

HOLISTIC STATISTICAL OPEN DATA INTEGRATION BASED ON INTEGER LINEAR PROGRAMMING

Alain Berro, Imen Megdiche, Olivier Teste

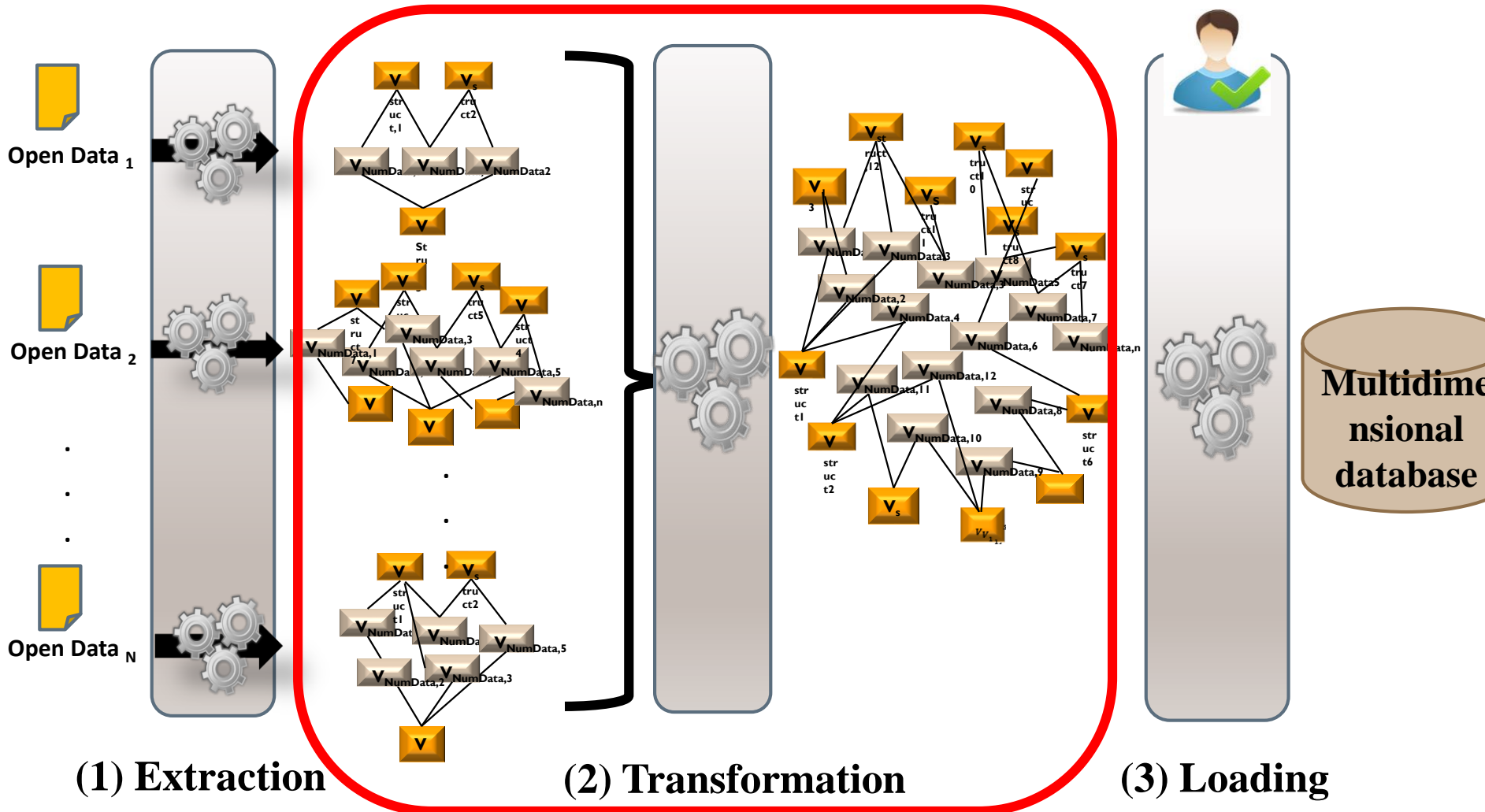


Outline

2

- 1. Objective
- 2. State of the art
- 3. A linear Program For Holistic Matching
- 4. Evaluation
 - ▣ Matching quality
 - ▣ Matching performance
- 5. Conclusion and perspectives

1. Objective



(1) Extraction

(2) Transformation

(3) Loading

(Berro et al., 2014)

1. Schema matching problem

4

SCHEMA 1

Library

Item

Isbn

Author

Title

Year

Author

First name

Last name

Borrowed Items

Item

Borrower

Borrower

First name

Last name

SCHEMA 2

Collection

Document

Identifier

Creator

Contributor

Publisher

Title

Year

Creator

Name

Name

First

Last

Publisher

Address

Name

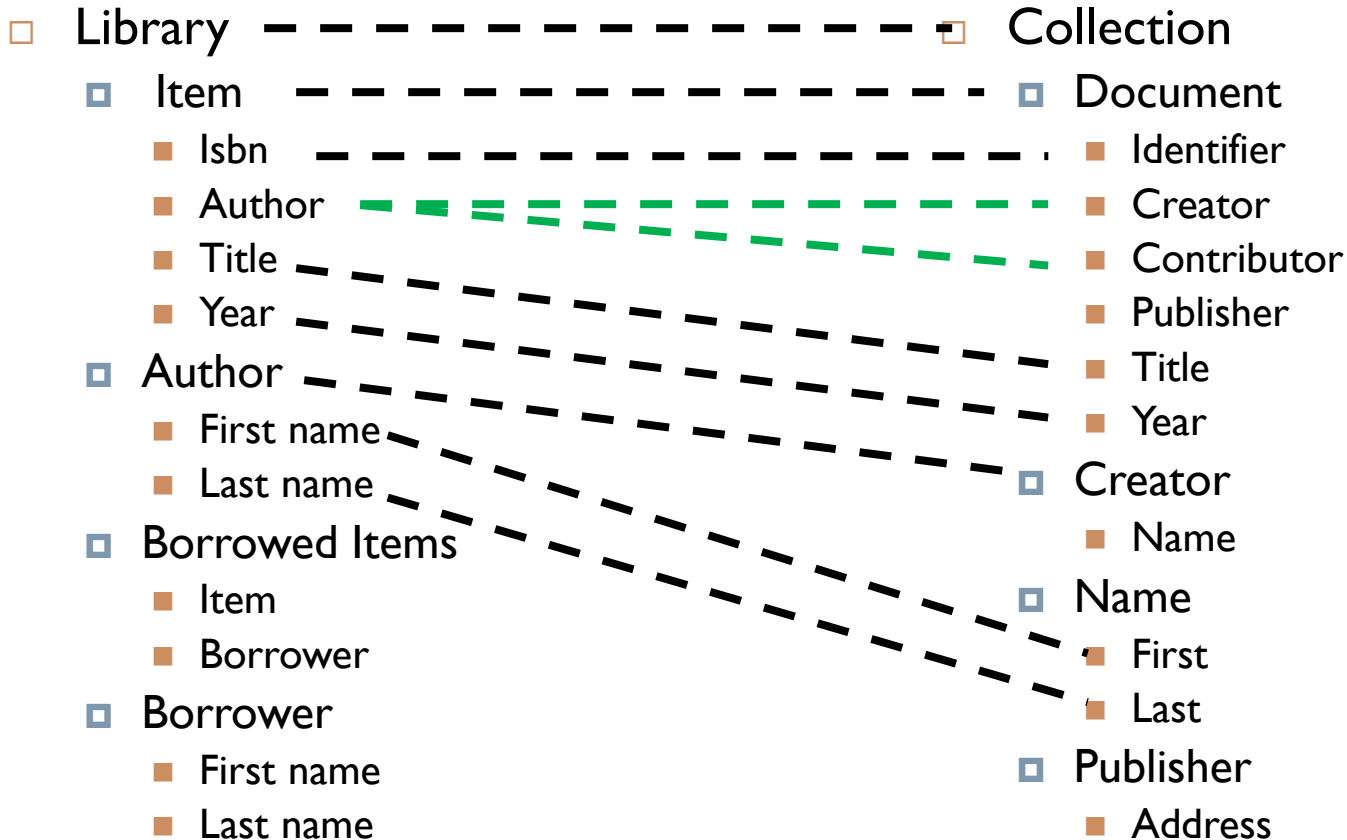


User 1 Solution

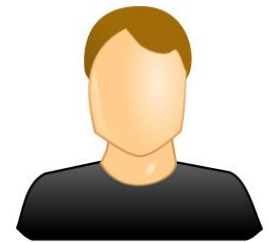
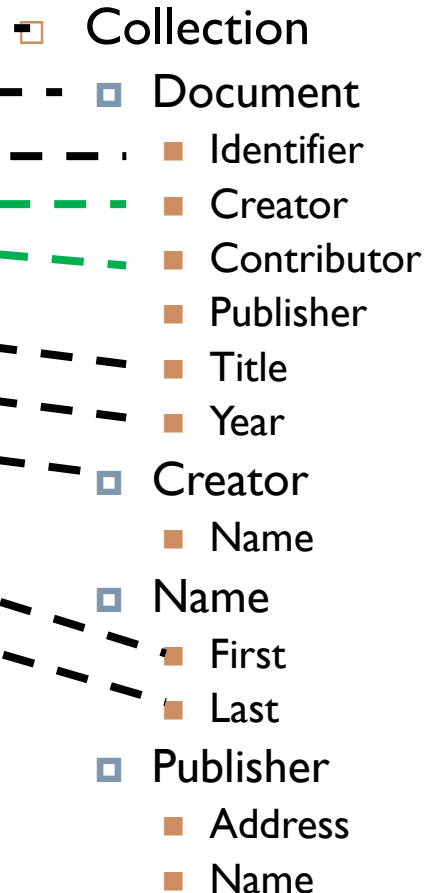
1. Schema matching problem

5

SCHEMA 1



SCHEMA 2



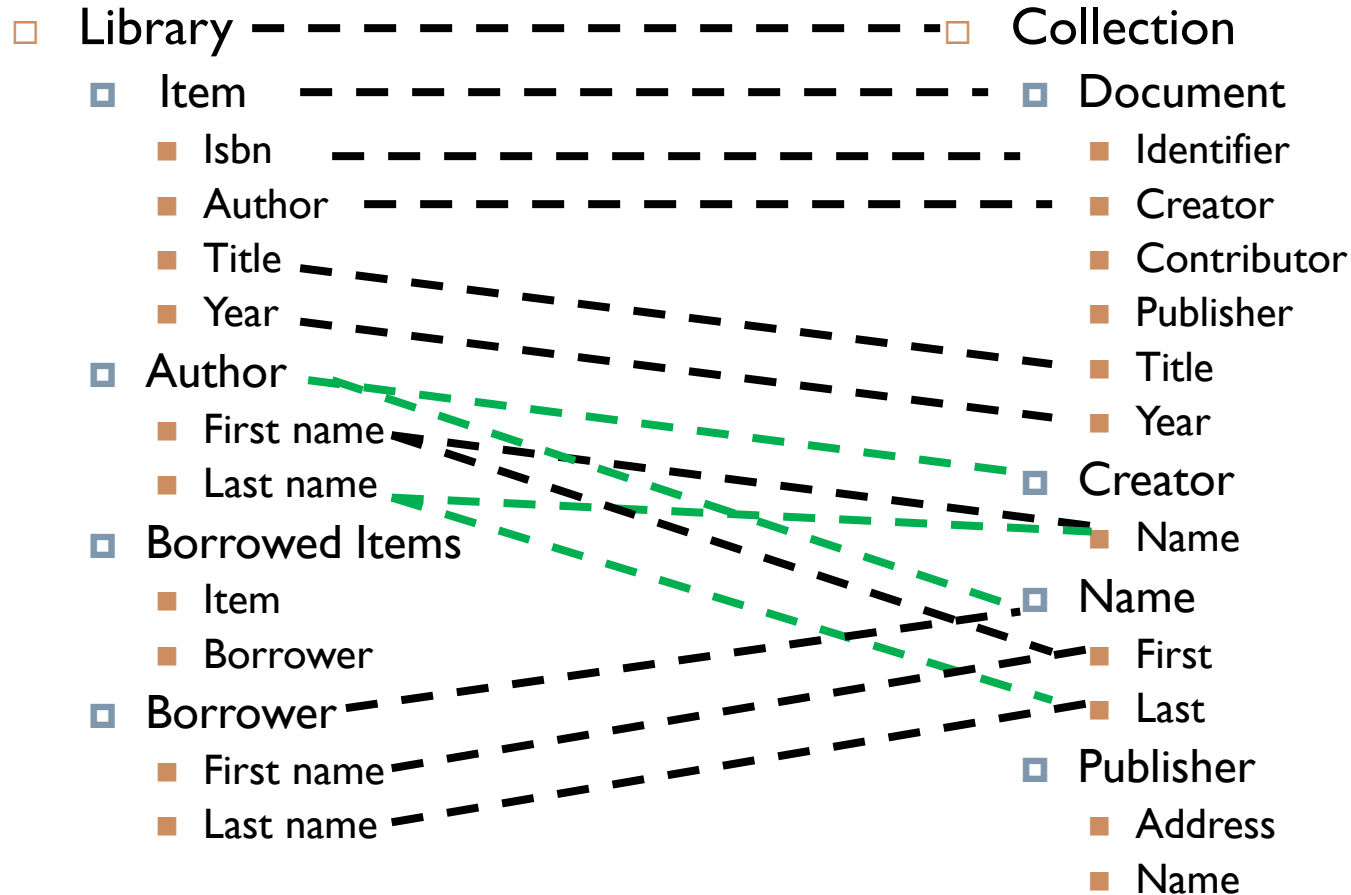
User 2 Solution

1. Schema matching problem

6

SCHEMA 1

SCHEMA 2

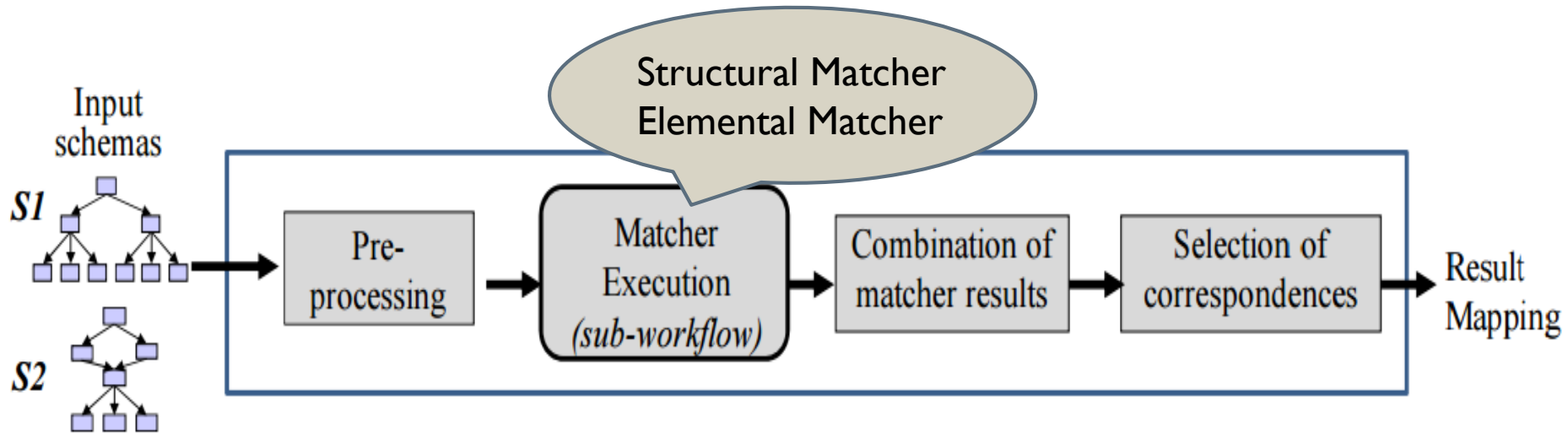


User 3 Solution

2. State of the art

Schema Matching Problem

7



General workflow for matching
(Bellahsene et al. 2011)

2. State of the art

Schema Matching Problem

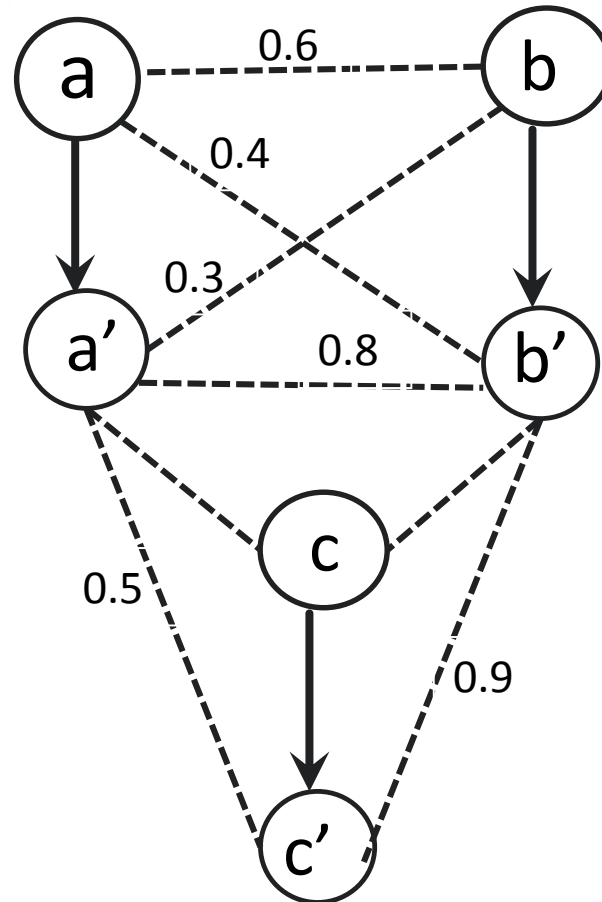
8

Approach	Type	Internal Model	Combinatorial Optimisation Reduction			Dependent to threshold
			Problem	Complexity	Used in	
COMA++	Pairwise	Directed acyclic graph	-	-	-	Yes
Similarity Flooding	Pairwise	Labelled graphs	Stable Marriage	Polynomial	Selection	Yes
BMatch	Pairwise	Tree	-	-	-	Yes
CODI	Pairwise	Labelled graphs	Max-Sat	NP-Hard	Structural matcher & Selection	Yes
OLA	Pairwise	Labelled graphs	Maximum weighted graph matching	Polynomial	Selection	Yes
DCM	Holistic	List of attributes	-	-	-	Yes
PORSCHE	Holistic	Tree	-	-	-	No
PLASMA	Holistic	Tree	-	-	-	Yes
LP4HM	Holistic	Directed acyclic graph	Maximum weighted graph matching	Polynomial	Structural matcher & Selection	No

3. A Linear Program For Holistic Matching (LP4HM)

9

Global optimal solution



Schema Matching
Find the best set of correspondences having 1:1 cardinality



Structural constraints

3. A Linear Program For Holistic Matching (LP4HM)

10

$$\max \sum_{i=1}^{N-1} \sum_{j=i+1}^N \sum_{k=1}^{n_i} \sum_{l=1}^{n_j} sim_{i_k, j_l} x_{i_k, j_l}$$

Objective function

$$s.t. \sum_{l=1}^{n_j} x_{i_k, j_l} \leq 1, \forall k \in [1, n_i] \\ \forall i \in [1, N-1] \forall j \in [i+1, N]$$

$$sim_{i_k, j_l} x_{i_k, j_l} \geq seuil x_{i_k, j_l} \\ \forall i \in [1, N-1] \forall j \in [i+1, N] \\ \forall k \in [1, n_i], \forall l \in [1, n_j]$$

Linear constraints

$$x_{i_k, j_l} \leq x_{i_{pred(k)}, j_{pred(l)}} \\ \forall i \in [1, N-1] \forall j \in [i+1, N] \\ \forall k \in [1, n_i], \forall l \in [1, n_j]$$

$$x_{i_k, j_l} + x_{i_{k'}, j_{l'}} - (dir_{i_k, k'} dir_{j_l, l'}) \leq 1 \\ \forall i \in [1, N-1] \forall j \in [i+1, N] \\ \forall k, k' \in [1, n_i], \forall l, l' \in [1, n_j]$$

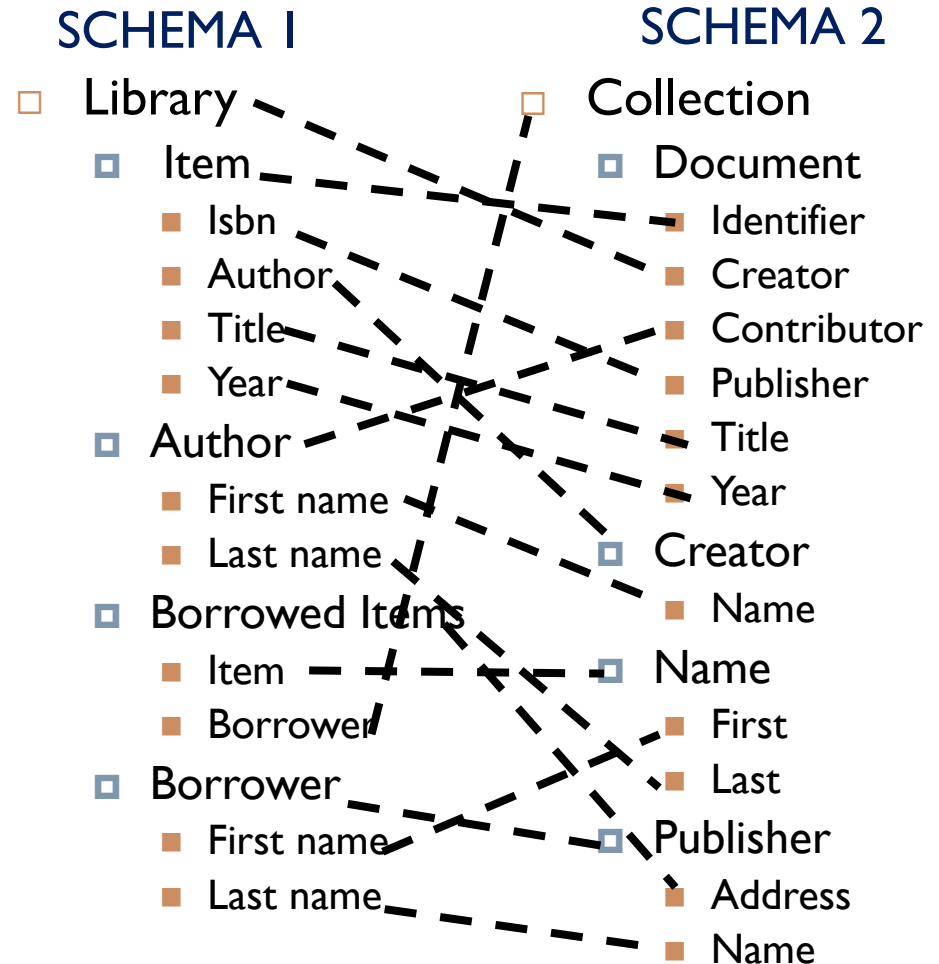
$$x_{i_k, j_l} \in \{0, 1\} \forall i \in [1, N-1] \forall j \in [i+1, N] \\ \forall k \in [1, n_i], \forall l \in [1, n_j]$$

Decision variables

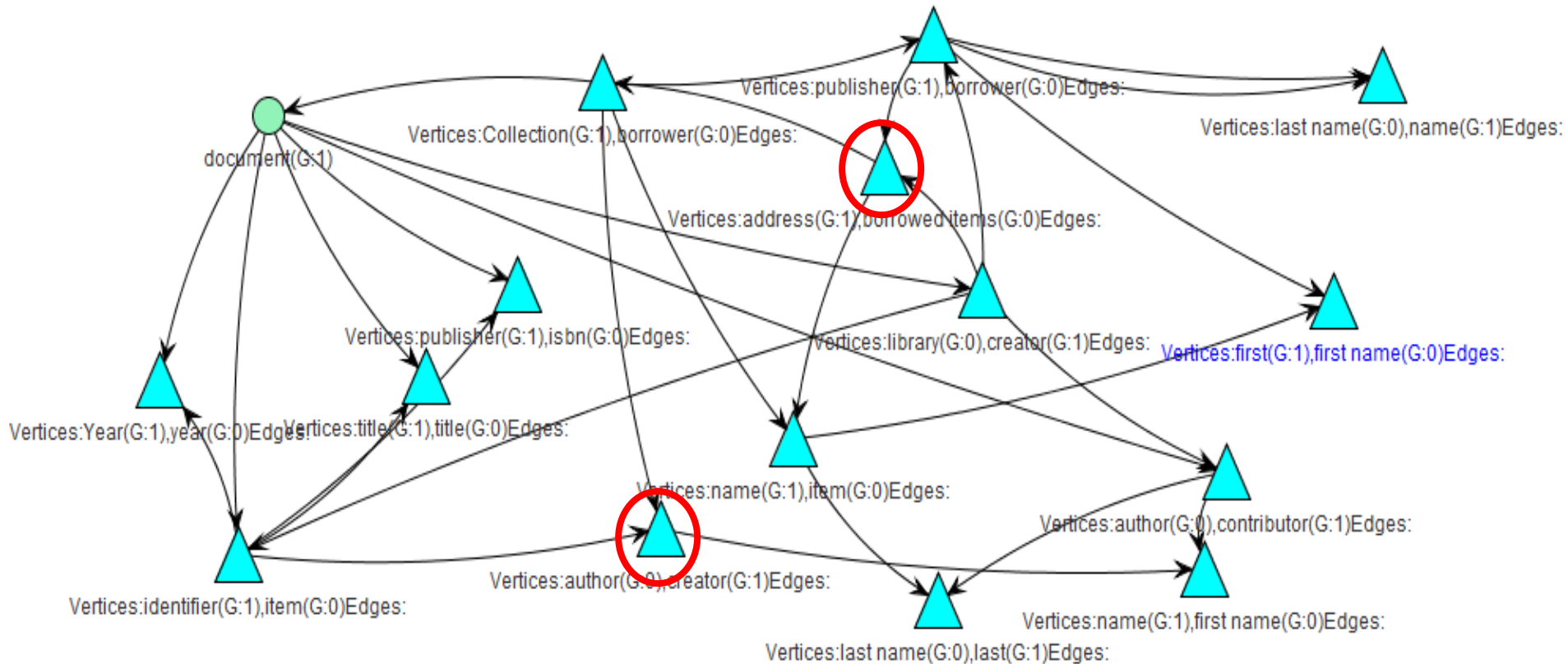
3. LP4HM : cardinality “1-1”

Each node in a graph G_i could match with at most one node in another graph G_j

$$\sum_{l=1}^{n_j} x_{ik,jl} \leq 1,$$



3. LP4HM : cardinality “1-1”

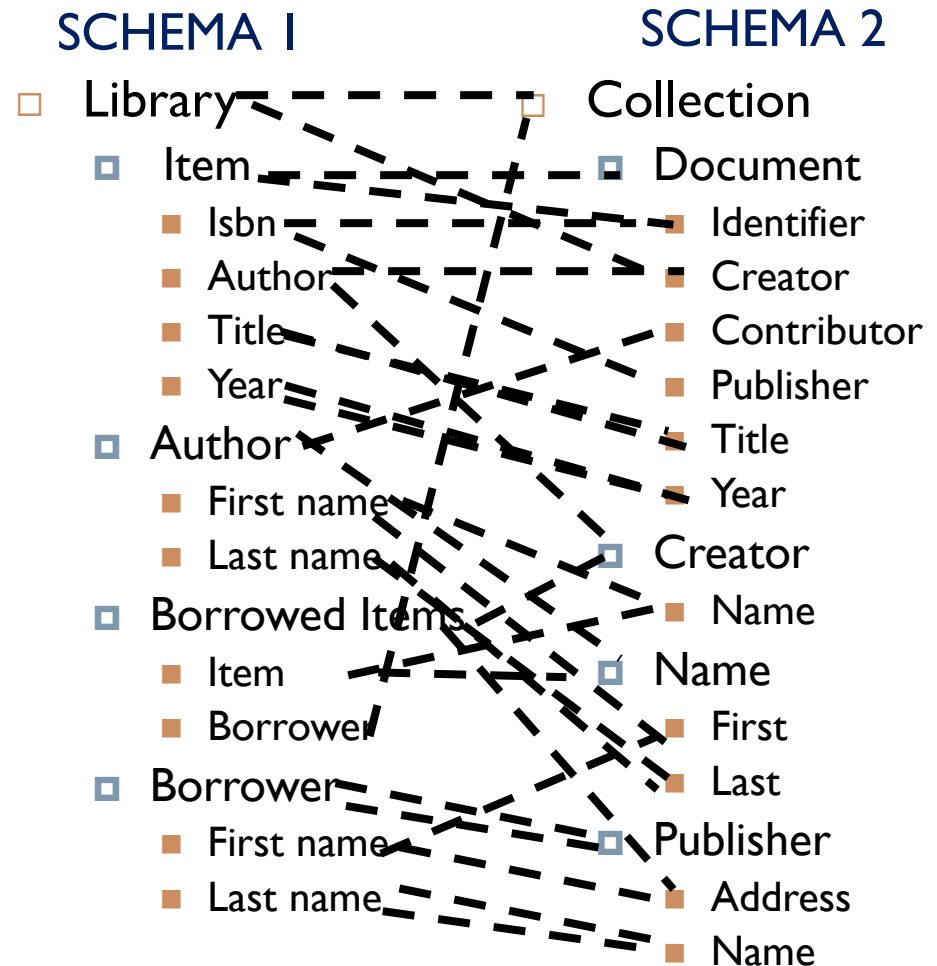


3. LP4HM : strict hierarchies

13

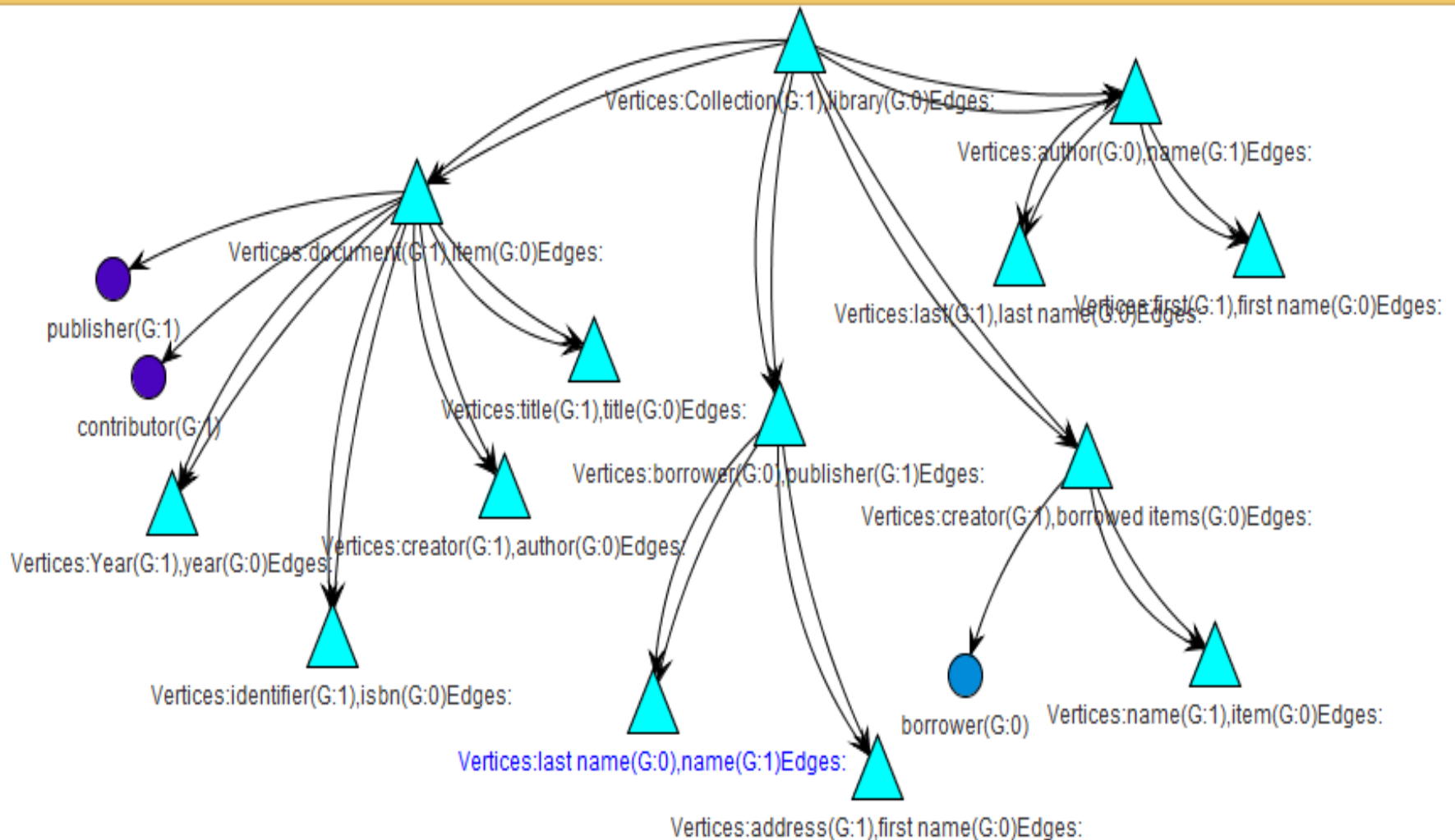
Each node should have at most one parent :
if children nodes match so
their predecessor have to
match too

$$x_{ik,jl} \leq x_{ipred(k),jpred(l)}$$



3. LP4HM : strict hierarchies

14



3. LP4HM : coherent structure

15

Prevent the generation of conflictual edges resulting from the matching of parents with childs' and childs with parents'

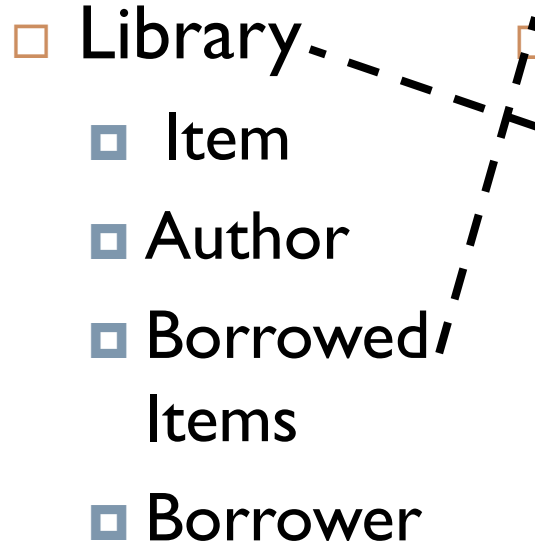
$$x_{i_k, j_l} + x_{i_{k'}, j_{l'}} - (dir_{i_k, k'} dir_{j_l, l'}) \leq 1$$

SCHEMA 1

- Library
 - Item
 - Author
 - Borrowed/Items
 - Borrower

SCHEMA 2

- Collection
 - Document
 - Creator
 - Name
 - Publisher



3. LP4HM : threshold

16

$$\begin{aligned} \max & \sum_{i=1}^{N-1} \sum_{j=i+1}^N \sum_{k=1}^{n_i} \sum_{l=1}^{n_j} sim_{i_k, j_l} \\ \text{s.t.} & \sum_{l=1}^{n_j} x_{i_k, j_l} \leq 1, \quad \forall k \in [1, n_i], \\ & \quad \quad \quad \forall i \in [1, N-1] \end{aligned}$$

For a given threshold, all correspondences should be above this threshold

$$\begin{aligned} sim_{i_k, j_l} x_{i_k, j_l} &\geq \text{seuil} x_{i_k, j_l} \\ &\forall i \in [1, N-1] \quad \forall j \in [i+1, N] \\ &\quad \forall k \in [1, n_i], \quad \forall l \in [1, n_j] \end{aligned}$$

$$\begin{aligned} x_{i_k, j_l} &\leq x_{i_{pred(k)}, j_{pred(l)}} \\ &\forall i \in [1, N-1] \quad \forall j \in [i+1, N] \\ &\quad \forall k \in [1, n_i], \quad \forall l \in [1, n_j] \end{aligned}$$

$$\begin{aligned} x_{i_k, j_l} + x_{i_{k'}, j_{l'}} - (dir_{i_k, k'} dir_{j_l, l'}) &\leq 1 \\ &\forall i \in [1, N-1] \quad \forall j \in [i+1, N] \\ &\quad \forall k, k' \in [1, n_i], \quad \forall l, l' \in [1, n_j] \end{aligned}$$

$$\begin{aligned} x_{i_k, j_l} &\in \{0, 1\} \\ &\forall i \in [1, N-1] \quad \forall j \in [i+1, N] \\ &\quad \forall k \in [1, n_i], \quad \forall l \in [1, n_j] \end{aligned}$$

3. LP4HM Relaxed

17

From binary to fractional decision variables (in $[0,1]$)

SCHEMA 1

Library

Item

Author

First name

Last name

Borrowed Items

Borrower

First name

Last name

SCHEMA 2

Collection

Document

Creator

Name

Name

First

Last

Publisher

Address

Name

4. Evaluation

Matching quality

18

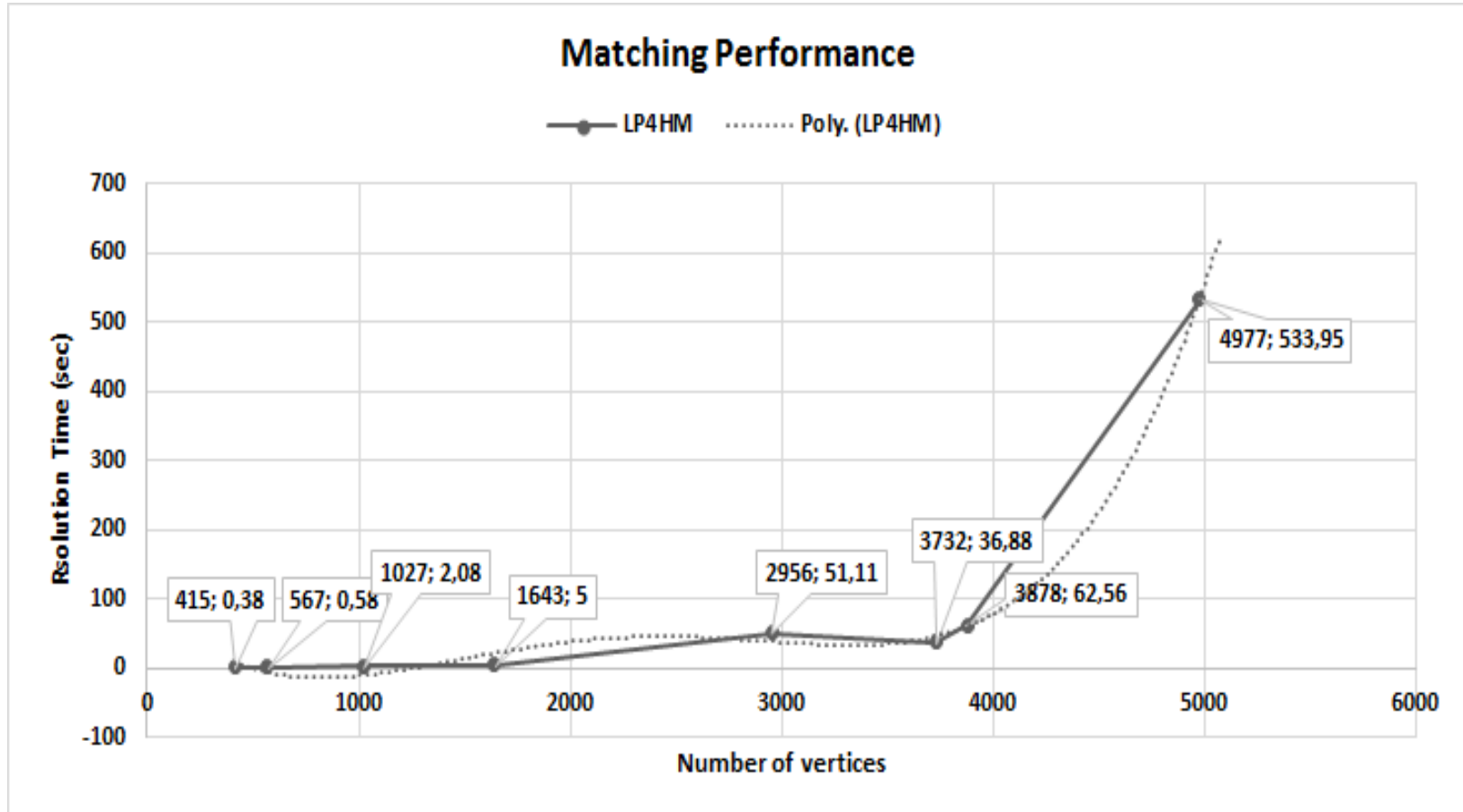
- Results on Similarity Flooding benchmark (7 users, 9 tasks)

	Precision	Recall	F-Measure	Accuracy	HSR
LP4HM	67%	58%	62%	30%	81%
LP4HM(Relaxed)	58%	66%	60%	23%	81%
COMA++	72%	50%	58%	32%	76%
Bmatch	22%	47%	28%	0%	69%
Similarity Flooding	81%	55%	65%	43%	80%

4. Evaluation

Matching performance

19



5. Conclusion and perspectives

20

- ❑ An approach to holistically integrate open data
 - A linear program focused on hierarchical structural constraints
 - Two strategies to resolve two types of cardinality problem

- ❑ Evaluation based on matching quality and performance
 - Competitive results without similarity threshold
 - A polynomial trendline for several input open data
 - A generalised approach for $N \geq 2$ graphs (holistic)

- ❑ Perspectives
 - Extend this model to handle labelled graphs
 - Experiments on larger datasets



**Thank you for you attention
Questions ?**

21

{berro,megdiche, teste}@irit.fr

IRIT, Toulouse, France