

# Tutorial on Semantic Schema Discovery: principles, methods and future research directions Part 3

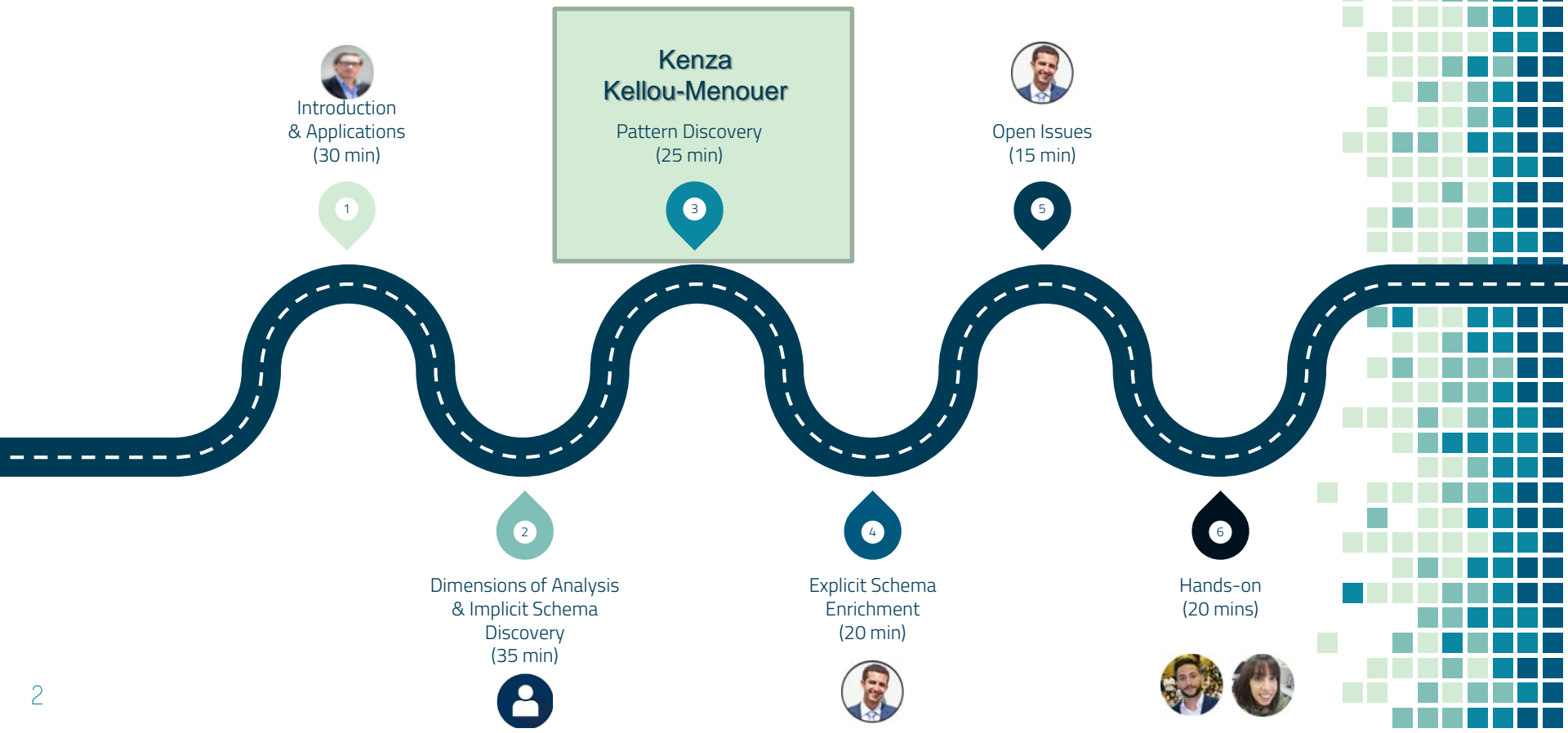
Kenza Kellou-Menouer, Nikolaos Kardoulakis, Georgia  
Troullinou, Zoubida Kedad, Dimitris Plexousakis,  
Haridimos Kondylakis



Équipes Traitement  
de l'Information  
et Systèmes

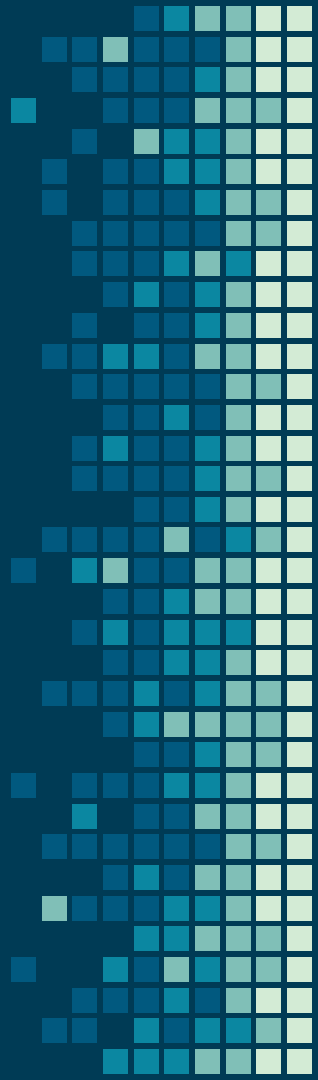


# ROADMAP





# Pattern Discovery Approaches



# Terminology

## Pattern:

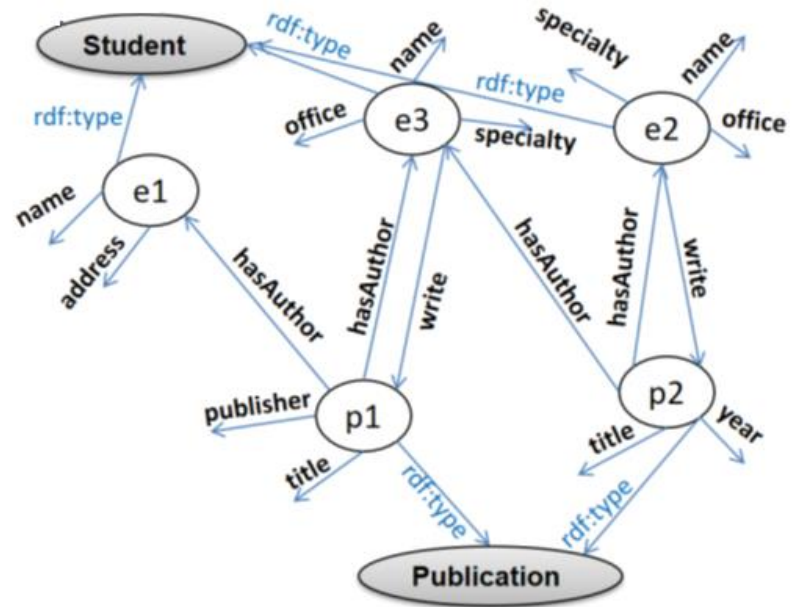
A set of properties  $V$  describing instances of a type.

e.g:  $V = \{\vec{\text{name}}, \vec{\text{address}}, \vec{\text{hasAuthor}}\}$

## Pattern Discovery Approaches:

Analyzing the co-occurrence relationships among the properties in the dataset.

- Identifying the possible set of properties that could describe an instance of a type.



# Terminology

## Pattern:

A set of properties  $V$  describing instances of a type:

$$- V = \{p_1, p_2, \dots, p_n\}$$

## Exact Pattern:

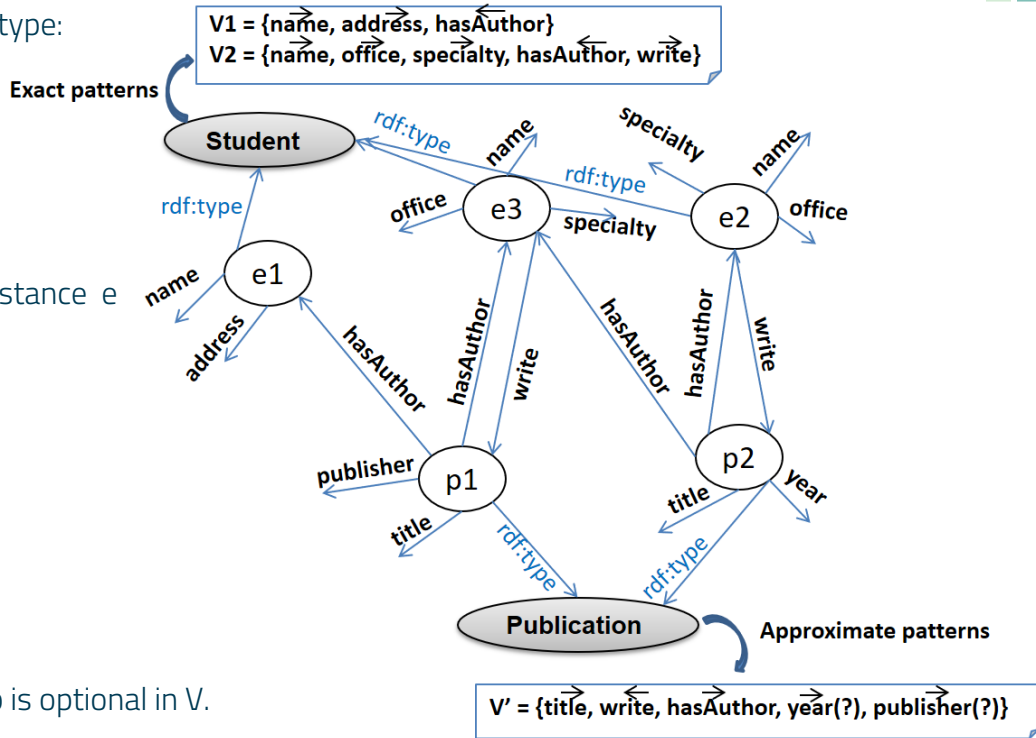
$V$  contains all the properties  $p$  describing an instance  $e$

- $\forall p \in V : p$  describes  $e$ ;
- $\forall p$  describing  $e : p \in V$ .

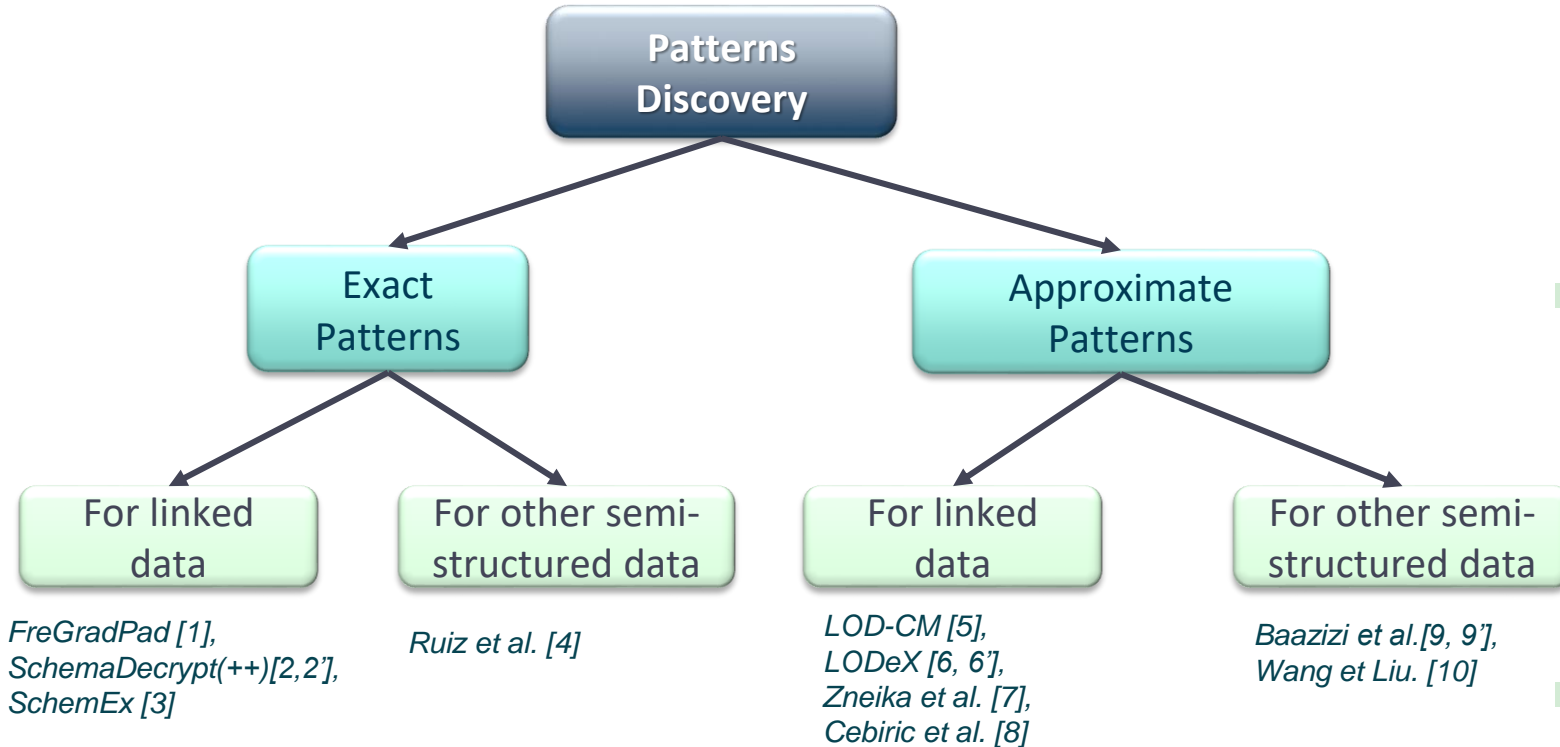
## Approximate Pattern:

Describe a set of similar instances  $T$

- $\forall p$  describing all the instances in  $T : p \in V$  ;
- If  $p$  describe only some instances in  $T$ , then  $p$  is optional in  $V$ .



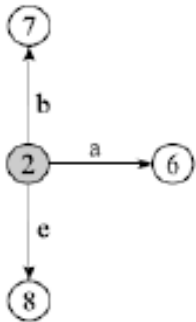
# Pattern discovery approaches classification



# Exact Pattern Discovery Approaches

## For linked data:

Approach	Inputs		Techniques	Outputs	Pros & Cons
	Data Description	Setting			
Belghaouti et al. <i>FreGraPad</i> (2016) [1]	<b>RDF data:</b> stream, outgoing properties	No required setting	Statistics by browsing and counting instances	Exact patterns with the frequency of each one	+ Incremental



**PHT**

Predicate	Ind.
a*	0
b*	1
c	2
d	3
e*	4

4 3 2 1 0

1	0	0	1	1
---	---	---	---	---

(19)

**Binary vector**

**GHT**

Graph	Freq.
.	.
.	.
.	.
19	1

**Graph n: set of triples**

# Exact Pattern Discovery Approaches

## For linked data:

Approach	Inputs		Techniques	Outputs	Pros & Cons
	Data Description	Setting			
Kellou-Menouer et Kedad <i>SchemaDecrypt++</i> (2017 - 2020) [2, 2']	RDF(S) data: remote source with restricted access, incoming and outgoing properties, part of schema ( <i>rdf:type</i> )	No required setting	Statistics on the properties of the instances. The set of properties could be automatically extracted relying on RDFS statements.	Exact patterns with the frequency of each one	+ On-line approach + Parallelized + Scalability can be by individual source restrictions



On-line access  
(through SPARQL endPoint)



Remote data sources  
(e.g: DBpedia)

Source restrictions:

- limited size of answer
- limited time to answer a query
- limited number of queries

[2] Kellou-Menouer, K., Kedad, Z.: On-line versioned schema inference for large semantic web data sources. In: Proceedings of the 29th International Conference on Scientific and Statistical Database Management, SSDBM 2017.

[2'] Kellou-Menouer, K., Kedad, Z.: SchemaDecrypt++: Parallel on-line versioned schema inference for large semantic web data sources. Information Systems Journal 2020.



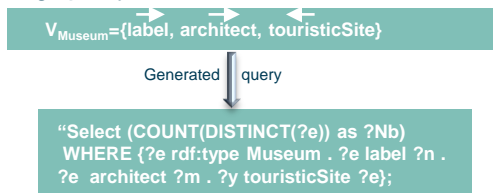
# Exact Pattern Discovery Approaches

## For linked data:

Approach	Inputs		Techniques	Outputs	Pros & Cons
	Data Description	Setting			
Kellou-Menouer et Kedad <i>SchemaDecrypt(++)</i> (2017 - 2020) [2, 2']	RDF(S) data: remote source with restricted access, incoming and outgoing properties, part of schema ( <i>rdf:type</i> )	No required setting	Statistics on the properties of the instances. The set of properties could be automatically extracted relying on RDFS statements.	Exact patterns with the frequency of each one	+ On-line approach + Parallelized + Scalability can be by individual source restrictions

### Basic Idea:

- The candidate versions for a class: all the combinations of its properties
- Each combination / candidate version is tested by sending the corresponding query to the source



- If the number of instances is positive, the candidate version is validated

### A combinatorial problem:

- $2^n$  queries to execute for a class with  $n$  properties
  - Average number of properties in DBpedia is 150:
    - $2^{150}$  queries to send
    - 15 queries by second<sup>(1)</sup>
- =>  $10^{36}$  years<sup>(2)</sup>

<sup>(1)</sup> The DBpedia online server can test a maximum of 15 queries per second.

<sup>(2)</sup>  $2^{150}/15$  seconds  $\approx 2^{146}$  seconds, knowing that 1 year  $\approx 31\,536\,000$  seconds  $\approx 2^{25}$  seconds,  $2^{146}$  seconds  $\approx 2^{121}$  years  $\approx 10^{36}$  years.

# Exact Pattern Discovery Approaches

## For linked data:

Approach	Inputs		Techniques	Outputs	Pros & Cons
	Data Description	Setting			
Kellou-Menouer et Kedad <i>SchemaDecrypt++</i> (2017 - 2020) [2, 2']	RDF(S) data: remote source with restricted access, incoming and outgoing properties, part of schema ( <i>rdf:type</i> )	No required setting	Statistics on the properties of the instances. The set of properties could be automatically extracted relying on RDFS statements.	Exact patterns with the frequency of each one	+ On-line approach + Parallelized + Scalability can be by individual source restrictions

## Strategy:

- Probabilistic profile: A property vector containing the properties of the instances of the class and their probabilities
  - Testing the most probable patterns at first
  - Reducing the search space
  - Defining a stopping criteria
- Properties co-occurrence rules
  - Pruning the search space
- Parallelization of candidate patterns exploration (SchemaDecrypt++)

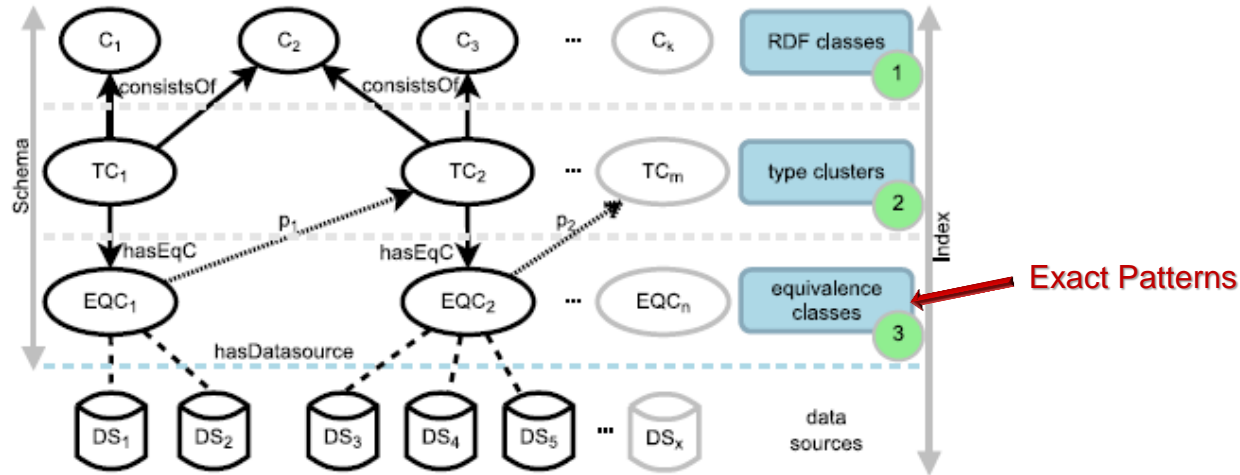
[2] Kellou-Menouer, K., Kedad, Z.: On-line versioned schema inference for large semantic web data sources. In: Proceedings of the 29th International Conference on Scientific and Statistical Database Management, SSDBM 2017.

[2'] Kellou-Menouer, K., Kedad, Z.: SchemaDecrypt++: Parallel on-line versioned schema inference for large semantic web data sources. Information Systems Journal 2020.

# Exact Pattern Discovery Approaches

## For linked data:

Approach	Inputs		Techniques	Outputs	Pros & Cons
	Data Description	Setting			
Konrath et al. <i>SchemEx</i> (2012) [3]	RDF data: distributed sources with unrestricted access, outgoing properties, part of schema ( <i>rdf:type</i> )	No required setting	Formal Method by browsing the data to build three layers	Exact patterns and an index with Three layers	+ Can be applied to RDF stream data + Web scale



# Exact Pattern Discovery Approaches

## For other semi-structured data:

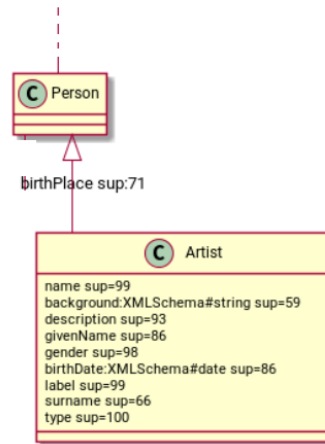
Approach	Inputs		Techniques	Outputs	Pros & Cons
	Data Description	Setting			
Ruiz et al. (2015) [4]	JSON data: outgoing properties, the kind of the object	No required setting	Formal Method by dividing the objects on the basis of the kind (type) with <i>Map-Reduce</i>	Exact patterns expressed in JSON	+ Provide tools for schema viewing and validation + Scalable

```
Entity Journal {  
  Version 1 {  
    issn: Tuple [String, String]  
    name: String  
    discipline: String  
  }  
  Version 2 {  
    issn: Tuple [String, String]  
    name: String  
    discipline: String  
    number: int  
  }  
}
```

# Approximate Pattern Discovery Approaches

## For linked data:

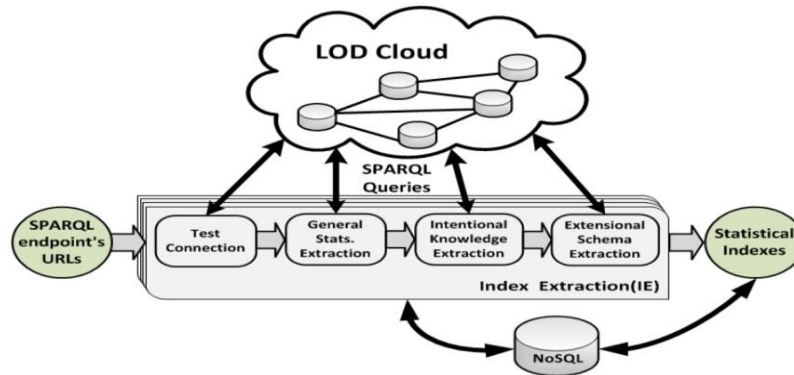
Approach	Inputs		Techniques	Outputs	Pros & Cons
	Data Description	Setting			
Issa S. et al. <i>LOD-CM</i> (2019) [5]	RDF data: outgoing properties, part of schema ( <i>rdf:type</i> )	Threshold for classifying a pattern as frequent	Statistics and Machine learning (FP-Growth) by mining maximal frequent patterns upon properties	Approximate patterns	+ Takes into account user quality expectations - Iterative process



# Approximate Pattern Discovery Approaches

## For linked data:

Approach	Inputs		Techniques	Outputs	Pros & Cons
	Data Description	Setting			
Benedetti et al. <i>LODeX</i> (2014-2015) [6, 6']	RDF data: <i>rdf:type</i> exploited if provided	No required setting	Statistics	Approximate patterns expressed on an RDF graph	+ Visual representation - Relies on online SPARQL endpoints with scalability restrictions



[6] Benedetti, F., Bergamaschi, S., Po, L.: Online index extraction from linked open data sources. Proceedings of the Second International Workshop on Linked Data for Information Extraction (LD4IE 2014) co-located with the 13th International Semantic Web Conference (ISWC 2014)

[6'] Benedetti, F., Bergamaschi, S., Po, L.: Exposing the underlying schema of LOD sources. In: IEEE/WIC/ACM International Conference on Web Intelligence and Intelligent Agent Technology, WI-IAT 2015, Singapore, December

# Approximate Pattern Discovery Approaches

## For linked data:

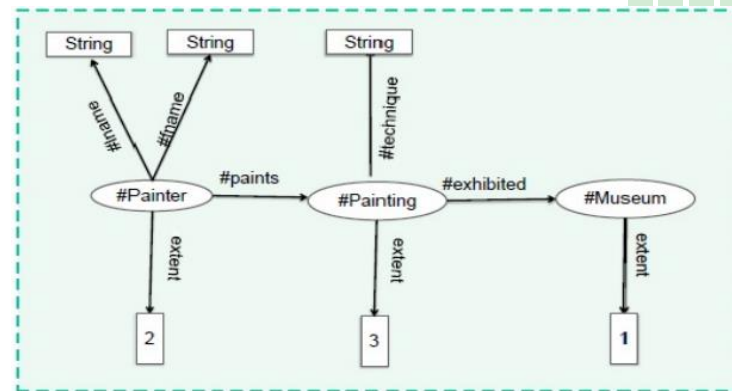
Approach	Inputs		Techniques	Outputs	Pros & Cons
	Data Description	Setting			
Zneika et al. (2016) [7]	RDF data: In/outgoing properties, part of schema ( <i>rdf:type</i> exploited if provided)	Settings for cost function	Machine learning using a pattern mining algorithm (PaNDA+)	Approximate patterns expressed on an RDF graph with the occurrence number of each pattern (using "extent" statement)	+ Detects the most representative approximate patterns - Requires dataset transformation into a binary matrix

	Painter(c)	Painting(c)	lname	fname	paints	exhibited	R_paints	R_exhibited	Museum(c)
Picasso	1	0	1	1	1	0	0	0	0
Rembrandt	1	0	1	1	1	0	0	0	0
Woman	0	1	0	0	0	0	1	0	0
Guernica	0	1	0	0	0	1	1	0	0
Abraham	0	1	0	0	0	1	1	0	0
museum.es	0	0	0	0	0	0	0	1	1

**Binary matrix**

ID	Pattern	Correspondence class
P1	Painting(c),exhibited, revers_paint	3
P2	Painter(c),paints, fname, lname	2
P3	Museum(c)	1

**Extracted patterns**

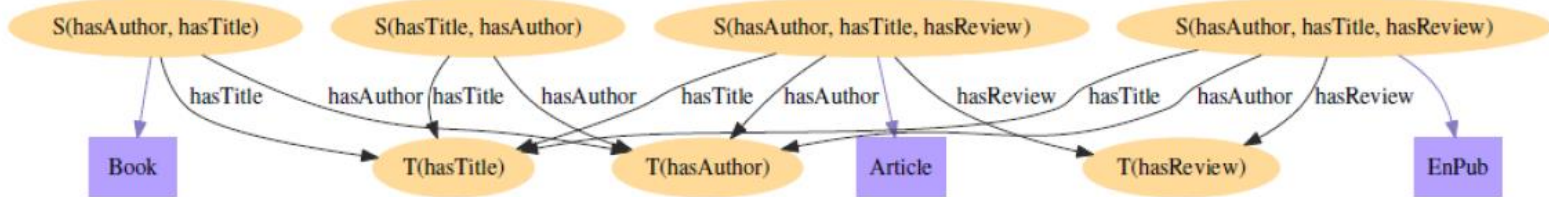


**RDF Summary graph for the set of patterns**

# Approximate Pattern Discovery Approaches

## For linked data:

Approach	Inputs		Techniques	Outputs	Pros & Cons
	Data Description	Setting			
Cebiric et al. (2015) [8]	RDFS data: outgoing properties, part of schema ( <i>rdf:type</i> , <i>rdf:domain</i> , <i>rdf:range</i> , <i>rdf:subClassOf</i> exploited if provided)	No required setting	Formal Methods by grouping properties according to the type of their subject in the triples, by considering <i>rdf:type</i> statements	Approximate patterns expressed on an RDF Query-oriented graph	<ul style="list-style-type: none"> <li>+ Considers both implicit and explicit triples</li> <li>- Computationally expensive <math>O( G ^5)</math></li> <li>- The optional properties are not marked in the produces schema</li> </ul>





# Approximate Pattern Discovery Approaches

## For other semi-structured data:

Approach	Inputs		Techniques	Outputs	Pros & Cons
	Data Description	Setting			
Baazizi et al. (2017-2019) [9, 9']	JSON data: outgoing properties, the kind of the object	Precision and conciseness level could be setting	Formal Method by a random division of data with <i>Map-Reduce</i>	Approximate patterns described by regular expressions	+ Spark-based scalable approach + Incremental

### Dataset:

1. Person { id : 12, name : John Smith, age : 14}
2. Office { id : 31, number : 3, address : 4 rue armengaud}
3. Person { id : 4, name : Kenza Kellou-Menouer}

### 1. Map: Structure inference

$T_1 = \{ \text{id} : \text{Number}, \text{name} : \text{String}, \text{age} : \text{Number} \}$

$T_2 = \{ \text{id} : \text{Number}, \text{number} : \text{Number}, \text{address} : \text{String} \}$

$T_3 = \{ \text{id} : \text{Number}, \text{name} : \text{String} \}$

### 2. Reduce: Type merging (Approximate Patterns)

$T_{1,3} = \{ \text{id} : \text{Number}, \text{name} : \text{String}, (\text{age} : \text{Number}) ? \}$

$T_2 = \{ \text{id} : \text{Number}, \text{number} : \text{Number}, \text{address} : \text{String} \}$

[9] Baazizi, M.A., Colazzo, D., Ghelli, G., Sartiani, C.: Parametric schema inference for massive JSON datasets. VLDB J. 28(4), 497–521 (2019)

[9'] Baazizi, M.A., Lahmar, H.B., Colazzo, D., Ghelli, G., Sartiani, C.: Schema inference for massive JSON datasets. In: Proceedings of the 20th International Conference on Extending Database Technology, EDBT 2017.



# Rewiew of Pattern Discovery Approaches

- A set of approaches proposed for:
  - Exact patterns discovery [1-4]
  - Approximate patterns discovery [5-10]
- In an approximate pattern, the co-occurrence of the properties is not clearly indicated:
  - Specifying the percentage of each property among the instances of a class [5]
  - Marking of the optional properties of a pattern [9, 9'], or not even [8, 10]
  - Including only the no optional properties in the pattern [7].
- Different data sources are processed:
  - Streaming data [1]
  - Remote data with constrained access [2, 2', 6, 6']
  - Distributed data [3]
  - Local data [4-7, 8-10]



Break.  
To be continued...

