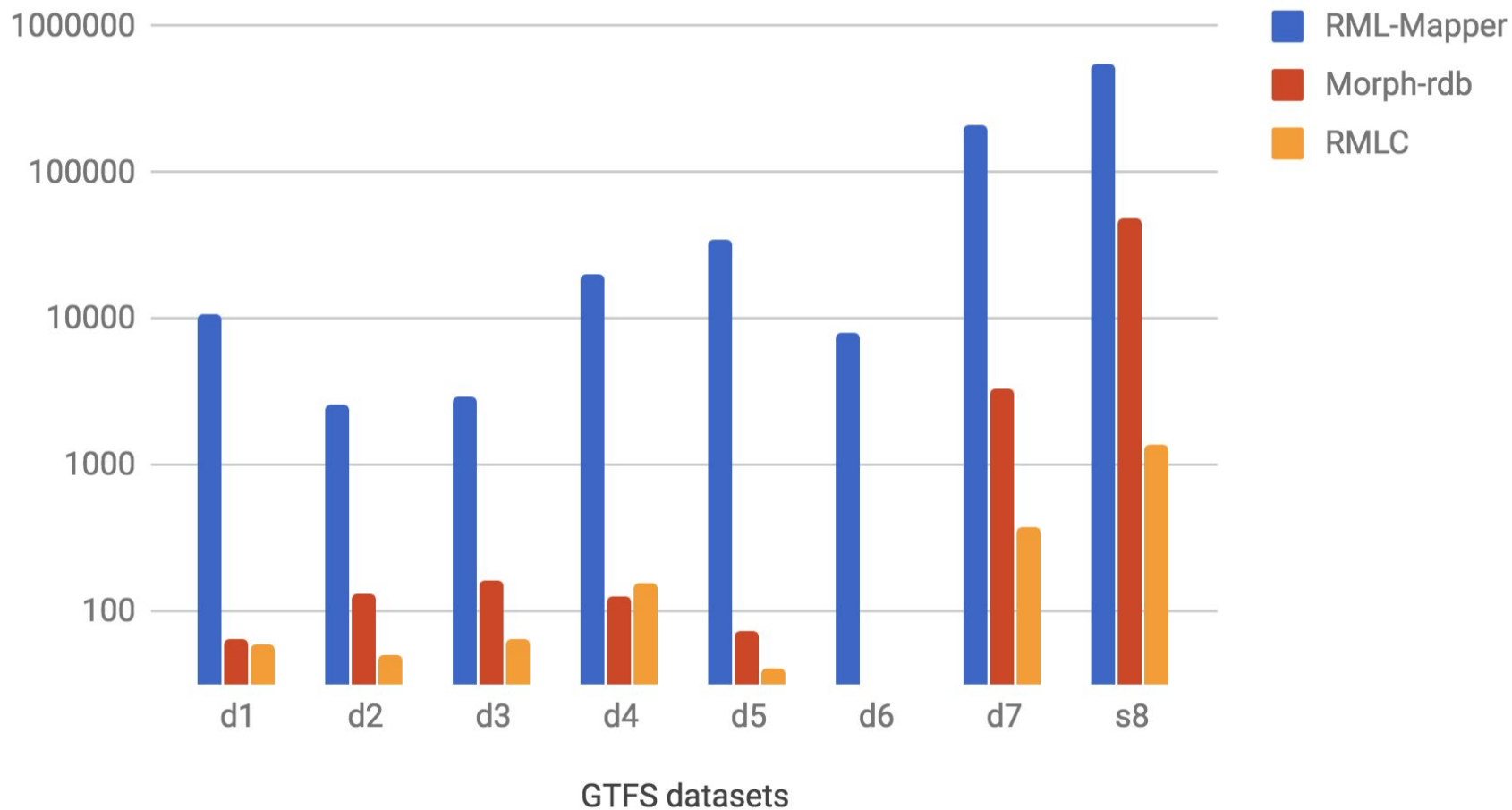

```
SELECT ?name ?birthday ?job
WHERE {
    ?name ex:birthday ?birthday.
    ?name ex:job ?job .
}
```

```
<#TriplesMap1>
....
rr:predicateObjectMap[
  rr:predicate ex:fullName;
  rr:objectMap [
    rr:parentTriplesMap <#TriplesMap2>;
    rr:joinCondition [
      rmlc:child [
        rmlc:functions "LOWER({FULL_NAME})";
      ];
      rmlc:parent [
        rmlc:functions "CONCAT({NAME},' ',REPLACE({SURNAME},' ',' '))";
      ];
    ];
  ];
];
.
```



```
SELECT name, birthday, table2.job FROM table1
INNER JOIN table2 ON table1.fullName = table2.fullName
```

GTFS to RDF materialization



SATET: Semantic Access for heTErogeneous Tabular data

Main Contributions:

- Discover implicit joins
- Apply transformation functions to individual columns
- Enriched database schema from mapping information
- Semantic preservation of R2RML

Future Work:

- Alignment with FnO → full specification (Possible collaboration)
- Alignment with RML (without FnO) for LD Generation from RDB/CSV
- Optimizations over generated SQL queries
- Query answering over SATET
- Applying to transport domain for linking potential datasets during a route planning creation
- Define the core for the mapping languages of SW (Possible collaboration)



Virtual Knowledge Graph Generation from heterogeneous data sources

**David Chaves-Fraga, Ontology Engineering Group
Universidad Politécnica de Madrid, Spain**

Freddy Priyatna, OEG-UPM

Oscar Corcho, OEG-UPM

✉ dchaves@fi.upm.es

🐦 [@dchavesf](https://twitter.com/dchavesf)

📅 24-25/10/2018

📍 [@imec/Ghent University](https://imec-ghent.com/)

