

# Spla: Open-Source Generalized Sparse Linear Algebra Framework with Vendor-Agnostic GPUs Accelerated Computations

Anonymous Author(s)

## ABSTRACT

Scalable high-performance graph analysis is an actual nontrivial challenge. 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 non-trivial task that is almost solved by GraphBLAST project. Though it is shown 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 a Spla library that aimed to solve these problems: it uses OpenCL to be portable and designed to utilize multiple GPGPUs. Preliminary evaluation shows that while further optimizations are required, the proposed solution demonstrates performance comparable with GraphBLAST on some tasks. Moreover, our solution on embedded GPU outperforms SuiteSparse:GrpahBLAS on the respective CPU on some graph analysis tasks.

## CCS CONCEPTS

• **Computer systems organization** → **Embedded systems**; *Redundancy*; Robotics; • **Networks** → Network reliability.

## KEYWORDS

graphs, algorithms, graph analysis, sparse linear algebra, GraphBLAS, GPGPU, OpenCL

### ACM Reference Format:

Anonymous Author(s). 2018. Spla: Open-Source Generalized Sparse Linear Algebra Framework with Vendor-Agnostic GPUs Accelerated Computations. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX)*. ACM, New York, NY, USA, 5 pages. <https://doi.org/XXXXXXX.XXXXXXX>

## 1 INTRODUCTION

Scalable high-performance graph analysis is an actual challenge. There is a big number of ways to attack this challenge [1] and the first promising idea is to utilize general-purpose graphic processing units (GPGPU). Such existing solutions, as CuSha [5] and Gunrock [6] show that utilization of GPUs can improve the performance of graph analysis, moreover it is shown that solutions

may be scaled to multi-GPU systems. But low flexibility and high complexity of API are problems of these solutions.

The second promising thing which provides a user-friendly API for high-performance graph analysis algorithms creation is a GraphBLAS API [4] which provides linear algebra based building blocks to create graph analysis algorithms. The idea of GraphBLAS is based on a well-known fact that linear algebra operations can be efficiently implemented on parallel hardware. Along with that, a graph can be natively represented using matrices: adjacency matrix, incidence matrix, etc. While reference CPU-based implementation of GraphBLAS, SuiteSparse:GraphBLAS [2], demonstrates good performance in real-world tasks, GPU-based implementation is challenging.

One of the challenges in this way is that real data are often sparse, thus underlying matrices and vectors are also sparse, and, as a result, classical dense data structures and respective algorithms are inefficient. So, it is necessary to use advanced data structures and procedures to implement sparse linear algebra, but the efficient implementation of them on GPU is hard due to the irregularity of workload and data access patterns. Though such well-known libraries as cuSPARSE show that sparse linear algebra operations can be efficiently implemented for GPGPU, 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 are almost all specified for operations over floats. The second problem is specific optimizations, such as masking fusion, which can not be natively implemented on top of existing kernels. Nevertheless, there is a number of implementations of GraphBLAS on GPGPU, such as GraphBLAST [8], GBTL [9], which show that GPGPUs utilization can improve the performance of GraphBLAS-based graph analysis solutions. But these solutions are not portable because they are based on Nvidia CUDA stack. Moreover, the scalability problem is not solved: all these solutions support only single-GPU, not multi-GPU computations.

To provide portable GPU implementation of GraphBLAS API we developed a *SPLA* library<sup>1</sup>. This library utilizes OpenCL for GPGPU computing to be portable across devices of different vendors. Moreover, it is initially designed to utilize multiple GPGPUs to be scalable. To sum up, the contribution of this work is the following.

- Design of portable GPU GraphBLAS implementation proposed. The design involves the utilization of multiple GPUS. Additionally, the proposed design is aimed to simplify library tuning and wrappers for different high-level platforms and languages creation.
- Subset of GraphBLAS API, including such operations as masking, matrix-matrix multiplication, matrix-matrix e-wise addition, is implemented. The current implementation is limited by COO and CSR matrix representation format and uses basic algorithms for some operations, but work

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).

Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2018 Association for Computing Machinery.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/XXXXXXX.XXXXXXX>

<sup>1</sup>Source code available at: <https://github.com/JetBrains-Research/spla>

in progress and more data formats will be supported and advanced algorithms will be implemented in the future.

- Preliminary evaluation on such algorithms as breadth-first search (BFS) and triangles counting (TC), and real-world graphs shows portability across different vendors and promising performance: for some problems Spla is comparable with GraphBLAST. Surprisingly, for some problems, the proposed solution on embedded Intel graphic card shows better performance than SuiteSparse:GraphBLAS on the respective CPU. At the same time, the evaluation shows that further optimization is required.

## 2 BACKGROUND OF STUDY

### 2.1 Related Work

Related work, existing solutions, systems.

### 2.2 GraphBLAS

GraphBLAS API. Discussion of drawbacks of current design and implementation.

### 2.3 GPU computations

Technologies, problems, challenges, architectures.

## 3 PROPOSED SOLUTION DESCRIPTION

Blah-blah-blah...

### 3.1 Design Principles

SPLA library is designed the way to maximize potential library performance, simplify its implementation and extensions, and to provided the end-user verbose, but effective interface allowing customization and precise control over operations execution. These ideas are captured in the following principles.

- *DAG-based expressions.* User constructs a computational expression from basic nodes and uses oriented edges to describe data dependencies between these nodes.
- *Automated hybrid-storage format.* Library uses internally specialized preprocessing to format data and automate its sharing between computational nodes.
- *Automated scheduling.* Computational work is automatically scheduled between available threads and nodes for execution. Scheduling order, dependencies, and granularity are defined from DAG expression, submitted by a user.
- *Customization of primitive types and operations.* Underlying primitives types and functions can be customized by user. The customization process does not require library re-compilation.
- *Exportable interface.* The library has a C++ interface with an automated reference-counting and with no-templates usage. It can be wrapped by C99 compatible API and exported to other languages, for example, in a form of a Python package.

### 3.2 Architecture Overview

Library general execution architecture is depicted in Fig. 1. As an input library accepts expression composed in the form of a

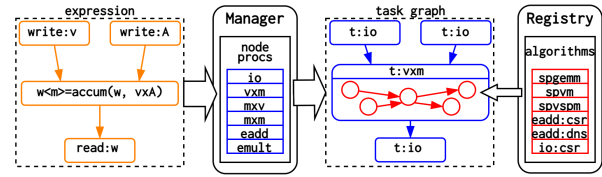


Figure 1: Library expression processing architecture.

DAG. Nodes represent fundamental operations, such as vector-matrix multiplication. Links describe dependencies between nodes. Expression execution is *asynchronous*. User can block and wait until its completion, or without blocking probe the expression until it is either *completed* or *aborted*.

Expression is transformed into a task graph. The task graph is submitted for execution to the task manager. Each task is processed by specialized *NodeProcessor*, capable of processing a particular node type. Each task, when executed, is split dynamically into a set of parallel sub-tasks. Each sub-task is processed by specialized *Algorithm*, which is capable of processing input blocks of matrices or vectors in particular storage formats with a specific set of options. *NodeProcessor* and *Algorithm* are selected at runtime from a registry using properties and arguments of the expression. Thus, it allows precise processing and optimization of edge-cases.

The granularity level of sub-tasks is defined by the structure of underlying processed primitives. The target device for execution is automatically assigned for the sub-task based on expression and node parameters. Currently, the fixed assignment is supported.

### 3.3 Containers

Library provides general *M-by-N Matrix*, *N Vector* and *Scalar* data containers. Underlying primitives types are specified by *Type* object. Primitives are stored in a hybrid storage in a form of two- or one-dimensional blocks' grid for matrices and vectors respectively. Each block is empty (not stored) or stores some data in any format. Blocks are immutable, they can be safely shared across computational units.

Currently, COO/CSR blocks for matrices and COO/Dense blocks for vectors are supported. Format choice is motivated by its simplicity and ease of implementation. Other formats, such as CSC, DCSR, ELL, etc., can be added to the library by the implementation of formats conversion or by the specialization of *Algorithm* for a specific format.

### 3.4 Algebraic Operations

Library supports all commonly used linear algebra operations, such as *mxm*, *vxm*, *eadd*, *reduce*, *transpose*. Other operations are coming soon since the library is still in development. Interface of operations is designed *similar* to GraphBLAS. It supports *masking*, *accum* of the result, *add* and *mult* user-functions specification, and *descriptor* object for additional operation tweaking.

### 3.5 Differences with GraphBLAS C API standard

To be clear...

## 4 IMPLEMENTATION DETAILS

### 4.1 Core

### 4.2 Linear Algebra Operations

### 4.3 Graph Algorithms

## 5 EVALUATION

For performance analysis of the proposed solution, we evaluated a few most common graph algorithms using real-world sparse matrix data. As a baseline for comparison we chose LAGraph [7] in connection with SuiteSparse [2] as a CPU tool, Gunrock [6] and GraphBLAST [8] as a Nvidia GPU tools. Also, we tested algorithms on several devices with distinct OpenCL vendors in order to validate the portability of the proposed solution. In general, these evaluation intentions are summarized in the following research questions.

- RQ1** What is the performance of the proposed solution relative to existing tools for both CPU and GPU analysis?
- RQ2** What is the portability of the proposed solution with respect to various device vendors and OpenCL runtimes?
- RQ3** What is the performance of the proposed solution in the CPU vs. integrated GPU comparison mode?

### 5.1 Evaluation Setup

**For evaluation of RQ1**, we use a PC with Ubuntu 20.04 installed, which has 3.40Hz Intel Core i7-6700 4-core CPU, DDR4 64Gb RAM, Intel HD Graphics 530 integrated GPU, and Nvidia GeForce GTX 1070 dedicated GPU with 8Gb on-board VRAM. **For evaluation of RQ2**, we use a PC with Ubuntu 22.04 installed, which has 4.70Hz AMD Ryzen 9 7900x 12-core CPU, DDR4 128 GB RAM, AMD GFX1036 integrated GPU, and either Intel Apc A770 flux dedicated GPU with 8GB on-board VRAM or AMD Radeon Vega Frontier Edition dedicated GPU with 16GB on-board VRAM. **For evaluation of RQ3**, the first PC with Intel CPU and embedded GPU and the second PC with AMD CPU and embedded GPU are used.

Programs using CUDA were compiled with GCC v8.4.0 and Nvidia NVCC v10.1. Release mode and maximum optimization level were enabled for all tested programs. Data loading time, preparation, format transformations, and host-device initial communications are excluded from time measurements. All tests are averaged across 10 runs. The deviation of measurements does not exceed the threshold of 10 percent. Additional warm-up run for each test execution is excluded from measurements.

### 5.2 Graph Algorithms

For preliminary study *breadth-first search* (BFS), *single-source shortest paths* (SSSP), *page rank* (PR) and *triangles counting* (TC) algorithms were chosen. Implementation of those algorithms is used from official examples packages of tested libraries with default parameters. Compared tools are allowed to make any optimizations as long as the result remains correct. The first vertex in the graph is chosen as a source for BFS and SSSP traversal for all tested tools and devices.

**Table 1: Dataset description.**

Graph	Vertices	Edges	Avg	Sd	Max
coAuthorsCit	227.3K	1.6M	7.2	10.6	1.4K
coPapersDBLP	540.5K	30.5M	56.4	66.2	3.3K
amazon2008	735.3K	7.0M	9.6	7.6	1.1K
hollywood2009	1.1M	112.8M	98.9	271.9	11.5K
comOrkut	3.1M	234.4M	76.3	154.8	33.3K
citPatents	3.8M	33.0M	8.8	10.5	793.0
socLiveJournal	4.8M	85.7M	17.7	52.0	20.3K
indochina2004	7.4M	302.0M	40.7	329.6	256.4K
belgiumosm	1.4M	3.1M	2.2	0.5	10.0
roadNetCA	2.0M	5.5M	2.8	1.0	12.0
rggn222s0	4.2M	60.7M	14.5	3.8	36.0
rggn223s0	8.4M	127.0M	15.1	3.9	40.0
roadcentral	14.1M	33.9M	2.4	0.9	8.0

### 5.3 Dataset

Thirteen matrices with graph data were selected from the Sparse Matrix Collection at University of Florida [3]. Information about graphs is summarized in Table 1. All datasets are converted to undirected graphs. Self-loops and duplicated edges are removed. For SSSP weights are initialized using pseudo-random generator with uniform  $[0, 1]$  distribution of floating-point values. Average, sd and max metrics relate to out degree property of the vertices. Graphs are divided roughly into two groups: relatively dense and relatively sparse. Graphs are sorted in the order of increasing number of vertices within each group.

### 5.4 Results Summary

Table 2 presents results of the evaluation and compares the performance of Spla against other Nvidia GPU tools and uses as a baseline LaGraph CPU tool. Table 3 presents result of the portability analysis of the proposed solution. It depicts runtime timings of the proposed solution on discrete GPUs of distinct vendors. Cell left empty with *none* if tested tool failed to analyze graph due to *out of memory* exception.

*RQ1. What is the performance of the proposed solution relative to existing tools for both CPU and GPU analysis?*

*RQ2. What is the portability of the proposed solution with respect to various device vendors and OpenCL runtimes?*

*RQ3. What is the performance of the proposed solution in the CPU vs. embedded GPU comparison mode?*

## 6 CONCLUSION

We presented a generalized sparse linear algebra framework with vendor-agnostic GPUs accelerated computations.

## REFERENCES

- [1] Miguel E. Coimbra, Alexandre P. Francisco, and Luís Veiga. 2021. An analysis of the graph processing landscape. *Journal of Big Data* 8, 1 (April 2021). <https://doi.org/10.1007/s13748-021-00100-0>

**Table 2: Performance comparison of the proposed solution. Time in milliseconds (lower is better).**

Dataset	GB	GR	LG	SP
BFS				
coAuthorsCit	5.0	1.9	6.3	6.9
coPapersDBLP	19.9	4.5	18.0	11.5
amazon2008	8.3	3.3	20.4	8.1
hollywood2009	64.3	20.3	23.4	20.3
belgiumosm	200.6	84.4	138.0	181.2
roadNetCA	116.3	32.4	168.2	101.7
comOrkut	none	205.0	40.6	53.2
citPatents	30.6	41.3	115.9	35.1
rggn222s0	367.3	95.9	1228.1	415.3
socLiveJournal	63.1	61.0	75.5	57.1
indochina2004	none	33.3	224.6	328.7
rggn223s0	615.3	146.2	2790.0	754.9
roadcentral	1383.4	243.8	1951.0	710.2
SSSP				
coAuthorsCit	14.7	2.1	38.9	10.3
coPapersDBLP	118.6	5.6	92.2	25.7
amazon2008	43.4	4.0	90.0	21.7
hollywood2009	404.3	24.6	227.7	57.5
belgiumosm	650.2	81.1	1359.8	240.9
roadNetCA	509.7	32.4	1149.3	147.9
comOrkut	none	219.0	806.5	241.0
citPatents	226.9	49.8	468.5	129.3
rggn222s0	21737.8	101.9	4808.8	865.4
socLiveJournal	346.4	69.2	518.0	189.5
indochina2004	none	40.8	821.9	596.6
rggn223s0	59015.7	161.1	11149.9	1654.8
roadcentral	13724.8	267.0	25703.4	1094.3
PR				
coAuthorsCit	1.6	10.0	24.3	3.2
coPapersDBLP	17.6	120.2	297.6	6.1
amazon2008	5.2	40.6	89.8	5.5
hollywood2009	62.9	559.5	1111.2	32.4
belgiumosm	4.4	22.9	167.6	9.4
roadNetCA	6.6	37.7	225.8	19.6
comOrkut	none	2333.6	5239.0	103.3
citPatents	27.0	686.1	1487.0	38.3
rggn222s0	45.2	320.0	563.5	26.6
socLiveJournal	none	445.9	2122.5	112.0
rggn223s0	none	662.7	1155.6	103.4
roadcentral	none	408.8	2899.9	172.0
TC				
coAuthorsCit	2.3	2.0	17.3	3.0
coPapersDBLP	105.2	5.3	520.8	128.4
amazon2008	11.2	3.9	73.9	10.8
roadNetCA	6.5	32.4	46.0	7.7
comOrkut	1776.9	218.0	23103.8	2522.0
citPatents	65.5	49.7	675.0	54.5
socLiveJournal	504.3	69.2	3886.7	437.8
rggn222s0	73.2	101.3	484.5	77.7
rggn223s0	151.4	158.9	1040.1	204.2
roadcentral	42.6	259.3	425.3	52.7

GraphBLAST (GB), Gunrock (GR), LaGraph (LG), Spla (SP).

**Table 3: Portability of the proposed solution. Time in milliseconds (lower is better).**

Dataset	Intel Arc	AMD Vega	Nvidia Gfx
BFS			
coAuthorsCit	12.8	8.3	6.9
coPapersDBLP	10.8	14.9	11.5
amazon2008	12.3	12.6	8.1
hollywood2009	15.3	26.7	20.3
belgiumosm	627.5	292.4	181.2
roadNetCA	265.5	259.8	101.7
comOrkut	33.2	63.6	53.2
citPatents	21.0	30.3	35.1
rggn222s0	825.3	1259.7	415.3
socLiveJournal	43.0	85.8	57.1
indochina2004	220.6	573.4	328.7
rggn223s0	1245.5	2519.6	754.9
roadcentral	1864.9	1680.8	710.2
SSSP			
coAuthorsCit	18.3	10.4	10.3
coPapersDBLP	22.9	27.7	25.7
amazon2008	23.4	22.2	21.7
hollywood2009	44.6	56.2	57.5
belgiumosm	1085.9	454.8	240.9
roadNetCA	447.3	422.5	147.9
comOrkut	79.7	111.5	241.0
citPatents	49.8	78.4	129.3
rggn222s0	1378.8	924.3	865.4
socLiveJournal	82.7	120.7	189.5
indochina2004	366.2	519.0	596.6
rggn223s0	1880.2	1201.4	1654.8
roadcentral	3176.3	2848.8	1094.3
PR			
coAuthorsCit	3.9	1.0	3.2
coPapersDBLP	5.7	6.1	6.1
amazon2008	25.2	4.0	5.5
hollywood2009	22.6	32.4	32.4
belgiumosm	10.2	7.1	9.4
roadNetCA	10.8	15.7	19.6
comOrkut	31.9	46.6	103.3
citPatents	12.3	21.3	38.3
rggn222s0	13.4	22.4	26.6
socLiveJournal	210.0	64.2	112.0
rggn223s0	38.6	57.2	103.4
roadcentral	57.9	89.6	172.0
TC			
coAuthorsCit	4.6	2.2	3.0
coPapersDBLP	57.6	106.2	128.4
amazon2008	6.9	8.5	10.8
roadNetCA	5.4	5.4	7.7
comOrkut	1533.5	3267.6	2522.0
citPatents	25.9	39.8	54.5
socLiveJournal	280.6	420.3	437.8
rggn222s0	21.0	57.8	77.7
rggn223s0	56.7	123.2	204.2
roadcentral	14.5	34.6	52.7

Distinct devices. Performance in not for comparison.

- 465 //doi.org/10.1186/s40537-021-00443-9
- 466 [2] Timothy A. Davis. 2019. Algorithm 1000: SuiteSparse:GraphBLAS: Graph Algo- 523
- 467 rithms in the Language of Sparse Linear Algebra. *ACM Trans. Math. Softw.* 45, 4, 524
- 468 Article 44 (Dec. 2019), 25 pages. <https://doi.org/10.1145/3322125> 525
- 469 [3] Timothy A. Davis and Yifan Hu. 2011. The University of Florida Sparse Matrix 526
- 470 Collection. *ACM Trans. Math. Softw.* 38, 1, Article 1 (dec 2011), 25 pages. <https://doi.org/10.1145/2049662.2049663> 527
- 471 [4] Jeremy Kepner, Peter Aaltonen, David Bader, Aydin Buluç, Franz Franchetti, 528
- 472 John Gilbert, Dylan Hutchison, Manoj Kumar, Andrew Lumsdaine, Henning 529
- 473 Meyerhenke, Scott McMillan, Carl Yang, John D. Owens, Marcin Zalewski, Tim- 530
- 474 othy Mattson, and Jose Moreira. 2016. Mathematical foundations of the Graph- 531
- 475 BLAS. In *2016 IEEE High Performance Extreme Computing Conference (HPEC)*. 1–9, 532
- 476 <https://doi.org/10.1109/HPEC.2016.7761646> 533
- 477 [5] Farzad Khorasani, Keval Vora, Rajiv Gupta, and Laxmi N. Bhuyan. 2014. CuSha: 534
- 478 Vertex-Centric Graph Processing on GPUs. In *Proceedings of the 23rd International 535*
- 479 *Symposium on High-Performance Parallel and Distributed Computing* (Vancouver, 536
- 480 BC, Canada) (HPDC '14). Association for Computing Machinery, New York, NY, 537
- 481 USA, 239–252. <https://doi.org/10.1145/2600212.2600227> 538
- 482 [6] Yuechao Pan, Yangzihao Wang, Yuduo Wu, Carl Yang, and John D. Owens. 2017. 539
- 483 Multi-GPU Graph Analytics. In *2017 IEEE International Parallel and Distributed 540*
- 484 *Processing Symposium (IPDPS)*. 479–490. <https://doi.org/10.1109/IPDPS.2017.117> 541
- 485 [7] Gábor Szárnyas, David A. Bader, Timothy A. Davis, James Kitchen, Timothy G. 542
- 486 Mattson, Scott McMillan, and Erik Welch. 2021. LAGraph: Linear Algebra, Network 543
- 487 Analysis Libraries, and the Study of Graph Algorithms. *arXiv:2104.01661 [cs.MS]* 544
- 488 [8] Carl Yang, Aydin Buluc, and John D. Owens. 2019. GraphBLAST: A 545
- 489 High-Performance Linear Algebra-based Graph Framework on the GPU. 546
- 490 *arXiv:1908.01407 [cs.DC]* 547
- 491 [9] Peter Zhang, Marcin Zalewski, Andrew Lumsdaine, Samantha Misurda, and Scott 548
- 492 McMillan. 2016. GBTL-CUDA: Graph Algorithms and Primitives for GPUs. In 549
- 493 *2016 IEEE International Parallel and Distributed Processing Symposium Workshops 550*
- 494 (IPDPSW). 912–920. <https://doi.org/10.1109/IPDPSW.2016.185> 551
- 495 Received 20 February 2007; revised 12 March 2009; accepted 5 June 2009 552
- 496 553
- 497 554
- 498 555
- 499 556
- 500 557
- 501 558
- 502 559
- 503 560
- 504 561
- 505 562
- 506 563
- 507 564
- 508 565
- 509 566
- 510 567
- 511 568
- 512 569
- 513 570
- 514 571
- 515 572
- 516 573
- 517 574
- 518 575
- 519 576
- 520 577
- 521 578
- 522 579
- 580