Graph Querying

Pierre Genevès CNRS





Reknown scientists born in Europe:

?x, ?y \leftarrow ?x hasWonPrize ?y, ?x wasBornIn/isLocatedIn⁺ Europe



Reknown scientists born in Europe:





Reknown scientists born in Europe:

conjunct conjunct



Reknown scientists born in Europe:

simple relation path relation



Recursion

- Essential mechanism to enable deep navigation in a graph
- ▶ Major feature for extracting valuable information from a linked data structure
- Navigational paths (regular path queries)
- PageRank, shortest paths, connected components, etc.

More Examples

All ancestors of Alice:

 $?y \leftarrow Alice (father|mother)^+ ?y$

Pairs of stops connected by tram lines A and B:

```
?x, ?y \leftarrow ?x nextStopA<sup>+</sup> ?y, ?x nextStopB<sup>+</sup> ?y
```

People with a certain skill that Bob knows, either directly or indirectly:

 $?x \leftarrow Bob knows^+ ?x:Person, ?x hasSkill skill1$

Analysis of drug interactions, etc.

Recursive graph query evaluation today

In practice: on large knowledge graphs (e.g. Yago, wikidata), certain recursive queries can be difficult to evaluate... or even hardly feasible

Recursive graph query evaluation today

In practice: on large knowledge graphs (e.g. Yago, wikidata), certain recursive queries can be difficult to evaluate... or even hardly feasible

> This depends on graph instance (size, topology) and **query shape**

Recursive graph query evaluation today

In practice: on large knowledge graphs (e.g. Yago, wikidata), certain recursive queries can be difficult to evaluate... or even hardly feasible

▶ This depends on graph instance (size, topology) and query shape



Recursive subqueries may result in intermediate results whose size is bigger than the overall graph size

N Red/Green* ?t

This query retrieves nodes reachable from node N through a given (recursive) path along edges































N Red/Green* ?t

Computing Green* might be very costly (or even unfeasible) on certain graphs.

N Red/Green* ?t

Computing Green* might be very costly (or even unfeasible) on certain graphs.

Retrieving the query results without computing the full relation $\ensuremath{\mathsf{Green}}^*$ is possible

N Red/Green* ?t

Computing Green* might be very costly (or even unfeasible) on certain graphs.

Retrieving the query results without computing the full relation Green* is possible

More generally

The evaluation strategy - the query execution plan - can have a huge impact on performance.

N Red/Green* ?t

Computing Green* might be very costly (or even unfeasible) on certain graphs.

Retrieving the query results without computing the full relation Green* is possible

More generally

The evaluation strategy - the query execution plan - can have a huge impact on performance.

Query answering might be feasible... or not. Even for graphs of moderate size.

N Red/Green* ?t

Computing Green* might be very costly (or even unfeasible) on certain graphs.

Retrieving the query results without computing the full relation Green* is possible

More generally

The evaluation strategy - the query execution plan - can have a huge impact on performance.

Query answering might be feasible ... or not. Even for graphs of moderate size.

Query execution plan is crucial.

Other examples

What if...

- the query contains more than one variable? ?x Red*/Green ?y
- more general forms of navigation? ?x (Red | Green)* N
- Several recursions? Which exploration should start first?



Fundamental Problem

Given a recursive graph query*, how to generate an efficient evaluation plan?

(*) Well-known recursive query language fragments include e.g.: RPQ, C2RPQ, UCRPQ, etc.

What is the state of the theory?

The Theory – Brief (and partial) Recap

1970: Codd's relational algebra (RA), established the domain of databases (multi \$G market in 2022) key ideas: separation of the query language (SQL) from the RA; tables (no recursion)

1979-1990s: attempts at extending RA with recursion, moderate success (limited forms of recursion); research on the Datalog (logic-programming view) side

1999: SQL supports recursive queries (but seen as optimization barriers for optimizers)

2000-2015: boom of NoSQL research (in particular for trees) key idea: preserve the native data structure (try not to split into tables)

Late 201x: a myriad of graph DB systems (many with poor/no support of recursion). Some with recursion increasingly inspired by Datalog or SQL.

From 2020: RA extended with a more general form of recursion (inspired from tree logics), with application to graphs.


Relational Data Model (Codd, 1970)

- Data organized in tables (relations)
- 1 database: 1 set of relations
- Type of a relation: set of (possibly typed) column names
- Data seen as tuples (or mappings of column names and values)

authId	name	birthDate
1	Edgar F. Codd	19/08/1923
2	Marie Curie	07/11/1867
3	Alain Turing	23/06/1912

. 1

articic			
arId	authId	title	
1	1	A relational model of data for large shared data banks	
2	1	Understanding relations	
3	2	Rayon semis par les composes de l'Uranium et du Thorium	
4	3	Computing machinery and intelligence	

article

Relational operators: filter

	author .				article		
	authId	name	birthDate		arId	authId	title
	1	Edgar F. Codd	19/08/1923		1	1	A relational model of data for large shared data banks
• name=Marie Curie	2	Marie Curie	07/11/1867		2	1	Understanding relations
	3	Alain Turing	23/06/1912		3	2	Rayon semis par les composes de l'Uranium et du Thorium
				I	4	3	Computing machinery and intelligence
L							
	authId	name	birthDate				
	2	Marie Curie	07/11/1867				

Relational operators: projection

		author	
	authId	name	birthDate
π	1	Edgar F. Codd	19/08/1923
name	2	Marie Curie	07/11/1867
	3	Alain Turing	23/06/1912



ar	ticle	
hId	title	

arId	authId	title
1	1	A relational model of data for large shared data banks
2	1	Understanding relations
3	2	Rayon semis par les composes de l'Uranium et du Thorium
4	3	Computing machinery and intelligence

Relational operators: natural join

author					ar	ticle
authId	name	birthDate		arId	authId	title
1	Edgar F. Codd	19/08/1923		1	1	A relational model of data for large shared data banks
2	Marie Curie	07/11/1867	\bowtie	2	1	Understanding relations
3	Alain Turing	23/06/1912		3	2	Rayon semis par les composes de l'Uranium et du Thorium
				4	3	Computing machinery and intelligence
			Π			

arId	authId	name	birthDate	title
1	1	Edgar F. Codd	19/08/1923	A relational model of data for large shared data banks
2	1	Edgar F. Codd	19/08/1923	Understanding relations
3	2	Marie Curie	07/11/1867	Rayon semis par les composes de l'Uranium et du Thorium

Other relational operators

• Column renaming $\rho_a^b(author)$: renames column *a* into column *b* in relation *author*

▶ Union ∪ of two relations (preferrably of same type!)

▶ ...

▶ No recursion in Codd's original relational algebra.

Graph Data

Possible representation: a table keeps track of source and target nodes connected by each relation $% \left({{{\left({{{\left({{{c}} \right)}} \right)}_{i}}}_{i}}} \right)$

For instance, for a social network:

src	trg
Anna	Tom
Tom	Ryan
Ryan	Jane

follows

Overview of some recent results in the area

Extended Relational Algebra (μ -RA, 2020)

Translation of Path Queries



Translation of Path Queries



Translation of Path Queries



 $Tr(A/B) = \widetilde{\pi}_m(\rho_t^m(Tr(A)) \bowtie \rho_s^m(Tr(B)))$

 $A^* = \text{Empty Path } or \underbrace{\text{Path of one or more edges labeled with A}}_{A^+ = A^*/A}$

$$A^* = \text{Empty Path or } \underbrace{\text{Path of one or more edges labeled with A}}_{A^+ = A^*/A}$$

$${\it Tr}(A^*) = {\it EmptyPath} \cup {\it Tr}(A^*)/A$$

$$A^* = \text{Empty Path or Path of one or more edges labeled with A}_{A^+ = A^*/A}$$

$$Tr(A^*) = EmptyPath \cup Tr(A^*)/A$$

 $= \mu \Big(X = EmptyPath \cup X/A \Big)$

$$A^* = \text{Empty Path } or \underbrace{\text{Path of one or more edges labeled with A}}_{A^+ = A^*/A}$$

$$Tr(A^*) = EmptyPath \cup Tr(A^*)/A$$

= $\mu \left(X = EmptyPath \cup X/A \right)$
= $\mu \left(X = \beta_t^s (AllNodes) \cup \tilde{\pi}_m \left(\rho_t^m (X) \bowtie \rho_s^m (Tr(A)) \right) \right)$

Algebraic transformations of fixpoints

Algebraic transformations of fixpoints

pushing filters?

$$\sigma_{\mathit{filter}}\left(\mu({\sf X}=arphi)
ight)\stackrel{?}{=}\mu({\sf X}=\sigma_{\mathit{filter}}\left(arphi
ight))$$

Algebraic transformations of fixpoints

pushing filters?

pushing joins?

$$\sigma_{\textit{filter}} \left(\mu(X = \varphi) \right) \stackrel{?}{=} \mu(X = \sigma_{\textit{filter}} \left(\varphi \right))$$

$$\psi \bowtie \mu(X = \varphi) \stackrel{?}{=} \mu(X = \psi \bowtie \varphi)$$

Algebraic transformations of fixpoints

pushing filters?

pushing joins?

pushing projections?

$$\sigma_{\textit{filter}} \left(\mu(X = \varphi) \right) \stackrel{?}{=} \mu(X = \sigma_{\textit{filter}} \left(\varphi \right))$$

$$\psi \bowtie \mu(X = \varphi) \stackrel{?}{=} \mu(X = \psi \bowtie \varphi)$$

$$\widetilde{\pi}_{p}\left(\mu(X=\varphi)\right)\stackrel{?}{=}\mu(X=\widetilde{\pi}_{p}\left(\varphi\right))$$

Algebraic transformations of fixpoints

pushing filters?

 $\sigma_{\textit{filter}} \left(\mu(X = \varphi) \right) \stackrel{?}{=} \mu(X = \sigma_{\textit{filter}} \left(\varphi \right))$

$$\psi \bowtie \mu(X = arphi) \stackrel{?}{=} \mu(X = \psi \bowtie arphi)$$

$$\widetilde{\pi}_{\rho}\left(\mu(X=\varphi)\right)\stackrel{?}{=}\mu(X=\widetilde{\pi}_{\rho}\left(\varphi\right))$$

pushing projections?

$$\mu(\boldsymbol{X}=\psi\cup\kappa) {\bowtie} \mu(\boldsymbol{X}=\varphi\cup\xi) \stackrel{?}{=} \mu(\boldsymbol{X}=\psi {\bowtie} \varphi\cup\xi\cup\kappa)$$

Algebraic transformations of fixpoints

pushing filters?

 $\sigma_{\textit{filter}} \left(\mu(X = \varphi) \right) \stackrel{?}{=} \mu(X = \sigma_{\textit{filter}} \left(\varphi \right))$

pushing joins?

pushing projections?

merging fixpoints?

$$\psi \bowtie \mu(X = \varphi) \stackrel{?}{=} \mu(X = \psi \bowtie \varphi)$$

$$\widetilde{\pi}_{p}\left(\mu(X=\varphi)\right)\stackrel{?}{=}\mu(X=\widetilde{\pi}_{p}\left(\varphi\right))$$

$$\mu(X = \psi \cup \kappa) \bowtie \mu(X = \varphi \cup \xi) \stackrel{?}{=} \mu(X = \psi \bowtie \varphi \cup \xi \cup \kappa)$$

All rules are semantics-preserving ? : decidable criteria [SIGMOD'20]

N Red/Green* ?t

N Red/Green* ?t

 $\widetilde{\pi}_{?s}\left(\sigma_{?s=N}\left(\operatorname{\mathsf{Red}}/\mu(X=\beta_{s}^{t}\left(\operatorname{\mathsf{AllNodes}}\right)\cup X/\operatorname{\mathsf{Green}}\right)\right)\right)$

N Red/Green* ?t

 $\widetilde{\pi}_{?s}(\sigma_{?s=N} (\text{Red}/\mu(X = \beta_s^t (\text{AllNodes}) \cup X/\text{Green})))$

 $\widetilde{\pi}_{?s}\left(\sigma_{?s=N}\left(\mu(X=\operatorname{\mathsf{Red}}/\beta_{s}^{t}\left(\operatorname{\mathsf{AllNodes}}\right)\cup X/\operatorname{\mathsf{Green}}\right)\right)\right)$

N Red/Green* ?t

 $\widetilde{\pi}_{?s}(\sigma_{?s=N} (\text{Red}/\mu(X = \beta_s^t (\text{AllNodes}) \cup X/\text{Green})))$

 $\widetilde{\pi}_{?s}(\sigma_{?s=N}(\mu(X = \text{Red}/\beta_s^t(\text{AllNodes}) \cup X/\text{Green})))$

 $\widetilde{\pi}_{?s}\left(\sigma_{?s=N}\left(\mu(X=\mathsf{Red}\cup X/\mathsf{Green})\right)\right)$

N Red/Green* ?t

 $\widetilde{\pi}_{?s}\left(\sigma_{?s=N}\left(\operatorname{\mathsf{Red}}/\mu(X=\beta_{s}^{t}\left(\operatorname{\mathsf{AllNodes}}\right)\cup X/\operatorname{\mathsf{Green}}\right)\right)\right)$

 $\widetilde{\pi}_{?s}\left(\sigma_{?s=N}\left(\mu(X=\operatorname{\mathsf{Red}}/\beta_{s}^{t}\left(\operatorname{\mathsf{AllNodes}}\right)\cup X/\operatorname{\mathsf{Green}}
ight)
ight)$

 $\widetilde{\pi}_{?s}\left(\sigma_{?s=N}\left(\mu(X=\mathsf{Red}\cup X/\mathsf{Green})\right)\right)$

 $\widetilde{\pi}_{?s}\left(\mu(X = \sigma_{?s=N} (\operatorname{\mathsf{Red}}) \cup X/\operatorname{\mathsf{Green}})\right)$

N Red/Green* ?t

 $\widetilde{\pi}_{?s}\left(\sigma_{?s=N}\left(\operatorname{\mathsf{Red}}/\mu(X=\beta_{s}^{t}\left(\operatorname{\mathsf{AllNodes}}\right)\cup X/\operatorname{\mathsf{Green}}\right)\right)\right)$

 $\widetilde{\pi}_{?s}\left(\sigma_{?s=N}\left(\mu(X=\operatorname{\mathsf{Red}}/\beta_{s}^{t}\left(\operatorname{\mathsf{AllNodes}}\right)\cup X/\operatorname{\mathsf{Green}}
ight)
ight)$

 $\widetilde{\pi}_{?s}(\sigma_{?s=N}(\mu(X = \mathsf{Red} \cup X/\mathsf{Green})))$

 $\widetilde{\pi}_{?s}\left(\mu(X = \sigma_{?s=N} (\operatorname{\mathsf{Red}}) \cup X/\operatorname{\mathsf{Green}})\right)$

 $\mu(X = \widetilde{\pi}_{?s} \left(\sigma_{?s=N} \left(\mathsf{Red} \right) \right) \cup X/\mathsf{Green})$

Sample query of the form: ?x Red⁺/Green⁺ ?y

plan 1 : unfold Green⁺ from right to left, ...

Sample query of the form: ?x Red⁺/Green⁺ ?y

plan 1 : unfold Green⁺ from right to left, ... plan 2 : start with Red⁺, left to right ...

Sample query of the form: ?x Red⁺/Green⁺ ?y

plan 1 : unfold Gr	een ⁺ from right to left,
plan 2 : start with	Red ⁺ , left to right
plan 3 : start with	Red ⁺ , right to left

Sample query of the form: ?x Red⁺/Green⁺ ?y

```
\begin{array}{rcl} \mathsf{plan 1} &:& \mathsf{unfold Green}^+ \text{ from right to left, }\dots \\ \mathsf{plan 2} &:& \mathsf{start with Red}^+, \mathsf{left to right }\dots \\ \mathsf{plan 3} &:& \mathsf{start with Red}^+, \mathsf{right to left }\dots \\ &\vdots && \vdots \\ \mathsf{plan i} &:& \mathsf{merged fixpoint:} \\ && \mu(X = \mathsf{Red}/\mathsf{Green} \cup \mathsf{Red}/X \cup X/\mathsf{Green}) \\ && \mathsf{start from Red}/\mathsf{Green}, \mathsf{then hop on left by Red}^*, \mathsf{and on the right by Green}^* \end{array}
```



1. A graph query is translated into $\mu\text{-RA}$



2. With transformation rules, we obtain a space of evaluation plans



3. An estimated most efficient plan is selected using a cost estimation and data statistics



 μ -RA plan space is richer compared to previous approaches (e.g. fixpoints can be merged)

Application for Graph Querying

In practice, on centralized systems

- Implementation on top of PostgreSQL (a popular relational database management system)
- Evaluation of queries on knowledge graphs:


Concluding Remarks

A glance at some hot research activity and recent results:

- Powerful recursive graph queries
- Extension of relational algebra with recursion [SIGMOD'20]
- Query plan enumeration (research results to be presented in 2024)
- Extension to neuro-symbolic graph queries...

References

Louis Jachiet, Pierre Genevès, Nils Gesbert and Nabil Layaïda,

On the Optimization of Recursive Relational Queries: Application to Graph Queries. In Proceedings of the ACM SIGMOD International Conference on Management of data, 2020 (SIGMOD'20)

Amela Fejza, Pierre Genevès, Nabil Layaïda, Sarah Chlyah
The μ-RA System for Recursive Path Queries over Graphs.
In Proceedings of the 32nd ACM International Conference on Information and Knowledge Management, 2023 (CIKM'23)

Amela Fejza, Pierre Genevès and Nabil Layaïda,

Efficient Enumeration of Recursive Plans in Transformation-based Query Optimizers. *Preprint, 2023: https://inria.hal.science/hal-03692274/file/rlqdag.pdf* (to appear in 2024).