

# A new Nested Graph Model for Data Integration

---

Giacomo Bergami

2/10/2017

1

## Introduction

- Graph data are a widely adopted.
- Currently, most of the No-SQL data representation is semistructured (JSON, XML).
- Graph operations allowing to integrate property graphs and semi-structured (**nested**) data are not provided.
  - An **intermediate data representation** allowing to represent both structured, semi-structured (nested) and non-structured data is missing.
  - A query language for this generalized data model is required. It should be able to sketch the data integration task and express both semistructured and graph queries.

2

## Related Works: Data Integration

- *Schema Oriented*, uniform data representation
  - A schema can be always extracted if missing (**generalization**).
  - **Schema alignment** allows to find data similarities.
  - All the data source schemas shall be aligned to the final **hub schema**.
- *Query Oriented*
  - The schema alignment can be represented as a query rewriting.
  - We can perform each query separately on the original sources first and then integrate the intermediate results (**LAV**) or translate and align the data first and then perform the query (**GAV**).
  - In both scenarios, part of the query has to be performed on top of the intermediate representation.

3

## Contributions

- An operation allowing to integrate different graphs was missing. Definition of the first class of binary operators for property graphs: **graph joins**.
  - Its definition pointed out the inefficiency of current graph query plans and data structures.
- Definition of a data model allowing to provide data integration between graph, semistructured and structured data.
  - **Future works**

4

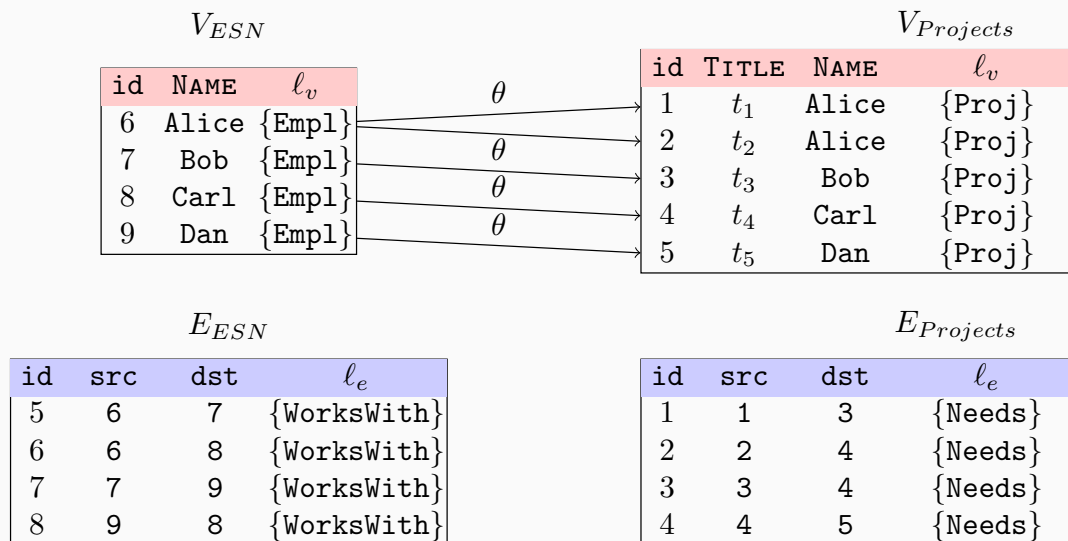
## Graph Joins: A query example

Consider an on-line service such as **ResearchGate** where researchers can {follow} each others' work, and a **citation graph**. Return the **paper graph** where **a paper cites another one iff. the first author of the first paper follows the first author of the second**.

- Binary operator returning just one graph.
- Vertex conditions are  $(\theta)$ -join conditions.
- A way to combine the edges is determined.

5

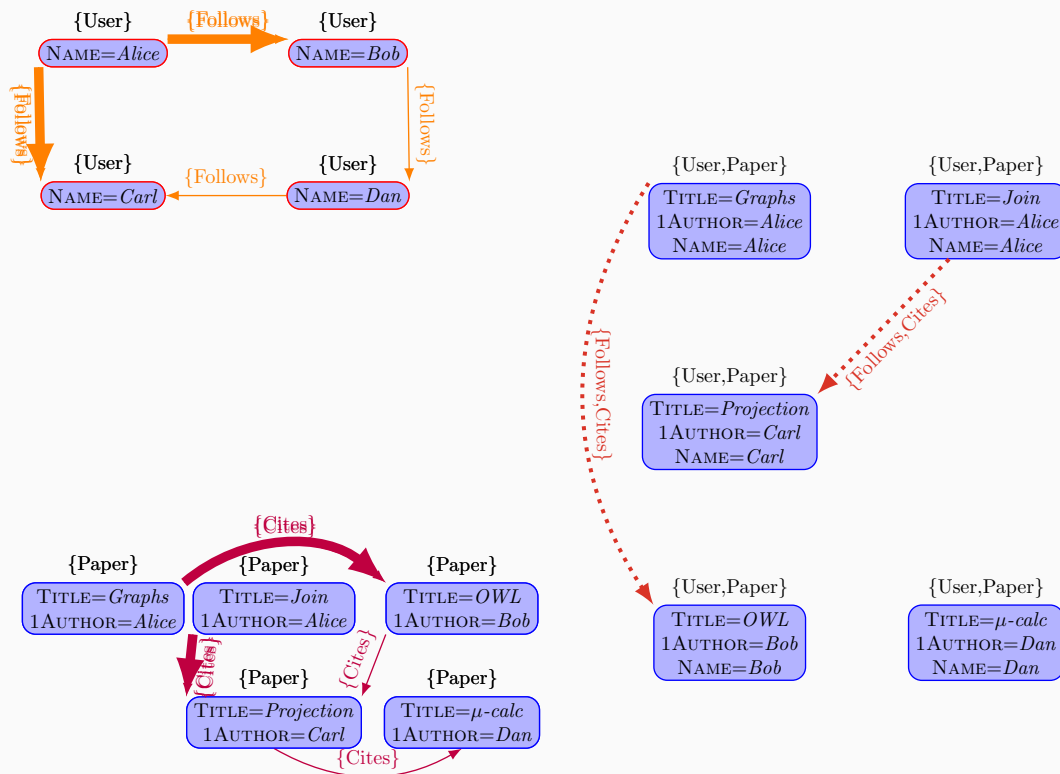
## Graph Joins: Relational Query Plan



**Figure 1:** A relational representation of the graph, where the vertices are  $\theta$ -joined.

6

## Graph Joins: The query result



7

## Graph Joins: Goals

1. The **data model** must *enhance the serialization* of both operands and graph result.
2. The **join definition** must be flexible enough to support further extensions (modularity, compositionality and properties preserving).
3. The **physical model** must allow a *quick access* to the data structures.

8

## Graph Joins: Physical Model

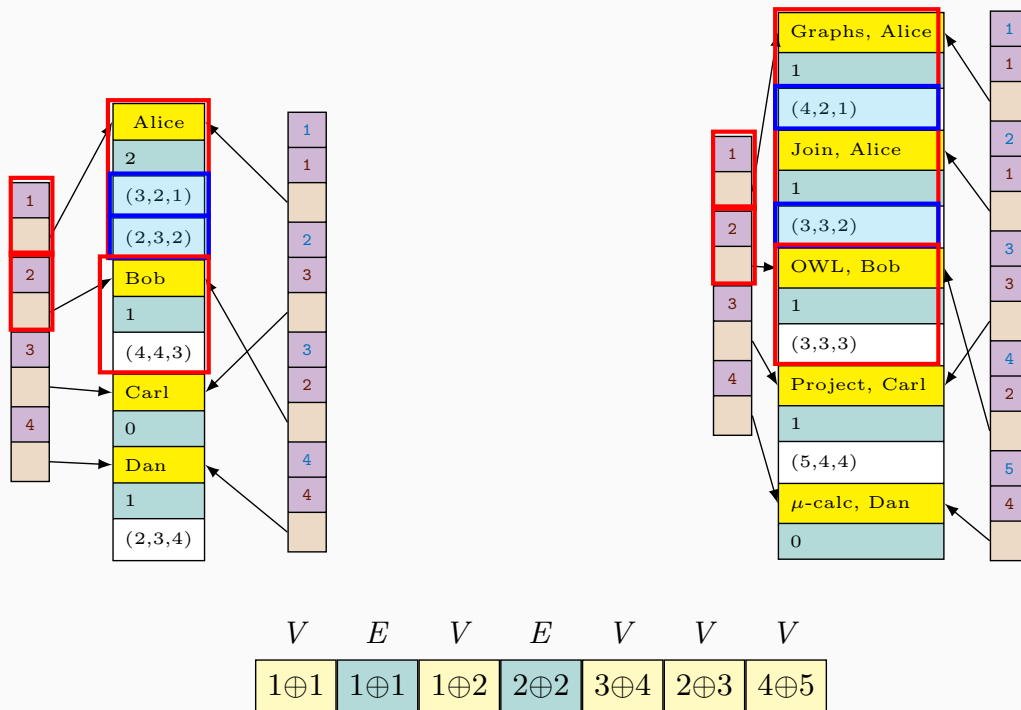


Figure 3: Representing Graph Operators and Output

9

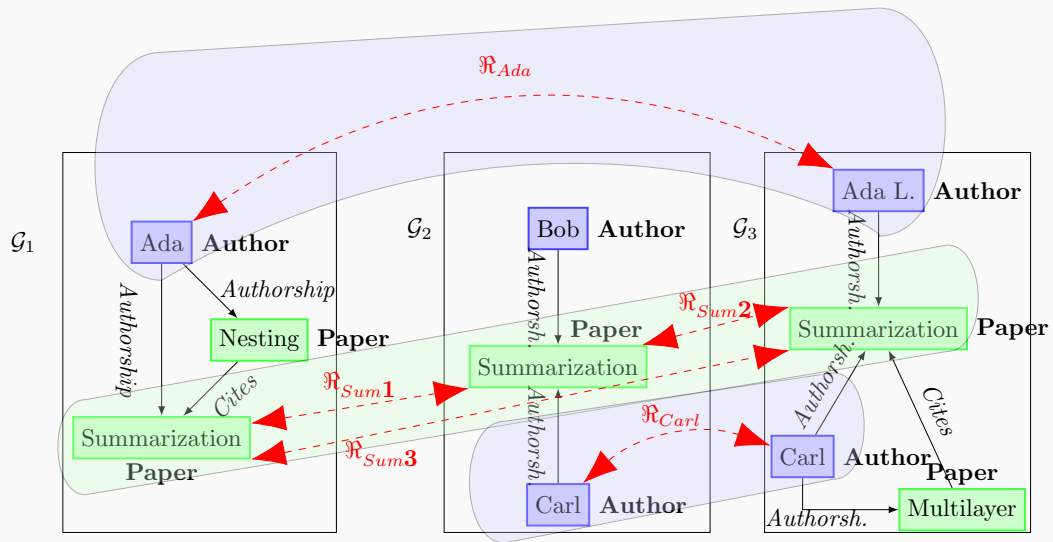
## Graph Joins: Benchmarks

Operands Size		Result		Join Time (C/C++) (ms)			Join Time (Java) (ms)	
Left ( $ V $ )	Right ( $ V $ )	Size ( $ V $ )	Size ( $ E $ )	Virtuoso	PostgreSQL	GCEA (C++)	Neo4J	GCEA (Java)
10	10	5	2	4.99	11.29	0.53	211.45	24.97
$10^2$	$10^2$	16	4	4.94	22.82	0.93	222.87	32.70
$10^3$	$10^3$	251	55	4.55	22.92	4.35	448.97	117.58
$10^4$	$10^4$	2734	680	117.71	12.00	183.90	40.42	3149.90
$10^5$	$10^5$	26803	7368	>4H	7150.74	411.78	241026.79	17178.49
$10^6$	$10^6$	151212	99558	>4H	99683.91	3966.72	>4H	178066.80

Figure 4: Conjunctive join execution on well known graph databases. Similar result were provided on graph database libraries.

10

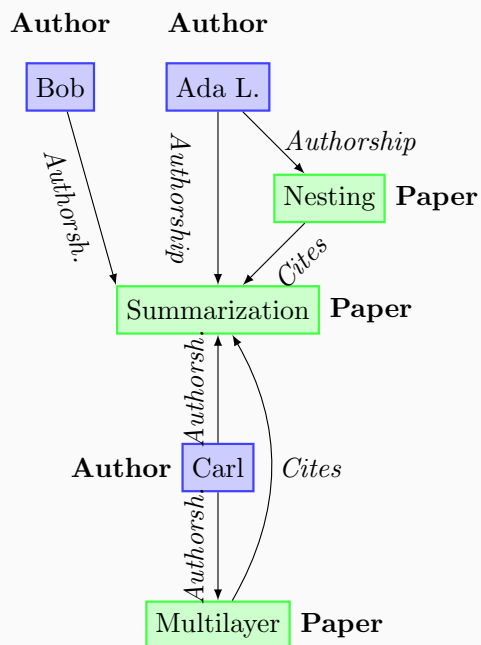
## Graph Joins: Limitations (1)



**Figure 5:** A data integration scenario, where the schema is already aligned

11

## Graph Joins: Limitations (2)



**Figure 6:** A data integration scenario, where the schema is already aligned

12

## Graph Joins: Limitations (3)

- Suppose to generalize the graph join by extending the  $\theta$ -join ( $\bowtie_{\theta}$ ) to a full  $\theta$ -join ( $\Join_{\theta}$ ) over the vertices.
  - This is possible within the previous definition.
- Suppose to use as a  $\theta$  some edges that remark the similarity values among the nodes.
  - After the pairwise application of the join, the  $\theta$  edges must be rewritten.
- In some data cleaning scenarios, we want to preserve the original non-aggregated information, while providing an integrated view as in the previous picture.
  - **A new data model is required**, generalizing statecharts and hypernodes.

13

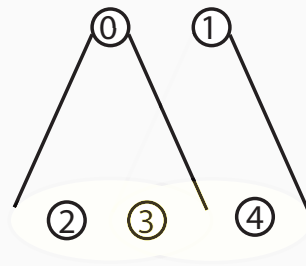
## Generalized Semistructured Model and Nested Graphs

$$GSM = (o, O, \ell, \xi, \phi)$$

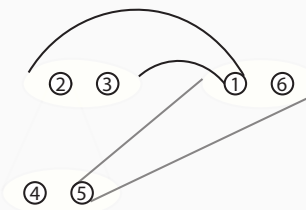
- $o$  is the **reference object**, designing the root of the GSM.
- $\phi$  is the object id containment function for each expression  $e$  ( $\phi(o', e)$ ).
- $O$  is the supset of all the objects contained by  $o$  ( $\bigcup_e \phi(o', e)$ )
- $\ell$  is the function associating to each object in  $O$  a list of labels.
- $\xi$  is the function associating to each object in  $O$  a list of expressions.

14

## GSM, generalizing semistructured data containment



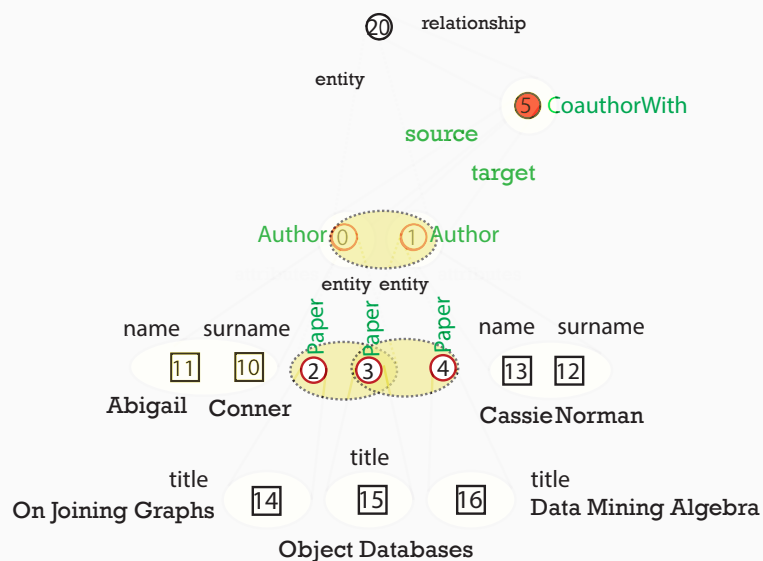
**Figure 7:** Allowing multiple containments of the same object.



**Figure 8:** Allowing multiple nesting and nestings with recursions.

15

## GSM, representing (nested) graphs



**Figure 9:** Representing a nested graph with two vertices and one edge.

16



## GSM, representing semistructured documents

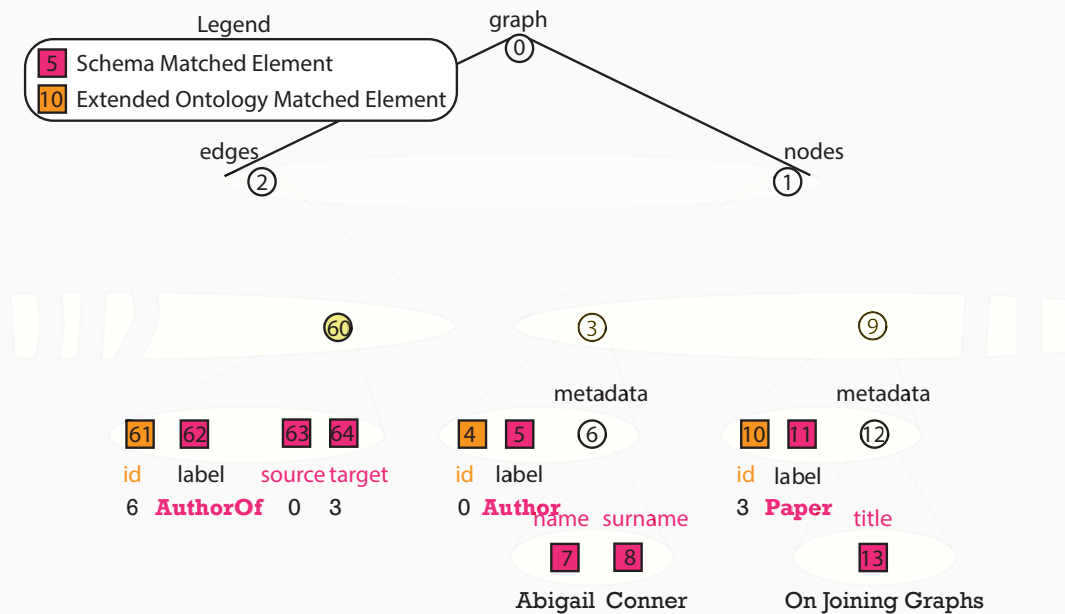


Figure 10: Another graph representation from a JSON document.

17

## GSM, aligning graphs with semistructured schemas (1)

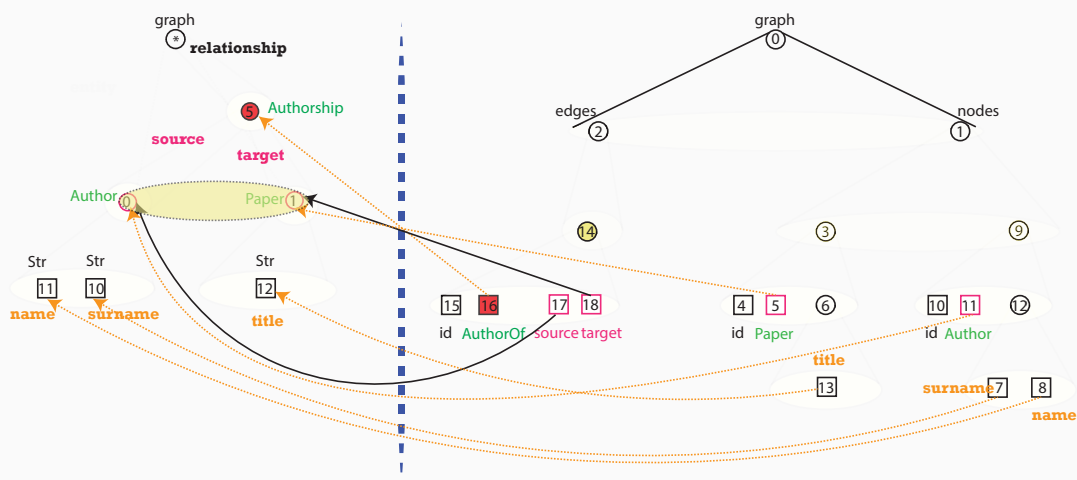
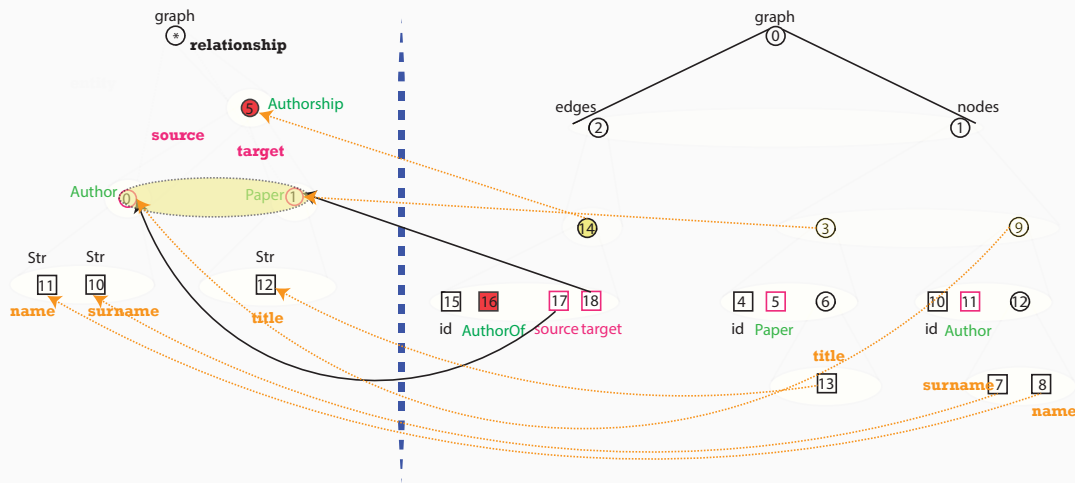


Figure 11: Schema alignment, using the usual techniques.

18

## GSM, aligning graphs with semistructured schemas (2)



**Figure 12:** Schema alignment, using the nesting information refinement.

19

## GSM, aligning graphs with semistructured schemas (3)

The previous data alignment process provides an interesting perspective for graph grammars:

- The schema (*on the right*) extracted from the original data defines the **matching graph**.
  - Given that the schema was extracted from the original data, we already know the result of such match.
- The hub schema (*on the left*) defines how the data on the right has to be rewritten **transformation graph**, using the edges as transformations.

Hereby, a generalization of the graph grammars is required. The operations required to do the matching and transformations are non trivial. Therefore, an algebra expressing such operations is required.

20

## paNGRAM Algebra

paNGRAM	COA Algebra	3W Algebra	Graph Languages
<b>Selection</b> , $\sigma_P$	Selection, $\sigma_P$ Projection, $\pi_P$	Selection, $\sigma_P$	Graph Traversal, $P$
<b>Map</b> , $\mu_{\ell', \xi', \phi'}$	Embedding, $\varepsilon_\phi$ Projection, $\pi_\phi$	Calc, $Calc_f$	Transform, Pattern Matching
<b>Creation</b> , $\kappa_{L,E}^\omega$	Embedding new elements	Regionize, $\kappa$	Transform,
<b>Fold</b> , $\text{fold}_{S,f}$	Expression using $\varsigma$	Loop operator, $\lambda$	Reduce (GrAlA)
<b>Disjoint Union</b> , $\sqcup$	$+\mu$ = Set operations	$+\mu$ = Set operations	Graph Binary Operators

- An operation for creating new elements is required because this algebra manipulates ids (already existing objects).
- All the other operations (such as graph unnesting  $\nu$ , splicing  $\varsigma$ , graph joins  $\bowtie$ ) are expressible through a composition of operators.

21

## Future Works

- In order to accomplish the data integration scenario, the alignment between the source schemas and the hub schema shall be defined.
- Algebraic expressions reproducing the single steps shall be defined.
- Check if the whole algorithmic plan can be implemented more efficiently than the single operators.

22