

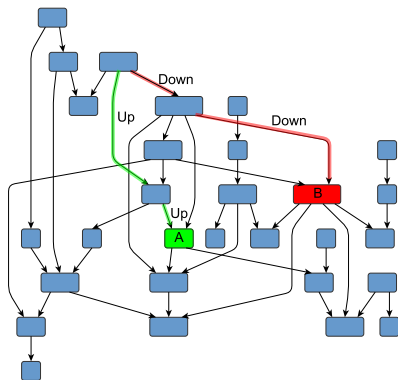
Context-Free Path Querying by Kronecker Product

Egor Orachev, Ilya Epelbaum,
Semyon Grigorev, **Rustam Azimov**

JetBrains Research, Programming Languages and Tools Lab
Saint Petersburg University

August 26, 2020

Context-Free Path Querying



Navigation through a graph

- Are nodes A and B on the same level of hierarchy?
- Is there a path of form $\text{Up}^n \text{Down}^n$?
- Find all paths of form $\text{Up}^n \text{Down}^n$ which start from the node A

- $\mathbb{G} = (\Sigma, N, P)$ — context-free grammar in normal form
 - ▶ $A \rightarrow BC$, where $A, B, C \in N$
 - ▶ $A \rightarrow x$, where $A \in N, x \in \Sigma \cup \{\varepsilon\}$
 - ▶ $L(\mathbb{G}, A) = \{\omega \mid A \Rightarrow^* \omega\}$

CFPQ: Query Semantics

- $\mathbb{G} = (\Sigma, N, P)$ — context-free grammar in normal form
 - ▶ $A \rightarrow BC$, where $A, B, C \in N$
 - ▶ $A \rightarrow x$, where $A \in N, x \in \Sigma \cup \{\varepsilon\}$
 - ▶ $L(\mathbb{G}, A) = \{\omega \mid A \Rightarrow^* \omega\}$
- $G = (V, E, L)$ — directed graph
 - ▶ $v \xrightarrow{I} u \in E$
 - ▶ $L \subseteq \Sigma$

CFPQ: Query Semantics

- $\mathbb{G} = (\Sigma, N, P)$ — context-free grammar in normal form
 - ▶ $A \rightarrow BC$, where $A, B, C \in N$
 - ▶ $A \rightarrow x$, where $A \in N, x \in \Sigma \cup \{\varepsilon\}$
 - ▶ $L(\mathbb{G}, A) = \{\omega \mid A \Rightarrow^* \omega\}$
- $G = (V, E, L)$ — directed graph
 - ▶ $v \xrightarrow{l} u \in E$
 - ▶ $L \subseteq \Sigma$
- $\omega(\pi) = \omega(v_0 \xrightarrow{l_0} v_1 \xrightarrow{l_1} \dots \xrightarrow{l_{n-2}} v_{n-1} \xrightarrow{l_{n-1}} v_n) = l_0 l_1 \dots l_{n-1}$

CFPQ: Query Semantics

- $\mathbb{G} = (\Sigma, N, P)$ — context-free grammar in normal form
 - ▶ $A \rightarrow BC$, where $A, B, C \in N$
 - ▶ $A \rightarrow x$, where $A \in N, x \in \Sigma \cup \{\varepsilon\}$
 - ▶ $L(\mathbb{G}, A) = \{\omega \mid A \Rightarrow^* \omega\}$
- $G = (V, E, L)$ — directed graph
 - ▶ $v \xrightarrow{l} u \in E$
 - ▶ $L \subseteq \Sigma$
- $\omega(\pi) = \omega(v_0 \xrightarrow{l_0} v_1 \xrightarrow{l_1} \dots \xrightarrow{l_{n-2}} v_{n-1} \xrightarrow{l_{n-1}} v_n) = l_0 l_1 \dots l_{n-1}$
- $R_A = \{(n, m) \mid \exists n \pi m, \text{ such that } \omega(\pi) \in L(\mathbb{G}, A)\}$

CFPQ: Original Matrix-Based Algorithm

Algorithm Context-free path querying algorithm

```
1: function EVALCFPQ( $D = (V, E, L), G = (\Sigma, N, P)$ )
2:    $n \leftarrow |V|$ 
3:    $T \leftarrow \{T^{A_i} \mid A_i \in N, T^{A_i} \text{ is a matrix } n \times n, T_{k,l}^{A_i} \leftarrow \text{false}\}$ 
4:   for all  $(i, x, j) \in E, A_k \mid A_k \rightarrow x \in P$  do  $T_{i,j}^{A_k} \leftarrow \text{true}$ 
5:   for all  $A_k \mid A_k \rightarrow \varepsilon \in P$  do
6:     for all  $i \in \{0, \dots, n-1\}$  do  $T_{i,i}^{A_k} \leftarrow \text{true}$ 
7:   while any matrix in  $T$  is changing do
8:     for  $A_i \rightarrow A_j A_k \in P$  do  $T^{A_i} \leftarrow T^{A_i} + (T^{A_j} \times T^{A_k})$ 
9:   return  $T$ 
```

Context-Free Path Querying: Grammar Transformation

- $\mathbb{G} = (\Sigma, N, P)$ — context-free grammar in general form
 - ▶ $A \rightarrow \alpha$, where $A \in N, \alpha \in (N \cup \Sigma \cup \{\varepsilon\})^*$

Context-Free Path Querying: Grammar Transformation

- $\mathbb{G} = (\Sigma, N, P)$ — context-free grammar in general form
 - ▶ $A \rightarrow \alpha$, where $A \in N, \alpha \in (N \cup \Sigma \cup \{\varepsilon\})^*$
- Every context-free grammar can be transformed to binary normal form

Context-Free Path Querying: Grammar Transformation

- $\mathbb{G} = (\Sigma, N, P)$ — context-free grammar in general form
 - ▶ $A \rightarrow \alpha$, where $A \in N, \alpha \in (N \cup \Sigma \cup \{\varepsilon\})^*$
- Every context-free grammar can be transformed to binary normal form
- The transformation takes time and can lead to a significant grammar size increase

Research Questions

- Can we create the matrix-based CFPQ algorithm that does not require grammar transformation?
- What matrix operations should be used?
- Does using matrix optimizations still significantly increases performance?

Recursive State Machines (RSM)

- RSM behaves as a set of finite state machines (FSM)
- Each FSM (box) works almost the same as a classical FSM, but it also handles additional recursive calls and employs an implicit call stack to call one component from another and then return execution flow back
- Any CFG can be easily encoded by an RSM with one box per nonterminal

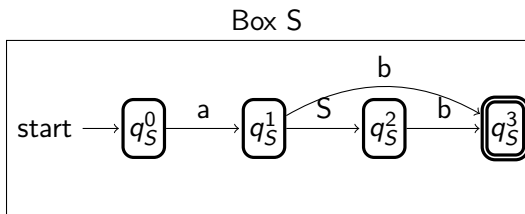


Figure: The RSM for grammar with rules $S \rightarrow aSb \mid ab$

Kronecker Product

- Kronecker product $A \otimes B$ for matrix A of size $m \times n$ and matrix B of size $p \times q$
 - ▶ Multiply each element of A and the matrix B
 - ▶ As a result we have $pm \times qn$ block matrix
- We need to intersect the FSMs from RSM generated by the context-free grammar and FSM generated by the graph
- Kronecker product can be used for constructing such intersections

Kronecker Product Based CFPQ Algorithm

Algorithm Kronecker product based CFPQ

```
1: function CONTEXTFREEPATHQUERYING( $G, \mathcal{G}$ )
2:    $R \leftarrow$  Recursive automata for  $G$ 
3:    $M_1, M_2 \leftarrow$  Adjacency matrices for  $R$  and  $\mathcal{G}$ 
4:   for  $s \in 0..dim(M_1) - 1$  do
5:     for  $i \in 0..dim(M_2) - 1$  do
6:        $M_2[i, i] \leftarrow M_2[i, i] \cup getNonterminals(R, s, s)$ 
7:   while Matrix  $M_2$  is changing do
8:      $M_3 \leftarrow M_1 \otimes M_2$  ▷ Evaluate Kronecker product
9:      $C_3 \leftarrow transitiveClosure(M_3)$ 
10:     $n \leftarrow dim(M_3)$  ▷ Matrix  $M_3$  size =  $n \times n$ 
11:    for  $i \in 0..n - 1$  do
12:      for  $j \in 0..n - 1$  do
13:        if  $C_3[i, j]$  then
14:           $s, f \leftarrow getStates(C_3, i, j)$ 
15:          if  $getNonterminals(R, s, f) \neq \emptyset$  then
16:             $x, y \leftarrow getCoordinates(C_3, i, j)$ 
17:             $M_2[x, y] \leftarrow M_2[x, y] \cup getNonterminals(R, s, f)$ 
18:  return  $M_2$ 
```

Kronecker Product Based CFPQ Algorithm: Technical Details

- We can use the sparse and block nature of the obtained matrices to apply wide class of optimizations
- We can operate over Boolean matrices
- We still can use existing high-performance math libraries for intersection if they provide a satisfying operation of element-wise multiplication

- **Kron** — implementation of the proposed algorithm using **SuiteSparse** C implementation of **GraphBLAS** API, which provides a set of sparse matrix operations

- **Kron** — implementation of the proposed algorithm using **SuiteSparse** C implementation of **GraphBLAS** API, which provides a set of sparse matrix operations
- We compare our implementation with **Orig** — the best CPU implementations of the original matrix-based algorithm using M4RI library with sparse matrix representation

- OS: Ubuntu 18.04
- CPU: Intel(R) Core(TM) i7-4790 CPU 3.60GHz
- RAM: DDR4 32 Gb

Evaluation results^{1 2}

	Graph	#V	#E	Kron	Orig		Graph	#V	#E	Kron	Orig
RDF	generations	129	351	0.04	0.03	RDF	core	1323	8684	0.28	0.12
	travel	131	397	0.05	0.05		pways	6238	37196	4.88	0.18
	skos	144	323	0.02	0.04	Worst case	WC ₁	64	65	0.03	0.04
	unv-bnch	179	413	0.05	0.04		WC ₂	128	129	0.16	0.23
	foaf	256	815	0.07	0.02		WC ₃	256	257	0.96	1.99
	atm-prim	291	685	0.24	0.02		WC ₄	512	513	7.14	23.21
	ppl_pets	337	834	0.18	0.03		WC ₅	1024	1025	121.99	528.52
	biomed	341	711	0.24	0.05	Full	F ₁	100	100	0.17	0.02
	pizza	671	2604	1.14	0.08		F ₂	200	200	1.04	0.03
	wine	733	2450	1.71	0.06		F ₃	500	500	18.86	0.03
	funding	778	1480	0.43	0.07		F ₄	1000	1000	554.22	0.07

¹Queries are based on the context-free grammars for nested parentheses

²Time is measured in seconds

Conclusion

- We show that in some cases our algorithm outperforms the original matrix-based algorithm

Conclusion

- We show that in some cases our algorithm outperforms the original matrix-based algorithm
- We still can use existing high-performance libraries for matrix operations

Conclusion

- We show that in some cases our algorithm outperforms the original matrix-based algorithm
- We still can use existing high-performance libraries for matrix operations
- The idea of the proposed algorithm looks viable

Conclusion

- We show that in some cases our algorithm outperforms the original matrix-based algorithm
- We still can use existing high-performance libraries for matrix operations
- The idea of the proposed algorithm looks viable
- Dataset is published: both graphs and queries
 - ▶ Link: https://github.com/JetBrains-Research/CFPQ_Data

- Improve our implementation to make it applicable for real-world graphs analysis

Future Research

- Improve our implementation to make it applicable for real-world graphs analysis
- Investigate such formal properties of the proposed algorithm as time and space complexity

- Improve our implementation to make it applicable for real-world graphs analysis
- Investigate such formal properties of the proposed algorithm as time and space complexity
- Analyze how the behavior depends on the query type and its form
 - ▶ Analyze regular path queries evaluation and context-free path queries in the form of extended context-free grammars (ECFG)

- Improve our implementation to make it applicable for real-world graphs analysis
- Investigate such formal properties of the proposed algorithm as time and space complexity
- Analyze how the behavior depends on the query type and its form
 - ▶ Analyze regular path queries evaluation and context-free path queries in the form of extended context-free grammars (ECFG)
- Compare our algorithm with the matrix-based one in cases when the size difference between Chomsky Normal Form and ECFG representation of the query is significant.

- Improve our implementation to make it applicable for real-world graphs analysis
- Investigate such formal properties of the proposed algorithm as time and space complexity
- Analyze how the behavior depends on the query type and its form
 - ▶ Analyze regular path queries evaluation and context-free path queries in the form of extended context-free grammars (ECFG)
- Compare our algorithm with the matrix-based one in cases when the size difference between Chomsky Normal Form and ECFG representation of the query is significant.
- Extend our algorithm to single-path and all-path query semantics

Contact Information

- Semyon Grigorev:
 - ▶ s.v.grigoriev@spbu.ru
 - ▶ Semen.Grigorev@jetbrains.com
- Rustam Azimov:
 - ▶ rustam.azimov19021995@gmail.com
 - ▶ Rustam.Azimov@jetbrains.com
- Egor Orachev: egor.orachev@gmail.com
- Ilya Epelbaum: iliyepelbaun@gmail.com

- Dataset: https://github.com/JetBrains-Research/CFPQ_Data

Thanks!