

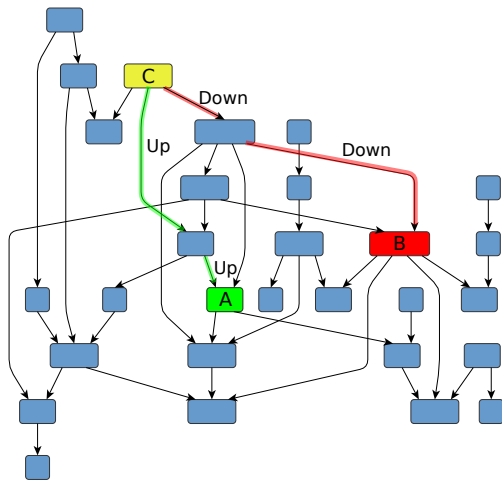
Multiple-Source Context-Free Path Querying in Terms of Linear Algebra

Arseniy Terekhov, Vlada Pogozhelskaya, Vadim Abzalov, Timur Zinnatulin, **Semyon Grigorev**

JetBrains Research, Programming Languages and Tools Lab
Saint Petersburg University

March 24, 2021

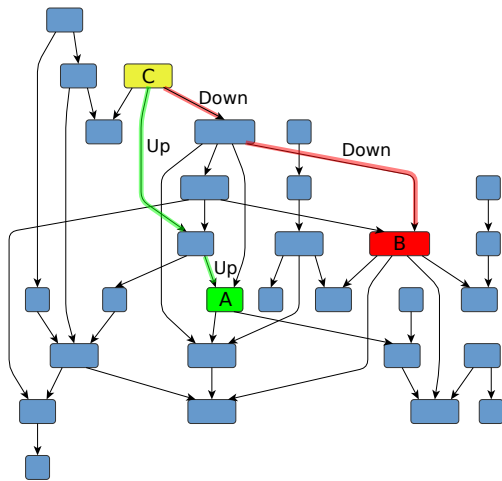
Formal Language Constrained Path Querying



Navigation through an edge-labeled graph

- **Path** specifies a **word** formed by labels of edges
- **Paths constraint** is a **language**: path is good if related word in the give language
- Constraints expressiveness is related to **formal languages classes**

Formal Language Constrained Path Querying



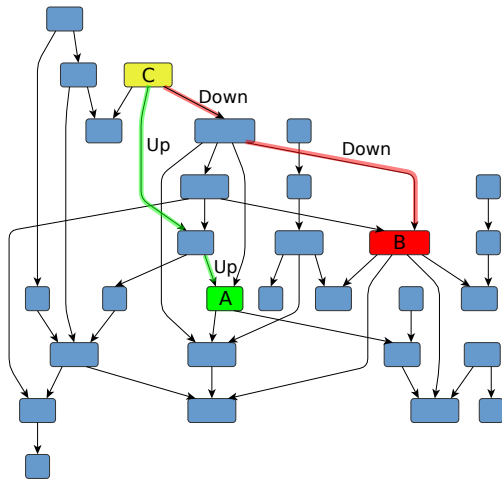
Navigation through an edge-labeled graph

- **Path** specifies a **word** formed by labels of edges
- **Paths constraint** is a **language**: path is good if related word in the give language
- Constraints expressiveness is related to **formal languages classes**

Regular path queries (RPQ)

- **Regular** languages are used as constraints
- Which nodes are reachable from **C** by arbitrary number of **Up** and **Down** edges?
- $\mathcal{L} = (\text{Up} \mid \text{Down})^*$

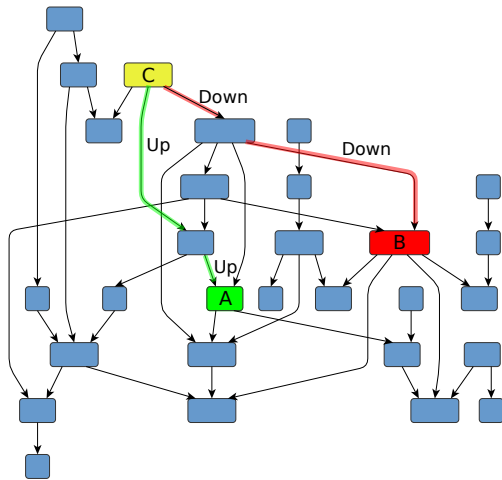
Context-Free Path Querying (CFPQ)



The constraint is a **context-free language**

- Are nodes A and B on the same level of hierarchy?
- Is there a path of form **Upⁿ Downⁿ**?
- Find all paths of form **Upⁿ Downⁿ** which start from the node A

Context-Free Path Querying (CFPQ)



The constraint is a **context-free language**

- Are nodes A and B on the same level of hierarchy?
- Is there a path of form **Upⁿ Downⁿ**?
- Find all paths of form **Upⁿ Downⁿ** which start from the node A

Applications

- Static code analysis [T. Reps, et al, 1995]
- Graph segmentation [H. Miao, et al, 2019]
- Biological data analysis [P. Sevon, et al, 2008]
- ...

The Problem

There is no support of of CFPQ in real-world graph analysis systems (graph databases)

¹Jochem Kuijpers, George Fletcher, Nikolay Yakovets, and Tobias Lindaaker. 2019. An Experimental Study of Context-Free Path Query Evaluation Methods.

²Arseniy Terekhov, Artyom Khoroshev, Rustam Azimov, and Semyon Grigorev. 2020. Context-Free Path Querying with Single-Path Semantics by Matrix Multiplication.

The Problem

There is no support of of CFPQ in real-world graph analysis systems (graph databases)

- J. Kuijpers, et al¹: existing algorithms are too slow to be practical (in the context of Neo4j)

¹Jochem Kuijpers, George Fletcher, Nikolay Yakovets, and Tobias Lindaaker. 2019. An Experimental Study of Context-Free Path Query Evaluation Methods.

²Arseniy Terekhov, Artyom Khoroshev, Rustam Azimov, and Semyon Grigorev. 2020. Context-Free Path Querying with Single-Path Semantics by Matrix Multiplication.

The Problem

There is no support of of CFPQ in real-world graph analysis systems (graph databases)

- J. Kuijpers, et al¹: existing algorithms are too slow to be practical (in the context of Neo4j)
- A. Terekhov, et al²: linear algebra based CFPQ algorithm can be performant enough

¹Jochem Kuijpers, George Fletcher, Nikolay Yakovets, and Tobias Lindaaker. 2019. An Experimental Study of Context-Free Path Query Evaluation Methods.

²Arseniy Terekhov, Artyom Khoroshev, Rustam Azimov, and Semyon Grigorev. 2020. Context-Free Path Querying with Single-Path Semantics by Matrix Multiplication.

The Problem

There is no support of of CFPQ in real-world graph analysis systems (graph databases)

- J. Kuijpers, et al¹: existing algorithms are too slow to be practical (in the context of Neo4j)
- A. Terekhov, et al²: linear algebra based CFPQ algorithm can be performant enough
- There is no full-stack support of CFPQ
 - ▶ Grammars instead of full-featured queries
 - ▶ Custom graph storage instead of real-world graph database

¹Jochem Kuijpers, George Fletcher, Nikolay Yakovets, and Tobias Lindaaker. 2019. An Experimental Study of Context-Free Path Query Evaluation Methods.

²Arseniy Terekhov, Artyom Khoroshev, Rustam Azimov, and Semyon Grigorev. 2020. Context-Free Path Querying with Single-Path Semantics by Matrix Multiplication.

- **Multiple-Source CFPQ** allows one to reduce computations in real-world cases

³Tobias Lindaaker, Path Patterns for Cypher, 2017,
<https://github.com/thobe/openCypher/blob/rpq/cip/1.accepted/CIP2017-02-06-Path-Patterns.adoc>

⁴!!!

⁵!!!

Proposed Solution

- **Multiple-Source CFPQ** allows one to reduce computations in real-world cases
- **Cypher** extended with **path patterns**³ allows one to express context-free constraints

³Tobias Lindaaker, Path Patterns for Cypher, 2017,
<https://github.com/thobe/openCypher/blob/rpq/cip/1.accepted/CIP2017-02-06-Path-Patterns.adoc>

⁴!!!

⁵!!!

Proposed Solution

- **Multiple-Source CFPQ** allows one to reduce computations in real-world cases
- **Cypher** extended with **path patterns**³ allows one to express context-free constraints
- **RedisGraph** database
 - ▶ Provides graph storage with matrix-based representation
 - ▶ Contains linear algebra based query engine (SuitSparse:GraphBLAS⁴ is used)
 - ▶ Allows one to use Cypher for querying (libcypher.parser⁵ is used)

³Tobias Lindaaker, Path Patterns for Cypher, 2017,

<https://github.com/thobe/openCypher/blob/rpq/cip/1.accepted/CIP2017-02-06-Path-Patterns.adoc>

⁴!!!

⁵!!!

Multiple-Source CFPQ

```
1: function MULTISRCCFPQ( $D = (V, E, \Sigma_V, \Sigma_E, \lambda_V, \lambda_E)$ ,  $G = (N, \Sigma, P, S)$ ,  $Src$ )
2:    $T \leftarrow \{T^A \mid A \in N, T^A[i, j] \leftarrow \text{false}, \text{ for all } i, j\}$ 
3:    $TSrc \leftarrow \{TSrc^A \mid A \in N, TSrc^A[i, j] \leftarrow \text{false}, \text{ for all } i, j\}$ 
4:   for all  $v \in Src$  do  $TSrc^S[v, v] \leftarrow \text{true}$ 
5:    $MSrc \leftarrow TSrc^S$ 
6:   for all  $A \rightarrow x \in P \mid x \in \Sigma_E$  do
7:     for all  $(v, to) \in E \mid x \in \lambda_E(v, to)$  do  $T^A[v, to] \leftarrow \text{true}$ 
8:   for all  $A \rightarrow x \in P \mid x \in \Sigma_V$  do
9:     for all  $v \in V \mid x \in \lambda_V(v)$  do  $T^A[v, v] \leftarrow \text{true}$ 
10:  while  $T$  or  $TSrc$  is changing do
11:    for all  $A \rightarrow BC \in P$  do
12:       $M \leftarrow TSrc^A * T^B$ 
13:       $T^A \leftarrow T^A + M * T^C$ 
14:       $TSrc^B \leftarrow TSrc^B + TSrc^A$ 
15:       $TSrc^C \leftarrow TSrc^C + \text{GETDST}(M)$ 
16:  return  $MSrc * T^S$ 
```

Multiple-Source CFPQ

```
1: function MULTISRCCFPQ( $(D = (V, E, \Sigma_V, \Sigma_E, \lambda_V, \lambda_E), G = (N, \Sigma, P, S), Src)$ )
2:    $T \leftarrow \{T^A \mid A \in N, T^A[i, j] \leftarrow false, \text{ for all } i, j\}$ 
3:    $TSrc \leftarrow \{TSrc^A \mid A \in N, TSrc^A[i, j] \leftarrow false, \text{ for all } i, j\}$ 
4:   for all  $v \in Src$  do  $TSrc^S[v, v] \leftarrow true$ 
5:    $MSrc \leftarrow TSrc^S$ 
6:   for all  $A \rightarrow x \in P \mid x \in \Sigma_E$  do
7:     for all  $(v, to) \in E \mid x \in \lambda_E(v, to)$  do  $T^A[v, to] \leftarrow true$ 
8:   for all  $A \rightarrow x \in P \mid x \in \Sigma_V$  do
9:     for all  $v \in V \mid x \in \lambda_V(v)$  do  $T^A[v, v] \leftarrow true$ 
10:  while  $T$  or  $TSrc$  is changing do
11:    for all  $A \rightarrow BC \in P$  do
12:       $M \leftarrow TSrc^A * T^B$ 
13:       $T^A \leftarrow T^A + M * T^C$ 
14:       $TSrc^B \leftarrow TSrc^B + TSrc^A$ 
15:       $TSrc^C \leftarrow TSrc^C + GETDST(M)$ 
16:  return  $MSrc * T^S$ 
```

Multiple-Source CFPQ

```
1: function MULTISRCCFPQ( $D = (V, E, \Sigma_V, \Sigma_E, \lambda_V, \lambda_E)$ ,  $G = (N, \Sigma, P, S)$ ,  $Src$ )
2:    $T \leftarrow \{T^A \mid A \in N, T^A[i, j] \leftarrow \text{false, for all } i, j\}$ 
3:    $TSrc \leftarrow \{TSrc^A \mid A \in N, TSrc^A[i, j] \leftarrow \text{false, for all } i, j\}$ 
4:   for all  $v \in Src$  do  $TSrc^S[v, v] \leftarrow \text{true}$ 
5:    $MSrc \leftarrow TSrc^S$ 
6:   for all  $A \rightarrow x \in P \mid x \in \Sigma_E$  do
7:     for all  $(v, to) \in E \mid x \in \lambda_E(v, to)$  do  $T^A[v, to] \leftarrow \text{true}$ 
8:   for all  $A \rightarrow x \in P \mid x \in \Sigma_V$  do
9:     for all  $v \in V \mid x \in \lambda_V(v)$  do  $T^A[v, v] \leftarrow \text{true}$ 
10:  while  $T$  or  $TSrc$  is changing do
11:    for all  $A \rightarrow BC \in P$  do
12:       $M \leftarrow TSrc^A * T^B$ 
13:       $T^A \leftarrow T^A + M * T^C$ 
14:       $TSrc^B \leftarrow TSrc^B + TSrc^A$ 
15:       $TSrc^C \leftarrow TSrc^C + \text{GETDST}(M)$ 
16:  return  $MSrc * T^S$ 
```

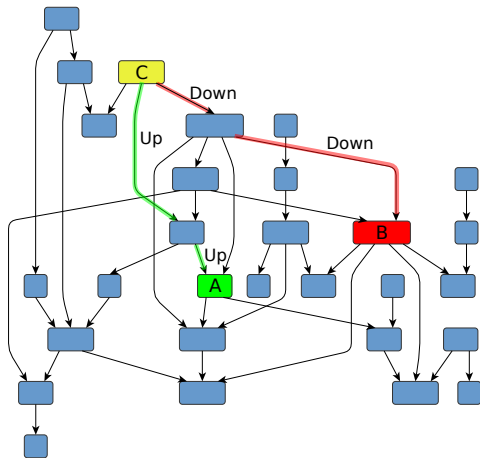
Multiple-Source CFPQ

```
1: function MULTISRCCFPQ( $D = (V, E, \Sigma_V, \Sigma_E, \lambda_V, \lambda_E)$ ,  $G = (N, \Sigma, P, S)$ ,  $Src$ )
2:    $T \leftarrow \{T^A \mid A \in N, T^A[i, j] \leftarrow \text{false, for all } i, j\}$ 
3:    $TSrc \leftarrow \{TSrc^A \mid A \in N, TSrc^A[i, j] \leftarrow \text{false, for all } i, j\}$ 
4:   for all  $v \in Src$  do  $TSrc^S[v, v] \leftarrow \text{true}$ 
5:    $MSrc \leftarrow TSrc^S$ 
6:   for all  $A \rightarrow x \in P \mid x \in \Sigma_E$  do
7:     for all  $(v, to) \in E \mid x \in \lambda_E(v, to)$  do  $T^A[v, to] \leftarrow \text{true}$ 
8:   for all  $A \rightarrow x \in P \mid x \in \Sigma_V$  do
9:     for all  $v \in V \mid x \in \lambda_V(v)$  do  $T^A[v, v] \leftarrow \text{true}$ 
10:   while  $T$  or  $TSrc$  is changing do
11:     for all  $A \rightarrow BC \in P$  do
12:        $M \leftarrow TSrc^A * T^B$ 
13:        $T^A \leftarrow T^A + M * T^C$ 
14:        $TSrc^B \leftarrow TSrc^B + TSrc^A$ 
15:        $TSrc^C \leftarrow TSrc^C + \text{GETDST}(M)$ 
16:   return  $MSrc * T^S$ 
```


Multiple-Source CFPQ

```
1: function MULTISRCCFPQ( $D = (V, E, \Sigma_V, \Sigma_E, \lambda_V, \lambda_E)$ ,  $G = (N, \Sigma, P, S)$ ,  $Src$ )
2:    $T \leftarrow \{T^A \mid A \in N, T^A[i, j] \leftarrow \text{false, for all } i, j\}$ 
3:    $TSrc \leftarrow \{TSrc^A \mid A \in N, TSrc^A[i, j] \leftarrow \text{false, for all } i, j\}$ 
4:   for all  $v \in Src$  do  $TSrc^S[v, v] \leftarrow \text{true}$ 
5:    $MSrc \leftarrow TSrc^S$ 
6:   for all  $A \rightarrow x \in P \mid x \in \Sigma_E$  do
7:     for all  $(v, to) \in E \mid x \in \lambda_E(v, to)$  do  $T^A[v, to] \leftarrow \text{true}$ 
8:   for all  $A \rightarrow x \in P \mid x \in \Sigma_V$  do
9:     for all  $v \in V \mid x \in \lambda_V(v)$  do  $T^A[v, v] \leftarrow \text{true}$ 
10:  while  $T$  or  $TSrc$  is changing do
11:    for all  $A \rightarrow BC \in P$  do
12:       $M \leftarrow TSrc^A * T^B$ 
13:       $T^A \leftarrow T^A + M * T^C$ 
14:       $TSrc^B \leftarrow TSrc^B + TSrc^A$ 
15:       $TSrc^C \leftarrow TSrc^C + \text{GETDST}(M)$ 
16:  return  $MSrc * T^S$ 
```

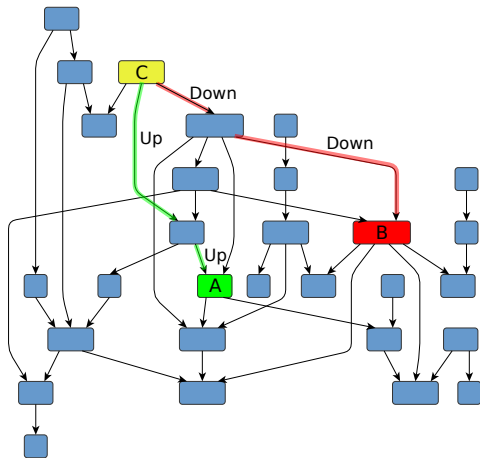
Cypher Extension



MATCH (a) - [(Down | Up)*] -> (b)

RETURN a.name, b.name

Cypher Extension



```
MATCH (a)-[(Down | Up)*]->(b)
```

```
RETURN a.name, b.name
```

Named path pattern

SameLvl $\rightarrow \overline{\text{Down}} \text{ SameLvl } \text{Down} \mid \epsilon$

```
PATH PATTERN SameLvl =  
  ()-/-<:Down [~SameLvl | ()] :Down>/->()
```

```
MATCH (a)-/-~SameLvl/->(b)
```

```
RETURN a.name, b.name
```

Implementation Details

- Linear algebra based multiple-source CFPQ is implemented as part of RedisGraph query engine
- Cypher parser is extended to support path patterns
- Path patterns are supported⁶ in RedisGreaph query execution workflow

⁶Partially. Full support is a nontrivial challenge: formal description of the extension is required

Evaluation Setup

- Ubuntu 18.04, Intel Core i7-6700 CPU, 3.4GHz, DDR4 64Gb RAM
- Graphs stored in RedisGraph with our extensions
- Queries are generated with template for given size of start set
- The union of all start sets is a V

Evaluation Setup

- Ubuntu 18.04, Intel Core i7-6700 CPU, 3.4GHz, DDR4 64Gb RAM
- Graphs stored in RedisGraph with our extensions
- Queries are generated with template for given size of start set
- The union of all start sets is a V

Graph	#V	#E	Q
core	1323	4342	g_1
pathways	6238	18 598	g_1
gohierarchy	45 007	980 218	g_1
enzyme	48 815	109 695	g_1
eclass_514en	239 111	523 727	g_1
geospecies	450 609	2 311 461	geo
go	272 770	534 311	g_1

Evaluation Setup

- Ubuntu 18.04, Intel Core i7-6700 CPU, 3.4GHz, DDR4 64Gb RAM
- Graphs stored in RedisGraph with our extensions
- Queries are generated with template for given size of start set
- The union of all start sets is a V

Graph	#V	#E	Q
core	1323	4342	g_1
pathways	6238	18 598	g_1
gohierarchy	45 007	980 218	g_1
enzyme	48 815	109 695	g_1
eclass_514en	239 111	523 727	g_1
geospecies	450 609	2 311 461	geo
go	272 770	534 311	g_1

PATH PATTERN S =

`()-/[<:SubClassOf [~S | ()] :SubClassOf] | [<:Type [~S | ()] :Type] /->()`

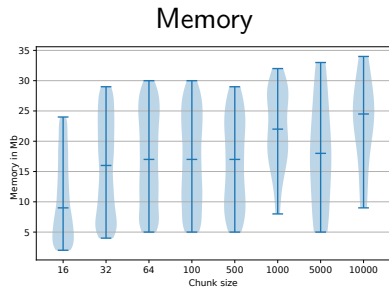
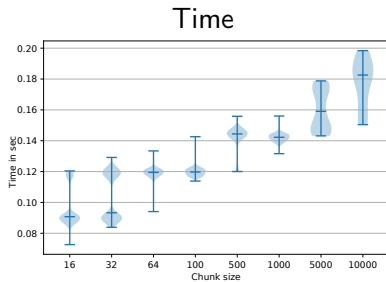
MATCH (src)-/[~S /->()

WHERE {id_from} <= src.id and src.id <= {id_to}

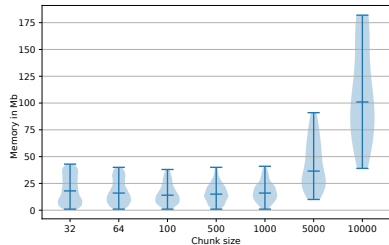
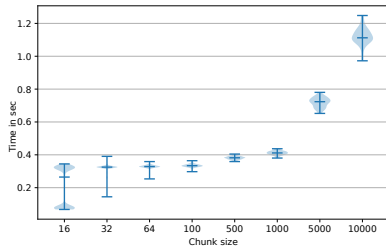
RETURN count(*)

Evaluation Results

eclass_514en
Query: g1



geospecies
Query: geo



- Full-stack support for CFPQ in real-world graph query language (Cypher) on the top of real-world graph database (RedisGraph)
 - ▶ No more ugly context-free grammars
 - ▶ No more custom graph formats and storages
- Reasonable performance of context-free path queries
 - ▶ Multiple-source scenario
 - ▶ Space-time ratio can be tuned
- Context-free path queries can be used in applications with well-established tools

- Mechanization of Cypher semantics in Coq
 - ▶ Including path patterns
 - ▶ Correctness of translation to linear algebra

⁷Egor Orachev, Ilya Epelbaum, R. Azimov and S. Grigorev. “Context-Free Path Querying by Kronecker Product.” ADBIS (2020).

- Mechanization of Cypher semantics in Coq
 - ▶ Including path patterns
 - ▶ Correctness of translation to linear algebra
- Integration of tensor-based CFPQ algorithm⁷ to RedisGraph
 - ▶ Allows one to construct paths, not only reachability facts
 - ▶ Should be modified to get multiple-source version

⁷Egor Orachev, Ilya Epelbaum, R. Azimov and S. Grigorev. “Context-Free Path Querying by Kronecker Product.” ADBIS (2020).

- Mechanization of Cypher semantics in Coq
 - ▶ Including path patterns
 - ▶ Correctness of translation to linear algebra
- Integration of tensor-based CFPQ algorithm⁷ to RedisGraph
 - ▶ Allows one to construct paths, not only reachability facts
 - ▶ Should be modified to get multiple-source version
- Detailed evaluation
 - ▶ More graphs and queries, including RPQs
 - ▶ Scalability of the solution
 - ▶ Comparison with other graph query engines

⁷Egor Orachev, Ilya Epelbaum, R. Azimov and S. Grigorev. “Context-Free Path Querying by Kronecker Product.” ADBIS (2020).

Contact Information

- Try it out (Docker image with all included):
<https://hub.docker.com/r/simpletondl/redisgraph>
- Sources of RedisGraph extended with CFPQ:
<https://github.com/YaccConstructor/RedisGraph>
- Sources of Cypher parser extended with path patterns:
<https://github.com/YaccConstructor/libcypher-parser>

Contact Information

- Try it out (Docker image with all included):
<https://hub.docker.com/r/simpletondl/redisgraph>
- Sources of RedisGraph extended with CFPQ:
<https://github.com/YaccConstructor/RedisGraph>
- Sources of Cypher parser extended with path patterns:
<https://github.com/YaccConstructor/libcypher-parser>

- Semyon Grigorev: s.v.grigoriev@spbu.ru
- Arseniy Terekhov: simpletondl@yandex.ru
- Vlada Pogozhelskaya: pogozhelskaya@gmail.com
- Vadim Abzalov: vadim.i.abzalov@gmail.com
- Timur Zinnatulin: teemychteemych@gmail.com

Contact Information

- Try it out (Docker image with all included):
<https://hub.docker.com/r/simpletondl/redisgraph>
- Sources of RedisGraph extended with CFPQ:
<https://github.com/YaccConstructor/RedisGraph>
- Sources of Cypher parser extended with path patterns:
<https://github.com/YaccConstructor/libcypher-parser>
- Semyon Grigorev: s.v.grigoriev@spbu.ru
- Arseniy Terekhov: simpletondl@yandex.ru
- Vlada Pogozhelskaya: pogozhelskaya@gmail.com
- Vadim Abzalov: vadim.i.abzalov@gmail.com
- Timur Zinnatulin: teemychteemych@gmail.com

Thanks!