

SPLA: Portable Multi-GPU Implementation of GraphBLAS API

1st Egor Orachev
Faculty of Mathematics and Mechanics
St. Petersburg State University,
Programming Languages and Tools Lab
JetBrains Research
St. Petersburg, Russia
egor.orachev@gmail.com
!!!OCRID

2nd Gleb
!!!
!!!
St. Petersburg, Russia
!!!email
!!!OCRID

3rd Semyon Grigorev
Faculty of Mathematics and Mechanics
St. Petersburg State University,
Programming Languages and Tools Lab
JetBrains Research
St. Petersburg, Russia
s.v.grigoriev@spbu.ru
semyon.grigorev@jetbrains.com
0000-0002-7966-0698

Abstract—Scalable high-performance graph analysis is an actual nontrivial challenge and usage of sparse linear algebra operations as building blocks for graph analysis algorithms, which is a core idea of GraphBLAS standard, is a promising way to attack it. While it is known that sparse linear algebra operations can be efficiently implemented on GPGPU, full GraphBLAS implementation on GPGPU is a nontrivial task that is almost solved by GraphBLAST project. Though it is shown by GraphBLAST that utilization of GPGPUs for GraphBLAS implementation significantly improves performance, portability and scalability problems are not solved yet: GraphBLAST uses Nvidia stack and utilizes only one GPGPU. In this work we propose an SPLA library that aimed to solve these problems: it uses OpenCL to be portable and designed to utilize multiple GPGPUS. Preliminary evaluation shows that !!!!

Index Terms—graphs, algorithms, graph analysis, sparse linear algebra, GraphBLAS, GPGPU, OpenCL

I. INTRODUCTION

Scalable high-performance graph analysis is an actual challenge. There is a big number of ways to attach this challenge [1]. And one of the promising ways is To utilize general purpose graphic processing units (GPGPU-s) for graph analysis. Existing solutions demonstrate good performance. Scalable. But existing solutions provide a very low level API for graph algorithms construction. As a result, Verbosity, so on

One of promising ways to high performance graph analysis is a GraphBLAS API¹ [?] provides a linear algebra based building blocks [2] to create graph analysis algorithms. The idea of GraphBLAS is based on is well-known fact that linear algebra operations can be efficiently implemented on parallel hardware. Within it, a graph can be natively represented using matrices: adjacency matrix, incidence matrix, etc. But real data often sparse, thus underlying matrices and vectors also sparse, and classical data structures and algorithms became inefficient. This fact leads to sparse linear algebra. Generic sparse linear algebra. SuiteSparse:GraphBLAS² [?].

Identify applicable funding agency here. If none, delete this.

¹GraphBLAS home page: <https://graphblas.org/>. Access date: 07.01.2022.

²!!!

Though such well-known libraries as cuSparse show that sparse linear algebra operations can be efficiently implemented for GPGPU-s, it is not so trivial to implement GraphBLAS on GPGPU. First of all, it requires *generic* sparse linear algebra, thus it is impossible just to reuse existing libraries which almost specified for operations over floats and doubles. The second problem is specific optimizations, such as maskings fusion, which can not be natively implemented on top of existing kernels. GraphBLAS on GPGPU. GraphBLAST: [3], GBTL [?], !!! . But these solutions are not portable (because they are based on Nvidia Cuda stack). Moreover, scalability problem is not solved: all these solutions support only single-GPU, not multi-GPU computations.

To provide portable multi-GPU implementation of GraphBLAS API we developed a *SPLA*³ library. Contribution

- Design
- Implemented
- Evaluation

II. SOLUTION DESCRIPTION

A. Design Principles

Multi-GPU

Portability across multiple devices and vendors.

C interface to simplify interaction with other languages.

Tasks for better load-balancing.

Extensibility: one can easily implement a specific algorithm for primitive operation.

....

B. Architecture Overview

Fig. 1. Architecture

C. Matrices and Vectors Representation

Block formats. Pros and cons.

Currently COO only.

³SPLA home page: <https://jetbrains-research.github.io/spla/>.

D. Algebraic Operations

Implemented operations and algorithms for it.

E. Implementation Details

OpenCL,

Used libraries (Boost.Compute⁴, taskflow⁵ [4]), etc

III. EVALUATION

What and how.

A. Graph Algorithms

BFS, SSSP, ... implemented

B. Dataset Description

SuiteSparse matrix collection⁶ [5]. Table with specific graphs.

C. Evaluation Setup

Hardware description, two GPU,

Single GPU

Two GPU

D. Evaluation Results

Single GPU

TABLE I
TABLE TYPE STYLES

Table Head	Table Column Head		
	Table column subhead	Subhead	Subhead
copy	More table copy ^a		

^aSample of a Table footnote.

Two GPU

TABLE II
TABLE TYPE STYLES

Table Head	Table Column Head		
	Table column subhead	Subhead	Subhead
copy	More table copy ^a		

^aSample of a Table footnote.

Analysis and conclusion.

IV. CONCLUSION

Work in progress!

And future work

Advanced algorithms for LA operations

More formats for sparse matrices/vectors

C interface

Python package

Graphalytics

!!!

⁴<https://github.com/boostorg/compute>.

⁵Taskflow project home page: <https://taskflow.github.io/>. Access date: 07.01.2022.

⁶<https://sparse.tamu.edu/>.

REFERENCES

- [1] M. E. Coimbra, A. P. Francisco, and L. Veiga, "An analysis of the graph processing landscape," *Journal of Big Data*, vol. 8, no. 1, Apr. 2021. [Online]. Available: <https://doi.org/10.1186/s40537-021-00443-9>
- [2] J. Kepner, P. Aaltonen, D. Bader, A. Buluç, 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, "Mathematical foundations of the graphblas," in *2016 IEEE High Performance Extreme Computing Conference (HPEC)*, 2016, pp. 1–9.
- [3] C. Yang, A. Buluc, and J. D. Owens, "Graphblast: A high-performance linear algebra-based graph framework on the gpu," 2019.
- [4] T.-W. Huang, D.-L. Lin, C.-X. Lin, and Y. Lin, "Taskflow: A lightweight parallel and heterogeneous task graph computing system," *IEEE Transactions on Parallel and Distributed Systems*, vol. 33, pp. 1303–1320, 2022.
- [5] T. A. Davis and Y. Hu, "The university of florida sparse matrix collection," *ACM Trans. Math. Softw.*, vol. 38, no. 1, dec 2011. [Online]. Available: <https://doi.org/10.1145/2049662.2049663>