

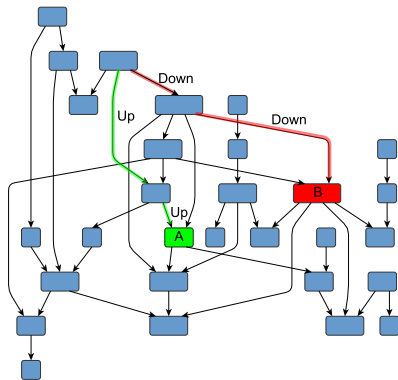
Context-Free Path Querying by Matrix Multiplication

Rustam Azimov, Semyon Grigorev
presents Kate Verbitskaia

JetBrains Research, Programming Languages and Tools Lab
Saint Petersburg University

June 10, 2018

Language-Constrained Path Filtering



Navigation through a graph

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

Language-Constrained Path Filtering: Relational 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 $A \in N, x \in \Sigma$
- $G = (V, E, L)$ — directed graph
 - ▶ $v \xrightarrow{l} u \in E$
 - ▶ $L \subseteq \Sigma$
- $\omega(\pi) = \omega(v_0 \xrightarrow{l_0} v_1 \xrightarrow{l_1} \dots \xrightarrow{l_{n-2}} v_{n-1} \xrightarrow{l_{n-1}} v_n) = l_0 l_1 \dots l_{n-1}$
- $R_A = \{(n, m) \mid \exists n\pi m, \text{ such that } \omega(\pi) \in L(\mathbb{G}, A)\}$

Regular language constraints

- Widely spread
 - ▶ Graph databases query languages (SPARQL, Cypher, PGQL)
 - ▶ Network analysis
- Still in active development
 - ▶ OpenCypher: <https://goo.gl/5h5a8P>
 - ▶ Scalability, huge graphs processing
 - ▶ Derivatives for graph querying: *Maurizio Nole and Carlo Sartiani*. Regular path queries on massive graphs. 2016

Context-free language constraints

- Graph databases and semantic networks (Context-Free Path Querying, CFPQ)
 - ▶ *Sevon P., Eronen L.* "Subgraph queries by context-free grammars." 2008
 - ▶ *Hellings J.* "Conjunctive context-free path queries." 2014
 - ▶ *Zhang X. et al.* "Context-free path queries on RDF graphs." 2016
- Static code analysis (Language Reachability Framework)
 - ▶ *Thomas Reps et al.* "Precise interprocedural dataflow analysis via graph reachability." 1995
 - ▶ *Qirun Zhang et al.* "Efficient subcubic alias analysis for C." 2014
 - ▶ *Dacong Yan et al.* "Demand-driven context-sensitive alias analysis for Java." 2011
 - ▶ *Jakob Rehof and Manuel Fahndrich.* "Type-base flow analysis: from polymorphic subtyping to CFL-reachability." 2001

Open Problems

- Development of efficient algorithms
- Effective utilization of GPGPU and parallel programming
- Lifting up the limitations on the input graph and the query language

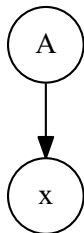
The algorithm

Algorithm 1 Context-free recognizer for graphs

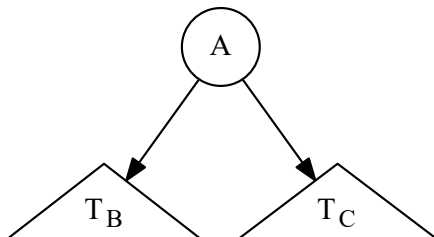
```
1: function CONTEXTFREEPATHQUERYING( $D, G$ )
2:    $n \leftarrow$  the number of nodes in  $D$ 
3:    $E \leftarrow$  the directed edge-relation from  $D$ 
4:    $P \leftarrow$  the set of production rules in  $G$ 
5:    $T \leftarrow$  the matrix  $n \times n$  in which each element is  $\emptyset$ 
6:   for all  $(i, x, j) \in E$  do ▷ Matrix initialization
7:      $T_{i,j} \leftarrow T_{i,j} \cup \{A \mid (A \rightarrow x) \in P\}$ 
8:   while matrix  $T$  is changing do
9:      $T \leftarrow T \cup (T \times T)$  ▷ Transitive closure  $T^{cf}$  calculation
10:  return  $T$ 
```

Derivation step

$$A \rightarrow x$$



$$A \rightarrow BC$$



Matrix multiplication

- Subset multiplication, $N_1, N_2 \subseteq N$
 - ▶ $N_1 \cdot N_2 = \{A \mid \exists B \in N_1, \exists C \in N_2 \text{ such that } (A \rightarrow BC) \in P\}$
- Subset addition: set-theoretic union.
- Matrix multiplication
 - ▶ Matrix of size $|V| \times |V|$
 - ▶ Subsets of N are elements
 - ▶ $c_{i,j} = \bigcup_{k=1}^n a_{i,k} \cdot b_{k,j}$

Transitive closure

- $a^{cf} = a^{(1)} \cup a^{(2)} \cup \dots$
- $a^{(1)} = a$
- $a^{(i)} = a^{(i-1)} \cup (a^{(i-1)} \times a^{(i-1)}), \quad i \geq 2$

Theorem

*Let $D = (V, E)$ be a graph and let $G = (N, \Sigma, P)$ be a grammar.
Then for any i, j and for any non-terminal $A \in N$, $A \in a_{i,j}^{cf}$ iff $(i, j) \in R_A$.*

Theorem

*Let $D = (V, E)$ be a graph and let $G = (N, \Sigma, P)$ be a grammar.
The Algorithm 1 terminates in a finite number of steps.*

Algorithm Complexity

Theorem

Let $D = (V, E)$ be a graph and let $G = (N, \Sigma, P)$ be a grammar. The Algorithm 1 calculates the transitive closure T^{cf} in $O(|V|^2|N|^3(BMM(|V|) + BMU(|V|)))$.

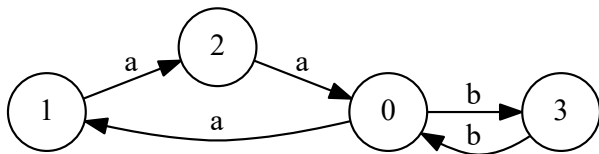
- $BMM(n)$ — number of elementary operations executed by the algorithm of multiplying two $n \times n$ Boolean matrices.
- $BMU(n)$ — number of elementary operations, executed by the matrix union operation of two $n \times n$ Boolean matrices

Algorithm Complexity: the Worst Case

Input graph: two cycles connected via a shared node

- first cycle has $2^k + 1$ edges labeled a
- second cycle has 2^k edges labeled b

$$\begin{array}{ccc} S & \rightarrow & a S b \\ | & & a b \end{array}$$



- dGPU (dense GPU): row-major matrix representation and a GPU for matrix operation calculation.
- sCPU (sparse CPU): CSR format for sparse matrix representation and a CPU for matrix operation calculation.
- sGPU (sparse GPU): CSR format for sparse matrix representation and a GPU for matrix operation calculation.

Evaluation: Same Generation Queries

Query 1 retrieves the concepts on the same layer

$$\begin{array}{lcl} S & \rightarrow & subClassOf^{-1} S subClassOf \\ & | & type^{-1} S type \\ & | & subClassOf^{-1} subClassOf \\ & | & type^{-1} type \end{array}$$

Query 2 retrieves concepts on the adjacent layers

$$\begin{array}{lcl} S & \rightarrow & B subClassOf \\ S & | & subClassOf \\ B & \rightarrow & subClassOf^{-1} B subClassOf \\ B & | & subClassOf^{-1} subClassOf \end{array}$$

Evaluation: Query 1

Ontology	V	E	#results	GLL(ms)	dGPU
skos	144	323	810	10	56
generations	129	351	2164	19	62
travel	131	397	2499	24	69
univ-bench	179	413	2540	25	83
atom-primitive	291	685	15454	255	19
biomedical-measure-primitive	341	711	15156	261	26
foaf	256	815	4118	39	15
people-pets	337	834	9472	89	39
funding	778	1480	17634	212	141
wine	733	2450	66572	819	204
pizza	671	2604	56195	697	110
g_1	6224	11840	141072	1926	—
g_2	5864	19600	532576	6246	—
g_3	5368	20832	449560	7014	—

Evaluation: Query 2

Ontology	V	E	#results	GLL(ms)	dGPU
skos	144	323	1	1	10
generations	129	351	0	1	9
travel	131	397	63	1	33
univ-bench	179	413	81	11	55
atom-primitive	291	685	122	66	36
biomedical-measure-primitive	341	711	2871	45	27
foaf	256	815	10	2	53
people-pets	337	834	37	3	14
funding	778	1480	1158	23	124
wine	733	2450	133	8	72
pizza	671	2604	1262	29	94
g_1	6224	11840	9264	167	—
g_2	5864	19600	1064	46	—
g_3	5368	20832	10096	393	—

Something about how this algorithm is matrix operations independent and great and stuff

Contact Information

- Semyon Grigorev: s.v.grigoriev@spbu.ru
- Rustam Azimov: st013567@student.spbu.ru
- YaccConstructor: <https://github.com/YaccConstructor>