

Parser-Combinators for Context-Free Path Querying*

Sophia Smolina
Electrotechnical University
St. Petersburg, Russia
sofysmol@gmail.com

Ilya Kirillov
Saint Petersburg State University
St. Petersburg, Russia
larst@affiliation.org

Ekaterina Verbitskaia
Saint Petersburg State University
St. Petersburg, Russia
webmaster@marysville-ohio.com

Semyon Grigorev
Saint Petersburg State University
St. Petersburg, Russia
semen.grigorev@jetbrains.com

ABSTRACT

Transparent intergration of domain-specific languages for graph-structured data access into general-purpose programming languages is an important for data-centric application development simplification. It is necessary to provide safety too (static errors checking, type chacking, etc) One of type of navigational queries is a context-free path queries which stands more and more popular. Context-free path querying required in some areas, theoretical research, but no languages (cfSPARQL) Grammars specification languages — combinators. We propose to use parser combinators technique to implement context-free path querying. We demonstrate library and show that it is applicable for realistic problems.

Abstract, abstract, abstract, abstract, abstract, abstract, abstract, abstract, abstract, abstract, Abstract, abstract, abstract, abstract, abstract, abstract, abstract, abstract, abstract, abstract,

CCS CONCEPTS

• **Information systems** → **Graph-based database models; Query languages for non-relational engines;** • **Software and its engineering** → *Functional languages;* • **Theory of computation** → *Grammars and context-free languages;*

KEYWORDS

Graph data bases, Language-constrained path problem, Context-Free path querying, Parser Combinators, Domain Specific Language, Generalized LL, GLL, Neo4j, Scala

ACM Reference Format:

Sophia Smolina, Ekaterina Verbitskaia, Ilya Kirillov, and Semyon Grigorev. 2018. Parser-Combinators for Context-Free Path Querying. In *Proceedings of Joint International Workshop on Graph Data Management Experiences & Systems (GRADES) and Network Data Analytics (NDA) 2018 (GRADES-NDA'18)*. ACM, New York, NY, USA, Article 4, 5 pages. https://doi.org/10.475/123_4

*This work is supported by grant from JetBrains Research.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for third-party components of this work must be honored. For all other uses, contact the owner/author(s).

GRADES-NDA'18, June 2018, Houston, Texas USA

© 2018 Copyright held by the owner/author(s).

ACM ISBN 123-4567-24-567/08/06...\$15.00

https://doi.org/10.475/123_4

1 INTRODUCTION

One of useful type of graph queries is language-constrained path queries [3]. There is a number of languages for graph traversing/-querying which support constraints formulated in terms of regular languages. For example SPARQL [17], Cypher ¹, and GremLin [21]. In this work we are focused on context-free path queries (CFPQ) which use context-free languages for constraints specification and used in such areas as bioinformatics [23], static code analysis [4, 16, 19, 27], RDF processing [26]. There are a lot of theoretical research and problem-specific solutions on CFPQ [2, 8, 9, 15, 20, 23, 25], but there is only one known graph query language with support of CF constraints: cfSPARQL [26].

When you develop a data-centric application, you want to use general purpose programming language and have an transparent and native access to data sources. The naive approach, named string-embedded DSLs, which means that you have a driver which can execute a string with query and return (possibly untyped) result have some serious problems. First of all, it is necessary to use special DSL which may require additional knowledge from developer. Moreover, string-embedded language is a source of errors and vulnerabilities which static detection is very hard [5]. These lead to creation such special techniques as Object Relationship Mapping (ORM) or Language Integrated Query (LINQ) [6, 14] which solve some problems, but have some difficulties with flexibility, for example with query decomposition and subquery reusing. In this paper we propose solution for natural transparent integration of CFPQs into general-purpose language.

One of the natural way to specify any language is specify its formal grammar which can be done by using special DSL based, for example, on EBNF-like notation [24]. The classical alternative way is a parser combinators technique which provide !!!Ekaterina, we need your help!!.

Unfortunately, classical combinators are based on LL(k) top-down parsing technique and can not handle left recursive and ambiguous grammars [?]. In [11] authors show that it is possible to use ideas of Generalized LL [22] (GLL) for creation parser combinators which can handle arbitrary context-free grammars. Meerkat ² parser combinators library is based on [11] and provide Shared Packed Parse Forest [18] (SPPF) — compact representation of parsing result. Furthermore, there is a solution which use GLL for

¹Cypher language web page: <https://neo4j.com/developer/cypher-query-language/>. Access date: 16.01.2018

²Meerkat project repository: <https://github.com/meerkat-parser/Meerkat>. Access date: 16.01.2018

CFPQ [7] and it is showed that SPPF can be used is finite structural representation of query result even when set of paths is infinite. It is important feature because such tasks as paths extraction, queries debuggibg and result processing [10] require appropriate representation of query result.

On the other hand, an idea to use combinators for graph traversing proposed in [12], but presented solution provide approximated handling of cycles in input graph and have a problem with left-recursive grammars. Authors pointed out that described idea is very close to classical parser combinators technique, but supported language class and restrictions are not discussed.

In this paper we show how to compose these ideas and present parser combinators for CFPQ which can handle arbitrary context-free grammar and provide structural representation of result. We make the following contributions in this paper.

- (1) We show that it is possible to create parser combinators for context-free path querying which supports both arbitrary context-free grammars and arbitrary graphs and provide finite structural representation of query result.
- (2) We provide the implementation of parser combinators library in Scala. This library provide integration with Neo4J graph data base. Source code available on GitHub: <https://github.com/YaccConstructor/Meerkat>.
- (3) We perform an evaluation on realistic data. Also we compare performance of our library with other GLL-based CFPQ tool and with Trails library. As a result we conclude that our solution provide good expressivity and performance to be applied for real-world problems.

2 PARSER-COMBITATORS FOR PATH QUERYING

Parser-Combinators is a way to describe context-free grammar in terms of functions and operations on them. Parser is a function which takes some input and returns either Success with result of parsing or Error in a case of failure. Such parsers are composable which makes it easy to create new parsers from existing ones. But classical parser combinators libraries can was made to work on a linear input, and Meerkat which we use as a base for our work too.

We extend input for Meerkat from linear to the graph one. It can be done by using an odservation which for (string or graph) parsing we need only to provide function for getting symbols follower by specified position. As far as Meerkat based on GLL, termination and correctness for arbitrary input graph and arbitrary context-free grammar is based on GLL generalization for CFPQ, which described in [7]. That allows us to adopt parser combinators technique for context-free path querying.

Let us introduce an example of same generation query by using Meerkat. Context-free grammar G presented in Fig. 1 and its Meerkat representation is in Fig. 2

Let's closely take a look at it. For every nonterminal in our CF grammar we create a val of Nonterminal type. Strings are implicitly converted to terminals. syn is a macro which creates new nonterminal and automatically assigns a name of our val to it. Inside a syn macro we've got a definition of nonterminal. It using two combinators \sim and $|$. The first one \sim says that after edge described by left operand we would like to have adjacence edge described by right

$$\begin{aligned} S &\rightarrow \text{subClassOf}^{-1} S \text{ subClassOf} \\ S &\rightarrow \text{type}^{-1} S \text{ type} \\ S &\rightarrow \text{subClassOf}^{-1} \text{subClassOf} \\ S &\rightarrow \text{type}^{-1} \text{type} \end{aligned}$$

Figure 1: Query 1 grammar

```
val S: Nonterminal = syn(
  "subclassof-1" ~ S ~ "subclassof" |
  "type-1" ~ S ~ "type" |
  "subclassof-1" ~ "subclassof" |
  "type-1" ~ "type")
```

Figure 2: Meerkat representation of Query 1

one. $|$ is an alternative combinator which have lower priority then \sim and used to describes possible alternatives of paths description.

The following table shows the some combinators available in Meerkat.

Combinator	Description
$a \sim b$	a and then b
$a b$	a or b
$a.?$	a or nothing
$a.*$	zero or more of a
$a.+$	at least one a

Table 1: Meerkat combinators

One of the most excited combinators feature is that queries using combinators can be used as first class values, which provide high ability for queries generalization and comositon. Let's take a look at Fig 4. sameGen creates a same generation query for every pair of brackets in brs. This generaization is independent from environment (graph and other parsers), and we can use it for creation of different queries. For example, the application of sameGen which one can see in Fig. 3 builds a query which is equivalents to the query presented in Fig. 2.

```
val query1 = syn( sameGen( List(
  ("subclassof-1", "subclassof")
  , ("type-1", "type"))) )
```

Figure 3: Application of sameGen for buildin the query 1

And let's fially parse graph from Fig. 5 and get all paths described by same-generation grammar from it.

To do it let us use meerkat function parseGraphFromAllPositions(parser, graph) which applies given parser to given graph and gets from SPPF all pairs of nodes such that exists a path between them described by a parser.

```

def sameGen(brs) =
  bs.map { case (lbr, rbr) =>
    lbr ~ syn(sameGen(bs).?) ~ rbr }
  match {
    case x :: Nil => syn(x)
    case x :: y :: xs =>
      syn(xs.foldLeft(x | y)(_ | _))
  }

```

Figure 4: Generic function for same generations query building

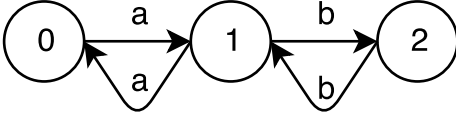


Figure 5: Example graph

The result for graph in Fig. 5 is $\{(1, 0), (1, 2), (0, 0), (2, 1), (2, 2), (0, 2), (0, 1), (1, 1)\}$, where (i, j) stands for the path from node with label i to the node with label j

Meerkat library consists of 3 main units (Fig. ??): the input data unit, the analysis unit, the visualization unit.

The input data unit can accept two main types: graphs and strings. It is necessary to implement the IGraph interface to analyze graphs with different implementations. There are two implementation of IGraph at the moment: Neo4JGraph and SimpleGraph. Neo4JGraph takes data from the Neo4j graph database. SimpleGraph allows you to describe the graph in the code of the program.

The analysis unit consists of parser and SPPFLookup. The parser gets the input data and the grammar, which is a input data query. The parser maps the input data to the input grammar and finds all the paths. The result is added to the SPPFLookup data structure, which contains all the paths in the input data corresponding to the grammar. All nodes in the SPPFLookup are in a single instance and re-used.

After the analysis is completed, the visualization unit gets SPPFLookup and SPPF query. The unit filters the path and then returns the file in dot format to the user.

3 EVALUATION

In this section we present evaluation of Meerkat graph querying library. We show its performance on a classical ontology graphs for in memory graph and for Neo4j database, show application on may-alias static code analysis problem, and compare with Trails [12] library for graph traversals.

All tests are performed on a machine running Fedora 27 with quad-core Intel Core i7 2.5 GHz CPU with 8 GB of memory.

3.1 Ontology querying

One of well-known graph querying problems is a queries for ontologies [1]. We use Meerkat to evaluate it on some popular

ontologies presented as RDF files from paper [26]. We convert RDF files to a labeled directed graph like the following: for every RDF triple $(subject, predicate, object)$ we create two edges $(subject, predicate, object)$ and $(object, predicate^{-1}, subject)$. On those graphs we apply two queries from the paper [7] which grammars are in Fig. 3, and Fig. 6

```

val query2 = syn(
  sameGen(List(("subclassof-1", "subclassof")))
  ~ "subclassof")

```

Figure 6: Query 2 grammar

The queries applied in two following ways.

- Convert RDF files to a graph input for meerkat and then directly parse on query 1 and query 2
- Convert RDF files to a Neo4j database and then parse this database on given queries

Table 2 shows experimental results of those two approaches over the testing RDF files where *number of results* is a number of pairs of nodes (v_1, v_2) such that exists S-path from v_1 to v_2 .

The performance is about 2 times slower than in [7] and shows the same results. If compare the performance of in memory graph querying and database querying, the second one is slower in about 2 – 4 times.

3.2 Static code analysis

Alias analysis is one of the fundamental static analysis problems [13]. Alias analysis checks may-alias relations between code expressions and can be formulated as a Context-Free language (CFL) reachability problem [19]. In that case program represented as Program Expression Graph (PEG) [27]. Vertices in PEG are program expressions and edges are relations between them. In a case of analysing C source code there is two kind of edges **D**-edge and **A**-edge.

- Pointer dereference edge (**D**-edge). For each pointer dereference $*e$ there is a directed D-edge from e to $*e$.
- Pointer assignment edge (**A**-edge). For each assignment $*e_1 = e_2$ there is a directed A-edge from e_2 to $*e_1$

Also, for the sake of simplicity, there are edges labeled by \bar{D} and \bar{A} which corresponds to reversed D-edge and A-edge, respectively.

The grammar for may-alias problem from [27] presented in Fig. 7. It consists of two nonterminals **M** and **V**. It allows us to make two kind of queries for each of nonterminals **M** and **V**.

- **M** production shows that two l-value expression are memory aliases i.e. may stands for the same memory location.
- **V** shows that two expression are value aliases i.e. may evaluate to the same pointer value.

We made **M** and **V** queries on the code some open-source C projects. The results are presented on the Table 3

It may be useful for tools development.

3.3 Comparison with Trails

Trails [12] is a Scala graph combinator library. It provides traversers for describing paths in graphs in terms of parser combinators and

Ontology	#tripples	Query 1				Query 2			
		#results	In memory graph (ms)	DB query time (ms)	Trails (ms)	#results	In memory graph (ms)	DB query time (ms)	Trails (ms)
atom-primitive	425	15454	174	236	2849	122	49	56	453
biomedical-mesure-primitive	459	15156	328	398	3715	2871	36	52	60
foaf	631	4118	23	42	432	10	1	2	1
funding	1086	17634	151	175	367	1158	18	23	76
generations	273	2164	9	27	9	0	0	0	0
people_pets	640	9472	68	87	75	37	2	3	2
pizza	1980	56195	711	792	7764	1262	44	56	905
skos	252	810	4	29	6	1	0	1	0
travel	277	2499	23	93	34	63	2	2	1
univ-bench	293	2540	19	74	31	81	2	3	2
wine	1839	66572	578	736	3156	133	5	7	4

Table 2: Evaluation results for In Memory Graph and Graph DB

$$M \rightarrow \bar{D} V D$$

$$V \rightarrow (M? \bar{A})^* M? (A M?)^*$$

Figure 7: Context-Free grammar for the may-alias problem

```

val M = syn("nd" ~ V ~ "d")
val V = syn((M.? ~ "na").* ~ M.? ~ ("a" ~ M.?).*)

```

Figure 8: Meerkat representation of may-alias problem grammar

Program	Code Size (KLOC)	Count of aliases		Time (ms)
		M aliases	V aliases	
wc-5.0	0.5K	0	174	350
pr-5.0	1.7K	13	1131	532
ls-5.0	2.8K	52	5682	436
bzip2-1.0.6	5.8K	9	813	834
gzip-1.8	31K	120	4567	1585

Table 3: Running may-alias queries on Meerkat on some C open-source projects

allows to get results as a stream (maybe infinite) of all possible paths described by composition of basic traversals. Trails as well as Meerkat support parsing in memory graphs, so we compare performance of Trails and Meerkat on the ontology queries which are described above. The result of comparison are in table 2. Trails gives the same results as Meerkat (column *results* in table 2) but slower than Meerkat.

4 CONCLUSION

We propose a native way to integrate language for language-constrained path querying into general purpose programming language. We implement it and show that our implementation can be applied for real problems. Arbitrary context-free grammars for

querying. This approach is language-independent and may be implemented for closely all general-purpose programming languages.

Code is available on GitHub:

Future work is

Technical improvements: sppf as a set of paths, combinators for vertices information processing, etc

SPPF utilization for debugging and results processing

Attributed grammars processing to provide mechanism for semantics calculation, filtration, etc

REFERENCES

- [1] Serge Abiteboul, Richard Hull, and Victor Vianu. 1995. Foundations of Databases. (1995).
- [2] Pablo Barceló, Gaelle Fontaine, and Anthony Widjaja Lin. 2013. Expressive Path Queries on Graphs with Data. In *International Conference on Logic for Programming Artificial Intelligence and Reasoning*. Springer, 71–85.
- [3] Chris Barrett, Riko Jacob, and Madhav Marathe. 2000. Formal-language-constrained path problems. *SIAM J. Comput.* 30, 3 (2000), 809–837.
- [4] Osbert Bastani, Saswat Anand, and Alex Aiken. 2015. Specification inference using context-free language reachability. In *ACM SIGPLAN Notices*, Vol. 50. ACM, 553–566.
- [5] Tevfik Bultan, Fang Yu, Muath Alkhalaf, and Abdulkali Aydin. 2018. String Analysis for Software Verification and Security. (2018).
- [6] James Cheney, Sam Lindley, and Philip Wadler. 2013. A Practical Theory of Language-integrated Query. *SIGPLAN Not.* 48, 9 (Sept. 2013), 403–416. <https://doi.org/10.1145/2544174.2500586>
- [7] Semyon Grigorev and Anastasiya Ragozina. 2017. Context-free Path Querying with Structural Representation of Result. In *Proceedings of the 13th Central & Eastern European Software Engineering Conference in Russia (CEE-SECR '17)*. ACM, New York, NY, USA, Article 10, 7 pages. <https://doi.org/10.1145/3166094.3166104>
- [8] Jelle Hellings. 2014. Conjunctive context-free path queries. (2014).
- [9] Jelle Hellings. 2015. Path Results for Context-free Grammar Queries on Graphs. *CoRR abs/1502.02242* (2015). <http://arxiv.org/abs/1502.02242>
- [10] Piotr Hofman and Wim Martens. 2015. Separability by short subsequences and subwords. In *LIPICs-Leibniz International Proceedings in Informatics*, Vol. 31. Schloss Dagstuhl-Leibniz-Zentrum fuer Informatik.
- [11] Anastasia Izmaylova, Ali Afroozeh, and Tijs van der Storm. 2016. Practical, General Parser Combinators. In *Proceedings of the 2016 ACM SIGPLAN Workshop on Partial Evaluation and Program Manipulation (PEPM '16)*. ACM, New York, NY, USA, 1–12. <https://doi.org/10.1145/2847538.2847539>
- [12] Daniel Kröni and Raphael Schweizer. 2013. Parsing Graphs: Applying Parser Combinators to Graph Traversals. In *Proceedings of the 4th Workshop on Scala (SCALA '13)*. ACM, New York, NY, USA, Article 7, 4 pages. <https://doi.org/10.1145/2489837.2489844>
- [13] Thomas J. Marlowe, William G. Landi, Barbara G. Ryder, Jong-Deok Choi, Michael G. Burke, and Paul Carini. 1993. Pointer-induced Aliasing: A Clarification. *SIGPLAN Not.* 28, 9 (Sept. 1993), 67–70. <https://doi.org/10.1145/165364.165387>
- [14] Erik Meijer, Brian Beckman, and Gavin Bierman. 2006. LINQ: Reconciling Object, Relations and XML in the .NET Framework. In *Proceedings of the 2006 ACM*

- SIGMOD International Conference on Management of Data (SIGMOD '06)*. ACM, New York, NY, USA, 706–706. <https://doi.org/10.1145/1142473.1142552>
- [15] A. Mendelzon and P. Wood. 1995. Finding Regular Simple Paths in Graph Databases. *SIAM J. Computing* 24, 6 (1995), 1235–1258.
 - [16] Polyvios Pratikakis, Jeffrey S Foster, and Michael Hicks. 2006. Existential label flow inference via CFL reachability. In *SAS*, Vol. 6. Springer, 88–106.
 - [17] Eric Prud, Andy Seaborne, et al. 2006. SPARQL query language for RDF. (2006).
 - [18] Joan Gerard Rekers. 1992. *Parser generation for interactive environments*. Ph.D. Dissertation. Universiteit van Amsterdam.
 - [19] Thomas Reps. 1997. Program Analysis via Graph Reachability. In *Proceedings of the 1997 International Symposium on Logic Programming (ILPS '97)*. MIT Press, Cambridge, MA, USA, 5–19. <http://dl.acm.org/citation.cfm?id=271338.271343>
 - [20] Juan L Reutter, Miguel Romero, and Moshe Y Vardi. 2015. Regular queries on graph databases. *Theory of Computing Systems* (2015), 1–53.
 - [21] Marko A Rodriguez. 2015. The gremlin graph traversal machine and language (invited talk). In *Proceedings of the 15th Symposium on Database Programming Languages*. ACM, 1–10.
 - [22] Elizabeth Scott and Adrian Johnstone. 2010. GLL parsing. *Electronic Notes in Theoretical Computer Science* 253, 7 (2010), 177–189.
 - [23] Petteri Sevon and Lauri Eronen. 2008. Subgraph queries by context-free grammars. *Journal of Integrative Bioinformatics* 5, 2 (2008), 100.
 - [24] Niklaus Wirth. 1996. Extended Backus-Naur Form (EBNF). *ISO/IEC 14977* (1996), 2996.
 - [25] Mihalios Yannakakis. 1990. Graph-theoretic methods in database theory. In *Proceedings of the ninth ACM SIGACT-SIGMOD-SIGART symposium on Principles of database systems*. ACM, 230–242.
 - [26] Xiaowang Zhang, Zhiyong Feng, Xin Wang, Guozheng Rao, and Wenrui Wu. 2016. Context-free path queries on RDF graphs. In *International Semantic Web Conference*. Springer, 632–648.
 - [27] Xin Zheng and Radu Rugina. 2008. Demand-driven Alias Analysis for C. In *Proceedings of the 35th Annual ACM SIGPLAN-SIGACT Symposium on Principles of Programming Languages (POPL '08)*. ACM, New York, NY, USA, 197–208. <https://doi.org/10.1145/1328438.1328464>