# Multiple-Source Context-Free Path Querying in Terms of Linear Algebra

Arseniy Terekhov
simpletondl@yandex.ru
Saint Petersburg State University
St. Petersburg, Russia

Vlada Pogozhelskaya
pogozhelskaya@gmail.com
Saint Petersburg State University
St. Petersburg, Russia

Vadim Abzalov
vadim.i.abzalov@gmail.com
Saint Petersburg State University
St. Petersburg, Russia

Timur Zinnatulin
!!!@!!!
Saint Petersburg State University
St. Petersburg, Russia

Semyon Grigorev
s.v.grigoriev@spbu.ru
semyon.grigorev@jetbrains.com
Saint Petersburg State University
St. Petersburg, Russia
JetBrains Research
St. Petersburg, Russia

## ABSTRACT

A long time ago in a galaxy far far away... Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract. Abstract is very abstract.

## 1 INTRODUCTION

Language-constrained path querying [2] is a way to find paths in edge-labeled graphs when constraints are formulated in terms of language which restrict words formed by paths: the word formed by path's labels concatenation should be in the specified language. This way is very natural for navigational queries in graph databases, and one of the most popular languages which are used for constraints is a regular language. But in some cases, regular languages are not expressive enough, as a result, context-free languages gain popularity. Constraints in the form of context-free languages, or context-free path querying (CFPQ), can be used for RDF analysis [11], biological data analysis [9], static code analysis [8, 12], and in other areas.

Big amount of research done on CFPQ, a number of CFPQ algorithms were proposed, but the application of context-free constraints for real-world data analysis faced with some problems problem. The first problem is a bad performance of proposed algorithms on real-world data, as was shown by Jochem Kuijpers et al. [5]. Moreover, there are no graph databases with full-stack support of CFPQ, the main effort was made in algorithms and their theoretical properties research. This fact hinders research of problems reducible to CFPQ, thus it hinders the development of new solutions for some problems. For example, recently graph segmentation in data provenance analysis was reduced to CFPQ [6], but authors faced the problem during the

evaluation of the proposed approach: no one graph database support CFPQ.

In [1] Rustam Azimov propose a matrix-based algorithm for CFPQ. This algoritm is one of promising way to solve the first problem and provide appropriate solution for real-world data analysis, as was shown by Nikita Mishim et al. in [7] and Arseniy Terekhov et al. in [10]. But this algorithm always computes information (reachability facts or single path which satisfies constraints) for all pairs of vertices in the graph, namely it solves *all-pairs* problem. It is unreasonable for some real-world scenarios when one can provide a relatively small set of start vertices or even single start vertex.

While all-pairs context-free path querying is a classical problem that investigates in a number of works, there is no, in our knowledge, solutions for single-source and multiple-source CFPQ. In this work we propose a matrix-based *multiple-source* (and *single-source* as a partial case) CFPQ algorithm.

Also, we provide full-stack support of CFPQ for the Redis-Graph[1] [3] graph database. We implement a Cypher query language extension[2] that allows one to express context-free constraints, and extend the RedisGraph to support this extension. In our knowledge, it is the first full-stack implementation of CFPQ.

To summarize, we make the following contribution in this paper.

(1) We modify Azimov's matrix-based CFPQ algorithm and provide a multiple-source matrix-based CFPQ algorithm. As a partial case, it is possible to use our algorithm in a single-source scenario. Our modification still based on linear algebra, hance it is simple to implementation and allows one to use high-performance libraries for implmementation.
(2) We evalute the proposed algorithm. Our evaluation shows that !!!
(3) We provide full-stack support of CFPQ by extending the RedisGraph graph database. To do it, we extend Cypher with syntax allows one to express context-free constraints,

---

[1]RedisGraph graph database Web-page: https://redislabs.com/redis-enterprise/redis-graph/. Access date: 19.07.2020.
[2]Proposal which describes path patterns specification syntax for Cypher query language: https://github.com/thobe/openCypher/blob/rpq/cip/1.accepted/CIP2017-02-06-Path-Patterns.adoc. The proposed syntax allows one to specify context-free constraints. Access date: 19.07.2020.

implement the proposed algorithm in a RedisGraph back-end, and support new syntax in the RedisGraph query execution engine. Finally, evaluate the poposed solution.

## 2 PRELIMINARIES

In this section we introduce common definitions in graph theory and formal language theory which will be used in this paper. Also, we provide brief description of Azimov's algorithm which is used as a base of our solution.

### 2.1 Graphs

In this work we use edge-labelled digraph as a data model and define it as follows.

*Definition 2.1.* *Labeled directed graph* is a triple $D = (V, E, \lambda)$, where

- $V$ is a set of vertices. For simplicity, we assume that the vertices are natural numbers.
- $E \subseteq V \times V$ is a set of edges
- $\lambda : (E \cup V) \to 2^\Sigma$ is a function that maps edges and vertices to their labels from the label set $\Sigma$.

An example of the graph is presented in figure 1. Here the set of labels $\Sigma = \{A, B, C, D, X\}$.
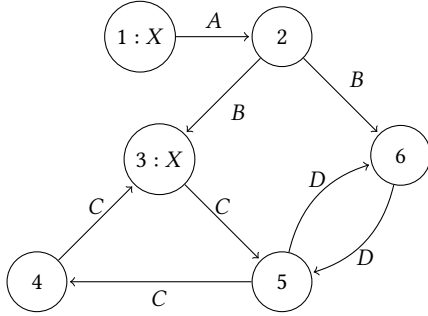


**Figure 1: The example of input graph $\mathcal{G}$**

We use adjacency matrix decomposed to a set of a boolean matrix as a representation of the graph.

*Definition 2.2.* An adjacency matrix $M$ of the labelled graph $\mathcal{G} = (V, E, \lambda)$ is a square $|V| \times |V|$ matrix, such that $M[i, j] = \lambda((i, j))$.

Adjacency matrix $M$ of the graph $\mathcal{G}$ is

$$M = \begin{pmatrix} . & \{A\} & . & . & . & . \\ . & . & \{B\} & . & . & \{B\} \\ . & . & . & . & \{C\} & . \\ . & . & \{C\} & . & . & . \\ . & . & . & \{C\} & . & \{D\} \\ . & . & . & . & \{D\} & . \end{pmatrix}.$$

*Definition 2.3.* Boolean decomposition of adjacency matrix $M$ of graph $\mathcal{G} = (V, E, \lambda)$ is set of Boolean matrix

$$\mathcal{M} = \{M^l \mid l \in \Sigma, M^l[i, j] = 1 \iff l \in M[i, j]\}.$$

Matrix $M$ can be represented as a set of two Boolean matrices $M^a$ and $M^b$ where

$$M^A = \begin{pmatrix} . & 1 & . & . & . & . \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ . & . & . & . & . & . \end{pmatrix}, M^B = \begin{pmatrix} . & . & . & . & . & . \\ . & . & 1 & . & . & 1 \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ . & . & . & . & . & . \end{pmatrix} \quad (1)$$

$$M^C = \begin{pmatrix} . & . & . & . & . & . \\ . & . & . & . & . & . \\ . & . & . & . & 1 & . \\ . & . & 1 & . & . & . \\ . & . & . & 1 & . & . \\ . & . & . & . & . & . \end{pmatrix}, M^D = \begin{pmatrix} . & . & . & . & . & . \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ . & . & . & . & . & . \\ . & . & . & . & . & 1 \\ . & . & . & . & 1 & . \end{pmatrix}. \quad (2)$$

*Definition 2.4.* An vertex label matrix $N$ of the labelled graph $(G) = (V, E, \lambda)$ is a square $|V| \times |V|$ matrix, such that $N[i, i] = \lambda(i)$ and $N[i, j] = \emptyset$ for $i \neq j$.

*Definition 2.5.* Boolean decomposition of vertex label matrix $N$ of graph $\mathcal{G} = (V, E, \lambda)$ is set of Boolean matrix

$$\mathcal{N} = \{N^l \mid l \in \Sigma, N^l[i, j] = 1 \iff l \in N[i, j]\}.$$

### 2.2 Languages

We formulate constraints in terms of context-free languages, for this reason there are following definitions.

*Definition 2.6.* *Context-free grammar* is a 4-tuple $G = (N, \Sigma, P, S)$, where

- $N$ is a finite set of nonterminals
- $\Sigma$ is a finite set of terminals
- $P$ is a finite set of productions of the following forms: $A \to \alpha$, $A \in N$, $\alpha \in (N \cup \Sigma)^*$
- $S$ is a starting nonterminal

*Definition 2.7.* *Context-free language* is a language generated by a context-free grammar $G$:

$$L(G) = \{w \in \Sigma^* \mid S \xRightarrow[G]{*} w\}$$

Where $S \xRightarrow[G]{*} w$ denotes that a string $w$ can be generated from a starting non-terminal $S$ using some sequence of production rules from $P$.

*Definition 2.8.* Context-free grammar $G = (N, \Sigma, P, S)$ is said to be in *Chomsky normal form* if all productions in $P$ are in one of the following forms:

- $A \to BC$, $A \in N$, $B, C \in N \setminus S$
- $A \to a$, $A \in N$, $a \in \Sigma$
- $S \to \varepsilon$, $\varepsilon$ is an empty string

Since matrix-based CFPQ algorithms processes grammars only in Chomsky normal form, it should be noted that every context-free grammar can be transformed into an equivalent one in this form.

*Definition 2.9.* Context-free grammar $G = (N, \Sigma, P, S)$ is said to be in *weak Chomsky normal form* if all productions in $P$ are in one of the following forms:

- $A \to BC$, $A, B, C \in N$
- $A \to a$, $A \in N$, $a \in \Sigma$
- $A \to \varepsilon$, $A \in N$

In other words, weak Chomsky normal form differs from Chomsky normal form in the following:

- $\varepsilon$ can be derived from any non-terminal
- $S$ can be at a right part of productions

For example, let's consider the following context-free grammar, which generates the language $L(G) = \{A^n B^n, n \in \mathbb{N}\}$: $G = (N, \Sigma, P, S), \ N = \{S\}, \ \Sigma = \{a, b\}$ and productions:

$$S \to aSb$$
$$S \to ab$$

After transformation to weak Chomsky normal form the resulting grammar:

$$S \to AB \quad S \to AS_1 \quad S_1 \to SB$$
$$A \to a \quad B \to b$$

These productions itself are the grammar that has the same result as original grammar.

We use a context-free grammar in the weak Chomsky normal form without a starting non-terminal, which will be specified in the path queries for the graph. Also we omit the rules of the form $A \to \varepsilon$ for the reason that they correspond to trivial paths, which are more convenient to consider separately.

*Definition 2.10. Context-free relation* is a relation $R_A \subseteq V \times V$ for edge-labeled graph $D = (V, E)$, context-free grammar $G = (N, \ \Sigma, \ P)$ and fixed non-terminal $A$:

$$R_A = \{(n, m) \mid \exists n \pi m \ (l(\pi) \in L(G_A))\}$$

Where $l(\pi)$ is a word obtained by concatenating the labels along the path $\pi$.

Finally, in these notations context-free path querying problem is the problem of finding context-free relations in which the language is specified by a context-free grammar.

## 2.3 Matrix-Based Algorithm

Let $G = (N, \Sigma, P)$ be the input grammar, $D = (V, E)$ be the input edge-labeled graph and language $L$ over alphabet $\Sigma$. For the context-free path query evaluation, we need to provide context-free relations $R_A \subseteq V \times V$ for every $A \in N$. The matrix-based algorithm for CFPQ can be expressed in terms of operations over Boolean matrices (see listing 1) which is an advantage for implementation.

---

**Algorithm 1** Context-free path querying algorithm

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

---

This CFPQ algorithm allows efficiently apply GPGPU techniques, but it solves all-pairs problem and takes unreasonable amount of memory in scenarios in which we want to find paths from a relatively small set of vertices, since it calculates a lot of redundant information.

## 3 MATRIX-BASED MULTIPLE-SOURCE CFPQ ALGORITHM

In this section we introduce two versions of multiple-source matrix-based CFPQ algorithm. This algorithm is a modification of Azimov's matrix-based algorithm for CFPQ and its idea is that we cut off those vertices from which we are not interested in paths.

Let $D = (V, E)$ be the input graph, $G = (N, \Sigma, P)$ be the input context-free grammar and $Src$ be the input set of vertices. For the multiple-source context-free path query evaluation for every $A \in Src$ we need to find all context-free relations $R_A$, i.e. all node pairs $(n, m)$ such that $\exists n \pi m \ (l(\pi) \in L(G_A))$. In order to solve the

---

**Algorithm 2** Multiple-source context-free path querying algorithm

1: **function** MULTISRCCFPQ($D = (V, E), G = (N, \Sigma, P, S), Src$)
2:     $T \leftarrow \{T^A \mid A \in N, T^A \leftarrow \emptyset\}$   ▷ Matrix in which every element is $\emptyset$
3:     $TSrc \leftarrow \{TSrc^A \mid A \in N \setminus S, TSrc^A \leftarrow \emptyset\}$   ▷ Matrix for input vertices in which every element is $\emptyset$
4:     **for all** $v \in Src$ **do**     ▷ Input matrix initialization
5:         $TSrc^S_{v,v} \leftarrow true$
6:     **for all** $A \to x \in P$ **do**     ▷ Simple rules initialization
7:         **for all** $(v, x, to) \in E$ **do**
8:             $T^A_{v,to} \leftarrow true$
9:     **while** $T$ or $TSrc$ is changing **do**     ▷ Algorithm's body
10:         **for all** $A \to BC \in P$ **do**
11:             $M \leftarrow TSrc^A * T^B$
12:             $T^A \leftarrow T^A + M * T^C$
13:             $TSrc^B \leftarrow TSrc^B + TSrc^A$
14:             $TSrc^C \leftarrow TSrc^C + \text{GETDST}(M)$
15:     **return** $T$
16:
17: **function** GETDST($M$)
18:     $A \leftarrow \emptyset$
19:     **for all** $(v, to) \in V^2 \mid M_{v,to} = true$ **do**
20:         $A_{to,to} \leftarrow true$
21:     **return** $A$

---

single-source and multiple-source CFPQ problem Azimov's algorithm was modified: operations of Boolean matrix multiplication $T_A = T_A + T_B \cdot T_C$ for each $A \to BC \in R$ represented in line **8** of Algorithm 1 was supplemented with one more matrix multiplication $T_A = T_A + (TSrc^A \cdot T_B) \cdot T_C$ for each $A \to BC \in R$ which saves only vertices we are interested in, where $TSrc^A$ — matrix of vertices to calculate the paths from. It is represented in lines **11-13** of the Algorithm 2. Also, after every iteration of while loop this is nessesary to update the set of vertices paths from which we need to calculate. To do this, the function **getDst**, represented in lines **17-21**, is called at line **14**. Thus, the modified algorithm does not calculate the paths from all vertices in case of query to calculate the paths small set of vertices.

We proposed the variant of the algorithm that can calculate the paths from a certain set of vertices, however there are such scenarios when queries are partially or completely repeated. In such cases it would be useful to add data caching to improve the performance. The problem is that every time we want to find all paths from the certain set of vertices, the Algorithm 2 calculates everything from scratch. Since recalculating might take the significant amount of time, we modified multiple-source

CFPQ algorithm to specify it for such scenarios. This version stores all the vertices the paths from which have already been calculated in cash *index*, which is used to filter such vertices in line **3** of Algorithm 3. Thus, modified algorithm calculates paths from the particular vertex only once.

---

**Algorithm 3** Optimized multiple-source context-free path querying algorithm

1: **function** MultiSrcCFPQSmart($index = (D, G, T, TSrc), Src$)
2: $\quad TNewSrc \leftarrow \{TNewSrc^A \mid A \in N \setminus S, TNewSrc^A \leftarrow \emptyset\}$
3: $\quad$ **for all** $v \in Src \mid index.TSrc_{v,v} = false$ **do**
4: $\quad\quad TNewSrc^S_{v,v} \leftarrow true$
5: $\quad$ **while** $index.T$ or $TNewSrc$ is changing **do**
6: $\quad\quad$ **for all** $A \rightarrow BC \in P$ **do**
7: $\quad\quad\quad M \leftarrow TNewSrc^A * index.T^B$
8: $\quad\quad\quad index.T^A \leftarrow index.T^A + M * index.T^C$
9: $\quad\quad\quad TNewSrc^B \leftarrow TNewSrc^B + TNewSrc^A \setminus index.TSrc^B$
10: $\quad\quad\quad TNewSrc^C \leftarrow TNewSrc^C + \text{getDst}(M) \setminus index.TSrc^C$

---

## 3.1 Implementation Details

All of the above versions have been implemented[3] using GraphBLAS framework that allows you to represent graphs as matrices and work with it in terms of linear algebra. For convenience, all the code is written in Python using pygraphblas[4], which is Python wrapper around GraphBLAS API and based on SuiteSparse:GraphBLAS[5] [4] — the full implementation of GraphBLAS standart. This library is specialized for working with sparse matrices, which most often appear in real graphs. Also, it should be noted that, despite the fact that the function **getDst** does not seem to be expressed in terms of linear algebra, the implementation used the function **reduce_vector** from pygraphblas that reduces matrix to a vector, with which further work takes place.

## 3.2 Algorithm Evaluation

We evaluate both described version of multiple-source algorithm on real-world graphs. For evaluation, we use a PC with Ubuntu 20.04 installed. It has Intel core !!! CPU, !!!GHz, and DDR3 32Gb RAM. As far as we evaluate only algorithm execution time, we store each graph fully in RAM as its adjacency matrix in sparse format. Note, that graph loading time is not included in the result time of evaluation.

For evaluation we use graphs and queries from CFPQ_Data dataset[6] Detailed information on graphs which we select for evaluation is provided in table 1. We use classical same-generation queries $G_1$ (eq. 3) and $G_2$ (eq. 4) which are used in other works for CFPQ evaluation. Also we use *Geo* (eq. 5) query which was provided by J. Kuijpers et. al [5] for *geospecies* RDF. Note that in queries we use $\overline{x}$ notation to denote inverse of $x$ relation and

---

---

**Table 1: Graphs for CFPQ evaluation: BT for broaderTransitive**

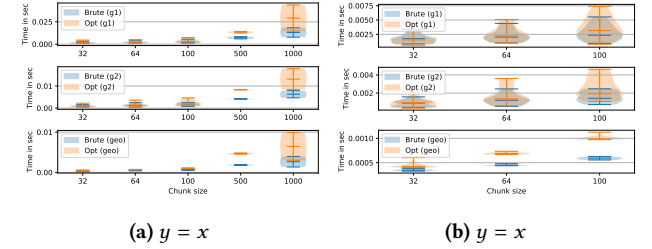| Graph | #V | #E | #subCalssOf | #type | #BT |
|---|---|---|---|---|---|
| core | 120 926 | 484 646 | | | — |
| eclass_514en | 120 926 | 484 646 | | | — |
| enzyme | 358 434 | 144 9711 | | | — |
| geospecies | 596 760 | 2 416 513 | | | |
| go | 1 188 340 | 4 820 728 | | | — |
| go-hierarchy | 1 780 956 | 7 228 358 | | | — |
| pathways | 120 926 | 484 646 | | | — |
| taxonomy | 2 308 385 | 9 369 511 | | | — |



**(a)** $y = x$      **(b)** $y = x$

**Figure 2: Single path extraction**

respective edge.

$$S \rightarrow \overline{subClassOf}\ S\ subClasOf \mid \overline{type}\ S\ type \\ \mid \overline{subClassOf}\ subClasOf \mid \overline{type}\ type \tag{3}$$

$$S \rightarrow \overline{subClassOf}\ S\ subClasOf \mid subClassOf \tag{4}$$

$$S \rightarrow broaderTransitive\ S\ \overline{broaderTransitive} \\ \mid broaderTransitive\ \overline{broaderTransitive} \tag{5}$$

Our main goal is to compare behavior of two proposed versions of the algorithm. To do it we measure query execution time for both versions for different sizes of star vertex set. Namely, for each graph we split all vertices into disjoint subsets of fixed size. After that, for each subset we evaluate queries using the given subset as a set of start vertices.

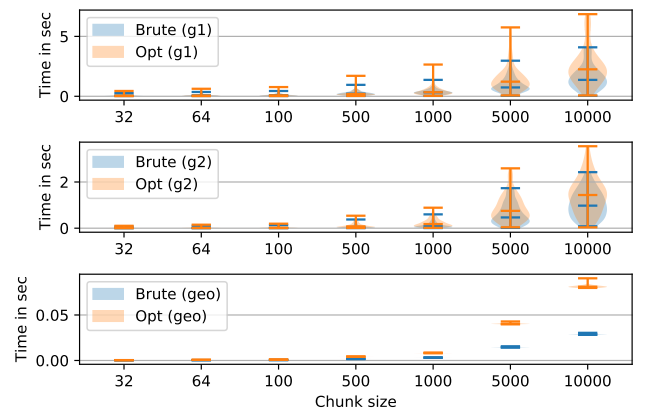Results of evaluation is presented in figures !!!!.



**Figure 3: Example of a parametric plot** $(\sin(x), \cos(x), x)$

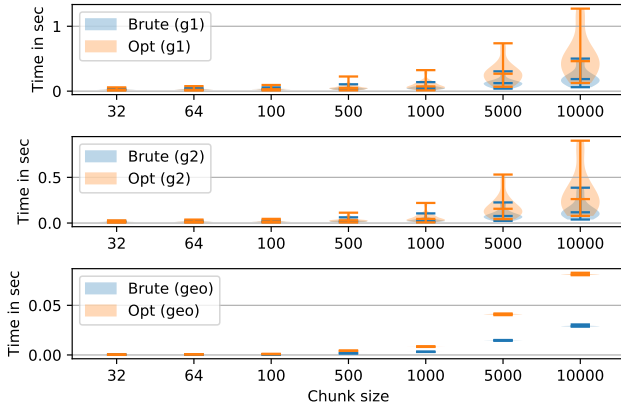We can see, that .... As a result, we select !!! to integrate into RedisGraph!!!

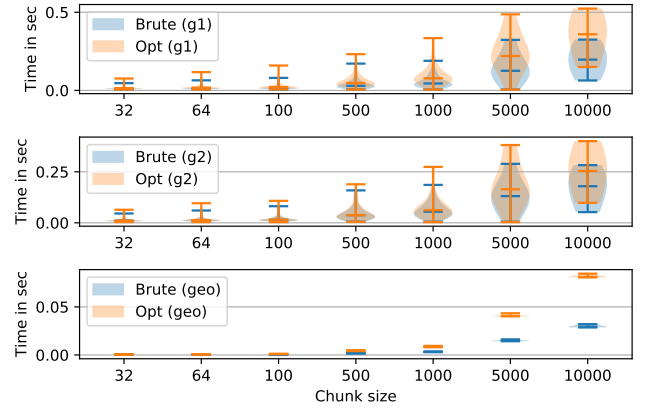**Figure 4: Example of a parametric plot** ($\sin(x), \cos(x), x$)



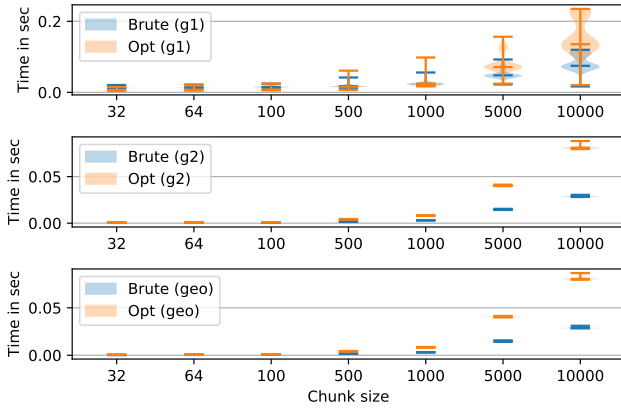**Figure 5: Example of a parametric plot** ($\sin(x), \cos(x), x$)



**Figure 6: Example of a parametric plot** ($\sin(x), \cos(x), x$)



**Figure 7: Example of a parametric plot** ($\sin(x), \cos(x), x$)



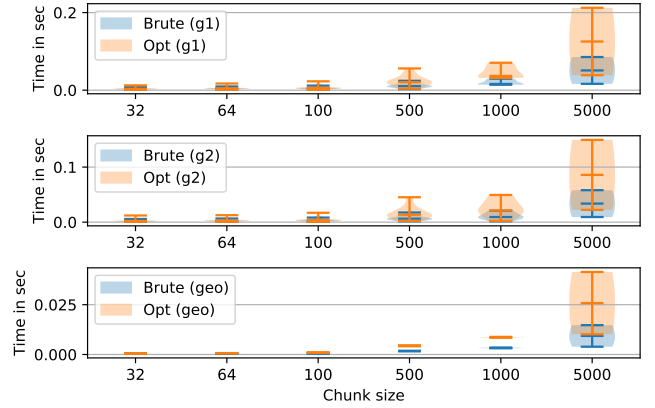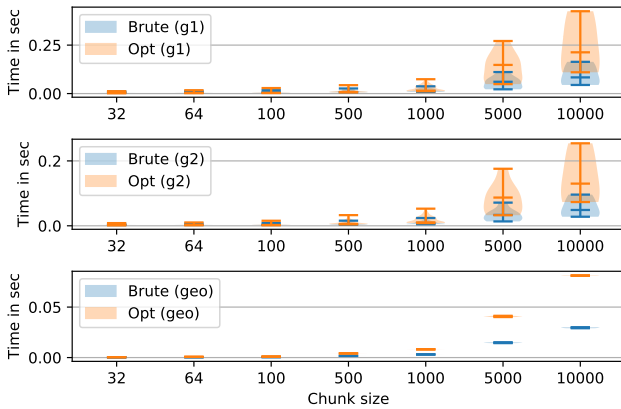**Figure 8: Example of a parametric plot** ($\sin(x), \cos(x), x$)

## 4 CFPQ FULL-STACK SUPPORT

In order to provide full-stack support of CFPQ it is necessatry to choose an appropriatr graph database. It was shown by Arseniy Terekhov et al. in [10] that matrix-based algorithm can be naturally integrated into RedisGraph graph database because both, the algorithm and the database, operates over matrix representation of graphs. Moreover, RedisGraph supports Cypher as a

query language and there is a proposal which describes Cypher extension which allows one to specify context-free constraints. Thus we choose RedisGraph as a base for our solution.

### 4.1 Cypher Extending

The first what we should do is to extend Cypher to be able to express context-free constraints. There is a description of the respective Cypher syntax extension[7], proposed by Tobias Lindaaker, but this syntax does not implement yet in Cypher parsers.

This extension introduces path patterns, which are a more powerful alternative to relationship patterns. Path patterns allow you to express regular constrains over basic patterns such as relationship and node patterns. Just like relationship patterns they can be specified in the MATCH clause between the node patterns.

---

**Listing 4** Example of using a simple path pattern

1: MATCH (v)-/ [:A (:X) :B] | [:C (:Y) :D] /->(to)
2: RETURN v, to

---

The listing 4 provides an example of query in extended syntax with a simple path pattern. In this example there are relationship

---

[7]Formal syntax specification: https://github.com/thobe/openCypher/blob/rpq/cip/1.accepted/CIP2017-02-06-Path-Patterns.adoc#11-syntax. Access date: 19.07.2020.

patterns :A, :B, :C :D and node patterns (:X), (:Y). The square brackets are used for grouping parts of the pattern. The | symbol denotes alternative between corresponding paths and the whitespace denotes sequence of paths. So the result of executing the query on the graph $D$ will be the following set of vertex pairs:

$$\{(v, to) : \exists \pi = (v, r_1, u, r_2, to) \in Paths(D) :$$

$$\begin{bmatrix} t(r_1) = A, l(u) = X, t(r_2) = B \\ t(r_1) = C, l(u) = Y, l(r_2) = D \end{bmatrix} \}$$

Main feature which allows one to specify context-free constraints is a *named path patterns*: one can specify a name for path pattern and after that use it in other patterns, or in the same pattern. Using this feature, structure of query is pretty similar to context-free grammar in the Extended Backus–Naur Form.

---

**Listing 5** Example of using a named path pattern

---
1: PATH PATTERN S = ()-/ [:A ~S :B] | [:A :B] /->()
2: MATCH (v)-/ ~S /->(to)
3: RETURN v, to

---

The listing 5 shows an example of using named path patterns. They can be defined in the PATH PATTERN clause and referenced within any other path pattern. In order to explain the semantics of the query let's consider contest-free grammar $G = (N, \Sigma, P, S)$ with $N = \{S\}$, $\Sigma = \{A, B\}$ and $P = \{S \rightarrow AB, S \rightarrow ASB\}$. Then $L(G) = \{A^n B^n : n \in \mathbb{N}\}$ specifies restrictions on the path labels and query result on the graph $D$ will be as follows:

$$\{(v, to) : \exists \pi = (v, r_1, u_1, ..., r_n, to) \in Paths(D) :$$

$$t(r_1)t(r_2)...t(r_n) \in L(G)\}$$

Thus this Cypher extension allows one express more complex queries including context-free path queries. RedisGraph database supports subset of Cypher language and uses `libcypher-parser`[8] library to parse queries. We extend this library by intoducing new syntax proposed [7]. We implement[9] full extension, not only part which is necessary for simple CFPQ.

## 4.2 RedisGraph Intro (TODO: move to introduction)

Named path patterns described in subsection 4.1 allows one to specify context-free constrains on the paths. In order to support the execution of these types of queries we need to extend backend of the RedisGraph database and integrate a suitable CFPQ algorithm into it.

There are quite a few algorithms that solve CFPQ problem ??, but their running time makes them unsuitable for practical use ??. Recent studies ?? have shown that one can achieve high performance through the use of matrix-based algorithms. These studies were conducted to analyze the performance of the Rustam Azimov algorithm described in ?? and have shown that it is acceptable for practical application.

Using the Rustam Azimov algorithm one can only find paths between all pairs of vertexes at once and in some cases it is quite wasteful. Queries to graph databases can be specified so that when they are executed, it is required to find paths from

[8]The `libcypher-parser` is an open-source parser library for Cypher query language. GitHub repository of the project: https://github.com/cleishm/libcypher-parser. Access date: 19.07.2020.
[9]The modified libsypher-pareser library with support of syntax for path patterns: https://github.com/YaccConstructor/libcypher-parser. Access date: 19.07.2020.

a given set of initial vertices. This set can be quite small due to the different filtering specified in the query. For example in the listing 6 path pattern -/ ~S /-> follows pattern (v)-[r]->(u) . The WHERE clause specifies some arbitrary predicate p(v, r, u) which also fixes a set of initial vertexes for a paths that must satisfy path pattern S . Depending on this predicate, this set of vertexes can have different sizes and for proper practical use the running time of the CFPQ algorithm should be sensitive to this.

---

**Listing 6** ...

---
1: PATH PATTERN S = ()-/ :A [~S | ()] :B /->()
2: MATCH (v)-[r]->(u)-/ ~S /->(to)
3: WHERE p(v, r, u)
4: RETURN to

---

The Multi-Source algorithm described in ?? is sensitive to the initial set of vertices and is therefore well suited for graph database query scenarios. In addition, it is based on matrix operations and works with graphs as sparse matrices, so it is suitable for integration in RedisGraph.

## 4.3 RedisGraph extension

This section describes the implementation of support for executing queries with the extended syntax in the RedisGraph.
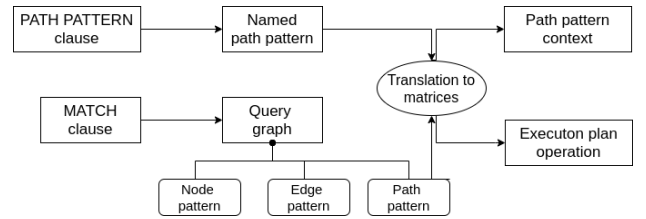


**Figure 9: Extension diagram for building a query execution plan**

*4.3.1 Execution plan building.* In the RedisGraph the main part of processing a query is building its execution plan. Execution plan consists of operations that perform basic processing such as filtering, pattern matching, aggregation and result construction. The diagram of its construction !!! is shown in Figure 9. It can be divided into two parts – processing named and unnamed path patterns, which are described below.

Let's consider the part that associated with unnamed path patterns. Unnamed path patterns relates to the pattern matching operations and is very similar to relationship patterns from the original Cypher. All pattern matching operations are derived from the MATCH clause that consists of relationship patterns and node patterns. In the first stage of processing, these patterns turn into an intermediate representation – the *query graph*. The nodes and edges of the *query graph* corresponds to node and relationship patterns. We extended query graph to be able to contain path patterns. Thus the query graph edges can be either relationship or path patterns, which are stored in a more convenient intermediate representation other than AST.

At the second stage, the query graph is translated into algebraic expressions over matrices. To do this, RedisGraph first linearizes the query graph, and then splits it into small paths. After that, each path is translated into a single algebraic expression. Its operands are a matrices specifying the type of edges and

labels of vertices. To support path patterns we first extended the split processing so that each path template corresponds to exactly one path after query graph splitting. After that we implemented translation of the path patterns into an algebraic expressions. To do this we needed to extend the matrix operands to support references to named paths patterns in algebraic expressions. Finally we have developed the semantics of path patterns in terms of algebraic expressions over matrices.

After obtaining algebraic expressions they are used to construct execution plan operations. Each operation is derived from a single algebraic expression that is involved in the further execution of the corresponding operation. We created a new *CFPQTraverse* operation for expressions that correspond to path patterns. During the query execution this operation performs path pattern matching and solves context-free path reachability problem if necessary. This completes the part of the query execution plan building which concerns unnamed path patterns.

Another processing that occurs during the execution plan construction and was supported by us is related to named path patterns. They are processed independently of the unnamed path patterns found in MATCH clause and don't produce execution plan operations.

All named path patterns are collected from PATH PATTERN clauses. Then they translated into algebraic expressions and stored in the corresponding global context of the query – *path pattern context*. This storage provides mapping of the path pattern name to its algebraic expression and can be used both when building an execution plan and during its execution.

Thus after execution plan building we receive *CFPQTraverse* operations that correspond to unnamed path patterns in MATCH clause and *path pattern context* that stores all named path patterns from PATH PATTERN clauses. Therefore we can proceed to the stage of execution plan evaluation.
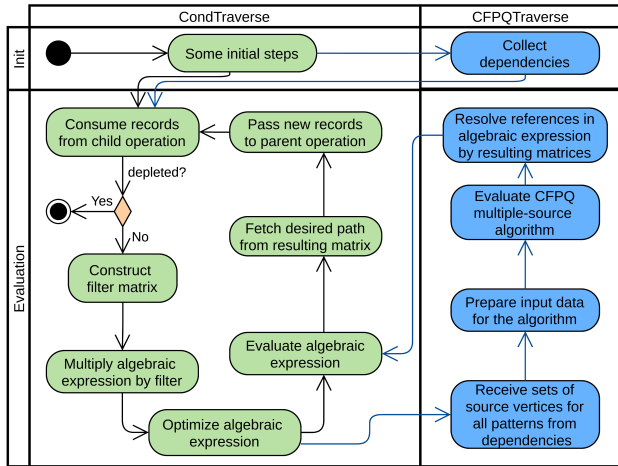


**Figure 10: CFPQTraverse and CondTraverse evaluation**

*4.3.2 Execution plan evaluating.* The remaining part of query processing is evaluation its execution plan. This section describes how the CFPQTraverse operation is performed.

Let's first consider the structure of the execution plan operations. Operations have parent-child relationships, so they are formed into a tree. Each operation can consume a record from a child operation, process it and produce another one for the parent. Records contain information necessary for the parent

operation, as well as everything to restore the response, such as identifiers of accumulated vertices and edges.

The CFPQTraverse operation is based on CondTraverse operation that already exists in the RedisGraph and performs a patterns matching. The activity diagram of this operations is shown in Figure 10 and described below. Actions that corresponds to CondTraverse operation are highlighted in green, actions of the CFPQTraverse operation that extend CondTraverse are highlighted in blue.

The CondTraverse works as follows. At first it consumes several records from the child operation and accumulates them in the buffer. Here each record corresponds to the path that built by the child operation and contains information about the destination vertex of the path. The task of the CondTraverse is to continue the path from this vertex in such way that the resulting path satisfies pattern corresponding to this operation. To do this CondTraverse uses the algebraic expression obtained in the previous step. The resulting matrix of this expression represents all pairs of vertices between which there is a path satisfying the pattern. In order to find paths that start from given sources vertices CondTraverse uses a filter matrix. This matrix is constructed from the destination vertices retrieved from the record buffer and multiplied to the right by algebraic expression of operation. Then the resulting expression is estimated to get the desired paths and to pass them to the parent operation by producing new records.

The CFPQTraverse operation is arranged in the same way as CondTraverse but performs some additional work. Since each CFPQTraverse corresponds to path pattern, its algebraic expression may contain references to named path patterns. Therefore all named path patterns that the algebraic expression depends on must be processed. For this they are extracted from *path pattern context* and stored in the *set of operation dependencies* during its initialization. In this case, dependencies are extracted recursively, so that references inside named path patterns are also extracted.

The CFPQTraverse execution stage starts the same way as CondTraverse. First filter matrix is constructed from record buffer and embedded in the algebraic expression. Then for each reference we need to determinate the set of source vertices. This can be done during algebraic expression evaluation which we extended for this purpose. After that we have everything to run *multiple-source* CFPQ algorithm to resolve all dependencies. This algorithm is slightly different from the one described in section 3 and is a generalization of it. It receives the *set of operation dependencies* and sets of source vertices. After running this algorithm a matrix is obtained for each named path pattern. This matrices represent a set of pairs of vertices between which there is a path that satisfies the pattern. Then all references in the algebraic expression are replaced with the resulting matrices and the algebraic expression is evaluated. Finally as well as CondTraverse, CFPQTraverse extracts desired paths from resulting matrix and passes them to parent operation.

## 4.4 Evaluation

Small basic evalustion on real-world graph (geo?). In order to show, that performance is reasonable.

Regular quries. Comparison with other DB?

## 5 CONCLUSION

In this paper we propose a number of multiple-source modifications of Azimov's CFPQ algorithm. Evaluation of the proposed modifications on the real-world examples shows that !!!!

Finally, we provide the full-stack support of CFPQ. For our solution we implement corresponding Cypher extension as a part of `libcypher-parser`, integrate the proposed algorithm into RedisGraph, and extend RedisGraph execution plan builder to support extended Cypher queres. We demonstrate, that our solution allows one evaluate not only context-free queryes, but also regular one.

In the future, it is necessary to provide formal translation of Cypher to linear algebra, or find a maximal subset of Cypher which can be translated to linear algebra. There is a number of work on a subset of SPARQL to linear algebra translation, such as [? ], but they are very limited. Deep investigation of this topic helps one to realize limits and restrictions of linear algebra utilization for graph databases. Moreover, it helps to improve existing solutions.

We show that evaluation of regular queryes is possible in practice by using CFPQ algorithm, as far as regular queries is a partial cas of the context-free one. But it seems, that the proposed solution is not optimal. For real-world solutions it is important to provide an optimal unified algorithm for both RPQ and CFPQ. One of possible way to solve this provlem is to use tensor-based algorithm [? ].

Another imprtant task is to compare non-linear-algebra-based approaches to multiple-source CFPQ with the proposed solution. In [? ] Johem Kujpers et.al. shows that all-pairs CFPQ algorithms implemented in Neo4j demonstrate unreasonable performance on real-world data for Neo4j. At the same time, Arseniy Terekhov et.al. shows that matrix-based all-pairs CFPQ algortihm implemented in appropriate linear algebra based graph database (RedisGraph) demonstrates good performance. But in the case of multiple-source scenario, when a number of sources is relatively small, non-linear-algebra-based solutions can be better, because such solutions naturally handle small reqired subgraph.

## REFERENCES

[1] Rustam Azimov and Semyon Grigorev. 2018. Context-free Path Querying by Matrix Multiplication. In *Proceedings of the 1st ACM SIGMOD Joint International Workshop on Graph Data Management Experiences & Systems (GRADES) and Network Data Analytics (NDA) (GRADES-NDA '18)*. ACM, New York, NY, USA, Article 5, 10 pages. https://doi.org/10.1145/3210259.3210264

[2] C. Barrett, R. Jacob, and M. Marathe. 2000. Formal-Language-Constrained Path Problems. *SIAM J. Comput.* 30, 3 (2000), 809–837. https://doi.org/10.1137/S0097539798337716 arXiv:https://doi.org/10.1137/S0097539798337716

[3] P. Cailliau, T. Davis, V. Gadepally, J. Kepner, R. Lipman, J. Lovitz, and K. Ouaknine. 2019. RedisGraph GraphBLAS Enabled Graph Database. In *2019 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*. 285–286.

[4] Timothy Davis. 2019. Algorithm 1000: SuiteSparse:GraphBLAS: Graph Algorithms in the Language of Sparse Linear Algebra. *ACM Trans. Math. Software* (2019). https://doi.org/10.1145/3322125

[5] Jochem Kuijpers, George Fletcher, Nikolay Yakovets, and Tobias Lindaaker. 2019. An Experimental Study of Context-Free Path Query Evaluation Methods. In *Proceedings of the 31st International Conference on Scientific and Statistical Database Management (SSDBM '19)*. ACM, New York, NY, USA, 121–132. https://doi.org/10.1145/3335783.3335791

[6] H. Miao and A. Deshpande. 2019. Understanding Data Science Lifecycle Provenance via Graph Segmentation and Summarization. In *2019 IEEE 35th International Conference on Data Engineering (ICDE)*. 1710–1713.

[7] Nikita Mishin, Iaroslav Sokolov, Egor Spirin, Vladimir Kutuev, Egor Nemchinov, Sergey Gorbatyuk, and Semyon Grigorev. 2019. Evaluation of the Context-Free Path Querying Algorithm Based on Matrix Multiplication. In *Proceedings of the 2Nd Joint International Workshop on Graph Data Management Experiences & Systems (GRADES) and Network Data Analytics (NDA) (GRADES-NDA'19)*. ACM, New York, NY, USA, Article 12, 5 pages. https://doi.org/10.1145/3327964.3328503

[8] Jakob Rehof and Manuel Fähndrich. 2001. Type-Base Flow Analysis: From Polymorphic Subtyping to CFL-Reachability. *SIGPLAN Not.* 36, 3 (Jan. 2001), 54–66. https://doi.org/10.1145/373243.360208

[9] Petteri Sevon and Lauri Eronen. 2008. Subgraph Queries by Context-free Grammars. *Journal of Integrative Bioinformatics* 5, 2 (2008), 157 – 172. https://doi.org/10.1515/jib-2008-100

[10] Arseniy Terekhov, Artyom Khoroshev, Rustam Azimov, and Semyon Grigorev. 2020. Context-Free Path Querying with Single-Path Semantics by Matrix Multiplication. In *Proceedings of the 3rd Joint International Workshop on Graph Data Management Experiences & Systems (GRADES) and Network Data Analytics (NDA) (GRADES-NDA'20)*. Association for Computing Machinery, New York, NY, USA, Article 5, 12 pages. https://doi.org/10.1145/3398682.3399163

[11] Xiaowang Zhang, Zhiyong Feng, Xin Wang, Guozheng Rao, and Wenrui Wu. 2016. Context-Free Path Queries on RDF Graphs. In *The Semantic Web – ISWC 2016*, Paul Groth, Elena Simperl, Alasdair Gray, Marta Sabou, Markus Krötzsch, Freddy Lecue, Fabian Flöck, and Yolanda Gil (Eds.). Springer International Publishing, Cham, 632–648.

[12] 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