

Generalized LL parsing for context-free constrained path search problem

Semyon Grigorev
Saint Petersburg State University
7/9 Universitetskaya nab.
St. Petersburg, 199034 Russia
semen.grigorev@jetbrains.com

Anastasiya Ragozina
Saint Petersburg State University
7/9 Universitetskaya nab.
St. Petersburg, 199034 Russia
ragozina.anastasiya@gmail.com

ABSTRACT

Aaaabstract is very abstract.... word1 word2 word3 word4
word5 word6 word7 word8 word9 word10 word1 word2
word3 word4 word5 word6 word7 word8 word9 word10
word1 word2 word3 word4 word5 word6 word7 word8
word9 word10 word1 word2 word3 word4 word5 word6
word7 word8 word9 word10 word1 word2 word3 word4
word5 word6 word7 word8 word9 word10 word1 word2
word3 word4 word5 word6 word7 word8 word9 word10
word1 word2 word3 word4 word5 word6 word7 word8 word9
word10 word1 word2 word3 word4 word5 word6 word7
word8 word9 word10 word1 word2 word3 word4 word5
word6 word7 word8 word9 word10 word1 word2 word3
word4 word5 word6 word7 word8 word9 word10

1. INTRODUCTION

Graph data model and graph data bases are very popular in many different areas such as bioinformatic, semantic web, social networks etc. Extraction of paths satisfying specific constraints may be useful for graph structured data investigation and for relations between data items detection. Path querying with constraints formulated in terms of formal grammars is a specific problem named formal language constrained path problem [3] and research in this area is still actual [8].

Classical parsing techniques can be used to solve formal language constrained path problem thus we can use “graph parsing” and it may be required not only in graph data base querying but also in other different areas: formal verification, and string-embedded language processing, for example.

Classical solution in DB area use such parsing algorithms as CYK or Earley. In string-embedded languages analysis (RN)GLR is used. It has better time complexity.

Graph parsing can be also used in string-embedde languages processing. Regular approximation for value set of string variable can by represened as directed graph of related fite automata. In orded to check corectness or safety (sql injections)... all generated strings (all paths from start states to final states) are correct w.r.t some context-free grammar. For example grammar of one of SQL dialects. GLR-based for string-embedded SQL checking [2, 4]. Solution based on RNGLR [11] for relaxed parsing of string-embedded languages [20] which allow to find all path between two specified vertices.

despite of the fact There is set of solutions in DB area but !!! Query result exploration is a challenge [6]. Classical parsing allow to construct derivation tree which is structural representation of parsed sentence. Structural represenatation of

query result can be useful for its exploration and qury debugging. In this paper, we propose graph parsing technique which allow to construct structural representation of query result with reletion to grammar query or derivation of result.

So, we propose algorithm based top-down parsing algorithm. LL is more natural and so on.

2. PRELIMINARIES

In this work we are focused on parsing algorithm, and not on the data representation, and de assume that full input graph can be located in RAM memory by the optimal for our algorithm way.

Also we need to introduce some definitions.

- Context-free grammar $G = (N, \Sigma, P, S)$ where N is a set of nonterminal symbols, Σ is a set of terminal symbols, $S \in N$ is a start nonterminal, and P is a productions set.
- $\mathcal{L}(G)$ is a language specified by grammar G .
- Directed graph $M = (V, E, L)$ where V — vertices set, $L \subseteq \Sigma$ — edge labels set, $E \subseteq V \times L \times V$. We assume that there are no parallel edges with equal labebs: for every $e_1 = (v_1, l_1, v_2) \in E, e_2 = (u_1, l_2, u_2) \in E$ if $v_1 = u_1$ and $v_2 = u_2$ then $l_1 \neq l_2$.
- Helper function for edge’s tag calculation $tag : E \rightarrow L$.

$$tag(e = (v_1, l, v_2), e \in E) = l$$

- Concatenation operation $\oplus : L^+ \times L^+ \rightarrow L^+$.
- Path p in graph M .

$$p = (v_0, l_0, v_1), (v_1, l_1, v_2), \dots, (v_{n-1}, l_{n-1}, v_n) \\ = e_0, e_1, \dots, e_{n-1}$$

where $v_i \in V, e_i \in E, l_i \in L, |p| = n, n \geq 1$.

- Set of paths $P = \{p : p \text{ path in } M\}$
- Helper function for string produced by path calculation $\Omega : P \rightarrow L^+$.

$$\Omega(p = e_0, e_1, \dots, e_{n-1}, p \in P) = \\ tag(e_0) \oplus \dots \oplus tag(e_{n-1})$$

As a result we can define that context-free language constrained path querying means that we get query as grammar G and result of this query is a set of paths

$$P = \{p | \Omega(p) \in \mathcal{L}(G)\}.$$

For some graphs and some queries P can be infinite set, and it can not be explicitly represented. In order to solve this problem, in this paper, we will construct not explicit representation of P but compact data structure which store all elements of P in finite space and allow to extract any of them. In this point our solution is slightly similar to subgraph quering proposed in article [16], but we are also construct derivation forest for result subgraph.

3. EXAMPLE

As a motivation of context-free constraints importance let we introduce the next example. Let we have graph $M = (\{0; 1; 2; 3\}, E, \{A; B\})$ presented in figure 1 where labels represent next relations:

- A — “is friend of” (v_0Av_1 means v_0 is friend of v_1);
- B — “has friend” (v_0Bv_1 means v_0 has friend v_1).

Suppose for each $n \geq 1$ we want to find all n -th generation friends with a common ancestor. In the other worlds, we wath to find all paths p , such that $\Omega(p) \in \{AB; AAB; AAAB; \dots\}$ or $\Omega(p) = A^n B^n$ where $n \geq 1$. This constraint can not be specified with regular language as far as $L = \{A^n B^n; n \geq 1\}$ is not regular but context free. Required language can be specified by grammar G_1 presented in picture 2 where $N = \{s; middle\}$, $\Sigma = \{A; B\}$, and $S = s$.

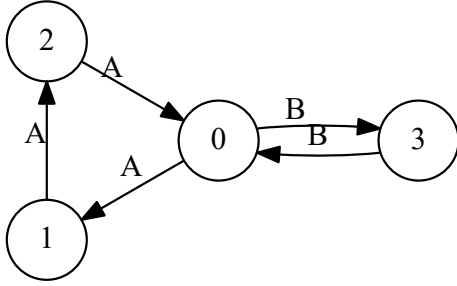


Figure 1: Input graph M

0: $s = L s R$
 1: $s = middle$
 2: $middle = L R$

Figure 2: Grammar G_1 for language $L = \{L^n R^n; n \geq 1\}$

Result is infinite for this query, and we can not... Also we want to know, who is common accesor. Futher we show ho to solve it.

4. ALGORITHM FOR GRAPH PARSING

We propose a context-free language constrained path problem solution which allow to find all paths satisfied specified arbitrary context-free grammar and to construct implicit representation of result. Finite representation of result set with structure related to specified grammar may be useful not only for results understanding and processing but also for query debugging especially for complex queries.

Our solution is based on generalized LL (GLL) [12, 1] parsing algorithm which allow to process ambiguous context-free grammars. Complexity is $O(n^3)$ in worst case and linear for unambiguous grammars, that better then complexity of CYK and Earley which used as base in other solutions (for example [5], [16]). This fact allow to demonstarte better performance on linear subgraphs and unambiguous grammars. Also it is not necessary to transform input grammar to CNF which required for CYK which allow to avoid grmmar size decreasing. It is important because real performance of parsing algorithm is sensitive to grammar size.

4.1 Generalized LL Parsing Algorithm

Generalized LL (GLL) is generalized top-down parsing algorithm which handle all context-free grammars (including left recursive) with worst-case cubic time complexity and linear for LL grammars. GLL is native for grammar, can be simple created blah-blah-blah.

GLL use descriptors for parsing states specification. Each descriptor contains full specification of process state enough to start parsing from state stored.

Descriptor is a triple (L, s, j) where L is a line label, s is a stack and j is a position in the input.

allows to restore parsing

Graph structured stack (GSS) [18] for multiple stack combining to prevent duplication. In GLL each GSS node is pair of position in input and grammar slot. Grammar slot is a !!!

4.2 Shared packed parse forest

Shared Packed Parse Forest (SPPF) is a spetial data structure for derivation forest compact representation which allow to reuse common nodes and subtrees. As a result multiple derivation trees, which csn be produced in case of ambiguos grammar, can be compressed in one SPPF with optimal reusing of common parts. Binarized form of SPPF proposed in [15] and it allow to achive worst-case cubic space complexity. GLL can use SPPF [13] for results representation achive cubic space complexity with binarised version.

Let we present an example of SPPF for ambiguos grammar G_0 (pic 3).

0: $s = NUM$
 1: $s = LBR s RBR$
 2: $s = s s$

Figure 3: Grammar G_0

Here N is token for number, L and R are tokens for '(' and ')' respectively.

Let we parse the sentence (1) (2) (3). There are two diferent leftmost derivations of this sentence in grammar G_0 (\rightarrow^n denote an application of production with nimber n):

1. $s \xrightarrow{2} ss \xrightarrow{2} sss \xrightarrow{1} LsRss \xrightarrow{0} LNRss \xrightarrow{1} LNRLsRs \xrightarrow{1} LNRLsRs \xrightarrow{0} LNRLNRs \xrightarrow{1} LNRLNRsR \xrightarrow{0} LNRLNRsR$
2. $s \xrightarrow{2} ss \xrightarrow{1} LsRs \xrightarrow{0} LNRs \xrightarrow{2} LNRss \xrightarrow{1} LNRLsRs \xrightarrow{1} LNRLsRs \xrightarrow{0} LNRLNRs \xrightarrow{1} LNRLNRsR \xrightarrow{0} LNRLNRsR$

As far as there are tho different derivations, SPPF should contains 2 different trees and it is presented in figure 4: re-

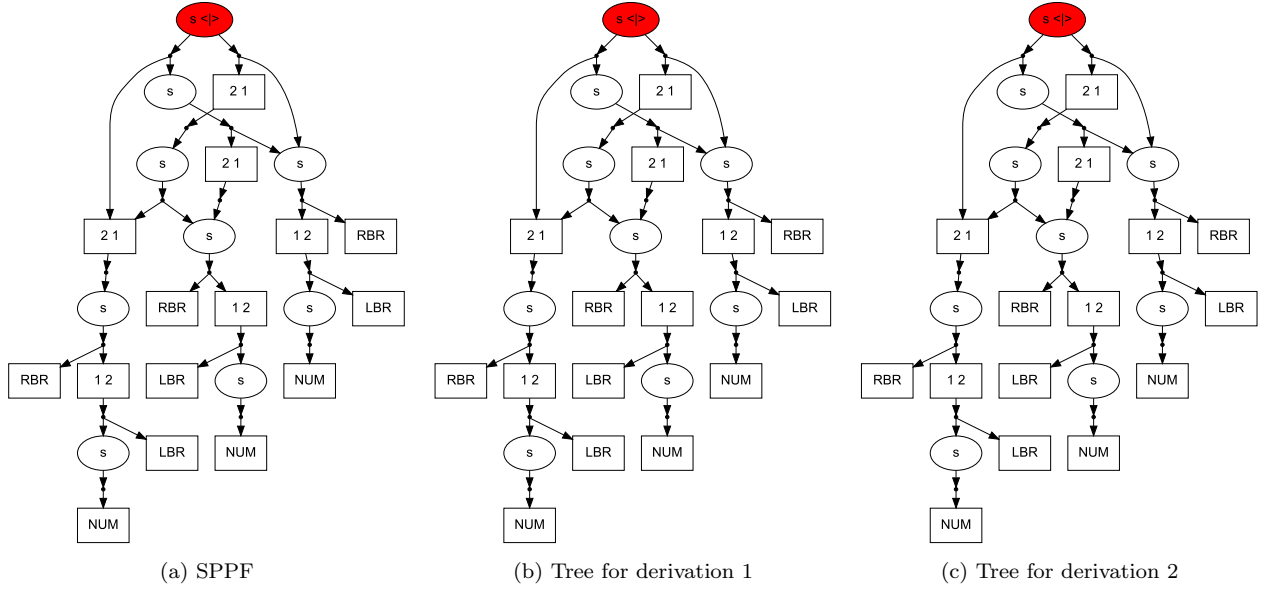


Figure 4: SPPF for sentence (1) (2) (3) and grammar G_0

sult SPPF 4a and trees for derivation 1 4b and derivation 2 4b respectively.

Binarised SPPF is a graph where !!! and each node has one of four types and one node marked as 'root' — node for start nonterminal.

- terminal node
- nonterminal node
- intermediate node
- packed node

Further we will remove redudant intermediate and packed nodes from SPPF to simplify it and decrease size of structure.

4.3 GLL-based graph parsing

$\mathbb{P} : G, M \rightarrow SPPF$

In order to use GLL for graph parsing we need only use graph verticea as position in input. As far as we work with context-free languages it is not important how this descriptor was created, and so descriptors management and other basic mechanisms of original algorithm can be reused “as is”. We can merge it if thea are equal.

We implement some optimizations: [1]

We also use binarised SPPF for result representation which allow to simplify query debugging and result exploration. In our case more then one root may be specified. For example, look at picture!!!! We

Binarized SPPF is at most cubic in terms of result size. Any path can be extracted in the linear time.

4.4 Complexity

Time complexity estimation in terms of input graph and grammar size is pretty similar to estimation of GLL complexity provided in [13].

LEMMA 1. For any descriptor (L, u, i, w) either $w = \$$ or w has extension (j, i) where u has index j .

PROOF. Proof of this lemma is the same as provided for riginal GLL in [13] because main function used for descriptor creation are the same as original one. \square

THEOREM 1. The GSS generated by GLL-based graph parsing algorithm for grammar G on input graph $M = (V, E, L)$ has at most $O(|V|)$ vertices and $O(|V|^2)$ edges.

PROOF. Proof the same as the proof of **Theorem 2** from [13].

\square

THEOREM 2. The SPPF generated by GLL-based graph parsing algorithm on input graph $M = (V, E, L)$ has at most $O(|V|^3 + |E|)$ vertices and edges.

PROOF. Let we estimate number of nodes of each type.

- Terminal nodes. Each of them has label of form (T, v_0, v_1) , and such lable can be created only if there is such $e \in E$ that $e = (v_0, T, v_1)$. Note, that there are no duplicate edges. Hence tere are at most $|E|$ terminal nodes.
- ϵ nodes labled with (ϵ, v, v) , hence there are at most $|E|$ of these.
- Nonterminal nodes have label of form (N, v_0, v_1) , so there are at most $O(|V|^2)$ of these.
- Intermideate nodes have label of form (t, v_0, v_1) , where t is grammar slot, so there are at most $O(|V|^2)$ of these.
- Packed nodes are children of intermediate or nonterminal nodes and have label of form (t, v) where t is a grammar slot $N : \alpha \cdot \beta$. There are at most $O(|V|^2)$ parents for packed nodes and each of them can have at most $O(|V|)$ children.

As a result there are at most $O(|V|^3 + |E|)$ nodes in SPPF.

The packed nodes have at most two children so there are at most $O(|V|^3 + |E|)$ edges with source in packed node. Nonterminal and intermediate nodes have at most $O(|V|)$ children and all of them are packed nodes. Thus there are at most $O(|V|^3)$ edges with source in nonterminal or intermediate nodes. As a result there are at most $O(|V|^3 + |E|)$ edges in SPPF.

□

THEOREM 3. *The space complexity of GLL-based graph parsing algorithm for graph $M = (V, E, L)$ is at most $O(|V|^3 + |E|)$.*

PROOF. From theorems 1 and 2 we have that space required for main data structures is at most $O(|V|^3 + |E|)$.

□

THEOREM 4. *The runtime complexity of GLL-based graph parsing algorithm for graph $M = (V, E, L)$ is at most*

$$O\left(|V|^3 * \max_{v \in V} (deg^+(v))\right).$$

PROOF. From Lemma 1 we get that there are at most $O(|V|^2)$ descriptors. Complexity of all functions are the same as in proof of **Theorem 4** from [13] except *processing* function where we should process not one next input token, but all outgoing edges. Thus for each descriptor we should examine at most

$$\max_{v \in V} (deg^+(v))$$

edges where $deg^+(v)$ is outdegree of vertex v .

So, worst-case complexity of proposed algorithm is

$$O\left(V^3 * \max_{v \in V} (deg^+(v))\right).$$

□

From theorem (4) we can get estimations for linear input and for LL grammars: for any $v \in V$ $deg^+(v) \leq 1$, so $\max(deg^+(v)) = 1$ and we get $O(|V|^3)$. For LL grammars and linear input complexity should be $O(|V|)$ for the same reason as for original GLL.

As discussed in [7] achieving of theoretical complexity required special datastructures which can be irrational for practice implementation and it is necessary to find balance between performance, software complexity, and hardware resources. As a result in practice we can get slightly worse performance than theoretical estimation.

Note that result SPPF contains only paths matched specified query, so result SPPF size is $O(|V'|^3 + |E'|)$ where $M' = (V', E', L')$ is a subgraph of input graph M which contains only matched paths. Also note that each specific path can be explored with linear SPPF traversal.

4.5 Example

In details, main function input is graph M , set of start vertices $V_s \subseteq V$, set of final vertices $V_f \subseteq V$, grammar G_1 . Output is Shared Packed Parse Forest (SPPF) [10] — finite data structure which contains all derivation trees for all paths in M , $\Omega(p) \in L(G_1)$ and allows to reconstruct any of

paths implicitly. As far as we can specify sets of start and final vertices, our solution can find all paths in graph, all paths from specified vertex, all paths between specified vertices. Also SPPF represents a structure of paths in terms of derivation which allow to get more useful information about result.

Let us introduce the next example. Grammar G_1 is a query and we want to find all paths in graph M (presented in picture 1) matched this query. Result SPPF for this input is presented in picture 5. Note that presented version does not contain obsolete nodes.

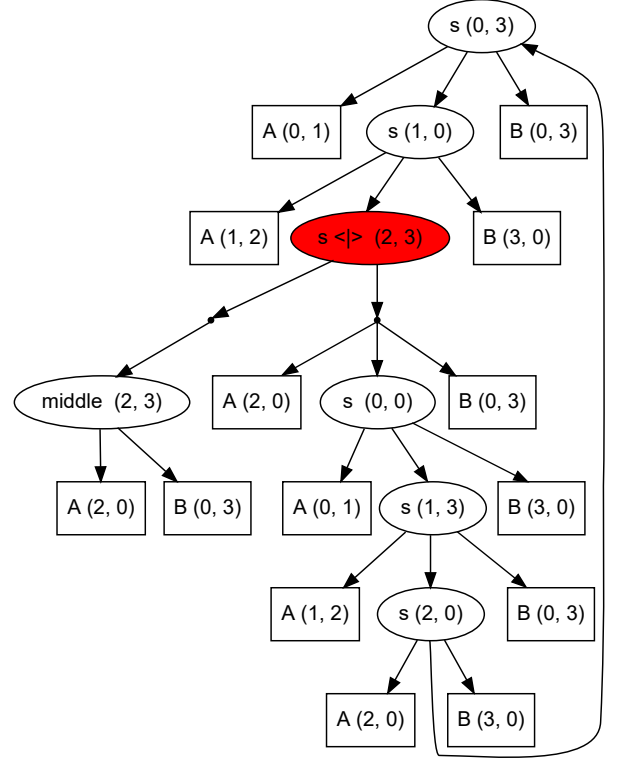


Figure 5: Result SPPF for input graph M (pic. 1) and query G_1 (pic. 2)

We use next markers for nodes which similar to original SPPF but have some additional information in order to relation with graph.

- Node with rectangle shape labeled with (v_0, T, v_1) is terminal node. Each terminal node corresponds with edge in the input graph: for each node with label (v_0, T, v_1) there is $e \in E : e = (v_0, T, v_1)$. Duplication of terminal nodes is only for figure simplification.
- Node with oval shape labeled with (v_0, nt, v_1) is non-terminal node. This node denotes that there is at least one path p from vertex v_0 to vertex v_1 in input graph M such that $nt \Rightarrow_G^* \Omega(p)$. All paths matched this condition can be extracted from SPPF by left-to-right top-down graph traversal started from respective node.
- Filled node with oval shape labeled with $(<|>, v_0, nt, v_1)$ is nonterminal node denoting that there are more than one path p from v_0 to v_1 such that $nt \Rightarrow_G^* \Omega(p)$.

- Node with dot shape is used for representation of derivation variants. Subgraph with root in one such node is one variant of derivation. Parent of such nodes is always node with label $(\langle \rangle (v_0, nt, v_1))$.
- v_0 and v_1 are left and right extensions of node respectively.

As an example of derivation structure usage we can find 'middle' of any path in example above simply by finding corresponded nonterminal *middle* in SPPF. So we can found that there is only one common ancestor for all results and it is vertex with $id = 0$.

Extensions stored in nodes allow to check whether path from u to v exists and extract it. To extract specified path we need only travers SPPF which can be done in linear time. Let for example we want to find path satisfying specified constraints from vertex 0. To do this we should find vertices with label $(0, s, -)$ in SPPF. There are two vertices with required label: $(0, s, 0)$ and $(0, s, 3)$. In our example there is cycle in SPPF so there are **at least** two different paths:

$p_0 = \{(0, A, 1); (1, A, 2); (2, A, 0); (0, B, 3); (3, B, 0); (0, B, 3)\}$
and

$p_1 = \{(0, A, 1); (1, A, 2); (2, A, 0); (0, A, 1); (1, A, 2); (2, A, 0);$
 $(0, B, 3); (3, B, 0); (0, B, 3); (3, B, 0); (0, B, 3); (3, B, 0)\}.$

5. EVALUATION

We use two grammars for balanced brackets in order to investigate performance relations with grammar ambiguity: ambiguous grammar G_0 3 and unambiguous grammar G_2 6.

0: $s = L s R s$
1: $s = \text{eps}$

Figure 6: Unambiguous grammar G_2 for balanced brackets

For input we use complete graphs where for each terminal symbol there is edge between two vertices labeled with it. Task we solve in performed experiments is find all paths from all vertices to all vertices satisfied specified query.

All tests were performed on a PC with following characteristics:

- OS Name: Microsoft Windows 10 Pro
- System Type: x64-based PC
- CPU: Intel(R) Core(TM) i7-4790 CPU @ 3.60GHz, 3601 Mhz, 4 Core(s), 4 Logical Processor(s)
- RAM: 32 GB

Performance measurement results presented in figure 7. As far as for complete graph $M = (V, E, L)$

$$\max_{v \in V} (\deg^+(v)) = (|V| - 1) * |\Sigma|$$

where Σ is terminals of input grammar, we should get time complexity at most $O(|V|^4)$ and space complexity at most $O(|V|^3)$. All two curves can be fit with polynomial function of degree 4 to a high level of confidence with R^2 .

Also we present SPPF size in terms of nodes 8 for both G_0 and G_2 grammars. All two curves are cubic to a high level of confidence with $R^2 = 1$.

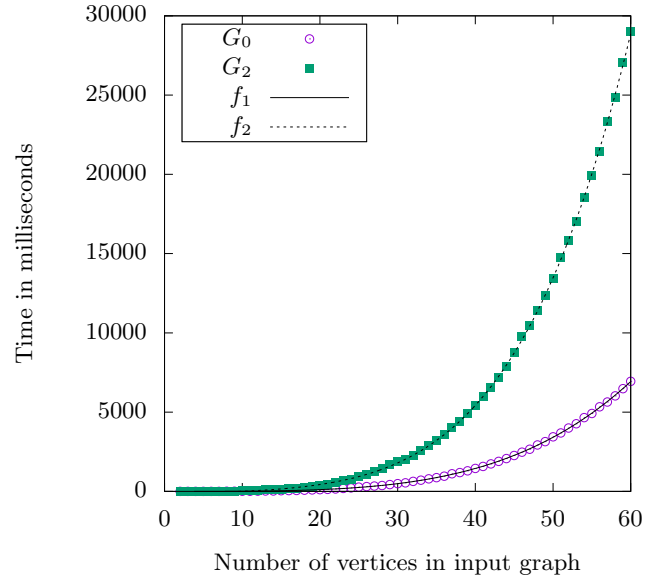


Figure 7: Performance on complete graphs for grammars G_0 and G_2

$$f_1(x) = 0.000495989 * x^4 + 0.001252184 * x^3 + 0.068491746 * x^2 - 0.306749160 * x; R^2 = 0.99996$$

$$f_2(x) = 0.003368883 * x^4 - 0.114919298 * x^3 + 3.161793404 * x^2 - 22.549491142 * x; R^2 = 0.99995$$

6. CONCLUSION AND FUTURE WORK

We propose GLL-based algorithm for context-free path querying which construct finite structural representation of all paths satisfying given constraint. Provided data structure can be useful for result investigation and processing, and query debugging. Presented algorithm implemented in F# [17] and available on GitHub: <https://github.com/YaccConstructor/YaccConstructor>.

In order to estimate practical value of proposed algorithm we should perform evaluation on real dataset and real queries. One of possible application of our algorithm is metagenomical assembly quering, and we are working on this topic.

Also we are working on performance improvement by implementation of recently proposed modifications in original GLL algorithm [14]. One of direction of our research is generalization of grammar factorization proposed in [14] which may be useful for regular query processing.

We are working on utilisation of GPGPU and multicore CPU power for graph parsing problem with Valiant [19] algorithm modification proposed by Alexander Okhotin [9]. One of possible benefit is ability to process more expressive queries because modification proposed by Alexander Okhotin extended to support boolean grammars.

7. REFERENCES

- [1] A. Afrozeh and A. Izmaylova. Faster, practical gll parsing. In *International Conference on Compiler Construction*, pages 89–108. Springer, 2015.
- [2] A. Annamaa, A. Breslav, J. Kabanov, and V. Vene. An interactive tool for analyzing embedded sql queries.

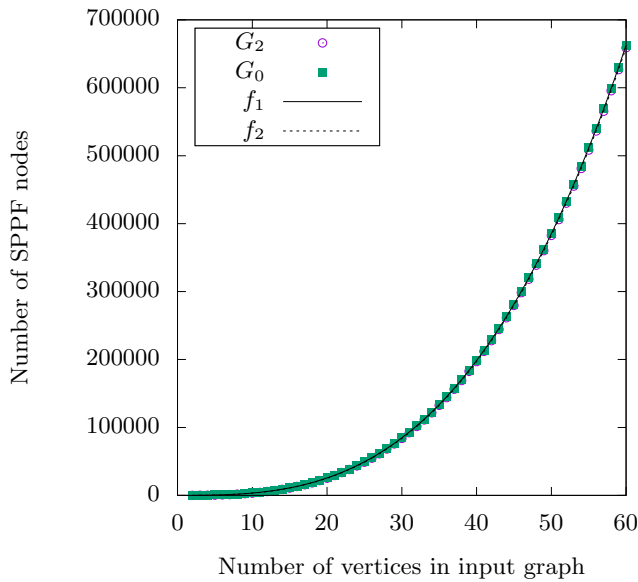


Figure 8: SPPF size on complete graph for grammars G_0 and G_2 on complete graphs

$$f_1(x) = 3.000047 * x^3 + 3.994579 * x^2 + 4.191568 * x; R^2 = 1$$

$$f_2(x) = 3.000050 * x^3 + 2.994338 * x^2 + 4.196472 * x; R^2 = 1$$

In *Asian Symposium on Programming Languages and Systems*, pages 131–138. Springer, 2010.

- [3] C. Barrett, R. Jacob, and M. Marathe. Formal-language-constrained path problems. *SIAM Journal on Computing*, 30(3):809–837, 2000.
- [4] A. Breslav, A. Annamalai, and V. Vene. Using abstract lexical analysis and parsing to detect errors in string-embedded dsl statements. In *Proceedings of the 22nd Nordic Workshop on Programming Theory*, pages 20–22, 2010.
- [5] J. Hellings. Conjunctive context-free path queries. 2014.
- [6] P. Hofman and W. Martens. Separability by short subsequences and subwords. In *LIPICs-Leibniz International Proceedings in Informatics*, volume 31. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik, 2015.
- [7] A. Johnstone and E. Scott. Modelling gll parser implementations. In *International Conference on Software Language Engineering*, pages 42–61. Springer Berlin Heidelberg, 2010.
- [8] J. A. Miller, L. Ramaswamy, K. J. Kochut, and A. Fard. Research directions for big data graph analytics. In *2015 IEEE International Congress on Big Data*, pages 785–794. IEEE, 2015.
- [9] A. Okhotin. Parsing by matrix multiplication generalized to boolean grammars. *Theoretical Computer Science*, 516:101–120, 2014.
- [10] J. G. Rekers. *Parser generation for interactive environments*. PhD thesis, Citeseer, 1992.
- [11] E. Scott and A. Johnstone. Right nulled glr parsers. *ACM Transactions on Programming Languages and Systems (TOPLAS)*, 28(4):577–618, 2006.
- [12] E. Scott and A. Johnstone. Gll parsing. *Electronic Notes in Theoretical Computer Science*, 253(7):177–189, 2010.
- [13] E. Scott and A. Johnstone. Gll parse-tree generation. *Science of Computer Programming*, 78(10):1828–1844, 2013.
- [14] E. Scott and A. Johnstone. Structuring the gll parsing algorithm for performance. *Science of Computer Programming*, 125:1–22, 2016.
- [15] E. Scott, A. Johnstone, and R. Economopoulos. Brnglr: a cubic tomita-style glr parsing algorithm. *Acta informatica*, 44(6):427–461, 2007.
- [16] P. Sevon and L. Eronen. Subgraph queries by context-free grammars. *Journal of Integrative Bioinformatics*, 5(2):100, 2008.
- [17] D. Syme, A. Granicz, and A. Cisternino. *Expert F# 3.0*. Springer, 2012.
- [18] M. Tomita. An efficient context-free parsing algorithm for natural languages. In *Proceedings of the 9th International Joint Conference on Artificial Intelligence - Volume 2, IJCAI’85*, pages 756–764, San Francisco, CA, USA, 1985. Morgan Kaufmann Publishers Inc.
- [19] L. G. Valiant. General context-free recognition in less than cubic time. *Journal of computer and system sciences*, 10(2):308–315, 1975.
- [20] E. Verbitskaia, S. Grigorev, and D. Avdyukhin. Relaxed parsing of regular approximations of string-embedded languages. In *International Andrei Ershov Memorial Conference on Perspectives of System Informatics*, pages 291–302. Springer International Publishing, 2015.