

Context-Free Path Querying: Obstacles on the Way to Adoption

Semyon Grigorev

JetBrains Research, Programming Languages and Tools Lab
St. Petersburg State University

16.07.2021

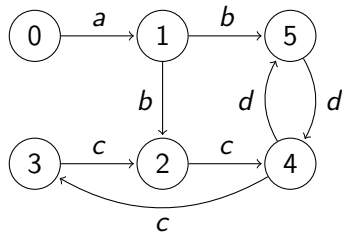
Who we are?

• !!!

• !!!

Formal language constrained path querying

Navigational queries in edge-labelled graph

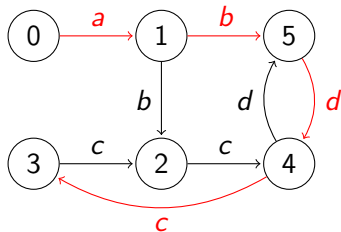


- Path to find:

$$0 \xrightarrow{a} v_0 \xrightarrow{b} v_1 \underbrace{\xrightarrow{d} v_2 \xrightarrow{c} v_3 \dots v_k}_{c \text{ or } d \text{ in arbitrary order}} \xrightarrow{c} v$$

Formal language constrained path querying

Navigational queries in edge-labelled graph

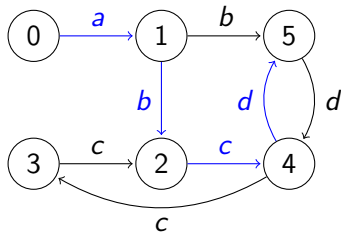


- Path to find:

$$0 \xrightarrow{a} v_0 \xrightarrow{b} v_1 \underbrace{\xrightarrow{d} v_2 \xrightarrow{c} v_3 \dots v_k}_{c \text{ or } d \text{ in arbitrary order}} \xrightarrow{c} v$$

Formal language constrained path querying

Navigational queries in edge-labelled graph

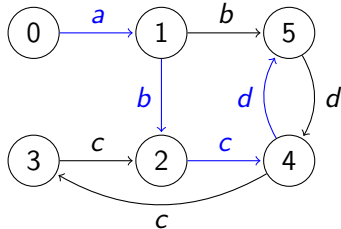


- Path to find:

$$0 \xrightarrow{a} v_0 \xrightarrow{b} v_1 \underbrace{\xrightarrow{d} v_2 \xrightarrow{c} v_3 \dots v_k}_{c \text{ or } d \text{ in arbitrary order}} \xrightarrow{c} v$$

Formal language constrained path querying

Navigational queries in edge-labelled graph



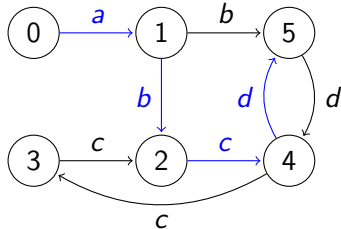
- $w(v_0 \xrightarrow{l_0} v_1 \xrightarrow{l_1} \dots \xrightarrow{l_{k-1}} v_k) = l_0 l_1 \dots l_{k-1}$
- $Q = \{(v_i, v_j) \mid \exists \pi = v_i \rightarrow \dots \rightarrow v_j; w(\pi) \in \mathcal{L}\}$,
where \mathcal{L} — formal language
 - ✓ Regular $(ab(c \mid d)^*)$
 - ⚙️ **Context-Free** $(a^n b^n)$
 - 🕒 Multiple Context-Free $(a^n c^m b^n d^m)$

- Path to find:

$$0 \xrightarrow{a} v_0 \xrightarrow{b} v_1 \underbrace{\xrightarrow{d} v_2 \xrightarrow{c} v_3 \dots v_k}_{c \text{ or } d \text{ in arbitrary order}} \xrightarrow{c} v$$

Formal language constrained path querying

Navigational queries in edge-labelled graph



- Path to find:

$$0 \xrightarrow{a} v_0 \xrightarrow{b} v_1 \xrightarrow{d} v_2 \xrightarrow{c} v_3 \dots v_k \xrightarrow{c} v$$

$\underbrace{\hspace{10em}}_{c \text{ or } d \text{ in arbitrary order}}$

- $w(v_0 \xrightarrow{l_0} v_1 \xrightarrow{l_1} \dots \xrightarrow{l_{k-1}} v_k) = l_0 l_1 \dots l_{k-1}$
- $Q = \{(v_i, v_j) \mid \exists \pi = v_i \rightarrow \dots \rightarrow v_j; w(\pi) \in \mathcal{L}\}$,
where \mathcal{L} — formal language
 - ✓ Regular $(ab(c \mid d)^*)$
 - ⚙️ **Context-Free** $(a^n b^n)$
 - ⌚ Multiple Context-Free $(a^n c^m b^n d^m)$

Variations:

- All-pairs
- Multiple source
- Reachability
- All paths
- ...

Applications of Context-Free Path Querying

Hierarchy analysis: variations of the *same-generation queries* is essence of CFPQ

¹Thomas Reps. 1997. “Program Analysis via Graph Reachability”.

²Mihalis Yannakakis. 1990. “Graph-theoretic Methods in Database Theory”.

Applications of Context-Free Path Querying

Hierarchy analysis: variations of the *same-generation queries* is essence of CFPQ

Graph databases

- Introduced by M. Yannakakis in 1990¹
- Biological data analysis
- Data provenance analysis
- ...

¹Thomas Reps. 1997. "Program Analysis via Graph Reachability".

²Mihalis Yannakakis. 1990. "Graph-theoretic Methods in Database Theory".

Applications of Context-Free Path Querying

Hierarchy analysis: variations of the *same-generation queries* is essence of CFPQ

Graph databases

- Introduced by M. Yannakakis in 1990¹
- Biological data analysis
- Data provenance analysis
- ...

Static code analysis

- Introduced by T. Reps in 1997²
- Interprocedural points-to analysis
- Interprocedural alias analysis
- Type inference related tasks
- ...

¹Thomas Reps. 1997. "Program Analysis via Graph Reachability".

²Mihalis Yannakakis. 1990. "Graph-theoretic Methods in Database Theory".

Applications of Context-Free Path Querying

Hierarchy analysis: variations of the *same-generation queries* is essence of CFPQ

Graph databases

- Introduced by M. Yannakakis in 1990¹
- Biological data analysis
- Data provenance analysis
- ...

Static code analysis

- Introduced by T. Reps in 1997²
- Interprocedural points-to analysis
- Interprocedural alias analysis
- Type inference related tasks
- ...

Graph databases

¹Thomas Reps. 1997. "Program Analysis via Graph Reachability".

²Mihalis Yannakakis. 1990. "Graph-theoretic Methods in Database Theory".

- There is no unified infrastructure for solutions comparison
 - ▶ Data is spread over articles in different communities
 - ▶ There is a huge number of different subclasses of the problem
 - ★ all-pairs, single source, multiple source, ...
 - ★ reachability, single path, all path, ...
 - ▶ The first and only attempt to compare different algorithms: “An Experimental Study of Context-Free Path Query Evaluation Methods”³

³Jochem Kuijpers, George Fletcher, Nikolay Yakovets, and Tobias Lindaaker. 2019

⁴H. Miao and A. Deshpande, “Understanding Data Science Lifecycle Provenance via Graph Segmentation and Summarization”, 2019

- There is no unified infrastructure for solutions comparison
 - ▶ Data is spread over articles in different communities
 - ▶ There is a huge number of different subclasses of the problem
 - ★ all-pairs, single source, multiple source, ...
 - ★ reachability, single path, all path, ...
 - ▶ The first and only attempt to compare different algorithms: “An Experimental Study of Context-Free Path Query Evaluation Methods”³
 - ▶ Conclusion: “We conclude that state of the art solutions are not able to cope with large graphs as found in practice.”

³Jochem Kuijpers, George Fletcher, Nikolay Yakovets, and Tobias Lindaaker. 2019

⁴H. Miao and A. Deshpande, “Understanding Data Science Lifecycle Provenance via Graph Segmentation and Summarization”, 2019

Problems

- There is no unified infrastructure for solutions comparison
 - ▶ Data is spread over articles in different communities
 - ▶ There is a huge number of different subclasses of the problem
 - ★ all-pairs, single source, multiple source, ...
 - ★ reachability, single path, all path, ...
 - ▶ The first and only attempt to compare different algorithms: “An Experimental Study of Context-Free Path Query Evaluation Methods”³
 - ▶ Conclusion: “We conclude that state of the art solutions are not able to cope with large graphs as found in practice.”
- There is no support in real-world graph database
 - ▶ H. Miao and A. Deshpande: “Though the problem has been first studied in our community [40], there is little follow up and support in the context of modern graph databases ...”⁴

³Jochem Kuijpers, George Fletcher, Nikolay Yakovets, and Tobias Lindaaker. 2019

⁴H. Miao and A. Deshpande, “Understanding Data Science Lifecycle Provenance via Graph Segmentation and Summarization”, 2019

Our Results

- ✓ Collection of linear algebra based algorithms for CFPQ
 - ▶ SuiteSparse is utilized for sparse linear algebra subroutines
 - ▶ Published: https://github.com/JetBrains-Research/CFPQ_PyAlgo

⁵Arseniy Terekhov et al. 2019. "Multiple-Source Context-Free Path Querying in Terms of Linear Algebra".

Our Results

- ✓ Collection of linear algebra based algorithms for CFPQ
 - ▶ SuiteSparse is utilized for sparse linear algebra subroutines
 - ▶ Published: https://github.com/JetBrains-Research/CFPQ_PyAlgo
- ✓ Full-stack support of CFPQ⁵
 - ▶ On the top of RedisGraph
 - ▶ openCypher extended to support CFPQ

⁵Arseniy Terekhov et al. 2019. "Multiple-Source Context-Free Path Querying in Terms of Linear Algebra".

Our Results

- ✓ Collection of linear algebra based algorithms for CFPQ
 - ▶ SuiteSparse is utilized for sparse linear algebra subroutines
 - ▶ Published: https://github.com/JetBrains-Research/CFPQ_PyAlgo
- ✓ Full-stack support of CFPQ⁵
 - ▶ On the top of RedisGraph
 - ▶ openCypher extended to support CFPQ
- ⚙ Collecting of the dataset for CFPQ benchmarking is started
 - ▶ Synthetic graphs
 - ▶ Real-world graphs
 - ★ Static code analysis
 - ★ Biological data analysis
 - ▶ Published: https://github.com/JetBrains-Research/CFPQ_Data

⁵Arseniy Terekhov et al. 2019. "Multiple-Source Context-Free Path Querying in Terms of Linear Algebra".

Our Results Evaluation

- All-pairs reachability queries
 - *geospecies*, *taxonomy* — biological data
 - *crypto*, *drivers*, *fs* — points-to analysis
 - Time in seconds
- GPU: Geforce GTX 1070, 1.5GHz, 8Gb RAM, 1920 CUDA cores
 - CPU: Intel core i7-6700 CPU, 3.4GHz, DDR4 64Gb RAM

Graph	#V	#E	Neo4j ⁶	RedisGraph ⁷	Lin.al. CPU ⁸	Lin.al. GPU ⁹
geospecies	450 609	2 311 461	6 953.9	80.1	7.1	0.8
taxonomy	5 728 398	14 922 125	n.a.	⚙️	1.1	0.7
crypto	3 464 970	5 976 774	n.a.	⚙️	84.8	28.1
drivers	4 273 803	7 415 538	n.a.	⚙️	269.9	62.5
fs	4 177 416	7 218 746	n.a.	⚙️	165.1	47.7

⁶Jochem Kuijpers et al. 2019. “An Experimental Study of Context-Free Path Query Evaluation Methods”.

⁷Arseniy Terekhov et al. 2019. “Multiple-Source Context-Free Path Querying in Terms of Linear Algebra”.

⁸Standalone linear algebra based algorithm on CPU, using SuiteSparse.

⁹Standalone linear algebra based algorithm on GPU, using spbla.

- ⚙️ Benchmarking of linear algebra based algorithms
 - ▶ Comparison of different algorithms for different query semantics
 - ▶ Investigation of scalability on multicore machines
 - ▶ Estimation of performance on GPGPU
- ⚙️ Developing and evaluating GLL-based CFPQ algorithm for Neo4j
 - ▶ Multiple-source
 - ▶ All paths and reachability-only

- ⚙️ Benchmarking of linear algebra based algorithms
 - ▶ Comparison of different algorithms for different query semantics
 - ▶ Investigation of scalability on multicore machines
 - ▶ Estimation of performance on GPGPU
- ⚙️ Developing and evaluating GLL-based CFPQ algorithm for Neo4j
 - ▶ Multiple-source
 - ▶ All paths and reachability-only
- ⌚ Describing semantics of openCypher in terms of linear algebra (in Coq)
- ⌚ Utilizing multiple context-free languages as path constraints

Topics to discuss

- Benchmarks.
 - ▶ Algorithms
 - ▶ Query language. Semantics
 - ▶ Local and Global queries. Can start from RPQ.
 - ▶ Queries: templates and real-world queries.
- Competition on language constrained path querying.
 - ▶ Part of existing competition for graph processing systems.
 - ▶ Involve static analysis community.
- Graph database support
 - ▶ Algorithms
 - ▶ Query language. Semantics