



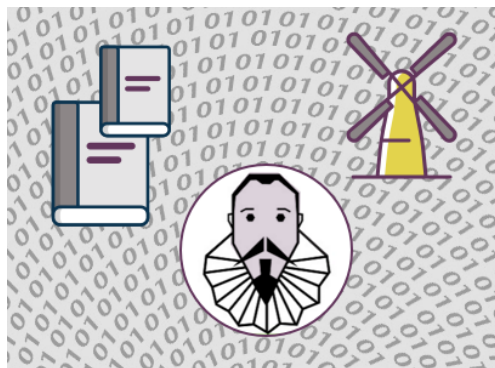
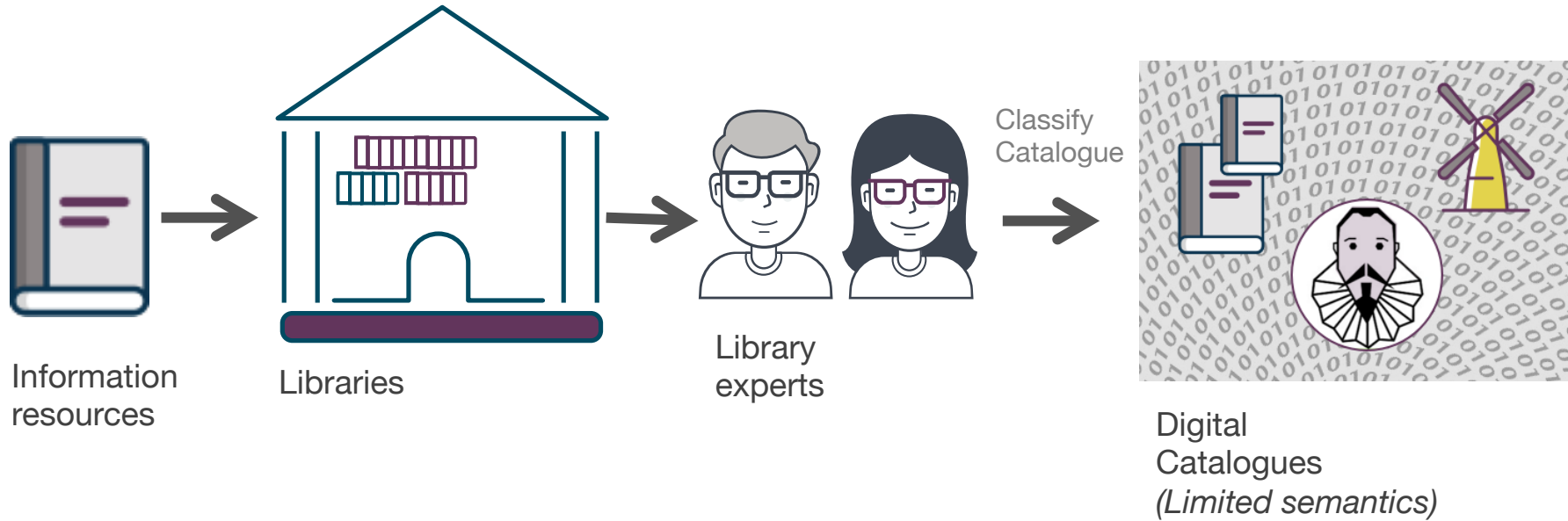
A FRAMEWORK FOR ONTOLOGY-BASED LIBRARY DATA GENERATION, ACCESS AND EXPLOITATION

Daniel Vila Suero

Advisors:

Prof. Dr. Asunción Gómez-Pérez and Dr. Jorge Gracia del Río

PhD in Artificial Intelligence,
Defense, Madrid, 27th of July 2016

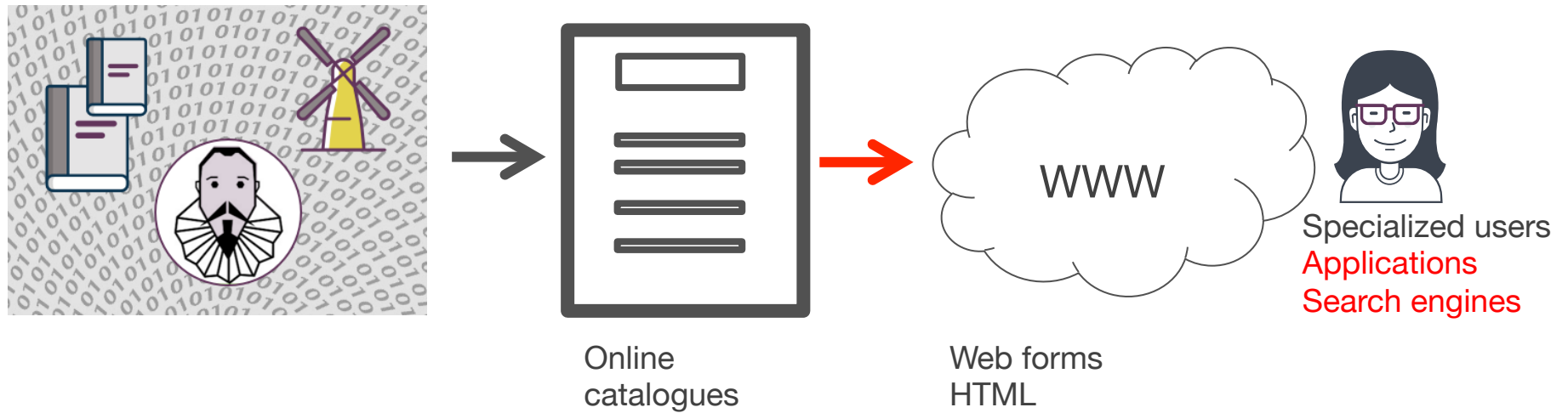


Library records describing **persons, concepts, works, and other entities.**

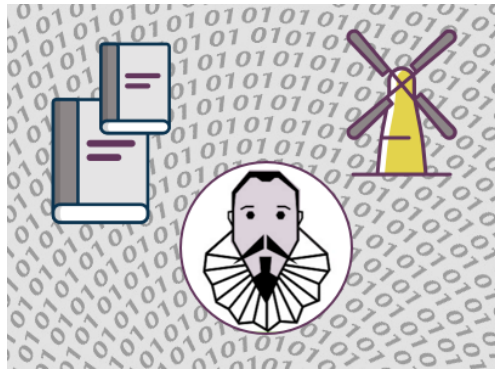
MARC 21 international standard describing **millions of records.**

(e.g., European Library +80 million records)

CATALOGUE RESOURCES ARE DISCONNECTED FROM THE WEB



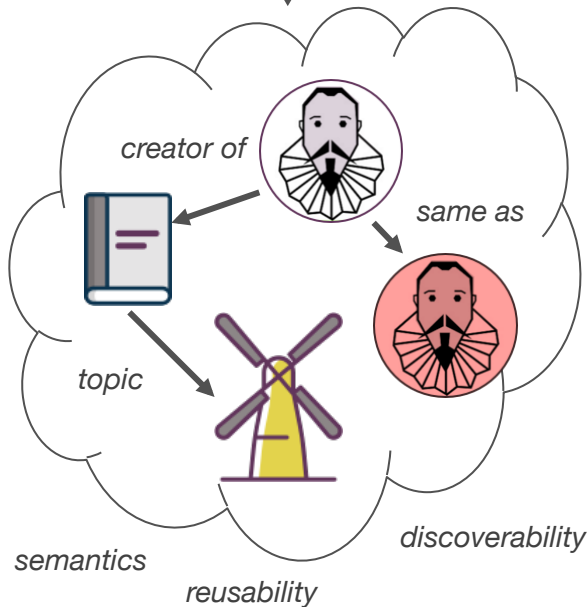
CATALOGUE RESOURCES ARE DISCONNECTED FROM THE WEB



Specialized users
Applications
Search engines



?



Online catalogues

Web forms
HTML

Semantic web and Linked Data technologies
(**Ontologies, RDF, HTTP URIs**):



My thesis

“**Creating and delivering library data** while providing a natural **extension** to the **collaborative sharing models** historically employed by libraries.”

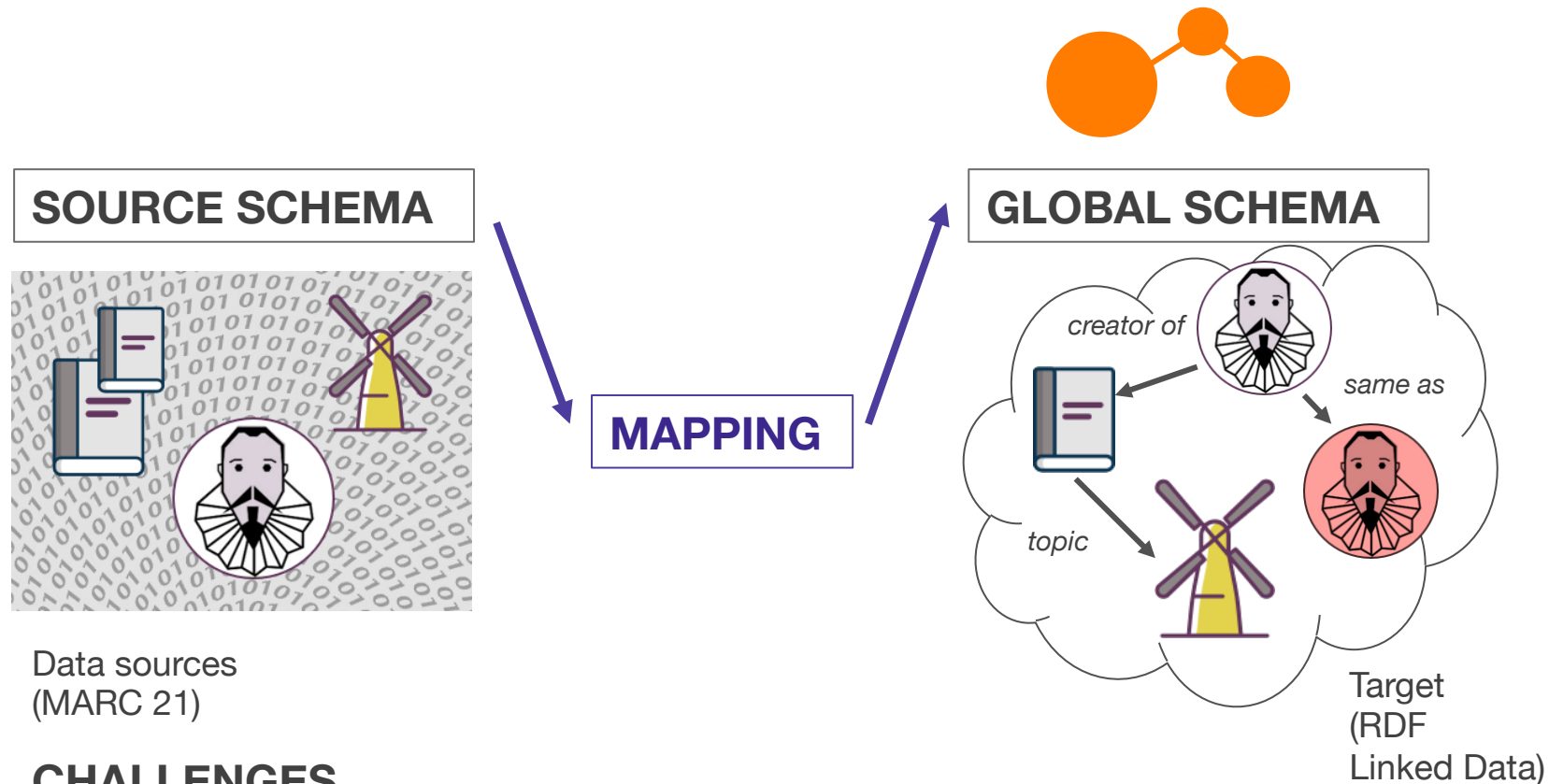
OUTLINE

- **Problem and state of the art**
- Research methodology and hypotheses
- Contributions
 - Mapping
 - Ontology development
 - Library applications
- Conclusion and future directions

RESEARCH PROBLEM

**PRINCIPLED TRANSFORMATION OF
LIBRARY CATALOGUES INTO
ONTOLOGY-BASED DATA
AND THEIR PUBLICATION AND CONSUMPTION
ON THE WEB.**

THE PROBLEM: GLOBAL-AS-VIEW ONTOLOGY-BASED DATA INTEGRATION



CHALLENGES

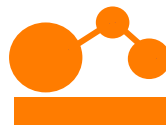
1. Library records present a **nested structure** (fields, subfields, etc.)
2. Library data sources lack **schema and query mechanisms**
3. Library standards and ontologies are **highly specialized**

RESEARCH PROBLEM AND AREAS

**PRINCIPLED TRANSFORMATION OF
LIBRARY CATALOGUES INTO
ONTOLOGY-BASED DATA
AND THEIR PUBLICATION AND CONSUMPTION
ON THE WEB.**



**TRANSFORMATION &
MAPPING LANGUAGES**



**ONTOLOGY
DEVELOPMENT**

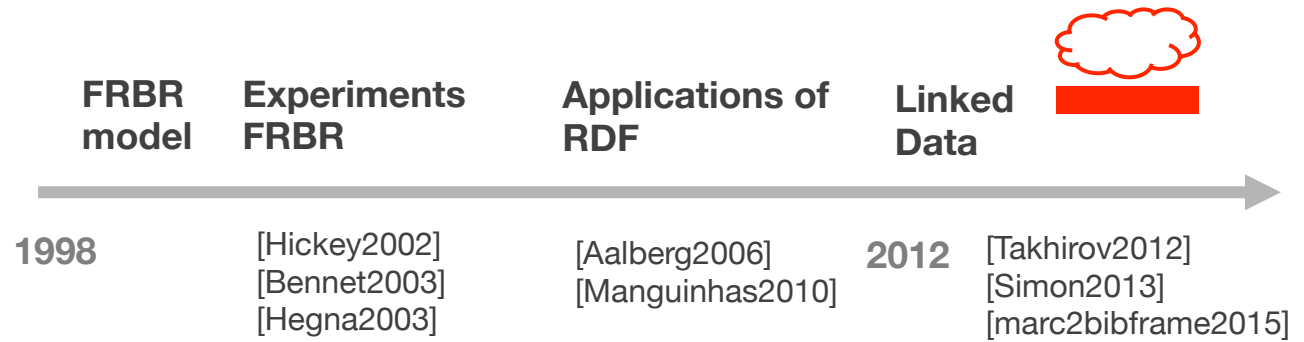


**LIBRARY
APPLICATIONS**

STATE OF THE ART



TRANSFORMATION METHODS



STATE OF THE ART



FRBR
model

Experiments
FRBR

Applications of
RDF

Linked
Data



1998

2012

TRANSFORMATION METHODS



Relational
databases into
RDF

R2RML
W3C standard

Extensions of R2RML
for non-RDB data sources

2004 [Barrasa04]
[Bizer04]

2014

RML [Dimou2014]
xR2RML [Michel2015]
KR2RML [Slepicka2015]

MAPPING LANGUAGES

STATE OF THE ART



FRBR model

Experiments FRBR

Applications of RDF

Linked Data



1998

2012

TRANSFORMATION METHODS



Relational databases into RDF

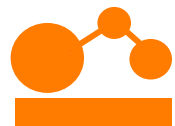
R2RML W3C standard

Extensions of R2RML for non-RDB data sources

2004

2014

MAPPING LANGUAGES



Initial methodologies

Ontology life-cycle models

Agile approaches

1995

[Grüniger1995]
[Fernández-López1997]

[Staab2001]
[Pinto2004]
Neon methodology
[Suarez-Figueroa2015]

[Auer2006]
[Presutti2012]

2015

ONTOLOGY DEVELOPMENT

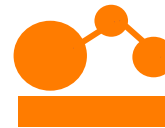
LIMITATION 1: MAPPING LANGUAGE AND METHODS



TRANSFORMATION



MAPPING LANGUAGES



ONTOLOGY
DEVELOPMENT



LIBRARY
LINKED DATA
APPLICATIONS

Lack of **interoperability** and reuse of mapping rules.

LIMITATION 1: MAPPING LANGUAGE AND METHODS



TRANSFORMATION

Lack of **interoperability** and reuse of mapping rules.



MAPPING LANGUAGES

Format-dependent (e.g., JSONPath)
[RML] [xR2RML]

Lack of **queries** over nested data
[KR2RML]



ONTOLOGY DEVELOPMENT



LIBRARY LINKED DATA APPLICATIONS

LIMITATION 2: FEEDBACK OF LIBRARY EXPERTS



TRANSFORMATION

Lack of **interoperability** and reuse of mapping rules.



MAPPING LANGUAGES

Format-dependent (e.g., JSONPath)



ONTOLOGY DEVELOPMENT

Lack of **queries** over nested data



LIBRARY LINKED DATA APPLICATIONS

Feedback of library experts not integrated within the transformation process.

[Takhirov2012]

LIMITATION 2: FEEDBACK OF LIBRARY EXPERTS



TRANSFORMATION

Lack of **interoperability** and reuse of mapping rules.

Feedback of library experts not integrated within the transformation process.
[Takhirov2012]



MAPPING LANGUAGES

Format-dependent (e.g., JSONPath)

Lack of **queries** over nested data



ONTOLOGY DEVELOPMENT

Limited **technical support for the active participation** of domain experts.
[NeOn]

Limited **methods for understanding similarities** among overlapping ontologies.
[Vandenbussche2014]



LIBRARY LINKED DATA APPLICATIONS

LIMITATION 3: **EXPERIMENTAL RESULTS****TRANSFORMATION**

Lack of **interoperability** and reuse of mapping rules.

Feedback of library experts not integrated within the transformation process.

**MAPPING LANGUAGES**

Format-dependent (e.g., JSONPath)

Lack of **queries** over nested data

**ONTOLOGY DEVELOPMENT**

Limited **technical support** for the **active participation** of domain experts.

Limited **methods for understanding similarities** among overlapping ontologies.

**LIBRARY LINKED DATA APPLICATIONS**

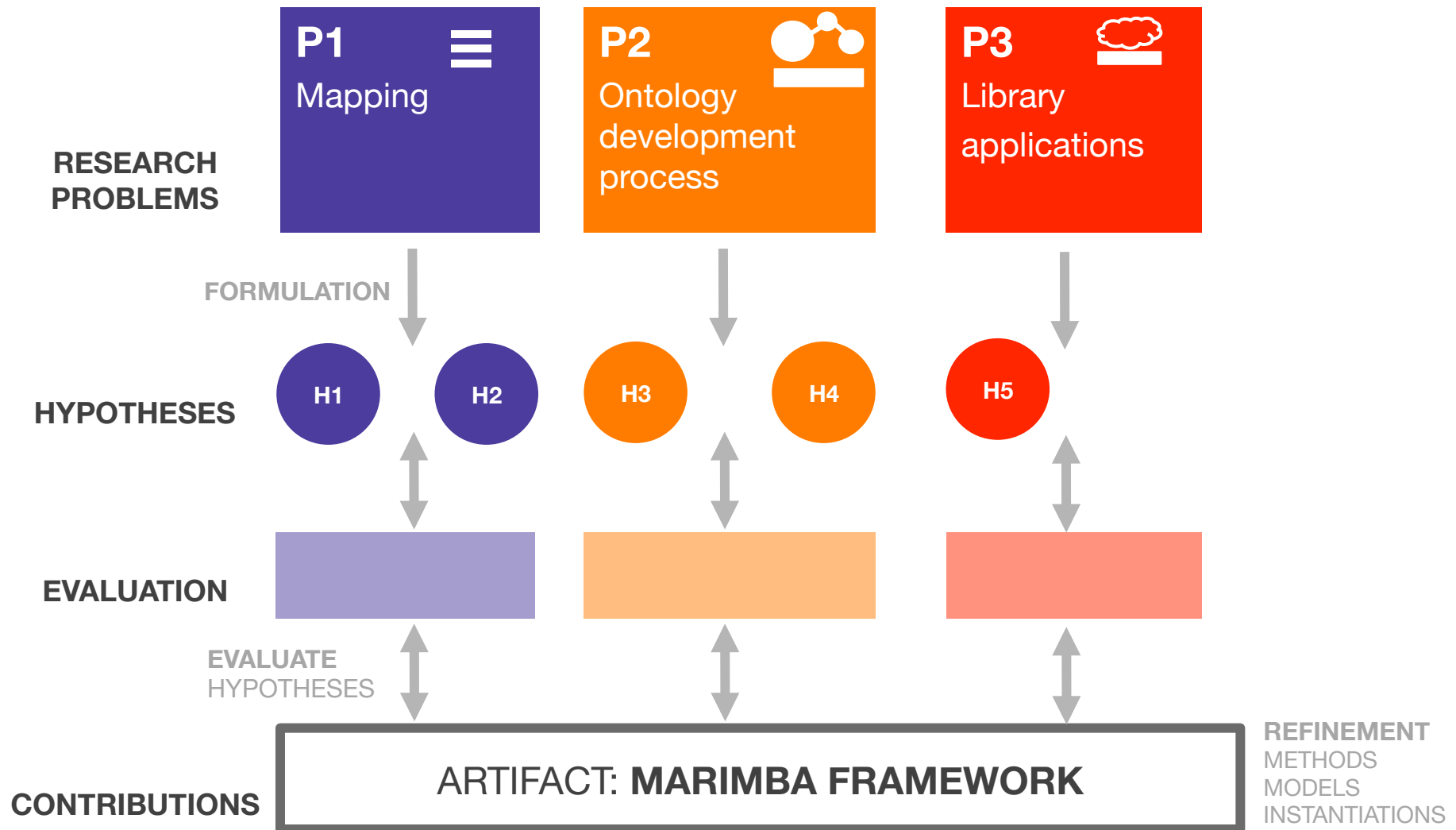
Lack of **evaluation and experimental results** of impact in end-user library applications
[Simon2013]

HYPOTHESES

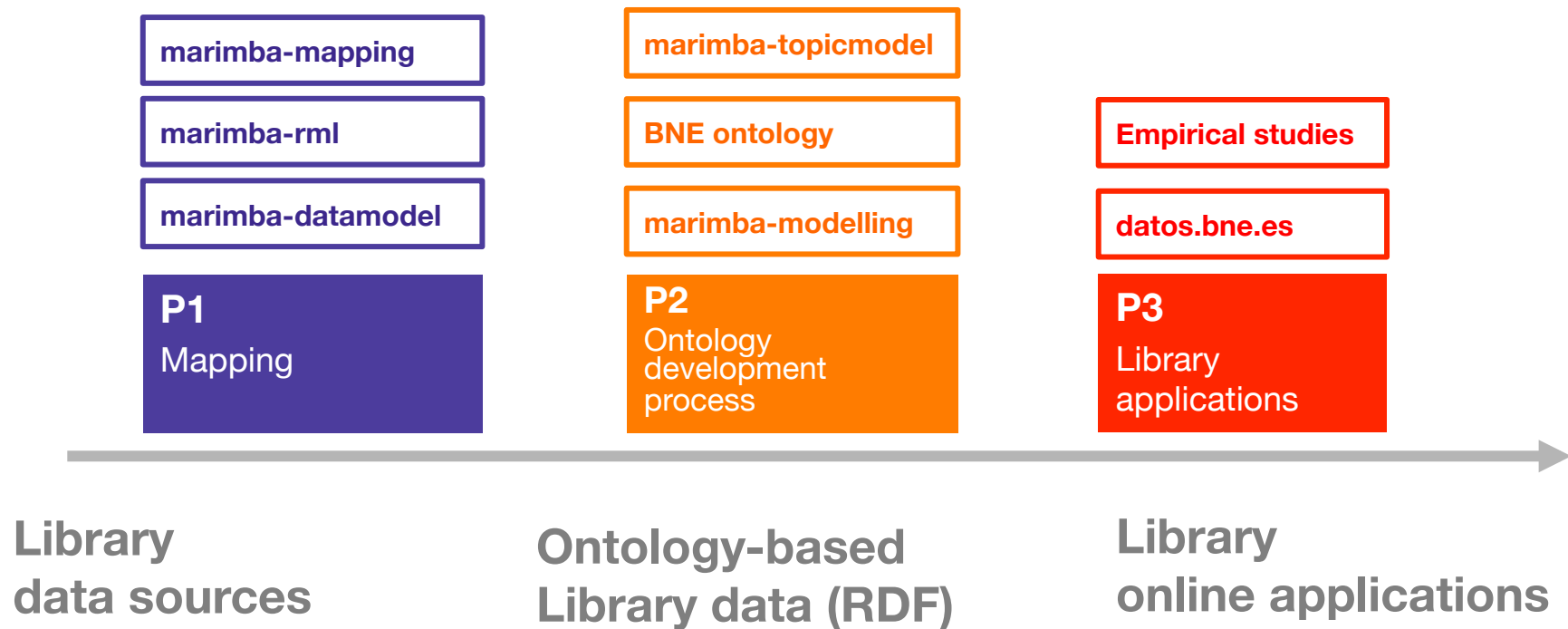
Research methodology, and hypotheses



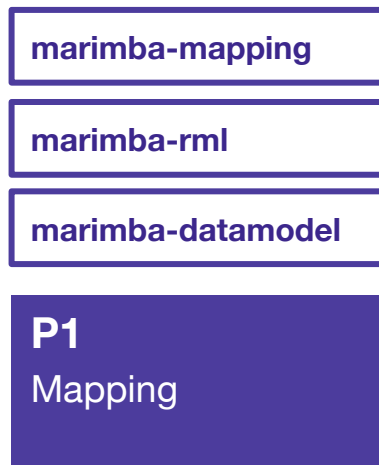
DESIGN-SCIENCE RESEARCH METHODOLOGY



MARIMBA FRAMEWORK: PROBLEMS AND CONTRIBUTIONS



P1: DEFINITION OF INTEROPERABLE MAPPING RULES

**H1**

The structure of MARC 21 records can be **fully represented** using a **nested relational model**.

H2

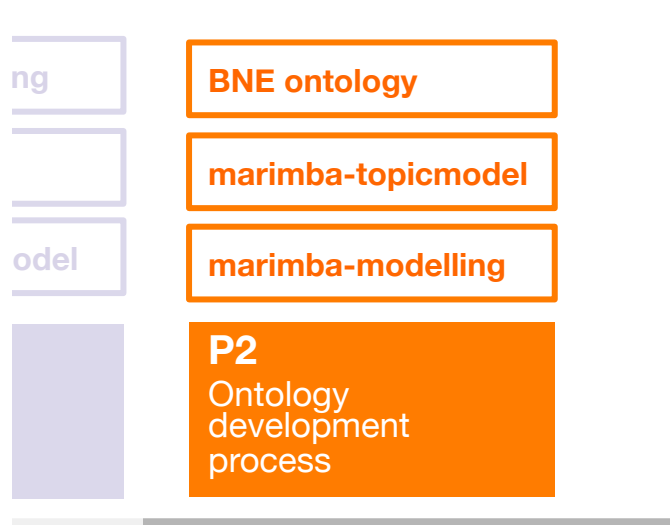
Minimal modifications to the **W3C R2RML** language to define mapping rules between **MARC 21 data sources into RDF**.

Library data sources

RESTRICTIONS:

No query optimization, only operational semantics of mapping language, authority and bibliographic formats, MARC 21 standard

P2: LIBRARY ONTOLOGY DEVELOPMENT



H3

Analytical data and the **feedback of library experts** can be used to develop a library **ontology** with sufficient quality

H4

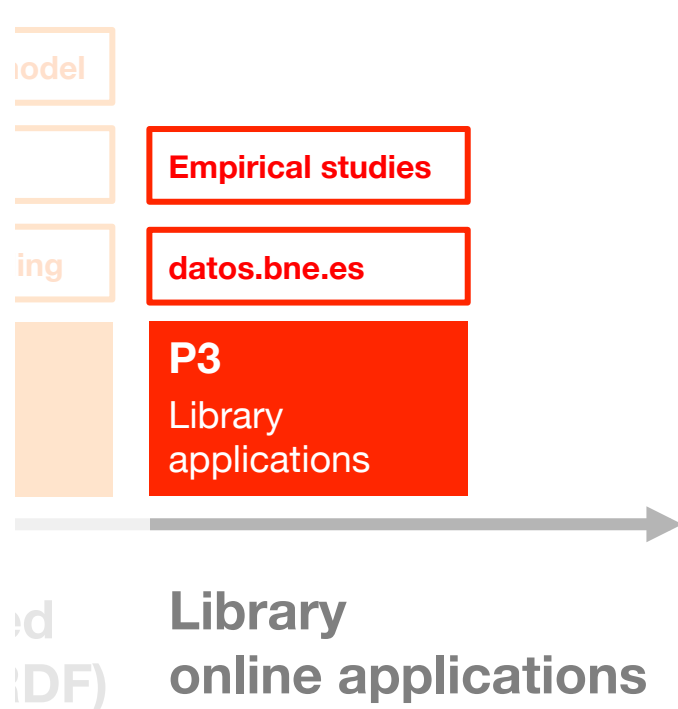
Probabilistic topic modelling techniques can produce **coherent descriptions of ontologies** and can perform **better than existing methods** used in ontology search.

Ontology-based
Library data (RDF)

ASSUMPTIONS:

Participation of library experts, available library ontologies and records

P3: APPLICATION OF SEMANTIC TECHNOLOGIES IN LIBRARY APPLICATIONS



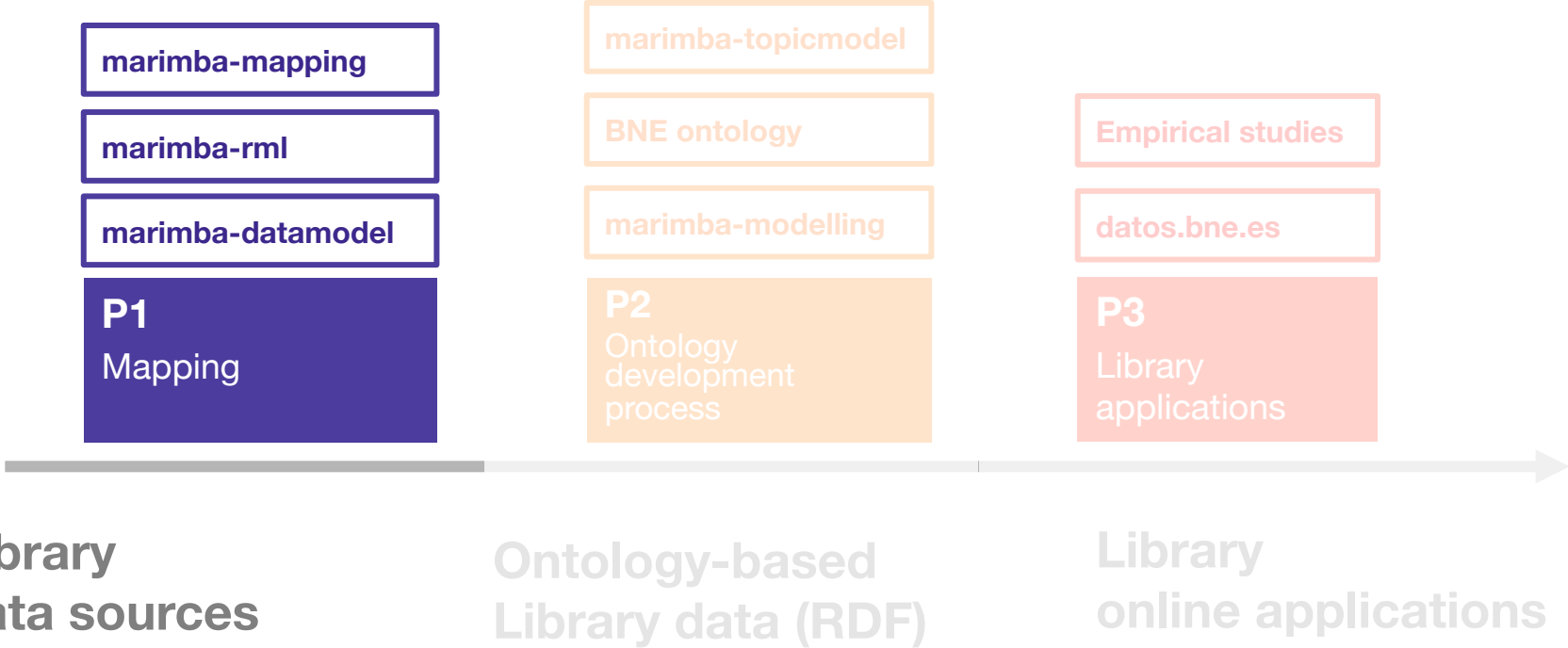
H5

The application of **semantic technologies** to end-user library applications can **increase user satisfaction** and **efficiency** for finding information

RESTRICTIONS:

Transformation in a batch process

OUTLINE



marimba-datamodel

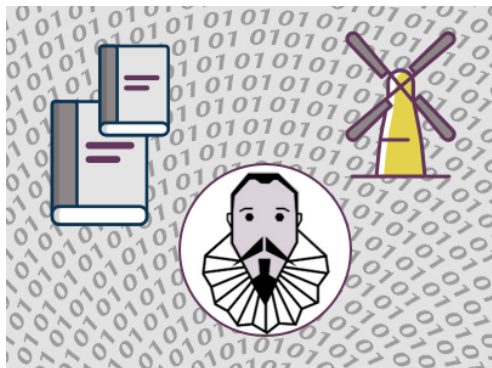
marimba-mapping

marimba-rml

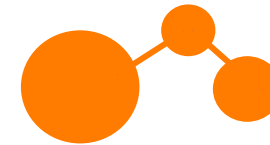
P1
Mapping

?

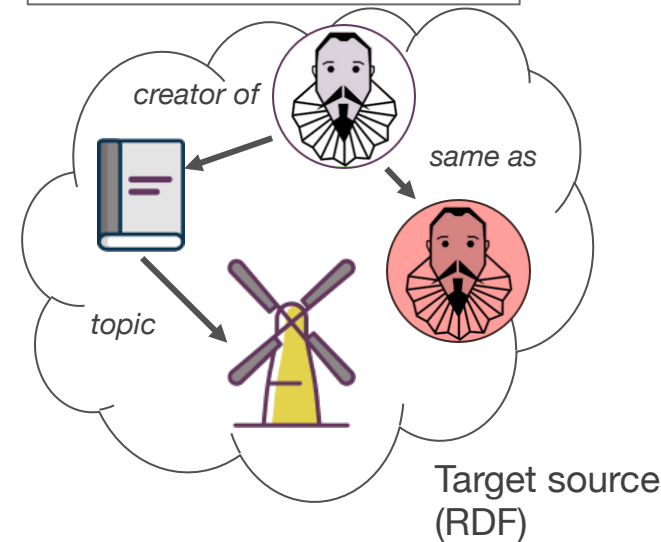
SOURCE SCHEMA



Data source
(MARC 21)



GLOBAL SCHEMA



Target source
(RDF)

marimba-mapping

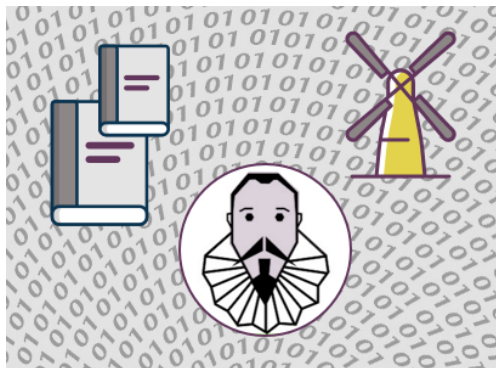
marimba-rml

P1
Mapping

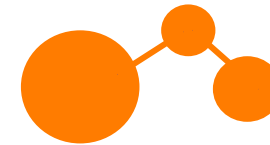
Schema extraction

marimba-datamodel

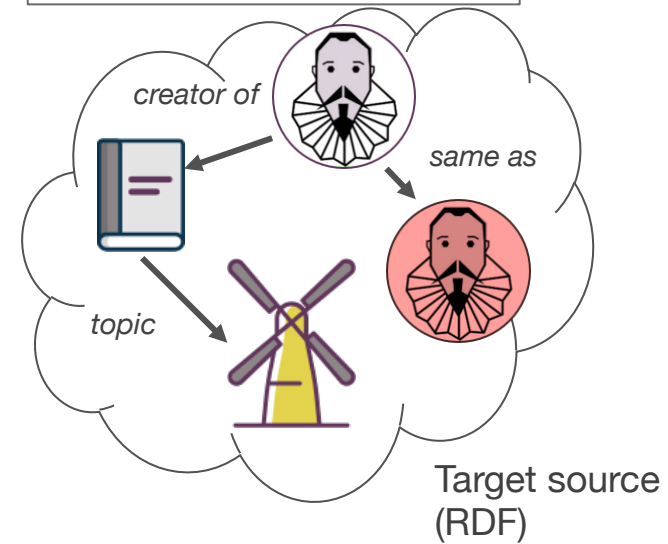
SOURCE SCHEMA



Data source
(MARC 21)



GLOBAL SCHEMA



marimba-datamodel

marimba-rml

marimba-mapping

- Study of **MARC 21 standard in Nested Relational Model** (Makinouchi [1977]).
- Attributes:
 - **Atomic:** e.g., Field 001
 - **Relation-valued:** e.g., Field 100, 008
- **Complete coverage** of MARC 21 standard elements

MARC 21	marimba-datamodel (NRM)	H1
<i>Format</i>	Relation scheme $R(S)$	
<i>Leader</i>	relation-valued attribute $R_{leader}(S_{leader}) \in S$ with $S_{leader} = (p00, p01, p02, \dots, p23)$	
<i>Control Field (fixed-length)</i>	relation-valued attribute $R_{cf}(S_{cf}) \in S$ with $S_{cf} = (p00, p01, p02, \dots, pN)$ for N positions	
<i>Control Field (single-valued)</i>	Atomic attribute $cf \in S$	
<i>Data Field</i>	relation-valued attribute $R_{df}(S_{df}) \in S$ with $S_{cf} = (i1, i2, sf1, \dots, sfN)$ with $i1$ and $i2$ indicators and $sf1..sfN$ the N subfields of the data field	
<i>Subfield</i>	Atomic attribute sf in a relation scheme of a data field $\in S_{df}$	
<i>Indicators</i>	Atomic attributes $i1$ and $i2$ in a relation scheme of a data field $\in S_{df}$	

marimba-datamodel

marimba-rml

marimba-mapping

Schema extraction

Novel **algorithm** for **schema extraction** of MARC 21 records.

- Support the **validation of queries to MARC 21 data**.
- Support the generation of **mapping templates for library experts**.
- **Complete representation of patterns in existing MARC 21 data sources**.

SCHEMA**TABLE bibliographic**

```
ITEM f001 UNIQUE NOT NULL /* C.field Unique ID */
```

```
ITEM TABLE leader /* Leader Info about language, etc.*/
```

```
ITEM p01
```

```
...
```

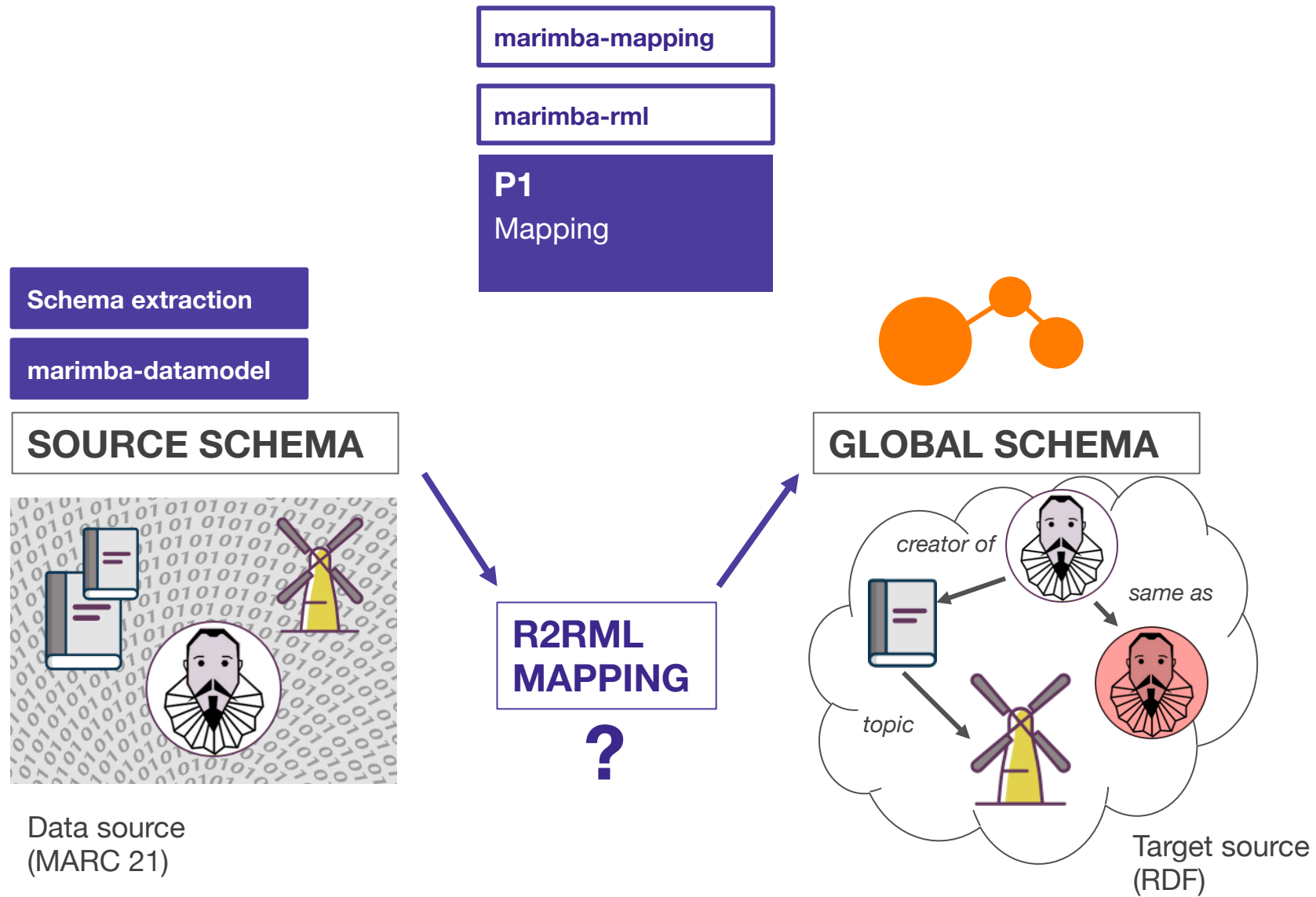
```
ITEM p23
```

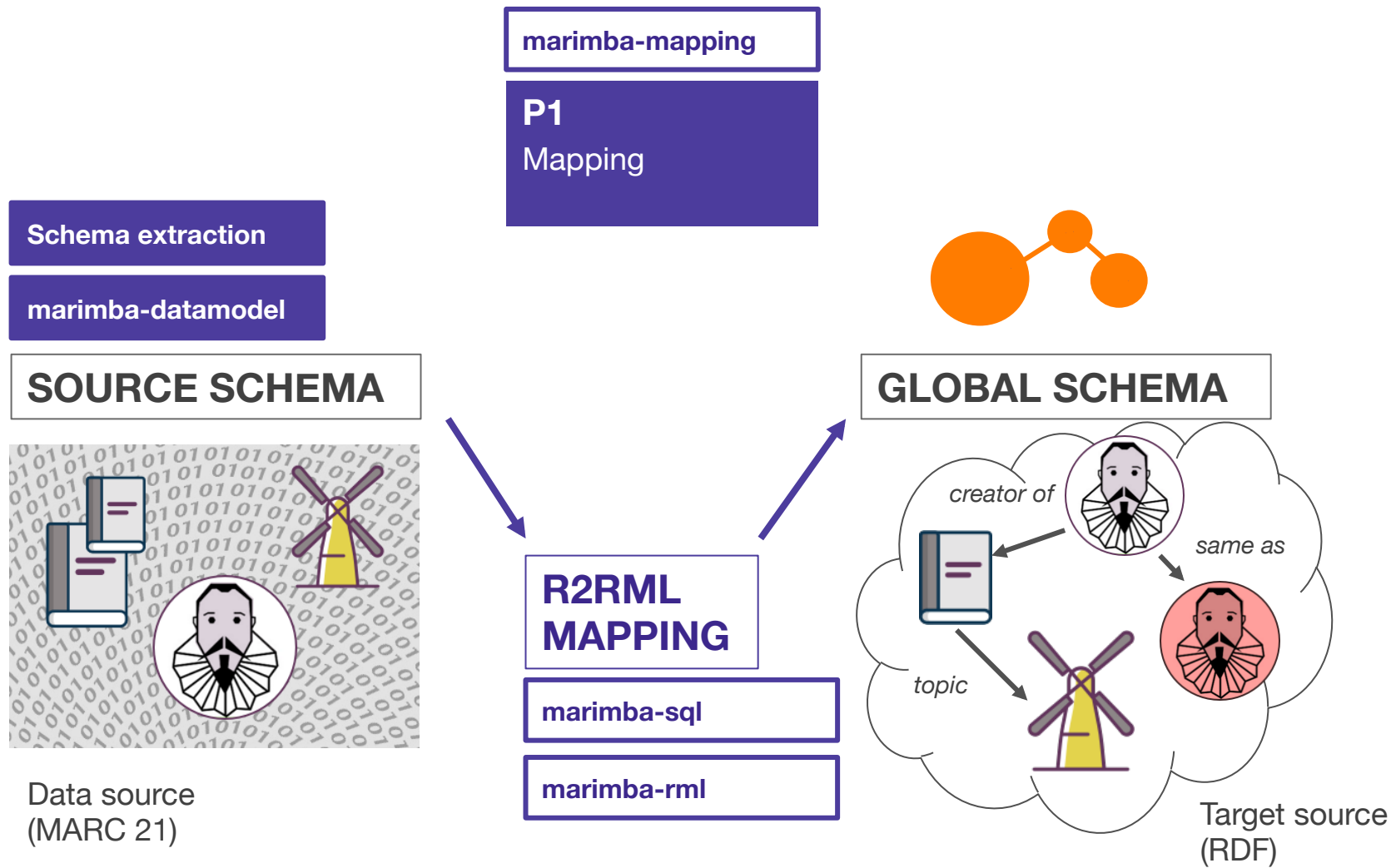
```
ITEM TABLE f100 /*Data field Author information */
```

```
ITEM a /* Subfield */
```

```
...
```

```
ITEM i0 /* Indicator */
```





marimba-datamodel

marimba-rml

marimba-mapping

marimba-sql

Minimal **query language** based on SQL/NF (Roth et al. [1987]):

- **Conciseness**
- **Orthogonality** of expressions → SFW nested expressions.

```
/* Selects bibliographic records of type Drawing*/  
bibliographic / * SELECT FROM * /  
WHERE  
EXISTS  
/ * NESTED EXPRESSION* /  
(leader WHERE p6 = "k" AND p7 != "s")  
AND  
/ * NESTED EXPRESSION * /  
EXISTS (f007 WHERE p2 = "d")
```



*BNF grammar provided in Annex I