

Angela Bonifati

Ugo Comignani

Emmanuel Coquery

Romuald Thion

angela.bonifati@univ-lyon1.fr

ugo.comignani@univ-lyon1.fr

emmanuel.coquery@univ-lyon1.fr

romuald.thion@univ-lyon1.fr

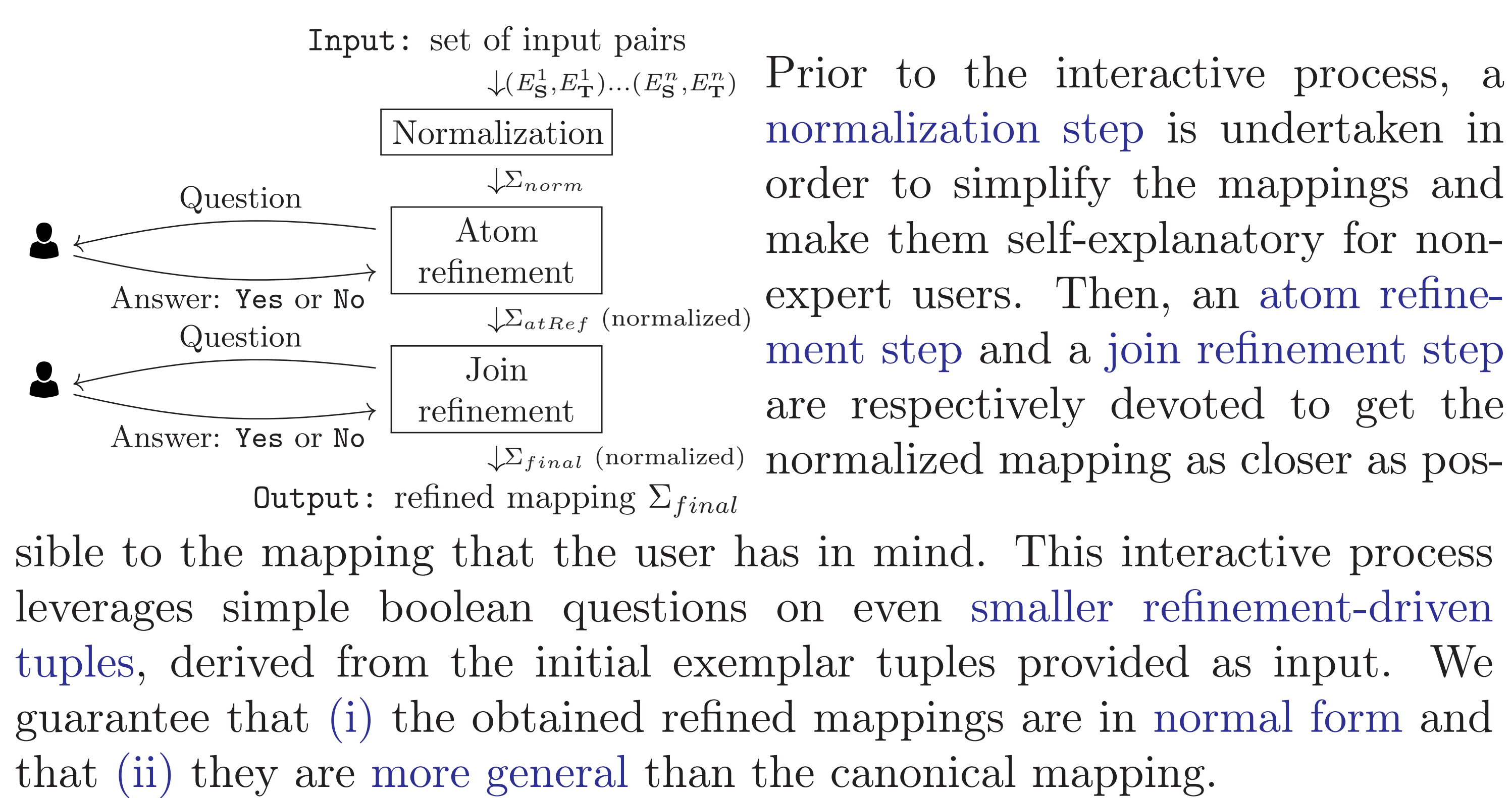
CNRS LIRIS and Université Claude Bernard Lyon 1

## 1. Motivations and goals

We propose an Interactive Mapping Specification process targeting non-expert users that:

- bootstraps with a set of **exemplar tuples** ( $E_S, E_T$ ), corresponding to a limited number of tuples provided as input;
- challenges the non-expert users with simple **boolean questions**, which are affordable for such users;
- is guaranteed to always produce a GLAV mapping  $\mathcal{M}'$  that **generalizes** a mapping  $\mathcal{M}$  in the user's mind, which is unknown beforehand.

## 2. Proposed framework



## 3. Specifying mappings using user exemplar tuples

User exemplar tuples:

Company			Flight		
IdCompany	Name	Town	Departure	Arrival	IdCompany
'C1'	'AA'	'Paris'	'Lyon'	'Paris'	'C1'
'C2'	'Ev'	'Lyon'	'Paris'	'Lyon'	'C2'

Travel Agency		
IdAgency	Name	Town
'A1'	'TC'	'L.A.'

Firm			Departure		Arrival	
Id	Name	Town	Town	IdFirm	Town	IdFirm
'Id1'	'AA'	'Paris'	'Lyon'	'Id1'	'Paris'	'Id1'
'Id2'	'Ev'	'Lyon'	'Paris'	'Id2'	'Lyon'	'Id2'
'Id3'	'TC'	'L.A.'				

Canonical mapping:

$m$  :  $Company(c1, aa, paris) \wedge Company(c2, ev, lyon) \wedge TravelAgency(a1, tc, la)$   
 $\wedge Flight(lyon, paris, c1) \wedge Flight(paris, lyon, c2)$   
 $\rightarrow \exists id1, id2, id3, Firm(id1, aa, paris) \wedge Departure(lyon, id1) \wedge Arrival(paris, id1)$   
 $\wedge Firm(id2, ev, lyon) \wedge Departure(paris, id2) \wedge Arrival(lyon, id2) \wedge Firm(id3, tc, la)$

Normalized mapping:

Before initiating the refinement steps, the canonical mapping is normalized [1] in order to separate unrelated connected components in the right-hand sides. This leads to the following split-reduced mapping:

$m_a$  :  $Company(c1, aa, paris) \wedge Company(c2, ev, lyon) \wedge TravelAgency(a1, tc, la)$   
 $\wedge Flight(lyon, paris, c1) \wedge Flight(paris, lyon, c2)$   
 $\rightarrow \exists id1, Firm(id1, aa, paris) \wedge Departure(lyon, id1) \wedge Arrival(paris, id1)$

$m_b$  :  $Company(c1, aa, paris) \wedge Company(c2, ev, lyon) \wedge TravelAgency(a1, tc, la)$   
 $\wedge Flight(lyon, paris, c1) \wedge Flight(paris, lyon, c2)$   
 $\rightarrow \exists id2, Firm(id2, ev, lyon) \wedge Departure(paris, id2) \wedge Arrival(lyon, id2)$

$m_c$  :  $Company(c1, aa, paris) \wedge Company(c2, ev, lyon) \wedge TravelAgency(a1, tc, la)$   
 $\wedge Flight(lyon, paris, c1) \wedge Flight(paris, lyon, c2)$   
 $\rightarrow \exists id3, Firm(id3, tc, la)$

During normalization, split reduction is followed by the deletion of logically equivalent tgds. In the above example, this leads to the removal of the tgd  $m_a$  (or, equivalently,  $m_b$ ) in order to obtain  $\Sigma_{norm}$ .

Final mapping after interactive specification:

$m_1$  :  $Company(c1, aa, paris_1) \wedge Flight(lyon, paris_2, c1)$   
 $\rightarrow \exists id1, Firm(id1, aa, paris_1) \wedge Departure(lyon, id1) \wedge Arrival(paris_2, id1)$

$m_2$  :  $TravelAgency(a1, tc, la) \rightarrow \exists id3, Firm(id3, tc, la)$

## 4. Atom refinement

Principle:

Remove superfluous atoms in the left-hand sides of normalized tgds.

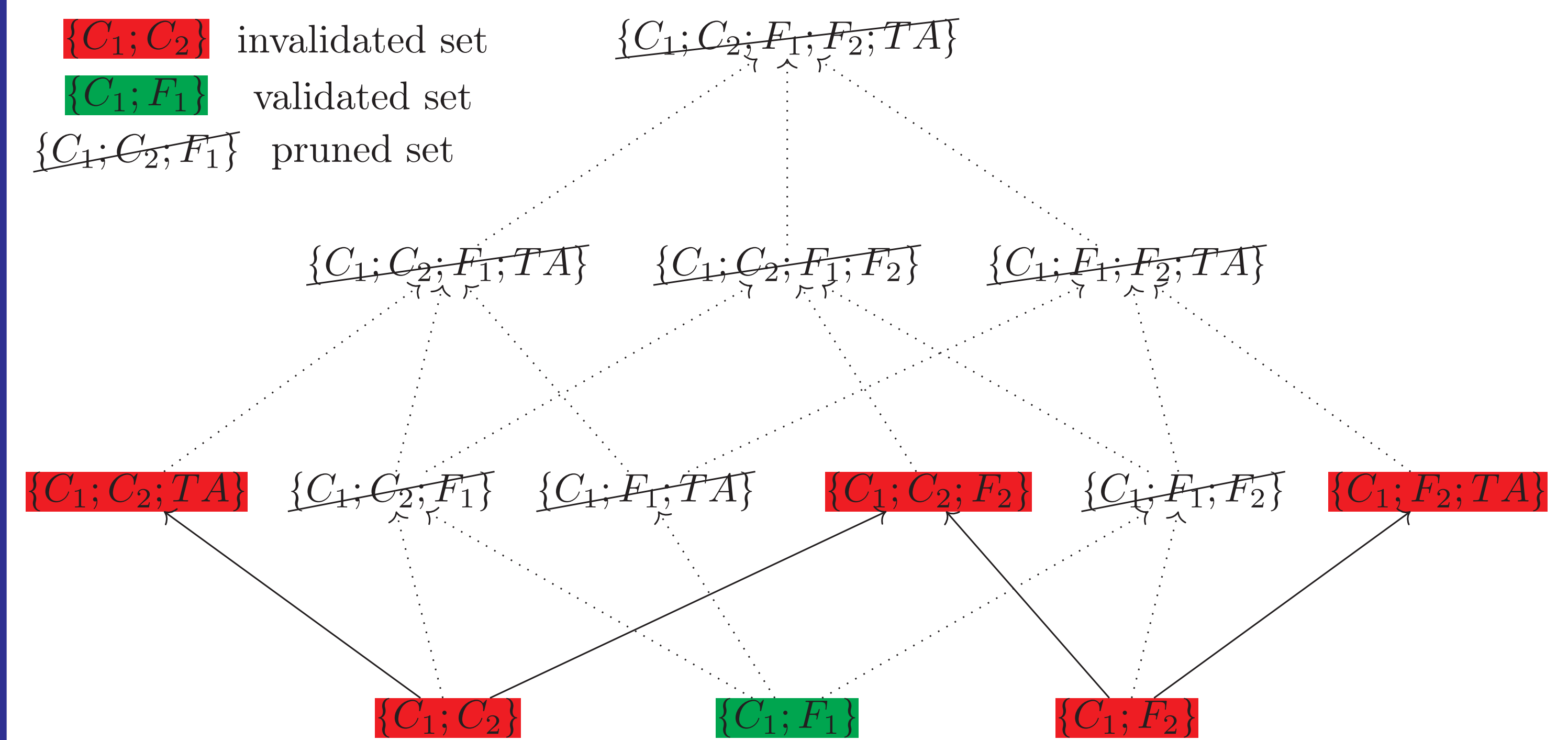
Exploration space:

Given a tgd, the possible sets of left-hand side atoms can be represented by a semi-lattice where the lower levels are the smallest valid sets of left-hand side atoms.

User feedback:

During semi-lattice exploration, the user is asked about the target tuples generated using only the subset of source tuples corresponding to the atoms in one node at a time.

Example with a bottom-up breadth-first approach:



As the sole set validated is  $\{C_1; F_1\}$ , the following tgd is generated:

$Company(c1, aa, paris) \wedge Flight(lyon, paris, c1)$   
 $\rightarrow \exists id1, Firm(id1, aa, paris) \wedge Departure(lyon, id1) \wedge Arrival(paris, id1)$

## 5. Join refinement

Principle:

For each tgd generated by atom refinement, identify redundant joins entailed by multiple occurrences of a given variable.

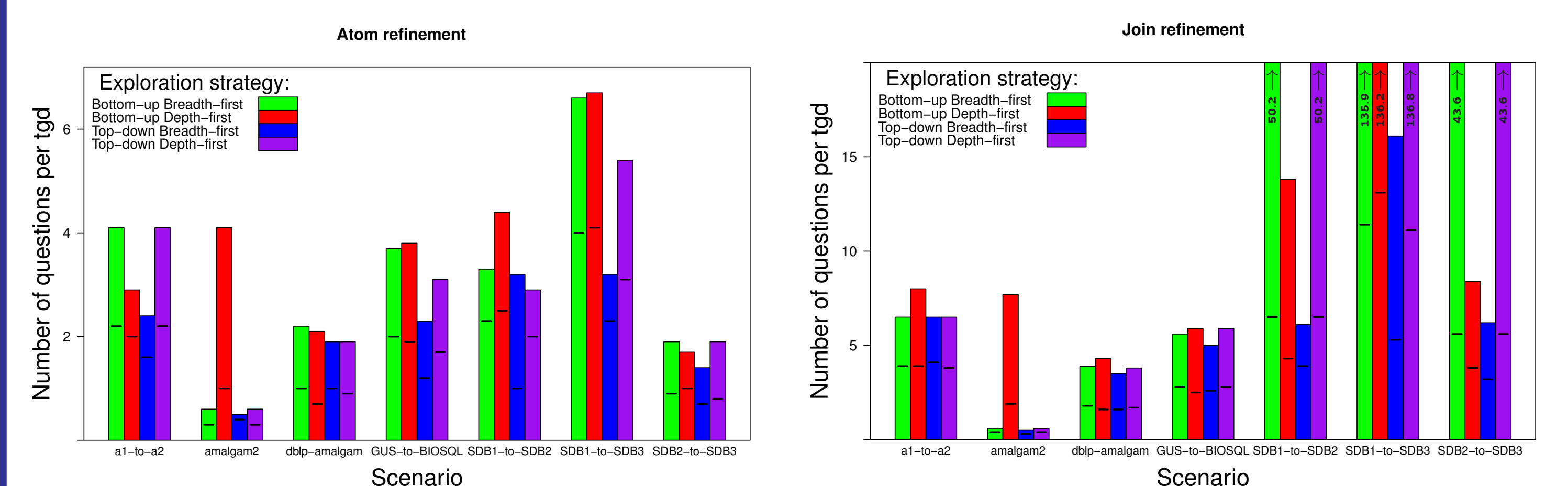
Exploration space:

As with atom refinement, we can build a semi-lattice of partitions representing possible joins between variable occurrences.

User feedback:

Similarly to atom refinement, the user is asked about the validity of small sets of tuples.

## 6. Experimental study



- We have used **seven real data integration scenarios** of the iBench benchmark [2].
- We have **simulated the user's ambiguities** by adding up to ten superfluous atoms or redundant joins to the initial exemplar tuples, derived from the above scenarios.
- We have tested **four exploration strategies** : bottom-up and top-down, each one with depth-first and breadth-first variations.

## References

- [1] G. Gottlob, R. Pichler, and V. Savenkov: *Normalization and optimization of schema mappings*. VLDB J., 20(2):277–302, 2011.
- [2] P.C. Arocena, B. Glavic, R. Ciucanu, and R.J. Miller: *The ibench integration metadata generator*. Proceedings of VLDB, 9(3):108–119, 2015.