

# Selection of representative subsets of link key candidates

Nacira Abbas, Alexandre Bazin, Jérôme David, Amedeo Napoli



Atelier Decade, le 28/06/2022, St Etienne

Partly funded by Elker ANR project (ANR-17-CE23-0007-01)  
Most of this work is from the PhD of Nacira Abbas (Loria, Nancy)

Data interlinking

Link keys

Link key candidates extraction (with FCA)

Link key candidate reduction and selection

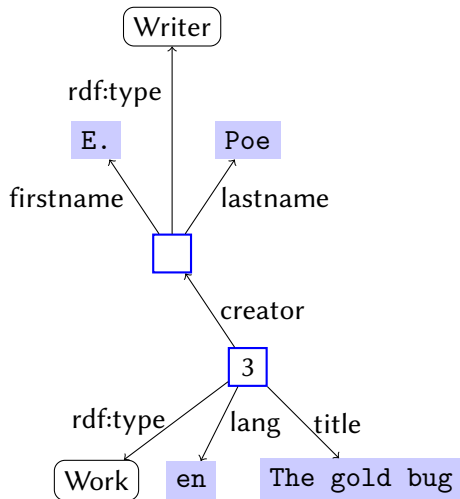
## Data interlinking

### Link keys

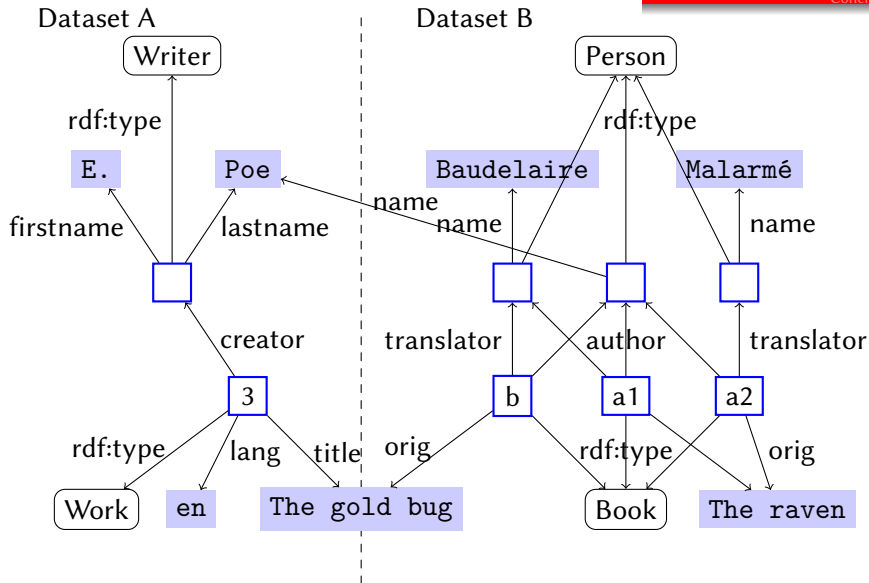
### Link key candidates extraction (with FCA)

### Link key candidate reduction and selection

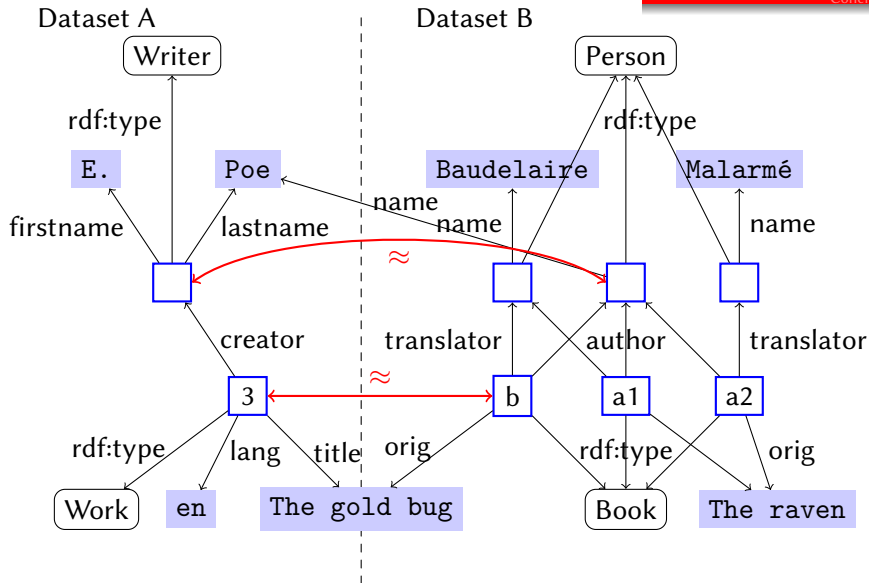
# The problem: RDF data interlinking



# The problem: RDF data interlinking



# The problem: RDF data interlinking



- ▶ Numerical specifications (Link Specifications)
  - ▶ Express or learn a similarity from RDF data
  - ▶ Generate links using frameworks such as SILK or LIMES
- ▶ NLP/IR based approaches
  - ▶ Change representation: from RDF space to VSM, or embedding spaces
  - ▶ Compute or learn a similarity on this new space
- ▶ Logical link specifications
  - ▶ Key-based: combine keys and alignments for deducing links
  - ▶ Link keys: cross dataset, generalization of keys without requiring alignment between properties or concepts

- ▶ Numerical specifications (Link Specifications)
  - ▶ Express or learn a similarity from RDF data
  - ▶ Generate links using frameworks such as SILK or LIMES
- ▶ NLP/IR based approaches
  - ▶ Change representation: from RDF space to VSM, or embedding spaces
  - ▶ Compute or learn a similarity on this new space
- ▶ Logical link specifications
  - ▶ Key-based: combine keys and alignments for deducing links
  - ▶ Link keys: cross dataset, generalization of keys without requiring alignment between properties or concepts



Data interlinking

Link keys

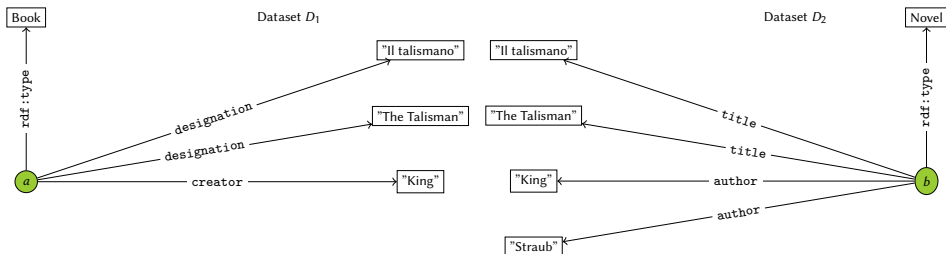
Link key candidates extraction (with FCA)

Link key candidate reduction and selection

# What is a link key?

two sets of property pairs and a pair of classes like

$$\underbrace{\{((\text{designation}, \text{title}))\}}_{\forall}, \underbrace{\{((\text{designation}, \text{title}), (\text{creator}, \text{author}))\}}_{\exists}, \underbrace{(\text{Book}, \text{Novel})}_{\text{Classes}}$$



On this example, the link key will generate the link  $(a, \text{owl:sameAs}, b)$

They may be several expressions having the form of link keys

$D$ (Employés)					$D'$ (Staff)				
id	prenom	datenaiss	poste	bât.	firstname	birthdate	position	building	id
$i_2$	Paul	1967	Dir.	B2	Paul		Dir.	B2	$z_2$
$i_3$	Mary	1963	Dir.	B1	Mary		Dir.	B1	$z_3$
$i_4$	John	1963	Pr.	B1	John		Pr.	B1	$z_4$
$i_6$	Bill	1980	Pr.	B1	William	1980	Pr.		$z_6$
$i_7$	Ana	1947	Dir.	B2	Ana	1947	Dir.		$z_7$
$i_8$	John	1967	Pr.	B2	John	1967	Pr.		$z_8$

Example of link key expressions:

- ▶  $k = \langle \{\}, \{\langle \text{datenaiss}, \text{birthdate} \rangle\}, \langle \text{Employee}, \text{Staff} \rangle \rangle$
- ▶  $h = \langle \{\langle \text{datenaiss}, \text{birthdate} \rangle\}, \{\langle \text{poste}, \text{position} \rangle\} \langle \text{Employee}, \text{Staff} \rangle \rangle$
- ▶  $l = \langle \{\langle \text{datenaiss}, \text{birthdate} \rangle, \langle \text{poste}, \text{position} \rangle\}, \{\langle \text{poste}, \text{position} \rangle\}, \langle \text{Employee}, \text{Staff} \rangle \rangle$

And generated links (if used as link keys):

- ▶  $L_k^{D,D'} = \{\langle i_7, z_7 \rangle, \langle i_8, z_8 \rangle, \langle i_6, z_6 \rangle, \langle i_2, z_8 \rangle\}$
- ▶  $L_l^{D,D'} = L_h^{D,D'} = \{\langle i_7, z_7 \rangle, \langle i_8, z_8 \rangle, \langle i_6, z_6 \rangle\}$

They may be several expressions having the form of link keys

$D$ (Employés)					$D'$ (Staff)				
id	prenom	datenaiss	poste	bât.	firstname	birthdate	position	building	id
$i_2$	Paul	1967	Dir.	B2	Paul		Dir.	B2	$z_2$
$i_3$	Mary	1963	Dir.	B1	Mary		Dir.	B1	$z_3$
$i_4$	John	1963	Pr.	B1	John		Pr.	B1	$z_4$
$i_6$	Bill	1980	Pr.	B1	William	1980	Pr.		$z_6$
$i_7$	Ana	1947	Dir.	B2	Ana	1947	Dir.		$z_7$
$i_8$	John	1967	Pr.	B2	John	1967	Pr.		$z_8$

Example of link key expressions:

- ▶  $k = \langle \{\}, \{\langle \text{datenaiss}, \text{birthdate} \rangle\}, \langle \text{Employe}, \text{Staff} \rangle \rangle$
- ▶  $h = \langle \{\langle \text{datenaiss}, \text{birthdate} \rangle\}, \{\langle \text{poste}, \text{position} \rangle\} \langle \text{Employe}, \text{Staff} \rangle \rangle$
- ▶  $l = \langle \{\langle \text{datenaiss}, \text{birthdate} \rangle, \langle \text{poste}, \text{position} \rangle\}, \{\langle \text{poste}, \text{position} \rangle\}, \langle \text{Employe}, \text{Staff} \rangle \rangle$

And generated links (if used as link keys):

- ▶  $L_k^{D,D'} = \{\langle i_7, z_7 \rangle, \langle i_8, z_8 \rangle, \langle i_6, z_6 \rangle, \langle i_2, z_8 \rangle\}$
- ▶  $L_l^{D,D'} = L_h^{D,D'} = \{\langle i_7, z_7 \rangle, \langle i_8, z_8 \rangle, \langle i_6, z_6 \rangle\}$

They may be several expressions having the form of link keys

$D$ (Employés)					$D'$ (Staff)				
id	prenom	datenaiss	poste	bât.	firstname	birthdate	position	building	id
$i_2$	Paul	1967	Dir.	B2	Paul		Dir.	B2	$z_2$
$i_3$	Mary	1963	Dir.	B1	Mary		Dir.	B1	$z_3$
$i_4$	John	1963	Pr.	B1	John		Pr.	B1	$z_4$
$i_6$	Bill	1980	Pr.	B1	William	1980	Pr.		$z_6$
$i_7$	Ana	1947	Dir.	B2	Ana	1947	Dir.		$z_7$
$i_8$	John	1967	Pr.	B2	John	1967	Pr.		$z_8$

Example of link key expressions:

- ▶  $k = \langle \{\}, \{\langle \text{datenaiss}, \text{birthdate} \rangle\}, \langle \text{Employe}, \text{Staff} \rangle \rangle$
- ▶  $h = \langle \{\langle \text{datenaiss}, \text{birthdate} \rangle\}, \{\langle \text{poste}, \text{position} \rangle\} \langle \text{Employe}, \text{Staff} \rangle \rangle$
- ▶  $l = \langle \{\langle \text{datenaiss}, \text{birthdate} \rangle, \langle \text{poste}, \text{position} \rangle\}, \{\langle \text{poste}, \text{position} \rangle\}, \langle \text{Employe}, \text{Staff} \rangle \rangle$

And generated links (if used as link keys):

- ▶  $L_k^{D,D'} = \{\langle i_7, z_7 \rangle, \langle i_8, z_8 \rangle, \langle i_6, z_6 \rangle, \langle i_2, z_8 \rangle\}$
- ▶  $L_l^{D,D'} = L_h^{D,D'} = \{\langle i_7, z_7 \rangle, \langle i_8, z_8 \rangle, \langle i_6, z_6 \rangle\}$

Data interlinking

Link keys

Link key candidates extraction (with FCA)

Link key candidate reduction and selection

## Problem: How to induce link keys from data ?

The number of set of pairs of properties is exponential

Our approach:

- ▶ restrict on link key expressions that would generate at least one link between the datasets
- ▶ consider only closed expressions : those which are maximal for a set of links

we call such expressions “link key candidates” (LKC) and we extract them using Formal Concept Analysis

# Formal context for candidate link key extraction

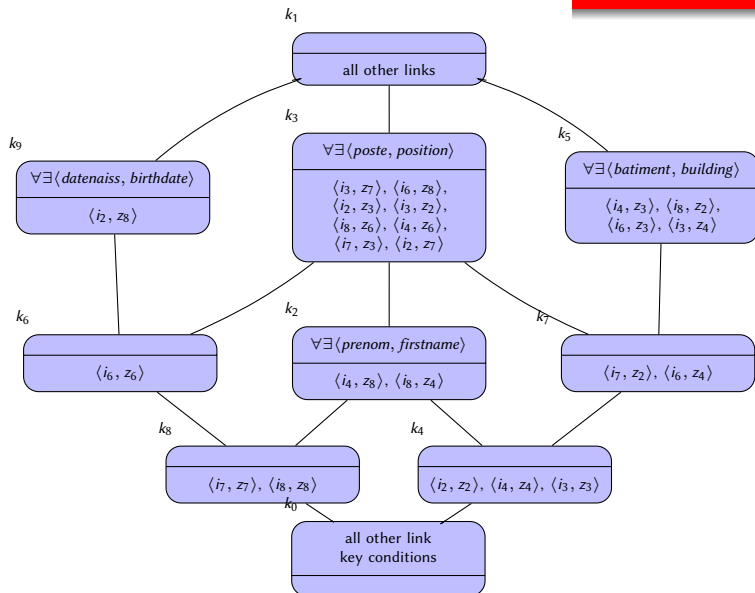
The formal context for link key candidates  $\langle G, M, I \rangle$  is:

$G \backslash M$	...	$\exists \langle p_i, p'_j \rangle$	...	...	$\forall \langle p_i, p'_j \rangle$	...
$\vdots$	$\ddots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$	$\ddots$
$\langle o, o' \rangle$	...	1 iff $p^D(o) \cap p'^{D'}(o') \neq \emptyset$	...	...	1 iff $p^D(o) = p'^{D'}(o')$	...
$\vdots$	$\ddots$	$\vdots$	$\ddots$	$\vdots$	$\vdots$	$\ddots$

- ▶  $G$ : the set of pairs of objects from each dataset
- ▶  $M$ : two sets of pairs of properties from each dataset
  - ▶  $\exists$ : if the objects share at least one value ( $\langle o, o' \rangle I \exists \langle p_i, p'_j \rangle$ )
  - ▶  $\forall$ : if the object have the same values ( $\langle o, o' \rangle I \forall \langle p_i, p'_j \rangle$ )



# Lattice of extracted link key candidates



FCA link key extraction is implemented in Linkex  
(<https://gitlab.inria.fr/moex/linkex>)

## Characteristics/functionalities:

- ▶ Fully unsupervised: only two RDF datasets as input.
- ▶ Normalization of textual content: lowercase, remove diacritics, tokenization, sort
- ▶ Can compute inverse and composition of properties
- ▶ Compute the class expression of concepts (i.e. covering instances)
- ▶ Different output formats: Alignment Format, GraphViz (dot), tabular file, etc.
- ▶ Implementation of diverse quality measures (discriminability, coverage, etc)

Data interlinking

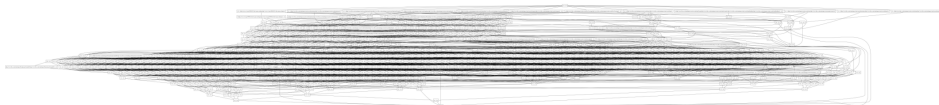
Link keys

Link key candidates extraction (with FCA)

Link key candidate reduction and selection

Many link key candidates can be extracted!

ex: OAEI Spimbench task, > 2k candidates



How to reduce the lattice and select the interesting/representative ones?

## 1. Reduce

- ▶ Identify redundant LKC according to owl:sameAs semantics
- ▶ Representative link key candidates based on clustering

## 2. Evaluate

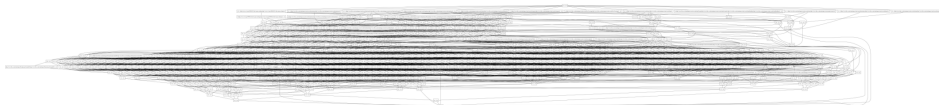
- ▶ Quality measures: discriminability and coverage

## 3. Combine: Disjunctions based on antichains

- ▶ Explore the antichains of the LKC lattice

Many link key candidates can be extracted!

ex: OAEI Spimbench task, > 2k candidates



How to reduce the lattice and select the interesting/representative ones?

## 1. Reduce

- ▶ Identify redundant LKC according to owl:sameAs semantics
- ▶ Representative link key candidates based on clustering

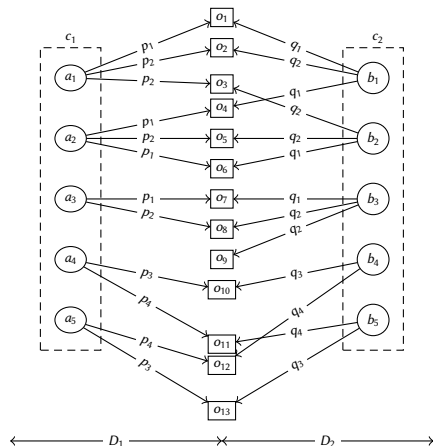
## 2. Evaluate

- ▶ Quality measures: discriminability and coverage

## 3. Combine: Disjunctions based on antichains

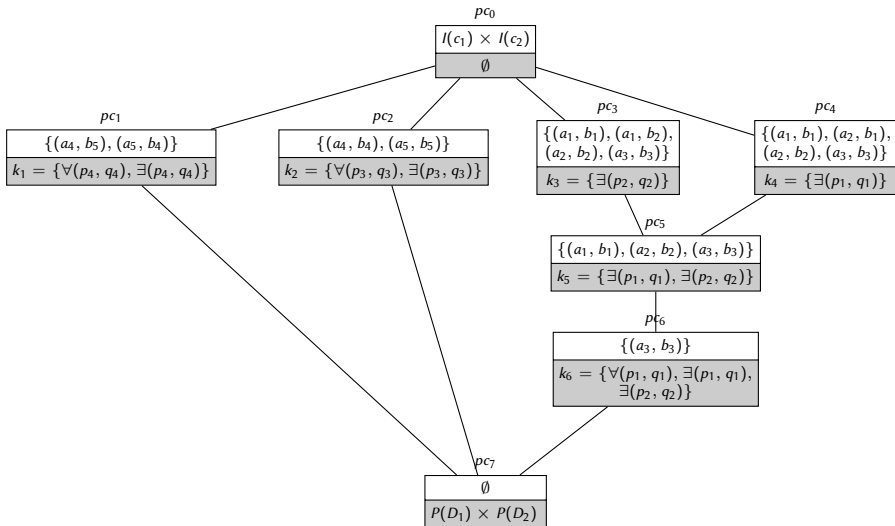
- ▶ Explore the antichains of the LKC lattice

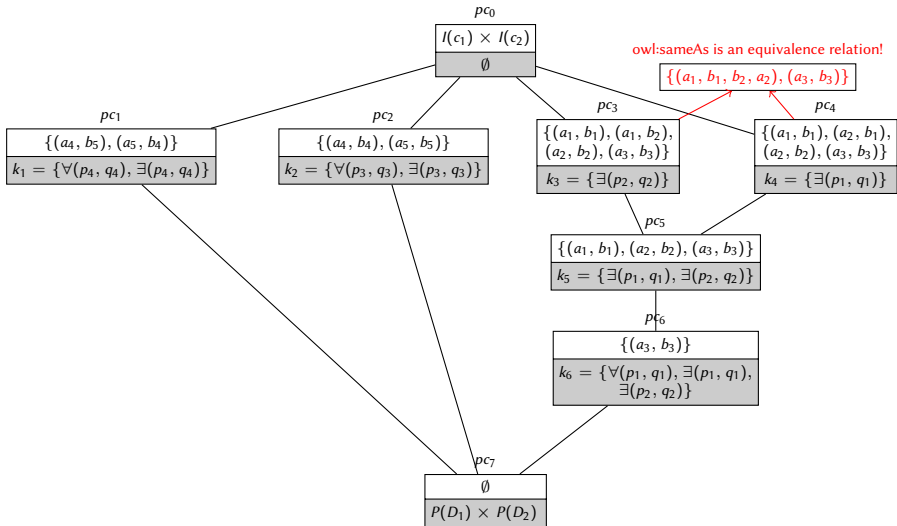
# Pattern structure for link key candidate discovery



Pattern structure:  $(G, (E, \sqcap), \delta)$

PS objects $G$	Descriptions ( $\delta$ ) $E$
$(a_1, b_1)$	$\{\exists(p_1, q_1), \exists(p_2, q_2)\}$
$(a_1, b_2)$	$\{\exists(p_2, q_2)\}$
$(a_2, b_1)$	$\{\exists(p_1, q_1)\}$
$(a_2, b_2)$	$\{\exists(p_1, q_1), \exists(p_2, q_2)\}$
$(a_3, b_3)$	$\{\forall(p_1, q_1), \exists(p_1, q_1), \exists(p_2, q_2)\}$
$(a_4, b_4)$	$\{\forall(p_3, q_3), \exists(p_3, q_3)\}$
$(a_4, b_5)$	$\{\forall(p_4, q_4), \exists(p_4, q_4)\}$
$(a_5, b_4)$	$\{\forall(p_4, q_4), \exists(p_4, q_4)\}$
$(a_5, b_5)$	$\{\forall(p_3, q_3), \exists(p_3, q_3)\}$

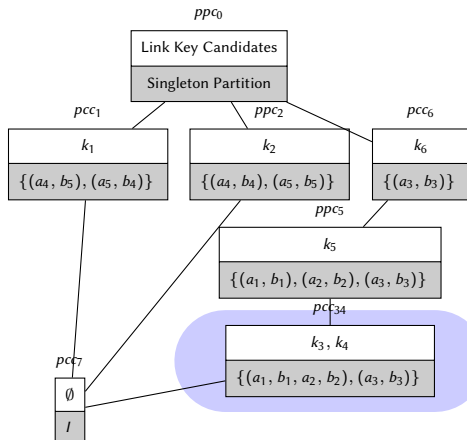






Idea: use Partition pattern structures to detect redundant LKC

Objects <i>LKC</i>	Descriptions <i>partitions induced by owl:sameAs</i>
$k_1$	$\{(a_4, b_5), (a_5, b_4)\}$
$k_2$	$\{(a_4, b_4), (a_5, b_5)\}$
$k_3$	$\{(a_1, b_1, a_2, b_2), (a_3, b_3)\}$
$k_4$	$\{(a_1, b_1, a_2, b_2), (a_3, b_3)\}$
$k_5$	$\{(a_1, b_1), (a_2, b_2), (a_3, b_3)\}$
$k_6$	$\{(a_3, b_3)\}$



Link key candidates having the same partition are in the same concept

Nacira Abbas, Alexandre Bazin, Jérôme David, Amedeo Napoli: A Study of the Discovery and Redundancy of Link Keys Between Two RDF Datasets Based on Partition Pattern Structures. CLA 2022: to appear

This works but ... this is not so useful in practice :- (

Interlinking task	datasets	#triple	#subj	#prop	#LKC	#NRLKC
Actor	db:Actor	94 606	5 807	16	2 198	<b>2 177 (↓ 1%)</b>
	yago:Actor	1 029 580	108 415	16		
Album	db:Album	594 144	85 002	5	44	44
	yago:Album	762 238	136 848	5		
Book	db:Book	247 372	29 846	7	82	82
	yago:Book	185 032	41 849	7		
Film	db:Film	1 369 600	82 099	9	18 718	<b>17 643 (↓ 5%)</b>
	yago:Film	1 067 084	123 822	9		
Mountain	db:Mountain	135 442	16 397	5	39	39
	yago:Mountain	233 562	32 874	5		
Museum	db:Museum	15 940	1 826	7	48	48
	yago:Museum	163 342	21 050	7		
Organization	db:Organization	4 487 205	183 665	17	1 425	1 425
	yago:Organization	4 410 854	430 071	17		
Scientist	db:Scientist	128 360	18 409	10	862	862
	yago:Scientist	671 266	92 828	18		
University	db:University	241 838	10 352	9	213	213
	yago:University	263 624	23 334	9		

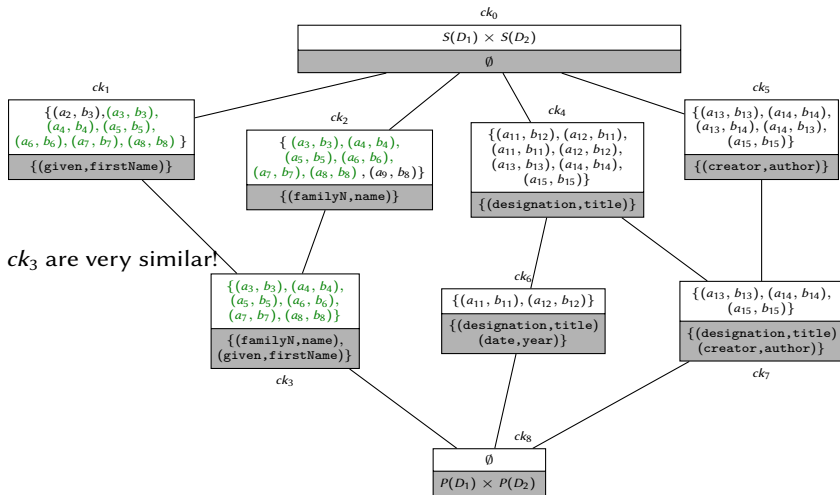
#LKC: number of link key candidates. #NRLKC: number of "non redundant" link key candidates.

## Datasets from

Danai Symeonidou, Luis Galárraga, Nathalie Pernelle, Fatiha Saïs, Fabian M. Suchanek: *VICKEY: Mining Conditional Keys on Knowledge Bases*. *ISWC (1) 2017: 661-677*

# Similarity between link keys

However, a lot of link key candidates generate almost the same links...

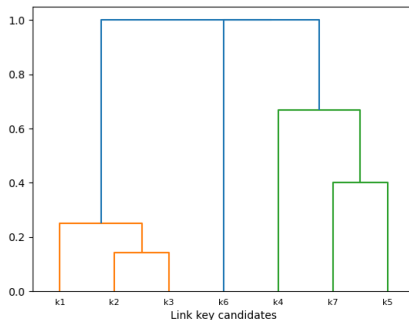
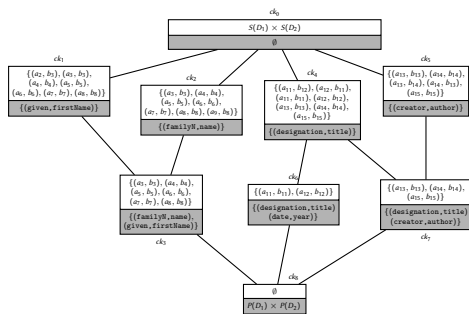


$ck_1, ck_2, ck_3$  are very similar!

We propose to select a subset of representative candidates.

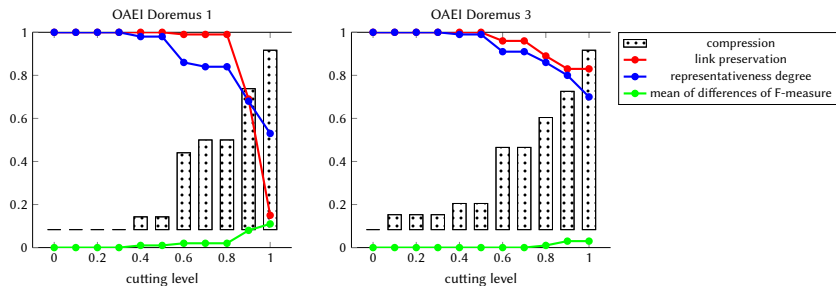
The procedure is as follows:

1. similarity is computed thanks to Jaccard index (between sets of links)
2. LKC are clustered with hierarchical agglomerative clustering
3. clusters are extracted from the resulting hierarchy
4. a representative LK is selected for each cluster
  - ▶ the LKC that minimizes the distance to the other



Cut at 0.5:

- ▶ 4 clusters:  $\{k_1, k_2, k_3\}$ ,  $\{k_4\}$  and  $\{k_6\}$
- ▶ representatives (core):  $\{k_3, k_4, k_6, k_7\}$
- ▶ This gives a compression ratio of 0.43 while preserving 77% of links.



- ▶ Almost 50% of LKC can be removed without loss
- ▶ Core LKC are good representative (and preserve F-Measure).

## The context:

- ▶ Link keys are symbolic and meaningful tool for interlinking data
- ▶ Fully unsupervised extraction from data with minimal input
- ▶ But... a lot of candidates can be discovered

## Contributions:

- ▶ A pattern structure for non redundant link key w.r.t. owl:sameAs
- ▶ A clustering based strategy for selecting representative subset of LKC

## Perspectives:

- ▶ Combine this approach with selecting disjunctions of LKC
- ▶ Generalize this clustering approach to FCA lattice reduction
  - ▶ Goal: reduce the # of concepts while preserving the order
  - ▶ adapt the similarity measure to FCA lattices
  - $\frac{|L(a \wedge b)|}{|L(a \vee b)|}$  instead of  $\frac{|L(a) \cap L(b)|}{|L(a) \cup L(b)|}$
  - ▶ design an optimised algorithm

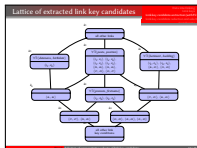
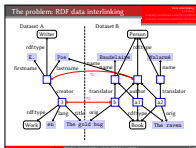
### Selection of representative subsets of link key candidates

Nicolas Abbat, Alexandre Bazin, Jérôme David, Arnaud Nagadi



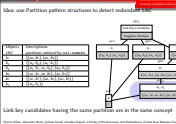
Abbat Nicolas, le 28/06/2022, 51 Etienne

Partly funded by Elkar ANR project (ANR-17-CE23-0007-01)  
Most of this work is from the PhD of Nicolas Abbat (Loria, Nancy)



### Partition pattern Structures lattice

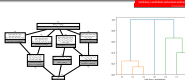
Idea: use Partition pattern structures to detect redundant link keys



Pattern	Number of link keys
all properties	1
all properties except rdfs:type	1
all properties except rdfs:type and owl:equivalentClass	1
all properties except rdfs:type and owl:equivalentProperty	1
all properties except rdfs:type, owl:equivalentClass, owl:equivalentProperty and owl:equivalentClass	1
all properties except rdfs:type, owl:equivalentClass, owl:equivalentProperty and owl:equivalentProperty	1
all properties except rdfs:type, owl:equivalentClass, owl:equivalentProperty and owl:equivalentProperty and owl:equivalentClass	1
all properties except rdfs:type, owl:equivalentClass, owl:equivalentProperty and owl:equivalentProperty and owl:equivalentProperty	1

Link key candidates having the same partition are in the same concept

### LKC clustering



Cut at 0.5:

- 4 clusters:  $\{k_1, k_2, k_3, k_4\}$  and  $\{k_5\}$
- representatives (rows):  $\{k_1, k_2, k_3, k_4, k_5\}$
- This gives a compression ratio of 0.43 while preserving 77% of links.



<https://moex.inria.fr>

Nacira . Abbas

Jerome . David @ inria . fr

Amedeo . Napoli @ loria . fr