

# Context-Free Path Querying Can be Fast if Done Properly

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

Artyom Khoroshev  
arthoroshev@gmail.com  
ITMO 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 recent study showed that the applicability of context-free path querying (CFPQ) algorithms integrated with Neo4j database is limited because of low performance and high memory consumption. In this work we implement a matrix-based CFPQ algorithm by using appropriate high-performance libraries for linear algebra and integrate it with RedisGraph graph database. Our evaluation shows that the provided implementation is, in some cases, up to 1000 times faster than the best Neo4j-based one.

## KEYWORDS

Context-free path querying, transitive closure, graph databases, linear algebra, context-free grammar, GPGPU, CUDA, boolean matrix, matrix multiplication

## 1 INTRODUCTION

Formal language constrained path querying, or formal language constrained path problem [4], is a graph analysis problem in which formal languages are used as constraints for navigational path queries. In this approach a path is viewed as a word constructed by concatenation of edge labels. Paths of interest are constrained with some formal language: a query should find only paths labeled by words from the language. The class of language constraints which is most widely spread is regular: it is used in various graph query languages and engines. Context-free path querying (CFPQ) [27], while being more expressive, is still at the early stage of development. Context-free constraints allow one to express such important class of queries as *same generation queries* [1] which cannot be expressed in terms of regular constraints.

Several algorithms for CFPQ based on such parsing techniques as (G)LL, (G)LR, and CYK were proposed recently [5, 6, 11, 13, 17, 20, 23, 25, 28]. Yet recent research by Jochem Kuijpers et.al. [16] shows that existing solutions are not applicable for real-world graph analysis because of significant running time and memory consumption. At the same time, Nikita Mishin et.al show in [18] that the matrix-based CFPQ algorithm demonstrates good performance on real-world data. A matrix-based algorithm proposed by Rustam Azimov [3] offloads the most critical computations onto boolean matrices multiplication. This algorithm is easy to implement and to employ modern massive-parallel hardware for CFPQ. The paper measures the performance

of the algorithm in isolation while J. Kuijpers provides the evaluation of the algorithms which are integrated with Neo4j<sup>1</sup> graph database. Also, in [16] the matrix-based algorithm is implemented as a simple single-thread Java program, while N. Mishin shows that to achieve the best performance, one should utilize high-performance matrix multiplication libraries which are highly parallel or utilize GPGPU better. Thus, it is required to evaluate a matrix-based algorithm which is integrated with a graph storage and makes use of performant libraries and hardware.

In this work we show that CFPQ in relational semantics (according to Hellings [12]) can be performant enough to be applicable to real-world graph analysis. We use RedisGraph<sup>2</sup> [7] graph database as a storage. This database uses adjacency matrices as a representation of a graph and GraphBLAS [15] for matrices manipulation. These facts allow us to integrate matrix-based CFPQ algorithm with RedisGraph with minimal effort. We make the following contributions in this paper.

- (1) We provide a number of implementations of the CFPQ algorithm which is based on matrix multiplication and uses RedisGraph as graph storage. The first implementation is CPU-based and utilizes SuiteSparse<sup>3</sup> [9] implementation of GraphBLAS API for matrices manipulation. The second implementation is GPGPU-based and includes both the existing implementation from [18] and our own CUSP<sup>4</sup>-based implementation. The source code is available on GitHub<sup>5</sup>.
- (2) We extend the dataset presented in [18] with new real-world and synthetic cases of CFPQ<sup>6</sup>.
- (3) We provide evaluation which shows that matrix-based CFPQ implementation for RedisGraph database is performant enough for real-world data analysis.

<sup>1</sup>Neo4j graph database web page: <https://neo4j.com/>. Access date: 12.11.2019.

<sup>2</sup>RedisGraph is a graph database which is based on Property Graph Model. Project web page: <https://oss.redislabs.com/redisgraph/>. Access date: 12.11.2019.

<sup>3</sup>SuiteSparse is a sparse matrix software which includes GraphBLAS API implementation. Project web page: <http://faculty.cse.tamu.edu/davis/suitesparse.html>. Access date: 12.11.2019.

<sup>4</sup>CUSP is an open source library for sparse matrix multiplication on GPGPU. Project site: <https://cusplibrary.github.io/>. Access date: 12.11.2019.

<sup>5</sup>Sources of matrix-based CFPQ algorithm for RedisGraph database: <https://github.com/YaccConstructor/RedisGraph>. Access date: 12.11.2019.

<sup>6</sup>The CFPQ\_Data dataset from CFPQ algorithms evaluation and comparison. GitHub page: [https://github.com/JetBrains-Research/CFPQ\\_Data](https://github.com/JetBrains-Research/CFPQ_Data). Access date: 12.11.2019.

## 2 THE MOTIVATING EXAMPLE

In this section, we formulate the problem of context-free path query evaluation, using a small graph and the classical *same-generation query* [1], which cannot be expressed using regular expressions.

Let us have a graph database or any other object, which can be represented as a graph. The same-generation query can be used for discovering a vertex similarity, for example, gene similarity [21]. For graph databases, the same-generation query is aimed at the finding all the nodes at the same hierarchy level. The language, formed by the paths between such nodes, is not regular and corresponds to the language of matching parentheses. Hence, the query is formulated as a context-free grammar.

For example, let us have a small double-cyclic graph (see Figure 1). One of the cycles has three edges, labeled with  $a$ , and the other has two edges, labeled with  $b$ . Both cycles are connected via a shared node 0.

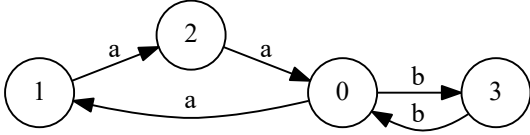


Figure 1: An example graph.

For this graph, we have a same-generation query, formulated as a context-free grammar, which generates a context-free language  $L = \{a^n b^n \mid n \geq 1\}$ .

The result of context-free path query evaluation w.r.t the relational query semantics for this example is a set of node pairs  $(m, n)$ , such that there is a path from the node  $m$  to the node  $n$ , whose labeling forms a word from the language  $L$ . For example, the node pair  $(0, 0)$  must be in this set, since there is a path from the node 0 to the node 0, whose labeling forms a string  $w = aaaaaabbbbb = a^6 b^6 \in L$ .

The result of context-free path query evaluation w.r.t the single-path query semantics also contains such a path for each node pair  $(m, n)$  returned after the context-free path query evaluation w.r.t the relational query semantics. For example, a path from the node 0 to the node 0, whose labeling forms a string  $w = a^6 b^6$  can be returned for the node pair  $(0, 0)$ .

## 3 PRELIMINARIES

Let  $\Sigma$  be a finite set of edge labels. Define an *edge-labeled directed graph* as a tuple  $D = (V, E)$  with a set of nodes  $V$  and a directed edge relation  $E \subseteq V \times \Sigma \times V$ .

A path  $\pi$  is a list of labeled edges  $[e_1, \dots, e_n]$  where  $e_i \in E$ . The concatenation of a path  $\pi_1$  with a path  $\pi_2$  we denote by  $\pi_1 + \pi_2$ .

For a path  $\pi$  in a graph  $D$ , we denote the unique word, obtained by concatenating the labels of the edges along the path  $\pi$  as  $l(\pi)$ . Also, we write  $n\pi m$  to indicate, that the path  $\pi$  starts at the node  $n \in V$  and ends at the node  $m \in V$ .

Following Hellings [12], we deviate from the usual definition of a context-free grammar in *Chomsky Normal Form* [?] by not including a special starting non-terminal, which will be specified in the path queries for the graph. Since

every context-free grammar can be transformed into an equivalent one in Chomsky Normal Form and checking, that empty string belongs to the language is trivial, it is sufficient to consider only grammars of the following type. A *context-free grammar* is a triple  $G = (N, \Sigma, P)$ , where  $N$  is a finite set of non-terminals,  $\Sigma$  is a finite set of terminals, and  $P$  is a finite set of productions of the following forms:

- $A \rightarrow BC$ , for  $A, B, C \in N$ ,
- $A \rightarrow x$ , for  $A \in N$  and  $x \in \Sigma \cup \{\varepsilon\}$ .

We use the conventional notation  $A \xrightarrow{*}_G w$  to denote, that a string  $w \in \Sigma^*$  can be derived from a non-terminal  $A$  by some sequence of production rule applications from  $P$  in grammar  $G$ . The *language* of a grammar  $G = (N, \Sigma, P)$  with respect to a start non-terminal  $S \in N$  is defined by

$$L(G_S) = \{w \in \Sigma^* \mid S \xrightarrow{*}_G w\}.$$

For a given graph  $D = (V, E)$  and a context-free grammar  $G = (N, \Sigma, P)$ , we define *context-free relations*  $R_A \subseteq V \times V$  for every  $A \in N$ , such that

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

For the context-free path query evaluation w.r.t. the single-path query semantics, we must provide such a path for each node pair from  $R_A$ . In order to do this, we introduce the

$$PathIndex = (left, right, middle, height, length)$$

— the elements of matrices which describe the found paths as concatenations of two smaller paths and help to restore each path and derivation tree for it at the end of evaluation. Here *left* and *right* stand for the indexes of starting and ending node in the founded path, *middle* — the index of intermediate node used in the concatenation of two smaller paths, *height* — the height of the derivation tree of the string corresponding to the founded path, and *length* is a length of founded path. When we do not find the path for some node pair  $i, j$ , we use the  $PathIndex = \perp = (0, 0, 0, 0, 0)$ .

Also, we will use the notation *proper matrix* which means that for every element of the matrix with indexes  $i, j$  it either  $PathIndex = (i, j, -, -, -)$  or  $\perp$ .

For proper matrices we use a binary operation  $\otimes$  defined for PathIndexes  $PI_1, PI_2$  which are not equal to  $\perp$  and with  $PI_1.right = PI_2.left$  as

$$PI_1 \otimes PI_2 = (PI_1.left, PI_2.right, PI_1.right, \max(PI_1.height, PI_2.height)+1, PI_1.length+PI_2.length).$$

If at least one of the operands is equal to  $\perp$  then  $PI_1 \otimes PI_2 = \perp$ .

For proper matrices we also use a binary operation  $\oplus$  defined for PathIndexes  $PI_1, PI_2$  with  $PI_1.left = PI_2.left$  and  $PI_1.right = PI_2.right$  as  $PI_1$  if  $PI_1.height \leq PI_2.height$  and  $PI_2$  otherwise.

Using  $\otimes$  as multiplication of PathIndexes, and  $\oplus$  as an addition, we can define a *matrix multiplication*,  $a \odot b = c$ , where  $a$  and  $b$  are matrices of a suitable size, that have PathIndexes as elements, as

$$c_{i,j} = \bigoplus_{k=1}^n a_{i,k} \otimes b_{k,j}.$$

Also, we use the element-wise  $+$  operation on matrices  $a$  and  $b$  with the same size:  $a + b = c$ , where  $c_{i,j} = a_{i,j} \oplus b_{i,j}$ .

## 4 MATRIX-BASED ALGORITHM FOR CFPQ

The matrix-based algorithm for CFPQ was proposed by Rustam Azimov [3]. This algorithm can be expressed in terms of operations over boolean matrices (see listing 1) which is an advantage for implementation.

**Listing 1** Context-free path quering algorithm

---

```

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

```

---

Here  $D = (V, E)$  is the input graph and  $G = (N, \Sigma, P)$  is the input grammar. For each matrix  $T^{A_k}$  indexed with a nonterminal  $A_k \in N$ , a cell holds a true value ( $T_{i,j}^{A_k} = \text{true}$ ) if and only if there exists  $i\pi j$  — a path in  $D$  such that  $A_k \xrightarrow{*}_G l(\pi)$ , where  $l(\pi)$  is a word formed by the labels along the path  $\pi$ . Thus, this algorithm solves the reachability problem, or, according to Hellings [12], implements relational query semantics.

The performance-critical part of the algorithm is boolean matrix multiplication, thus one can achieve better performance by using libraries which efficiently multiply boolean matrices. There is also the following optimization: if the matrices  $T^{A_j}$  and  $T^{A_k}$  have not changed at the previous iteration, then we can skip the update operation in line 7. Data in real-world problems is often sparse, thus employing libraries which manipulate sparse matrices improves running time even more.

## 5 MATRIX BASED CFPQ FOR SINGLE-PATH SEMANTICS

In this section, we propose the matrix-based algorithm for CFPQ w.r.t. the single-path query semantics (see listing 2). This algorithm constructs the set of matrices  $T$  with PathIndexes as elements.

**Listing 2** CFPQ algorithm w.r.t. single-path query semantics

---

```

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

```

---

After constructing the set of matrices  $T$  for every node pair  $i, j$  and nonterminal  $A$  we can extract a path  $i\pi j$  from  $i$  to  $j$  such that  $A \xrightarrow{*}_G l(\pi)$  if such path exists. We also

propose the algorithm (see listing 3) for extracting one of those paths which forms a string with minimal height of derivation tree. Our algorithm returns the empty path  $\perp$  only if  $i = j$  and  $A \rightarrow \varepsilon \in P$ . Note that if the PathIndex for given  $i, j, A$  is equal to  $\perp$  then our algorithm returns a special path  $\pi_\emptyset$  to denote that such a path does not exist.

**Listing 3** Path extraction algorithm

---

```

1: function EXTRACTPATH( $i, j, A, T = \{T^{A_i}\}, G = (N, \Sigma, P)$ )
2:    $index \leftarrow T_{i,j}^A$ 
3:   if  $index = \perp$  then
4:     return  $\pi_\emptyset$  ▷ Such a path does not exist
5:   if  $index.height = 1$  then
6:     if  $index.length = 0$  then
7:       return  $\perp$  ▷ Return an empty path
8:     for all  $x \mid (i, x, j) \in E$  do
9:       if  $A \rightarrow x \in P$  then
10:        return  $[(i, x, j)]$  ▷ Return a path of length one
11:   for all  $A \rightarrow BC \in P$  do
12:      $index_B \leftarrow T_{i, index.middle}^B$ 
13:      $index_C \leftarrow T_{index.middle, j}^C$ 
14:     if  $(index_B \neq \perp) \wedge (index_C \neq \perp)$  then
15:        $maxH \leftarrow \max(index_B.height, index_C.height)$ 
16:       if  $index.height = maxH + 1$  then
17:          $\pi_1 \leftarrow \text{EXTRACTPATH}(i, index.middle, B, T, G)$ 
18:          $\pi_2 \leftarrow \text{EXTRACTPATH}(index.middle, j, C, T, G)$ 
19:       return  $\pi_1 + \pi_2$ 

```

---

### 5.1 Correctness

Let  $T^{(p)} = \{T^{(p), A_i}\}$  be a constructed matrix  $T$  by the algorithm in listing 2 after  $p - 1$  loop iterations for  $p \geq 2$ , and  $T^{(1)} = \{T^{(1), A_i}\}$  be a constructed matrix  $T$  by this algorithm after initialization in lines 3-5. Note that the matrix  $T$  returned by this algorithm is equal to  $\sum_{p=1}^{\infty} T^{(p)}$ . Then the following lemma and theorem hold.

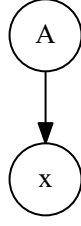
**LEMMA 5.1.** *Let  $D = (V, E)$  be a graph, let  $G = (N, \Sigma, P)$  be a grammar. Then for any  $i, j$  and for any non-terminal  $A \in N$ ,  $index = T_{i,j}^{(p), A}$  and  $index = (i, j, k, h, l) \neq \perp$  iff  $(i, j) \in R_A$  and  $i\pi j$ , such that there is a derivation tree of the minimal height  $h \leq p$  for the string  $l(\pi)$  of length  $l$  and a context-free grammar  $G_A = (N, \Sigma, P, A)$ .*

**PROOF.** (Proof by Induction)

**Base case:** Show that the lemma holds for  $p = 1$ . For any  $i, j$  and for any non-terminal  $A \in N$ ,  $(i, j, k, h, l) = T_{i,j}^{(1), A}$  iff there is either  $i\pi j$  of length 1 that consists of a unique edge  $e$  from the node  $i$  to the node  $j$  and  $(A \rightarrow x) \in P$ , where  $x = l(\pi)$ , or  $i = j$  and  $(A \rightarrow \varepsilon) \in P$ , where  $\varepsilon = l(\pi)$ . Therefore  $(i, j) \in R_A$  and there is a derivation tree of the minimal height  $h = p = 1$ , shown on Figure 2, for the string  $x$  and a context-free grammar  $G_A = (N, \Sigma, P, A)$ . Thus, it has been shown that the lemma holds for  $p = 1$ .

**Inductive step:** Assume that the lemma holds for any  $p \leq (q - 1)$  and show that it also holds for  $p = q$ , where  $q \geq 2$ .

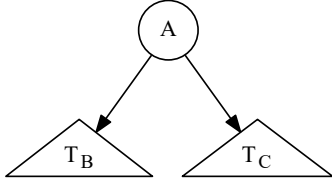
The index  $(i, j, k, h, l) = T_{i,j}^{(q), A}$  iff there exists a rule  $(A \rightarrow BC) \in P$  such that  $(i, j, k, h, l) = M_{i,j}$  where  $M = T^{(q-1), A} + (T^{(q-1), B} \odot T^{(q-1), C})$ .



**Figure 2: The derivation tree of the minimal height  $p = 1$  for the string  $x = l(\pi)$  where  $x \in \Sigma \cup \{\varepsilon\}$ .**

Let  $(i, j, k, h, l) = T_{i,j}^{(q-1),A}$ . By the inductive hypothesis,  $(i, j, k, h, l) = T_{i,j}^{(q-1),A}$  iff  $(i, j) \in R_A$  and there exists  $i\pi j$ , such that there is a derivation tree of the minimal height  $h \leq (q-1)$  for the string  $l(\pi)$  and a context-free grammar  $G_A = (N, \Sigma, P, A)$ . The statement of the lemma holds for  $p = q$  since the height  $h$  of this tree is also less than or equal to  $q$ .

Now, let  $(i, j, k, h, l) = (T^{(q-1),B} \odot T^{(q-1),C})_{i,j}$ . By the definition of the binary operation  $\odot$ ,  $(i, j, k, h, l) = (T^{(q-1),B} \odot T^{(q-1),C})_{i,j}$  iff there are  $r = k$ ,  $(i, r, -, h_1, l_1) = T_{i,r}^{(q-1),B}$  and  $(r, j, -, h_2, l_2) = T_{r,j}^{(q-1),C}$ , such that  $q = \max(h_1, h_2) + 1$ ,  $l = l_1 + l_2$ . Hence, by the inductive hypothesis, there are  $i\pi_1 r$  and  $r\pi_2 j$ , such that  $(i, r) \in R_B$  and  $(r, j) \in R_C$ , and there are the derivation trees  $T_B$  and  $T_C$  of minimal heights  $h_1 \leq (q-1)$  and  $h_2 \leq (p-1)$  for the strings  $w_1 = l(\pi_1)$ ,  $w_2 = l(\pi_2)$  and the context-free grammars  $G_B$ ,  $G_C$  respectively. Thus, the concatenation of paths  $\pi_1$  and  $\pi_2$  is  $i\pi j$ , where  $(i, j) \in R_A$  and there is a derivation tree of the minimal height  $h = 1 + \max(h_1, h_2)$ , shown on Figure 3, for the string  $w = l(\pi)$  of length  $l = l_1 + l_2$  and a context-free grammar  $G_A$ .



**Figure 3: The derivation tree of the minimal height  $h = 1 + \max(h_1, h_2)$  for the string  $w = l(\pi)$ , where  $T_B$  and  $T_C$  are the derivation trees for strings  $w_1$  and  $w_2$  respectively.**

The statement of the lemma holds for  $p = q$  since the minimal height  $h = 1 + \max(h_1, h_2) \leq q$ . This completes the proof of the lemma.  $\square$

**THEOREM 1.** 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  $A \in N$ ,  $index = T_{i,j}^A$  and  $index = (i, j, k, h, l) \neq \perp$  iff  $(i, j) \in R_A$  and  $i\pi j$ , such that there is a derivation tree of the minimal height  $h$  for the string  $l(\pi)$  of length  $l$  and a context-free grammar  $G_A = (N, \Sigma, P, A)$ .

**PROOF.** Since the matrix  $T = \sum_{p=1}^{\infty} T^{(p)}$  for any  $i, j$  and for any non-terminal  $A \in N$ ,  $index = T_{i,j}^A$  and  $index =$

$(i, j, k, h, l) \neq \perp$  iff there is  $p \geq 1$ , such that  $index \in T_{i,j}^{(p),A}$ . By the lemma 5.1,  $index = T_{i,j}^{(p),A}$  iff  $(i, j) \in R_A$  and  $i\pi j$ , such that there is a derivation tree of the minimal height  $h \leq p$  for the string  $l(\pi)$  of length  $l$  and a context-free grammar  $G_A = (N, \Sigma, P, A)$ . This completes the proof of the theorem.  $\square$

We can, therefore, determine whether  $(i, j) \in R_A$  by asking whether  $\perp = T_{i,j}^A$ . Also, we can extract such a path which forms a string with a derivation tree of minimal height by using our algorithm in listing 2. Thus, we show how the context-free path query evaluation w.r.t. the single-path semantics can be solved in terms of matrix operations.

Correctness of path extraction algorithm?

## 5.2 Complexity

Complexity analysis.

## 5.3 Example?

Example

## 6 IMPLEMENTATION

We showed that CFPQ can be naturally reduced to linear algebra. Linear algebra for graph problems is an actively developed area. One of the most important results is a GraphBLAS API which provides a way to operate over matrices and vectors over user-defined semirings.

Previous works show [3, 18] that existing linear algebra libraries utilization is the right way to achieve high-performance CFPQ implementation with minimal effort. But neither of these works provide an evaluation with data storage: algorithm execution time has been measured in isolation.

We provide a number of implementations of the matrix-based CFPQ algorithm. We use RedisGraph as storage and implement CFPQ as an extension by using the mechanism provided. Note that currently, we do not provide complete integration with the querying mechanism: one cannot use Cypher — a query language used in RedisGraph. Instead, a query should be provided explicitly as a file with grammar in Chomsky normal form. This is enough to evaluate querying algorithms and we plan to improve integration in the future to make our solution easier to use.

**CPU-based implementation (RG\_CPU)** uses SuteS-parse implementation of GraphBLAS, which is also used in RedisGraph, and a predefined boolean semiring. Thus we avoid data format issues: we use native RedisGraph representation of the adjacency matrix in our algorithm.

**GPGPU-based implementation** has four versions. The first one (**RG\_M4RI<sub>rel</sub>**) uses the Method of the Four Russians implemented in [18], the second one (**RG\_CUSP<sub>rel</sub>**) utilizes a CUSP [8] library for matrix operations, the third one (**RG\_SPARSE<sub>rel</sub>**) is our implementation based on the idea from this paper [19] and the fourth one (**RG\_SPARSE<sub>path</sub>**) is our implementation of single-path semantics. First two implementations require matrix format conversion but the last two does not.

The first (**RG\_M4RI<sub>rel</sub>**) implementation works with dense matrices so it cannot be applicable to huge graphs. Therefore, we choose the CUSP library as base solution that uses sparse matrices. CUSP is a C++ templated library

which allows us to multiply boolean matrices (that solve relational path semantic problem). But in fact, performing CFPQ with relational paths semantics on the largest graph using CUSP does not fit in GPU memory and this fact led us to develop an algorithm that would be more memory efficient.

The third (**RG\_SPARSE<sub>rel</sub>**) implementation utilizes low-latency on-chip shared memory for hash table of each row of the result matrix. For more details of algorithm see the original paper [19]. An original solution designed for single and double precision SpGEMM. Since we have a boolean matrix in CSR format, we can discard the array of values and optimize usage of shared memory. But boolean matrix multiplication is only one part of the algorithm, since we must effectively combine two boolean sparse matrices. We use merge path [10] algorithm to merge corresponding rows of the result matrix.

The fourth (**RG\_SPARSE<sub>path</sub>**) algorithm must perform SpGEMM on *PathIndex*. It is more complex task than boolean SpGEMM because domain of operation changed and it is no way to perform atomically  $\oplus$  defined for *PathIndex*s. But for every element of the matrix with indexes  $i, j$  it either  $PathIndex = (i, j, \rightarrow, \rightarrow, -)$  or  $\perp$ . Also note that *length* information does not matter in the algorithm and can be restored later, so only two elements are really important: *middle* and *height*. We can store two 4 bytes value into one 8 byte value and perform an atomic operation. In high four bytes we store the *height* and in the low four bytes we store the *middle*. For the value of  $\perp$  we use maximum unsigned integer value of 8 bytes in size. Now we can use *atomicMin* as  $\oplus$ . As a result the algorithm consist of two steps. On the first step we evaluate the **RG\_SPARSE<sub>rel</sub>** algorithm to determinate final structure of *proper matrix*. On the second step we fill matrices with  $\perp$  or edges and start multiplication.

## 7 DATASET DESCRIPTION

In our evaluation we use combined dataset which contains the following parts.

- CFPQ\_Data dataset which is provided in<sup>7</sup> [18] and contains both synthetic and real-world graphs and queries. Real-world data includes RDFs, synthetic cases include theoretical worst-case and random graphs.
- Dataset which is provided in [16]. Both Geospecies (RDF which contains information about biological hierarchy<sup>8</sup> and same generation query over *broader-Transitive* relation), and Synthetic (the set of graphs generated by using the Barabási-Albert model [2] of scale-free networks and same generation query), are integrated with CFPQ\_Data and used in our evaluation.
- It was shown in [18] that matrix-based algorithm is performant enough to handle bigger RDFs than those used in the initial datasets, such as [28]. So, we add a number of big RDFs to CFPQ\_Data and use them in our evaluation. New RDFs: *go-hierarchy*, *go*,

*enzyme*, *core*, *pathways* are from UniProt database<sup>9</sup>, and *eclass-514en* is from eClassOWL project<sup>10</sup>.

The variants of the *same generation query* [1] are used in almost all cases because it is an important example of real-world queries that are context-free but not regular. So, variations of the same generation query are used in our evaluation. All queries are added to the CFPQ\_Data dataset.

For RDFs ([**RDF**] dataset) we use two queries over *subClassOf* and *type* relations. The first query is the grammar  $G_1$ :

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

The second one is the grammar  $G_2$ :

$$s \rightarrow \text{subClassOf}^{-1} s \text{ subClassOf} \mid \text{subClassOf}$$

For geospecies and free scale graphs querying we use same-generation queries from the original paper.

## 8 EVALUATION AND DISCUSSION

We evaluate all the described implementations on all the datasets and the queries presented. We compare our implementations with [18] and [16]. We measure the full time of query execution including all overhead on data preparation. This way we can estimate the applicability of the matrix-based algorithm to real-world problems.

For evaluation, we use a PC with Ubuntu 18.04 installed. It has Intel core i7-6700 CPU, 3.4GHz, DDR4 32Gb RAM, and Geforce GTX 1070 GPGPU with 8Gb RAM.

The results of the evaluation are summarized in the tables below. We provide results only for a part of the collected dataset because of the page limit. Running time is measured in seconds, RAM memory consumption is measured in megabytes unless specified otherwise. Note that we provide results from the corresponding papers for all implementations except our own. The cell is left blank if the time limit is exceeded, or if there is not enough memory to allocate the data.

The results of the first dataset [**RDF**] are presented in table 1. We can see that the running time of both CPU and GPGPU versions is small even for graphs with a big number of vertices and edges. The relatively small number of edges of interest may be the reason for such behavior. We believe it is necessary to extend the dataset with new queries which involve more different types of edges. Also, we can see, that *m4ri* version which uses dense bit matrices requires more memory. Thus we recommend to use sparse matrices on GPGPU.

Geospecies dataset currently can be processed only by using CPU version and we compare our matrix-based CPU implementation with the result from [16] for *AnnGram<sub>rel</sub>* algorithm<sup>11</sup>. Fortunately, both algorithms calculate queries

<sup>7</sup>CFPQ\_Data dataset GitHub repository: [https://github.com/JetBrains-Research/CFPQ\\_Data](https://github.com/JetBrains-Research/CFPQ_Data). Access date: 12.11.2019.

<sup>8</sup><https://old.datahub.io/dataset/geospecies>. Access date: 12.11.2019.

<sup>9</sup>Protein sequences data base: <https://www.uniprot.org/>. RDFs with data are available here: [ftp://ftp.uniprot.org/pub/databases/uniprot/current\\_release/rdf](ftp://ftp.uniprot.org/pub/databases/uniprot/current_release/rdf). Access date: 12.11.2019

<sup>10</sup>eClassOWL project: <http://www.heppnetz.de/projects/eclassowl/>. eclass-514en file is available here: <http://www.ebusiness-unibw.org/ontologies/eclass/5.1.4/eclass.514en.owl>. Access date: 12.11.2019.

<sup>11</sup>Only *AnnGram* works correctly and fits limits, other implementations are faster, but either return an incorrect result or do not fit the memory.

**Table 1: RDFs querying results**

RDF					Query $G_1$			Query $G_2$		
Name	#V	#E	#type	#subClassOf	RG_CPU	RG_M4RI	RG_CUSP	RG_CPU	RG_M4RI	RG_CUSP
funding	778	1480	304	90	0.01	0.01	0.02	0.01	0.01	0.01
pizza	671	2604	365	259	0.01	0.01	0.02	0.01	0.01	0.01
wine	733	2450	485	126	0.01	0.01	0.02	0.01	0.01	0.01
core	1323	8684	1412	178	0.01	0.12	0.02	0.01	0.01	0.01
pathways	6238	37196	3118	3117	0.01	0.18	0.03	0.01	0.06	0.01
go-hierarchy	45007	1960436	0	490109	0.09	-	1.50	0.01	-	0.55
enzyme	48815	219390	14989	8163	0.02	61.23	0.10	0.01	6.97	0.02
eclass_514en	239111	1047454	72517	90962	0.06	-	0.39	0.01	-	0.10
go	272770	1068622	58483	90512	0.49	-	0.83	0.01	-	0.11

under relational semantics. The result is provided in the table 2.

**Table 2: Evaluation results on geospecies data**

RG_CPU		Neo4j_AnnGram <sub>rel</sub>	
Time	Memory (Gb)	Time	Memory (Gb)
6.8	6.83	6 953.9	29.17

As we can see, the matrix-based algorithm implemented for RedisGraph is more than 1000 times faster than the one based on annotated grammar implemented for Neo4j and uses more than 4 times less memory. We can conclude that the matrix-based algorithm is better than other CFPQ algorithms for query evaluation under a relational semantics for real-world data processing. CFPQ evaluation under other semantics (single path, all paths, etc) by using a matrix-based algorithm is a direction for future research.

The next is the [FreeScale] dataset. We compare our implementations with two implementations from [16] which evaluate queries under relational semantics: *Neo4j\_AnnGram<sub>rel</sub>* and *Neo4j\_Matrix*. The results are presented in table 3. The evaluation shows that sparsity of graphs (value of parameter  $p$ ) is important both for implementations which use sparse matrices and for implementations which use dense matrices. Note that the results for implementations for Neo4j are restored from graphics provided in [16]. So, values are not precise, but it is possible to compare implementations.

Evaluation shows that our CPU version is comparable with *Neo4j\_AnnGram<sub>rel</sub>* and for relatively dense graphs (each vertex has 10 connections) our implementation is faster. Moreover, while *Neo4j\_Matrix* exceeded limits on the biggest graph, our implementation works fine. This demonstrates the importance of using of appropriate libraries for matrix-based algorithm implementation. Also, we can see, that GPGPU version which utilizes sparse matrices is significantly faster than the other implementations. Note, that for GPGPU versions we include time required for data transferring and formats conversion.

Finally, we conclude that the matrix-based algorithm paired with a suitable database and employing appropriate libraries for linear algebra is a promising way to make CFPQ applicable for real-world data analysis. We show that SuiteSparse-based CPU implementation is performant enough to be comparable with GPGPU-based implementations on real-world data. It means that we can handle

more complex data. We can also see, that more complex queries should be added to the dataset to make it more representable.

## 9 CONCLUSION AND FUTURE WORK

We implemented a CPU and GPGPU based context-free path querying for RedisGraph and showed that CFPQ can be performant enough to analyze real-world data. However, our implementations are prototypes and we plan to provide full integration of CFPQ to RedisGraph. First of all, it is necessary to extend Cypher graph query language used in RedisGraph to support syntax for specification of context-free constraints. There is a proposal which describes such syntax extension<sup>12</sup> and we plan to support this syntax in libcypher-parser<sup>13</sup> used in RedisGraph.

Current version uses CUSP matrix multiplication library for GPGPU utilization, but it may be better to use GraphBLAST<sup>14</sup> [26] — Gunrock<sup>15</sup> [24] based implementation of GraphBLAS API for GPGPU. We plan to evaluate GraphBLAST based implementation of CFPQ and to investigate how multi-GPU support for GraphBLAST influences the performance of CFPQ in the case of processing huge real-world data.

Our implementations compute relational semantics of a query, but some problems require to find a path which satisfies the constraints. To the best of our knowledge, there is no matrix-based algorithm for single path or all path semantics, thus we see it as a direction for future research.

Another important open question is how to update the query results dynamically when data changes. The mechanism for result updating allows one to recalculate query faster and use the result as an index for other queries.

Also, further improvements of the dataset are required. For example, it is necessary to include real-world cases from the area of static code analysis [14, 22, 29].

<sup>12</sup>Proposal with path pattern syntax for openCypher: <https://github.com/thobe/openCypher/blob/rpq/cip/1.accepted/CIP2017-02-06-Path-Patterns.adoc>. It is shown that context-free constraints can be expressed with the proposed syntax. Access date: 12.11.2019

<sup>13</sup>Web page of libcypher-parser project: <http://cleishm.github.io/libcypher-parser/>. Access date: 12.11.2019

<sup>14</sup>GraphBLAST project: <https://github.com/gunrock/graphblast>. Access date: 12.11.2019.

<sup>15</sup>Gunrock project: <https://gunrock.github.io/docs/>. Access date: 12.11.2019.

**Table 3: Free scale graphs querying results**

Graph	RG_CPU		RG_m4ri		RG_CUSP		Neo4j_AnnGram <sub>rel</sub>		Neo4j_Matrix	
	Time	Mem	Time	Mem	Time	Mem	Time	Mem	Time	Mem
G(100,1)	↓ 0.01	↓ 0.01	↓ 0.01	0.10	0.01	2.00	↓ 0.02	0.08	0.20	0.03
G(100,3)	↓ 0.01	↓ 0.01	↓ 0.01	0.10	0.04	2.00	0.02	0.15	0.40	0.03
G(100,5)	↓ 0.01	↓ 0.01	↓ 0.01	0.10	0.05	2.00	0.03	0.21	0.40	0.03
G(100,10)	↓ 0.01	↓ 0.01	0.01	0.10	0.07	2.00	0.09	0.60	0.60	0.03
G(500,1)	↓ 0.01	↓ 0.01	↓ 0.01	2.00	0.01	2.00	↓ 0.02	0.20	20.00	0.60
G(500,3)	↓ 0.01	↓ 0.01	↓ 0.01	2.00	0.07	2.00	0.03	0.50	40.00	0.60
G(500,5)	↓ 0.01	0.17	↓ 0.01	2.00	0.10	2.00	0.10	1.10	50.00	0.60
G(500,10)	1.24	0.78	0.01	2.00	0.11	4.00	0.50	4.00	55.00	0.60
G(2500,1)	↓ 0.01	0.11	0.07	30.00	0.03	2.00	0.03	0.70	0.023	14.00
G(2500,3)	0.01	0.11	0.11	30.00	0.10	2.00	0.15	2.50	0.105	14.00
G(2500,5)	2.06	0.11	0.11	30.00	0.12	4.00	0.70	8.00	1.636	14.00
G(2500,10)	3.25	3.77	0.13	30.00	0.31	31.20	5.00	20.00	13.071	14.00
G(10000,1)	↓ 0.01	0.47	1.55	200.00	0.04	2.0	0.10	2.50	-	-
G(10000,3)	5.439	1.15	3.60	200.00	0.20	3.20	0.40	10.00	-	-
G(10000,5)	7.978	2.64	3.32	200.00	0.25	13.20	3.00	35.00	-	-
G(10000,10)	13.180	21.08	3.60	200.00	1.23	198.00	40.00	240.00	-	-

## ACKNOWLEDGMENTS

The research was supported by the Russian Science Foundation grant 18-11-00100 and a grant from JetBrains Research.

## REFERENCES

- [1] Serge Abiteboul, Richard Hull, and Victor Vianu. 1995. Foundations of Databases.
- [2] Réka Albert and Albert-lászló Barabási. [n.d.]. Statistical mechanics of complex networks. *Rev. Mod. Phys.* ([n.d.]), 2002.
- [3] 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>
- [4] Chris Barrett, Riko Jacob, and Madhav Marathe. 2000. Formal-language-constrained path problems. *SIAM J. Comput.* 30, 3 (2000), 809–837.
- [5] Phillip G Bradford. 2007. Quickest path distances on context-free labeled graphs. In *Appear in 6-th WSEAS Conference on Computational Intelligence, Man-Machine Systems and Cybernetics*. Citeseer.
- [6] Phillip G Bradford and Venkatesh Choppella. 2016. Fast point-to-point Dyck constrained shortest paths on a DAG. In *2016 IEEE 7th Annual Ubiquitous Computing, Electronics & Mobile Communication Conference (UEMCON)*. IEEE, 1–7.
- [7] P. Cailliau, T. Davis, V. Gadepally, J. Kepner, R. Lipman, J. Lovitz, and K. Ouaknine. 2019. RedisGraph Graph-BLAS Enabled Graph Database. In *2019 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*. 285–286. <https://doi.org/10.1109/IPDPSW.2019.00054>
- [8] Steven Dalton, Nathan Bell, Luke Olson, and Michael Garland. 2014. Cusp: Generic Parallel Algorithms for Sparse Matrix and Graph Computations. <http://cusplibrary.github.io/> Version 0.5.0.
- [9] Timothy A. Davis. 2018. Algorithm 9xx: SuiteSparse:GraphBLAS: graph algorithms in the language of sparse linear algebra.
- [10] Oded Green, Rob Mccoll, and David Bader. 2014. GPU merge path: a GPU merging algorithm. *Proceedings of the International Conference on Supercomputing*. <https://doi.org/10.1145/2304576.2304621>
- [11] 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>
- [12] Jelle Hellings. 2014. Conjunctive context-free path queries. In *Proceedings of ICDT'14*. 119–130.
- [13] Jelle Hellings. 2015. Querying for Paths in Graphs using Context-Free Path Queries. *arXiv preprint arXiv:1502.02242* (2015).
- [14] Nicholas Hollingum and Bernhard Scholz. 2017. Cauliflower: a Solver Generator for Context-Free Language Reachability. In *LPAR-21. 21st International Conference on Logic for Programming, Artificial Intelligence and Reasoning (EPIC Series in Computing)*, Thomas Eiter and David Sands (Eds.), Vol. 46. EasyChair, 171–180. <https://doi.org/10.29007/tbm7>
- [15] J. Kepner, P. Aaltonen, D. Bader, A. Buluc, F. Franchetti, J. Gilbert, D. Hutchison, M. Kumar, A. Lumsdaine, H. Meyerhenke, S. McMillan, C. Yang, J. D. Owens, M. Zalewski, T. Mattson, and J. Moreira. 2016. Mathematical foundations of the GraphBLAS. In *2016 IEEE High Performance Extreme Computing Conference (HPEC)*. 1–9. <https://doi.org/10.1109/HPEC.2016.7761646>
- [16] 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>
- [17] Ciro M. Medeiros, Martin A. Musicante, and Umberto S. Costa. 2018. Efficient Evaluation of Context-free Path Queries for Graph Databases. In *Proceedings of the 33rd Annual ACM Symposium on Applied Computing (SAC '18)*. ACM, New York, NY, USA, 1230–1237. <https://doi.org/10.1145/3167132.3167265>
- [18] 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>
- [19] Y. Nagasaka, A. Nukada, and S. Matsuoka. 2017. High-Performance and Memory-Saving Sparse General Matrix-Matrix Multiplication for NVIDIA Pascal GPU. In *2017 46th International Conference on Parallel Processing (ICPP)*. 101–110. <https://doi.org/10.1109/ICPP.2017.19>
- [20] Fred C. Santos, Umberto S. Costa, and Martin A. Musicante. 2018. A Bottom-Up Algorithm for Answering Context-Free Path Queries in Graph Databases. In *Web Engineering*, Tommi Mikkonen, Ralf Klamma, and Juan Hernández (Eds.). Springer International Publishing, Cham, 225–233.
- [21] Petteri Sevon and Lauri Eronen. 2008. Subgraph queries by context-free grammars. *Journal of Integrative Bioinformatics* 5, 2 (2008), 100.
- [22] Jyothi Vedurada and V Krishna Nandivada. [n.d.]. Batch Alias Analysis. ([n.d.]).

- [23] Ekaterina Verbitskaia, Ilya Kirillov, Ilya Nozkin, and Semyon Grigorev. 2018. Parser Combinators for Context-free Path Querying. In *Proceedings of the 9th ACM SIGPLAN International Symposium on Scala (Scala 2018)*. ACM, New York, NY, USA, 13–23. <https://doi.org/10.1145/3241653.3241655>
- [24] Yangzihao Wang, Yuechao Pan, Andrew Davidson, Yuduo Wu, Carl Yang, Leyuan Wang, Muhammad Osama, Chen-shan Yuan, Weitang Liu, Andy T. Riffel, and John D. Owens. 2017. Gunrock: GPU Graph Analytics. *ACM Trans. Parallel Comput.* 4, 1, Article 3 (Aug. 2017), 49 pages. <https://doi.org/10.1145/3108140>
- [25] Charles B. Ward, Nathan M. Wiegand, and Phillip G. Bradford. 2008. A Distributed Context-Free Language Constrained Shortest Path Algorithm. In *Proceedings of the 2008 37th International Conference on Parallel Processing (ICPP '08)*. IEEE Computer Society, Washington, DC, USA, 373–380. <https://doi.org/10.1109/ICPP.2008.67>
- [26] Carl Yang, Aydin Buluc, and John D. Owens. 2019. Graph-BLAST: A High-Performance Linear Algebra-based Graph Framework on the GPU. *arXiv:cs.DC/1908.01407*
- [27] Mihalīs Yannakakis. 1990. Graph-theoretic Methods in Database Theory. In *Proceedings of the Ninth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems (PODS '90)*. ACM, New York, NY, USA, 230–242. <https://doi.org/10.1145/298514.298576>
- [28] X. Zhang, Z. Feng, X. Wang, G. Rao, and W. Wu. 2016. Context-free path queries on RDF graphs. In *International Semantic Web Conference*. Springer, 632–648.
- [29] Xin Zheng and Radu Rugina. 2008. Demand-driven Alias Analysis for C. *SIGPLAN Not.* 43, 1 (Jan. 2008), 197–208. <https://doi.org/10.1145/1328897.1328464>