

Санкт-Петербургский государственный университет

Кафедра системного программирования

Группа 17Б.11-мм

*Орачев Егор Станиславович*

# Реализация алгоритма поиска путей в графовых базах данных через тензорное произведение на GPGPU

Отчёт по преддипломной практике

Научный руководитель:  
доцент кафедры информатики, к. ф.-м. н., С. В. Григорьев

Санкт-Петербург  
2021

Saint Petersburg State University

Software Engineering

Group 17B.11-mm

*Egor Orachyov*

# Context-Free path querying by tensor product for graph databases on GPGPU

Pre-Diploma practice report

Scientific supervisor:  
C.Sc., docent Semyon Grigorev

Saint Petersburg  
2021

# Оглавление

<b>Введение</b>	<b>4</b>
<b>1. Цель и задачи</b>	<b>6</b>
<b>2. Обзор предметной области</b>	<b>7</b>
2.1. Терминология . . . . .	7
2.2. Поиск путей с ограничениями . . . . .	7
2.3. Существующие решения . . . . .	8
2.4. Поиск путей с КС ограничениями через тензорное произведение . . . . .	9
2.5. Вычисления на GPGPU . . . . .	11
2.6. Библиотеки линейной алгебры для GPGPU . . . . .	12
<b>3. Архитектура библиотеки</b>	<b>13</b>
3.1. Структура . . . . .	13
3.2. Последовательность обработки операций . . . . .	16
<b>4. Детали реализации</b>	<b>18</b>
<b>5. Алгоритм поиска путей с КС ограничениями через тензорное произведение на GPGPU</b>	<b>19</b>
<b>6. Экспериментальное исследование</b>	<b>20</b>
6.1. Постановка экспериментов . . . . .	20
6.2. Результаты . . . . .	23
<b>7. Заключение</b>	<b>24</b>
<b>Список литературы</b>	<b>26</b>

# Введение

Все чаще современные системы аналитики и рекомендаций строятся на основе анализа данных, структурированных с использованием *графовой модели*. В данной модели основные сущности представляются вершинами графа, а отношения между сущностями — ориентированными ребрами с различными метками. Подобная модель позволяет относительно легко и практически в явном виде моделировать сложные иерархические структуры, которые не так просто представить, например, в классической *реляционной модели*. В качестве основных областей применения графовой модели можно выделить следующие: графовые базы данных [4], анализ RDF данных [6], биоинформатика [22] и статистический анализ кода [13].

Поскольку графовая модель используется для моделирования отношений между объектами, при решении прикладных задач возникает необходимость в выявлении неявных взаимоотношений между объектами. Для этого формируются запросы в специализированных программных средствах для управления графовыми базами данных. В качестве запроса можно использовать некоторый *шаблон* на путь в графе, который будет связывать объекты, т.е. выражать взаимосвязь между ними. В качестве такого шаблона можно использовать формальные грамматики, например, регулярные или контекстно-свободные (КС). Используя вычислительно более выразительные грамматики, можно формировать более сложные запросы и выявлять нестандартные и скрытые ранее взаимоотношения между объектами. Например, *same-generation queries* [1] могут быть выражены КС грамматиками, в отличие от регулярных.

Результатом запроса может быть множество пар объектов, между которыми существует путь в графе, удовлетворяющий заданным ограничениям. Также может возвращаться один экземпляр такого пути для каждой пары объектов или итератор всех путей, что зависит от семантики запроса. Поскольку один и тот же запрос может иметь разную семантику, требуются различные программные и алгоритмические сред-

ства для его выполнения.

Запросы с регулярными ограничениями изучены достаточно хорошо, языковая и программная поддержка выполнения подобных запросов присутствует в некоторых в современных графовых базах данных. Однако, полноценная поддержка запросов с КС ограничениями до сих пор не представлена. Существуют алгоритмы [3,6,7,15,17] для вычисления запросов с КС ограничениями, но потребуется еще время, прежде чем появиться полноценная высокопроизводительная реализация одного из алгоритмов, способная обрабатывать реальные графовые данные.

Работы Никиты Мишина и др. [11] и Арсения Терехова и др. [8] показывают, что реализация алгоритма Рустама Азимова [3], основанного на операциях линейной алгебры, с использованием GPGPU для выполнения наиболее вычислительно сложных частей алгоритма, дает *существенный* прирост в производительности.

Недавно представленный алгоритм [7] для вычисления запросов с КС ограничениями также полагается на операции линейной алгебры: тензорное произведение, матричное умножение и сложение в булевом полукольце. Данный алгоритм в сравнении с [8] позволяет выполнять запросы для всех ранее упомянутых семантик и потенциально поддерживает бóльшие по размеру КС запросы.

Для его реализации на GPGPU требуются высокопроизводительные библиотеки операций линейной алгебры. Подобные инструменты для работы со стандартными типами данных, такими как *float*, *double*, *int* и *long*, уже представлены. Однако библиотека, которая бы работала с разреженными данными и имела специализацию указанных ранее операций для булевых значений, еще не разработана.

Поэтому важной задачей является не только реализация перспективного алгоритма [7] на GPGPU, но и разработка библиотеки примитивов булевой алгебры, которая позволит реализовать этот и подобные алгоритмы на данной вычислительной платформе.

# 1. Цель и задачи

Целью данной работы является реализация алгоритма поиска путей в графовых базах данных через тензорное произведение на GPGPU. Для ее выполнения были поставлены следующие задачи.

- Разработка архитектуры библиотеки примитивов разреженной линейной булевой алгебры для вычислений на GPGPU.
- Реализация библиотеки в соответствии с разработанной архитектурой.
- Реализация алгоритма поиска путей с КС ограничениями через тензорное произведение с использованием разработанной библиотеки.
- Экспериментальное исследование полученных результатов.

## 2. Обзор предметной области

### 2.1. Терминология

*Ориентированный граф с метками*  $\mathcal{G} = \langle V, E, L \rangle$  это тройка объектов, где  $V$  конечное непустое множество вершин графа,  $E \subseteq V \times L \times V$  конечное множество ребер графа,  $L$  конечное множество меток графа. Здесь и далее будем считать, что вершины графа индексируются целыми числами, т.е.  $V = \{0 \dots |V| - 1\}$ .

Граф  $\mathcal{G} = \langle V, E, L \rangle$  можно представить в виде матрицы смежности  $M$  размером  $|V| \times |V|$ , где  $M[i, j] = \{l \mid (i, l, j) \in E\}$ . Используя булеву матричную декомпозицию, можно представить матрицу смежности в виде набора матриц  $\mathcal{M} = \{M^l \mid l \in L, M^l[i, j] = 1 \iff l \in M[i, j]\}$ .

Путь  $\pi$  в графе  $\mathcal{G} = \langle V, E, L \rangle$  это последовательность ребер  $e_0, e_1, \dots, e_{n-1}$ , где  $e_i = (v_i, l_i, u_i) \in E$  и для любых  $e_i, e_{i+1} : u_i = v_{i+1}$ . Путь между вершинами  $v$  и  $u$  будем обозначать как  $v\pi u$ . Слово, которое формирует путь  $\pi = (v_0, l_0, v_1), \dots, (v_{n-1}, l_{n-1}, v_n)$  будем обозначать как  $\omega(\pi) = l_0 \dots l_{n-1}$ , что является конкатенацией меток вдоль этого пути  $\pi$ .

*Контекстно-свободная (КС) грамматика*  $G = \langle \Sigma, N, P, S \rangle$  это четверка объектов, где  $\Sigma$  конечное множество терминалов или алфавит,  $N$  конечное множество нетерминалов,  $P$  конечное множество правил вывода вида  $A \rightarrow \gamma, \gamma \in (N \cup \Sigma)^*$ ,  $S \in N$  стартовый нетерминал. Вывод слова  $w$  в грамматике из нетерминала  $S$  применением одного или нескольких правил вывода обозначается как  $S \rightarrow_G^* w$ .

Язык  $L$  над конечным алфавитом символов  $\Sigma$  — множество слов, составленных из символов этого алфавита, т.е.  $L \subseteq \{w \mid w \in \Sigma^*\}$ . Язык, задаваемый грамматикой  $G$ , обозначим как  $L(G) = \{w \mid S \rightarrow_G^* w\}$ .

### 2.2. Поиск путей с ограничениями

При вычислении запроса на поиск путей в графе  $\mathcal{G} = \langle V, E, L \rangle$  в качестве ограничения выступает некоторый язык  $L$ , которому должны удовлетворять результирующие пути.

Поиск путей в графе с семантикой **достижимости** — это поиск всех

таких пар вершин  $(v, u)$ , что между ними существует путь  $v\pi u$  такой, что  $\omega(\pi) \in L$ . Результат запроса обозначается как  $R = \{(v, u) \mid \exists v\pi u : \omega(\pi) \in L\}$ .

Поиск путей в графе с семантикой **всех путей** — это поиск всех таких путей  $v\pi u$ , что  $\omega(\pi) \in L$ . Результат запроса обозначается как  $\Pi = \{v\pi u \mid v\pi u : \omega(\pi) \in L\}$ .

Необходимо отметить, что множество  $\Pi$  может быть бесконечным, поэтому в качестве результата запроса предполагается не всё множество в явном виде, а некоторый *итератор*, который позволит последовательно извлекать все пути.

Семантика **одного пути** является ослабленной формулировкой семантики всех путей, так как для получения результата достаточно найти всего один путь вида  $v\pi u : \omega(\pi) \in L$  для каждой пары  $(v, u) \in R$ .

Поскольку язык  $L$  может быть бесконечным, при составлении запросов используют не множество  $L$  в явном виде, а некоторое правило формирования слов этого языка. В качестве таких правил и выступают регулярные выражения или КС грамматики. При именовании запросов отталкиваются от типа правил, поэтому запросы именуются как регулярные или КС соответственно.

## 2.3. Существующие решения

Впервые проблема выполнения запросов с контекстно-свободными ограничениями была сформулирована в 1990 году в работе Михалиса Яннакакиса [27]. С того времени были представлены многие работы, в которых так или иначе предлагалось решение данной проблемы. Однако в недавнем исследовании Йохем Куиджперс и др. [12] на основе сравнения нескольких алгоритмов [3, 15, 23] для выполнения запросов с контекстно-свободными ограничениями заключили, что существующие алгоритмы неприменимы для анализа реальных данных в силу того, что обработка таких данных занимает значительное время. Стоит отметить, что алгоритмы, используемые в статье, были реализованы на языке программирования *Java* и исполнялись в среде *JVM* в однопо-



точном режиме, что не является сколь-угодно производительным решением.

Это подтверждают результаты работы Арсения Терехова и др. [8], в которой с использованием программных и аппаратных средств Nvidia Cuda был реализован алгоритм Рустама Азимова [3]. В данном алгоритме задача поиска путей с КС ограничениями была сведена к операциям линейной алгебры, что позволило использовать высокопроизводительные библиотеки для выполнения данных операций на GPGPU.

Алгоритм Рустама Азимова [8] способен выполнять запросы только в семантике одного пути. Поскольку в качестве формализма для представления грамматики КС запроса используется *ослабленная нормальная форма Хомского (ОНФХ)* [16], увеличение числа правил в исходной грамматике запроса может приводить к существенному разрастанию ОНФХ, что негативно влияет на время работы алгоритма.

## 2.4. Поиск путей с КС ограничениями через тензорное произведение

Представленный в работе Егора Орачева и др. [7] алгоритм для выполнения КС запросов использует операции линейной булевой алгебры: произведение Кронекера (частный случай тензорного произведения), матричное умножение и сложение. Данный алгоритм позволяет выполнять запросы в семантике достижимости и всех путей, а также он подлежит распараллеливанию на многоядерных системах, что делает его потенциально применимым для анализа реальных данных. Кроме этого, данный алгоритм использует в качестве формализма для представления запроса *рекурсивный автомат (РА)* [2], что потенциально может решить проблему разрастания исходной грамматики запроса.

Идея алгоритма состоит в *пересечении* РА и графа с использованием некоторой модификации классического алгоритма пересечения рекурсивных автоматов [16]. Пересечение выполняется с использованием произведения Кронекера, а множество рекурсивных вызовов учитывается с помощью транзитивного замыкания, что также выражается с

использованием матричных операций умножения и поэлементного сложения. Данный процесс итеративный, и он выполняется до тех пор, пока результат не достигнет фиксированной точки.

В листинге 1 представлен псевдокод алгоритма. Необходимо отметить, что алгоритм использует булеву матричную декомпозицию в строках **3** – **4** для представления матрицы переходов  $R$  и матрицы смежности графа, а также использует матричное умножение, сложение и произведение Кронекера в строках **14** – **16**.

Данный алгоритм является относительно простым в реализации, так как всю сложность выполнения он перекладывает на операции линейной алгебры, которые должны быть реализованы в сторонних высокопроизводительных библиотеках.

---

### Listing 1 Поиск путей через произведение Кронекера

---

```

1: function KRONECKERPRODUCTBASEDCFPQ( $G, \mathcal{G}$ )
2:    $R \leftarrow$  Рекурсивный автомат для грамматики  $G$ 
3:    $\mathcal{M}_1 \leftarrow$  Матрица переходов  $R$  в булевой форме
4:    $\mathcal{M}_2 \leftarrow$  Матрица смежности  $\mathcal{G}$  в булевой форме
5:    $C_3 \leftarrow$  Пустая матрица
6:   for  $s \in \{0, \dots, \dim(\mathcal{M}_1) - 1\}$  do
7:     for  $S \in \text{getNonterminals}(R, s, s)$  do
8:       for  $i \in \{0, \dots, \dim(\mathcal{M}_2) - 1\}$  do
9:          $M_2^S[i, i] \leftarrow \{1\}$ 
10:      end for
11:    end for
12:  end for
13:  while Матрица смежности  $\mathcal{M}_2$  изменяется do
14:     $\mathcal{M}_3 \leftarrow \mathcal{M}_1 \otimes \mathcal{M}_2$  ▷ Вычисление произведения Кронекера
15:     $M'_3 \leftarrow \bigvee_{M_3^a \in \mathcal{M}_3} M_3^a$  ▷ Слияние матриц в одну булеву матрицу достижимости
16:     $C_3 \leftarrow \text{transitiveClosure}(M'_3)$  ▷ Транзитивное замыкание для учета рекурсивных вызовов
17:     $n \times n \leftarrow \dim(\mathcal{M}_3)$ 
18:    for  $(i, j) \mid C[i, j] \neq 0$  do
19:       $s, f \leftarrow \text{getStates}(C_3, i, j)$ 
20:       $x, y \leftarrow \text{getCoordinates}(C_3, i, j)$ 
21:      for  $S \in \text{getNonterminals}(R, s, f)$  do
22:         $M_2^S[x, y] \leftarrow \{1\}$ 
23:      end for
24:    end for
25:  end while
26:  return  $\mathcal{M}_2, C_3$ 
27: end function

```

---

## 2.5. Вычисления на GPGPU

*GPGPU* (от англ. General-purpose computing on graphics processing units) — техника использования графического процессора видеокарты компьютера для осуществления неспециализированных вычислений, которые обычно проводит центральный процессор. Данная техника позволяет получить значительной прирост производительности, когда необходимо обрабатывать большие массивы данных с фиксированным набором команд по принципу *SIMD*.

Исторически видеокарты в первую очередь использовались как графические ускорители для создания высококачественной трехмерной графики в режиме реального времени. Позже стало ясно, что мощность графического процессора можно использовать не только для графических вычислений. Так появились программируемые вычислительные блоки (англ. compute shaders), которые позволяют выполнять на видеокарте неграфические вычисления.

На данный момент существует несколько промышленных стандартов для создания программ, использующих графический процессор, одними из которых являются Vulkan [26], OpenGL [25], DirectX [10] как API для графических и неспециализированных вычислительных задач, а также OpenCL [21], Nvidia Cuda [18] как API для неспециализированных вычислений.

В качестве GPGPU в этой работе используется Nvidia Cuda. В то время как OpenCL создавался как кросс-платформенный стандарт для программирования вычислений, Cuda API специфично только для видеокарт производства компании Nvidia, однако оно имеет более широкий набор инструментов как для написания, так и для отладки программ, а также собственный компилятор NVCC, который позволяет осуществлять кросс-компиляцию кода на языке Cuda, и прозрачно использовать его вместе с кодом на языке C/C++. Кроме этого, в данной работе используются результаты исследования Арсения Терехова и др. [8], в котором также использовалось Cuda API.

## 2.6. Библиотеки линейной алгебры для GPGPU

Для эффективной реализации алгоритмов [3, 7] требуются высокопроизводительные библиотеки операций линейной алгебры. Реальные графовые данные насчитывают порядка  $10^5 - 10^9$  вершин и являются сильно разреженными, т.е. количество ребер в графе сравнимо с количеством вершин, поэтому плотные матрицы не подходят для представления такого типа данных.

Библиотеки *cuSPARSE* [19] и *CUSP* [9] для платформы Nvidia и *clSPARSE* [5] для платформы OpenCL предоставляют функциональность для работы с разреженными данными, однако они имеют фокусируются на обработку численных данных и специализируются только на стандартных типах, таких как *float*, *double*, *int* и *long*. Для реализации алгоритмов [3, 7] требуются операции над разреженными булевыми матрицами, поэтому требуется специализация вышеуказанных библиотек для работы с булевыми примитивами. С одной стороны, библиотека *cuSPARSE* имеет закрытый исходный код, что делает невозможным ее модификацию, с другой стороны, библиотеки *CUSP* и *clSPARSE* имеют открытый исходный код и свободную лицензию, однако используемые ими алгоритмы умножения разреженных матриц *достаточно* требовательны к ресурсам памяти [8], что делает их неприменимым для обработки данных большого размера.

В работе Арсения Терехова и др. [8] была предпринята попытка самостоятельно реализовать алгоритм умножения разреженных матриц *Nsparse*, предложенный в работе Юсуке Нагасака и др. [20], и специализировать его для булевых значений. Данный алгоритм эксплуатирует возможности видеокарт Nvidia и за счет бóльшего количества шагов обработки позволяет снизить количество расходуемой видеопамяти. Эксперименты показали, что подобный подход позволяет не только снизить в разы количество расходуемой видеопамяти, но и снизить общее время работы алгоритма Рустама Азимова [3] по сравнению с его реализацией на *CUSP*.

## 3. Архитектура библиотеки

### 3.1. Структура

Архитектура библиотеки представлена на рис. 1. Структура библиотеки и ее конечная функциональность в основном продиктованы следующими высокоуровневыми требованиями, которые продиктованы как конечными вычислительными задачами на GPGPU, так и наличием существующей инфраструктуры для осуществления экспериментов [28].

- Поддержка вычислений на Cuda-девайсе.
- Поддержка вычислений на CPU.
- C-совместимое API для работы с библиотекой.
- Python-пакет для работы с примитивами и операциями библиотеки в управляемой высокоуровневой среде языка Python.
- Поддержка логирования, функций для отладки и прототипирования конечных пользовательских алгоритмов.

### Core

Класс **Library** поддерживает глобальное состояние библиотеки, осуществляет конфигурацию и инициализацию, выбор конкретного вычислительного бэкенда, первичную валидацию вызовов функций и входных данных пользователя, а также осуществляет хранение всех созданных объектов.

Класс **Matrix** является проху-классом, который осуществляет доступ к операциям конкретного вычислительного бэкенда, которой был выбран пользователем на этапе инициализации всей библиотеки. Данный подход позволяет не только динамически выбирать платформу вычислений, но также позволяет осуществлять дополнительную обработку

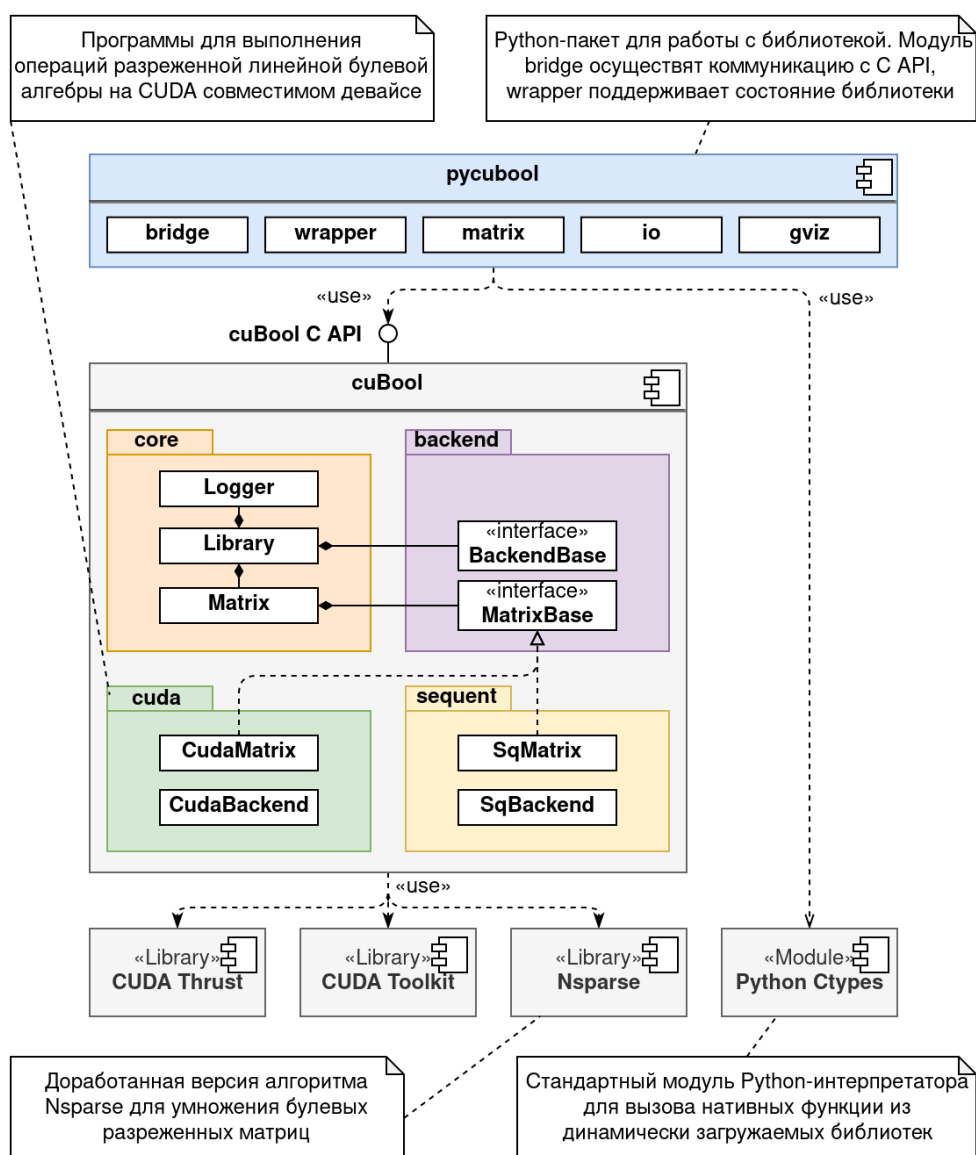


Рис. 1: Архитектура разработанной библиотеки

ошибок, а также поддерживать дополнительные операции над матрицами.

Класс **Logger** осуществляет логгирование в выбранный пользователем текстовый файл в процессе использования функций библиотеки, а также позволяет профилировать операций и также сохранять время их выполнения в текстовом виде.

## Backend

Интерфейс **MatrixBase** предоставляет набор основных функций и операций, которые каждый вычислительный бэкенд должен реализовывать для того, чтобы эти матрицы можно было использовать в **Core** непосредственно для вычислений.

Интерфейс **BackendBase** описывает базовый контракт, который должен предоставлять вычислительный бэкенд. Данный интерфейс включает в себя функции для создания и удаления матриц, специфичных для этого бэкенда, а также функции для корректной инициализации и завершения работы данного бэкенда.

## Cuda

Класс **CudaMatrix** реализует интерфейс **MatrixBase** и предоставляет операции, для работы с матрицами и осуществления вычислений на Cuda-девайсе. **CudaMatrix** поддерживает данные матрицы (не нулевые элементы) в видео-памяти, и использует Cuda *kernels* для осуществления вычислений на GPU. *Kernels* хранятся в виде функций вместе с исходным кодом модуля.

Данный вычислительный бэкенд выбирается по умолчанию, если в компьютере пользователя имеется Cuda-девайс. Однако пользователь может в таком случае все равно может выбрать **Sequent** вычисления.

## Sequent

Предоставляет реализацию класса матрицы и операций над ней для вычислений на CPU. Все вычисления осуществляются последовательно, в однопоточном режиме, не требуют дополнительных библиотек или компонентов.

Данный вычислительный бэкенд используется по умолчанию на устройствах без Cuda-девайса. Данный подход позволит использовать библиотеку всем пользователям без исключения. Также данный подход может быть удобен для прототипирования алгоритмов на локальном компью-

тере, чтобы позже запустить вычисления на высокопроизводительном сервере с поддержкой Cuda.

## Py cubool

Python-пакет предоставляет доступ к примитивам и операциям библиотеки в языковой среде Python. Модуль **matrix** предоставляет доступ к классу матрицы и основным операциям, доступным в C API. Модуль **bridge** осуществляет коммуникацию с библиотекой через механизмы вызова нативных методов. Модуль **wrapper** поддерживает глобальное состояние библиотеки в во время работы Python-интерпретатора. Модули *io* и **gviz** предоставляют доступ к операциям ввода/вывода данных, позволяют загружать или сохранять матрицы в виде текстовых на диск, а также экспортировать набор матриц в виде графа в формате GraphViz, что может быть полезно для отладки алгоритмов.

### 3.2. Последовательность обработки операций

На рис. 2 представлена последовательность обработки вычислительной операции над матрицей (матрицами) на Cuda-девайсе.

Пользовательский Python-код инициирует выполнение операции над матрицей или несколькими матрицами внутри инфраструктуры Python. Этот вызов обрабатывает **pycubool**-пакет, который осуществляет первичную базовую валидацию аргументом, осуществляет их запаковку и передачу в нативную функцию **cuBool C API**. На стороне реализации данного интерфейса, полученные аргументы приводятся к требуемому типу и передаются далее в модуль **Core**, который поддерживает состояние библиотеки, осуществляет валидацию аргументов, а также определяет допустимость выполнения операции. Далее вызов передается непосредственно вычислительному бэкенду **Cuda**, который осуществляет подготовку и непосредственный запуск вычислений на стороне **Nvidia GPU**.

Когда вычисление завершается, **Cuda**-бэкенд обновляет состояние матриц в соответствии с полученными результатами. Модуль **Core** осуще-



ствяет финальное логирование операции, а также сохраняет временные показатели выполнения вычислений в файл (опционально), и возвращает в качестве результата выполнения статус операции, либо возможное исключение, которое могло возникнуть на этапе обработки запроса вычислений. **cuBool C API** осуществляет финальную обработку исключения (если таковое возникло), и возвращает вызывающему числовой идентификатор статуса операции.

В результате выполнения операции **pucubool** уведомляет пользователя о потенциально возникших ошибках и возвращает управление из вызываемой функции. Обновленное состояние библиотеки находится в **Core**, а состояние матриц после выполнения операций хранится на стороне **Cuda**-бэкенда.

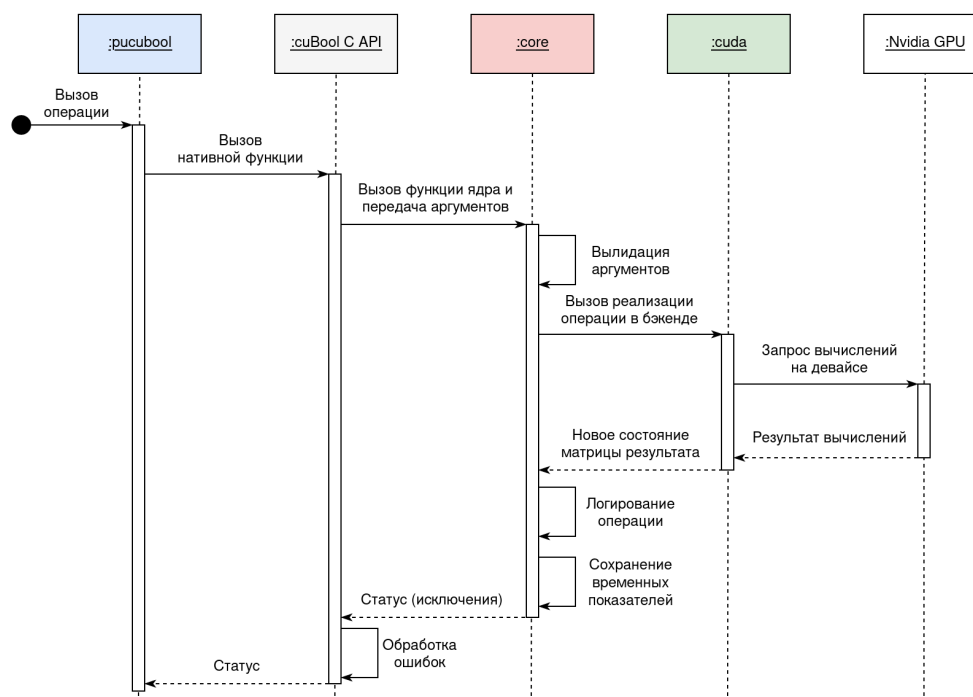


Рис. 2: Последовательность обработки вызова вычислительной матричной операции на Nvidia GPU

## 4. Детали реализации

## **5. Алгоритм поиска путей с КС ограничениями через тензорное произведение на GPGPU**

## 6. Экспериментальное исследование

### 6.1. Постановка экспериментов

Для экспериментов использовалась рабочая станция с процессором Intel Core i7-6790, тактовой частотой 3.40GHz, RAM DDR4 с объемом памяти 64Gb, видеокартой GeForce GTX 1070 с 8Gb VRAM, ОС под управлением Ubuntu 20.04.

#### Исследовательские вопросы

Для того, чтобы структурировать исследование, были сформулированы следующие вопросы.

- B1:** Какова производительность отдельных операций реализованной библиотеки примитивов разреженной линейной булевой алгебры на GPGPU по сравнению с существующими аналогами?
- B2:** Какова производительность реализованного алгоритма поиска путей через тензорное произведение на GPGPU по сравнению с существующими аналогами, также полагающимися на примитивы линейной алгебры?

B1 направлен на определение эффективности отдельных матричных операций в реализованной библиотеке. В качестве таких операций выступают *матричное умножение* и *матричное сложение* в булевом полукольце, как наиболее распространенные и критически важные операции в прикладных алгоритмах. Для сравнения производительности в этих операциях предлагается использовать популярные существующие библиотеки разреженной линейной алгебры для платформ Nvidia Cuda, OpenCL и CPU, что позволит составить наиболее полную картину. В качестве таких библиотек были выбраны CUSP и cuSPARSE для Nvidia Cuda, clSPARSE для OpenCL, и SuiteSparse для CPU. CUSP предоставляет реализацию операций, основанную на шаблонах для параметризации используемого типа данных, однако библиотека не делает каких-либо дополнительных оптимизаций конкретно

для булевых значений. `cuSPARSE` и `clSPRASE` предоставляют операции только для основных типов данных с плавающей запятой. Однако данное ограничение можно обойти, если интерпретировать ненулевые значения как *true*. Библиотека `SuiteSprase` является эталонной реализацией `GraphBLAS API` и имеет встроенное булево полукольцо для вычислений.

`B2` направлен на определение эффективности реализованного алгоритма и его сравнение с алгоритмом Рустама Азимова [8], который также полагается на операции линейной булевой алгебры. Данный алгоритм также реализован с использованием разработанного в данной работе Python-пакета, что делает сравнение корректным.

## Описание данных

Для замеров производительности отдельных операций реализованной библиотеки были выбраны 10 различных квадратных матриц из известной коллекции университета Флориды [24] для проверки эффективности алгоритмов, реализующих операции над разреженными матрицами. Информация о матрицах представлена в таблице 1. Для обозначения числа ненулевых элементов используется аббревиатура *Nnz* (англ. number of non-zero elements). В таблице приведено официальное название матрицы, количество строк (соответствует числу столбцов), а также количество ненулевых элементов в данной матрице и в производных от нее, полученных умножением матрицы самой на себя, что обозначается как степень  $M^2$ , и поэлементным сложением данной матрицы с собой также возведенной в степень, что обозначается как  $M + M^2$ . Вычисление данных артефактов имитирует шаг транзитивного замыкания, которое является элементом многих алгоритмов на графах, связанных с анализом достижимости. Эффективное вычисление этого шага во многом определяет производительность конечных алгоритмов.

Для замеров производительности алгоритмов поиска путей с КС ограничениями используется коллекция графовых данных лаборатории языковых инструментов JetBrains Research [14], которая использовалась в ряде работ [3, 7, 8, 11] для подобных экспериментов. Данная

Таблица 1: Разреженные матричные данные

№	Матрица $M$	Кол-во Строк	Nnz $M$	Nnz $M^2$	Nnz $M + M^2$
0	wing	62,032	243,088	714,200	917,178
1	luxembourg_osm	114,599	239,332	393,261	632,185
2	amazon0312	400,727	3,200,400	14,390,544	14,968,909
3	amazon-2008	735,323	5,158,388	25,366,745	26,402,678
4	web-Google	916,428	5,105,039	29,710,164	30,811,855
5	roadNet-PA	1,090,920	3,083,796	7,238,920	9,931,528
6	roadNet-TX	1,393,383	3,843,320	8,903,897	12,264,987
7	belgium_osm	1,441,295	3,099,940	5,323,073	8,408,599
8	roadNet-CA	1,971,281	5,533,214	12,908,450	17,743,342
9	netherlands_osm	2,216,688	4,882,476	8,755,758	13,626,132

коллекция содержит RDF данные, онтологии, графы программ для анализа указателей, а также ряд сгенерированных графов для анализа особых случаев.

Данные, необходимые для замеров, предварительно загружаются в RAM или VRAM в формате, требуемом для тестируемого инструмента. Время, необходимое на чтение данных с диска, их конвертацию, а также подготовку начального состояния входных матриц исключено из замеров.

## Метрики

Для ответа на поставленные исследовательские вопросы в качестве метрик производительности используется время, требуемое для выполнения операции, а также пиковое количество потребляемой видеопамяти на GPU в момент вычисления. Показатели времени усреднены по 10 запускам. Предварительно совершался не учитывающийся в замерах запуск, чтобы проинициализировать начальное состояние тестируемых библиотек. Показатели потребления видеопамяти получены с помощью инструмента *nvidia-smi*, который с точностью до 1 миллисекунды позволяет отслеживать количество потребляемой памяти процессом ОС на стороне видеокарты.

## 6.2. Результаты

1) B1: Какова производительность отдельных операций реализованной библиотеки примитивов разреженной линейной булевой алгебры на GPGPU по сравнению с существующими аналогами?

2) B2: Какова производительность реализованного алгоритма поиска путей через тензорное произведение на GPGPU по сравнению с существующими аналогами, также полагающимися на примитивы линейной алгебры?

## 7. Заключение

В рамках выполнения данной работы были получены следующие результаты:

- Спроектирована библиотека примитивов линейной булевой алгебры для работы с разреженными данными на GPGPU. Данная библиотека экспортирует C-совместимый интерфейс, имеет поддержку различных вычислительных модулей, а также предоставляет модуль для работы конечного пользователя с примитивами библиотеки в высокоуровневой среде вычислений с управляемыми ресурсами.
- Реализована библиотека `cuBool` в соответствии с разработанной архитектурой. Ядро библиотеки написано на языке C++, а математические операции, выполняющиеся на GPGPU, реализованы на языке CUDA C/C++. Библиотека предоставляет модуль CPU вычислений для компьютеров без Cuda девайсов. Также создан Python-пакет `ruscubool`, который позволяет использовать функциональность библиотеки в среде Python. Данный пакет доступен для скачивания через пакетный менеджер PyPI.
- Реализован алгоритм поиска путей с КС ограничениями через тензорное произведение с использованием Python пакета разработанной библиотеки. Данный алгоритм использует операции матричного умножения, сложения и произведение Кронекера в булевом полукольце, а также различные операции для манипуляций над значениями матриц. На вход алгоритм получает представление графа и КС грамматики в виде набора матриц, а на выходе — возвращает матрицу смежности графа достижимости, а также индекс, который позволяет восстанавливать все пути в графе, в соответствии с входной грамматикой.
- Выполнено экспериментальное исследование реализованной библиотеки с использованием синтетических и реальных данных из



коллекции Разреженных Матриц Университета Флориды. Замеры производительности операции матричного умножения показывают ускорение до 5 раз в сравнении с существующими аналогами при меньшем потреблении памяти. Матричное умножение сравнимо по времени с существующими аналогами, однако потребляет меньше памяти. Также выполнено экспериментальное исследование производительности реализованного алгоритма и его сравнение с существующими алгоритмами для выполнения КС запрос. В качестве данных для замеров использовалась коллекция RDF и синтетических данных Лаборатории Языковых Инструментов JetBrains Research. Замеры показали что (этого мы пока не знаем).

На основе результатов, полученных в данном исследовании, была написана статья, принятая на конференцию GrAPL 2021<sup>1</sup>.

Библиотека cuBool и Python-пакет для работы с данной библиотекой доступны для скачивания через следующие онлайн ресурсы: <https://github.com/JetBrains-Research/cuBool> и <https://test.pypi.org/project/pycubool/>.

---

<sup>1</sup>GrAPL 2021: Workshop on Graphs, Architectures, Programming, and Learning. Дата обращения: 1.04.2021. Сайт конференции: <https://hpc.pnl.gov/grapl/>.

## Список литературы

- [1] Abiteboul Serge, Hull Richard, Vianu Victor. Foundations of Databases. — 1995. — 01. — ISBN: [0-201-53771-0](#).
- [2] Analysis of Recursive State Machines / Rajeev Alur, Michael Benedikt, Kousha Etessami et al. // [ACM Trans. Program. Lang. Syst.](#) — 2005. — Jul. — Vol. 27, no. 4. — P. 786–818. — Access mode: <https://doi.org/10.1145/1075382.1075387>.
- [3] Azimov Rustam, Grigorev Semyon. [Context-free path querying by matrix multiplication](#). — 2018. — 06. — P. 1–10.
- [4] Barceló Baeza Pablo. [Querying Graph Databases](#) // Proceedings of the 32nd ACM SIGMOD-SIGACT-SIGAI Symposium on Principles of Database Systems. — PODS '13. — New York, NY, USA : Association for Computing Machinery, 2013. — P. 175–188. — Access mode: <https://doi.org/10.1145/2463664.2465216>.
- [5] [CLSPARSE: A Vendor-Optimized Open-Source Sparse BLAS Library](#) / Joseph L. Greathouse, Kent Knox, Jakub Poła et al. // Proceedings of the 4th International Workshop on OpenCL. — IWOCCL '16. — New York, NY, USA : Association for Computing Machinery, 2016. — Access mode: <https://doi.org/10.1145/2909437.2909442>.
- [6] Context-Free Path Queries on RDF Graphs / Xiaowang Zhang, Zhiyong Feng, Xin Wang et al. // CoRR. — 2015. — Vol. abs/1506.00743. — [1506.00743](#).
- [7] [Context-Free Path Querying by Kronecker Product](#) / Egor Orachev, Ilya Epelbaum, Rustam Azimov, Semyon Grigorev. — 2020. — 08. — P. 49–59. — ISBN: [978-3-030-54831-5](#).
- [8] [Context-Free Path Querying with Single-Path Semantics by Matrix Multiplication](#) / Arseniy Terekhov, Artyom Khoroshev, Rustam Azimov, Semyon Grigorev. — 2020. — 06. — P. 1–12.

- [9] Dalton Steven, Bell Nathan, Olson Luke, Garland Michael. Cusp: Generic Parallel Algorithms for Sparse Matrix and Graph Computations. — 2014. — Version 0.5.0. Access mode: <http://cusplibrary.github.io/>.
- [10] Direct3D 12 Graphics // Microsoft Online Documents. — 2018. — Access mode: <https://docs.microsoft.com/ru-ru/windows/win32/direct3d12/direct3d-12-graphics?redirectedfrom=MSDN> (online; accessed: 08.12.2020).
- [11] [Evaluation of the Context-Free Path Querying Algorithm Based on Matrix Multiplication](#) / Nikita Mishin, Iaroslav Sokolov, Egor Spirin et al. — 2019. — 06. — P. 1–5.
- [12] [An Experimental Study of Context-Free Path Query Evaluation Methods](#) / Jochem Kuijpers, George Fletcher, Nikolay Yakovets, Tobias Lindaaaker // Proceedings of the 31st International Conference on Scientific and Statistical Database Management. — SSDBM '19. — New York, NY, USA : ACM, 2019. — P. 121–132. — Access mode: <http://doi.acm.org/10.1145/3335783.3335791>.
- [13] Fast Algorithms for Dyck-CFL-Reachability with Applications to Alias Analysis / Qirun Zhang, Michael R. Lyu, Hao Yuan, Zhendong Su // [SIGPLAN Not.](#) — 2013. — Jun. — Vol. 48, no. 6. — P. 435–446. — Access mode: <https://doi.org/10.1145/2499370.2462159>.
- [14] Graphs and grammars for Context-Free Path Querying algorithms evaluation // Github. — 2021. — Access mode: [https://github.com/JetBrains-Research/CFPQ\\_Data](https://github.com/JetBrains-Research/CFPQ_Data) (online; accessed: 11.03.2021).
- [15] Hellings Jelle. Path Results for Context-free Grammar Queries on Graphs. — 2015. — 02.
- [16] Hopcroft John E., Motwani Rajeev, Ullman Jeffrey D. Introduction to Automata Theory, Languages, and Computation (3rd Edition). — USA : Addison-Wesley Longman Publishing Co., Inc., 2006. — ISBN: [0321455363](#).

- [17] Medeiros Ciro, Musicante Martin, Costa Umberto. An Algorithm for Context-Free Path Queries over Graph Databases. — 2020. — 04.
- [18] NVIDIA. CUDA Toolkit Documentation // NVIDIA Developer Zone. — 2020. — Access mode: <https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html> (online; accessed: 01.12.2020).
- [19] NVIDIA. cuSPARSE reference guide // NVIDIA Developer Zone. — 2020. — Access mode: <https://docs.nvidia.com/cuda/cusparse/index.html> (online; accessed: 09.12.2020).
- [20] Nagasaka Yusuke, Nukada Akira, Matsuoka Satoshi. [High-Performance and Memory-Saving Sparse General Matrix-Matrix Multiplication for NVIDIA Pascal GPU](#). — 2017. — 08. — P. 101–110.
- [21] OpenCL: Open Standard for Parallel Programming of Heterogeneous Systems // Khronos website. — 2020. — Access mode: <https://www.khronos.org/opencl/> (online; accessed: 08.12.2020).
- [22] Quantifying variances in comparative RNA secondary structure prediction / James Anderson, Adám Novák, Zsuzsanna Sükösd et al. // [BMC bioinformatics](#). — 2013. — 05. — Vol. 14. — P. 149.
- [23] Santos Fred, Costa Umberto, Musicante Martin. [A Bottom-Up Algorithm for Answering Context-Free Path Queries in Graph Databases](#). — 2018. — 01. — P. 225–233. — ISBN: 978-3-319-91661-3.
- [24] T. Davis. The SuiteSparse Matrix Collection (the University of Florida Sparse Matrix Collection). — 2020. — Access mode: <https://sparse.tamu.edu> (online; accessed: 09.03.2021).
- [25] The Khronos Working Group. OpenGL 4.4 Specification // Khronos Registry. — 2014. — Access mode: <https://www.khronos.org/registry/OpenGL/specs/gl/glspec44.core.pdf> (online; accessed: 08.12.2020).

- [26] The Khronos Working Group. Vulkan 1.1 API Specification // Khronos Registry. — 2019. — Access mode: <https://www.khronos.org/registry/vulkan/specs/1.1/html/vkspec.html> (online; accessed: 08.12.2020).
- [27] Yannakakis Mihalis. [Graph-Theoretic Methods in Database Theory](#) // Proceedings of the Ninth ACM SIGACT-SIGMOD-SIGART Symposium on Principles of Database Systems. — PODS '90. — New York, NY, USA : Association for Computing Machinery, 1990. — P. 230–242. — Access mode: <https://doi.org/10.1145/298514.298576>.
- [28] A collection of CFPQ algorithms implemented in PyGraph-BLAS // Github. — 2020. — Access mode: [https://github.com/JetBrains-Research/CFPQ\\_PyAlgo](https://github.com/JetBrains-Research/CFPQ_PyAlgo) (online; accessed: 16.12.2020).