

# SPbLA: The Library of GPGPU-powered Sparse Boolean Linear Algebra Operations

Egor Orachev<sup>1, 3</sup>, Maria Karpenko<sup>2</sup>, Pavel Alimov<sup>1</sup>, and Semyon Grigorev<sup>1, 3</sup>

<sup>1</sup> Saint Petersburg State University <sup>2</sup> ITMO University <sup>3</sup> JetBrains Research

DOI: [10.21105/joss.03743](https://doi.org/10.21105/joss.03743)

## Software

- [Review](#) ↗
- [Repository](#) ↗
- [Archive](#) ↗

Editor: [Nikoleta Glynatsi](#) ↗

## Reviewers:

- [@abb58](#)
- [@mlxd](#)

Submitted: 07 September 2021

Published: 21 September 2021

## License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

## Summary

SPbLA is a sparse Boolean linear algebra primitives and operations for GPGPU computations. It comes as stand-alone self-sufficient library with C API for high-performance computing with multiple backends for Nvidia Cuda, OpenCL and CPU-only platforms. The library has PyPI `pyspb1a` package ([Orachev et al., 2021](#)) for work within Python runtime. The primary library primitive is a sparse matrix of Boolean values. The library provides the most popular operations for matrix manipulation, such as construction from values, transpose, sub-matrix extraction, matrix-to-vector reduce, matrix-matrix element-wise addition, multiplication and Kronecker product.

## Statement of need

Sparse matrices are widely applicable in data analysis and GraphBLAS API provides a set of unified building linear algebra based blocks for reducing data analysis algorithms to sparse linear algebra operations. While GPGPU utilization for high-performance linear algebra is common, the high complexity of GPGPU programming makes the implementation of the complete set of sparse operations on GPGPU challenging. Thus, it is worth addressing this problem by focusing on a basic but still important case — sparse Boolean algebra.

The primary goal of the SPbLA is to provide a base for the implementation, testing and profiling high-performance algorithms for solving data analysis problems, such as RDF analysis ([X. Zhang et al., 2015](#)), RNA secondary structure analysis ([Anderson et al., 2013](#)), static code analysis (such as points-to or alias analysis) ([Q. Zhang et al., 2013](#)) and evaluation of regular and CFL-reachability queries ([Azimov & Grigorev, 2018](#); [Orachev et al., 2020](#)).

Thus, we can offload different language-constrained path querying related problems, and other problems that can be reduced to manipulation of Boolean matrices, to GPGPU uniformly.

Moreover, real data analysis leads to huge matrix processing that can not be efficiently handled on a single GPGPU. The creation of the library which supports multi-GPU and out-of-VRAM computations helps to create an efficient solution for a wide range of applied problems. The creation of such a solution is an open problem while ad-hoc solutions exist in specific areas. We propose an SPbLA as a base for such a solution.

Also, we hope, that the library is a small step to move the implementation of the fully-featured sparse linear algebra as specified in GraphBLAS forward multi-GPU computations.

## Related tools

GraphBLAS API (Kepner et al., 2016) is one of the first standards, that formalize the mathematical building blocks in the form of the programming interface for implementing algorithms in the language of the linear algebra. SuiteSparse (T. A. Davis, 2019) is a reference implementations of the GraphBLAS API for CPU computation. It is mature and fully featured library with number of bindings for other programming languages, such pygraphblas (Pygraphblas, 2021) for Python programming.

GPGPU's utilization for data analysis and for linear algebra operations is a promising way to high-performance data analysis because GPGPU is much more powerful in parallel data processing. However, GPGPU programming is still challenging. Best to our knowledge, there is no complete GraphBLAS API implementation for GPGPU computations, except GraphBLAST (Yang et al., 2019), which is currently in the active development. Some work is also done to move SuiteSparse forward GPGPU computations.

However, the sparsity of data introduces issues with load balancing, irregular data access, thus sparsity complicates the implementation of high-performance algorithms for sparse linear algebra on GPGPU. There is number of open-source libraries, which implement independently different sparse formats and operations. Thus, there is no single solid sparse linear algebra framework. Libraries such as cuSPARSE (Sparse Matrix Library in Cuda, n.d.), bhSPARSE (Liu & Vinter, 2015), c1SPARSE (Greathouse et al., 2016) and CUSP (Dalton et al., 2014) have limited type and operators customization features with major focus on numerical types only.

## Performance

We evaluate the applicability of the proposed library for some real-world matrix data. The experiment itself is designed as a computational tasks, which arises as stand-alone or intermediate step in the solving of practical problems. Results of the evaluation compared to CPU GraphBLAS implementation SuiteSparse and existing GPU sparse linear algebra libraries. The comparison is not entirely fair, since there are still no Boolean linear algebra libraries for GPU computations.

Machine for performance evaluation has the following configuration: PC with OS Ubuntu 20.04 installed, Intel Core i7-6700 3.4Hz CPU, 64Gb DDR4 RAM, GeForce GTX 1070 GPU with 8Gb VRAM.

For evaluation, we selected a number of square real-world matrices, widely applicable for sparse matrices benchmarks, from the Sparse Matrix Collection at University of Florida (T. Davis, n.d.). Information about matrices summarized bellow. Table contains matrix name, number of rows in the matrix (the same as number of columns), number of non-zero elements (Nnz) in the matrix, average and maximum nnz in row, nnz in the result matrix.

Matrix name	# Rows	Nnz M	Nnz/row	Max Nnz/row	Nnz $M^2$
amazon0312	400,727	3,200,440	7.9	10	14,390,544
amazon-2008	735,323	5,158,388	7.0	10	25,366,745
web-Google	916,428	5,105,039	5.5	456	29,710,164
roadNet-PA	1,090,920	3,083,796	2.8	9	7,238,920
roadNet-TX	1,393,383	3,843,320	2.7	12	8,903,897
roadNet-CA	1,971,281	5,533,214	2.8	12	12,908,450
netherlands_osm	2,216,688	4,882,476	2.2	7	8,755,758

Experiment is intended to measure the performance of matrix-matrix multiplication as  $M \times M$ .

Results of the evaluation presented in Figure 1 and Figure 2. SPbLA library shows the best performance among competitors for both OpenCL and Nvidia Cuda backends. CUSP and cuSPARSE show good performance as well. However, they have significant memory consumption in some cases, what can be a critical limitation in practical analysis tasks. SuiteSparse library on CPU has acceptable performance characteristics, and it is still a good alternative for CPU-only computations.

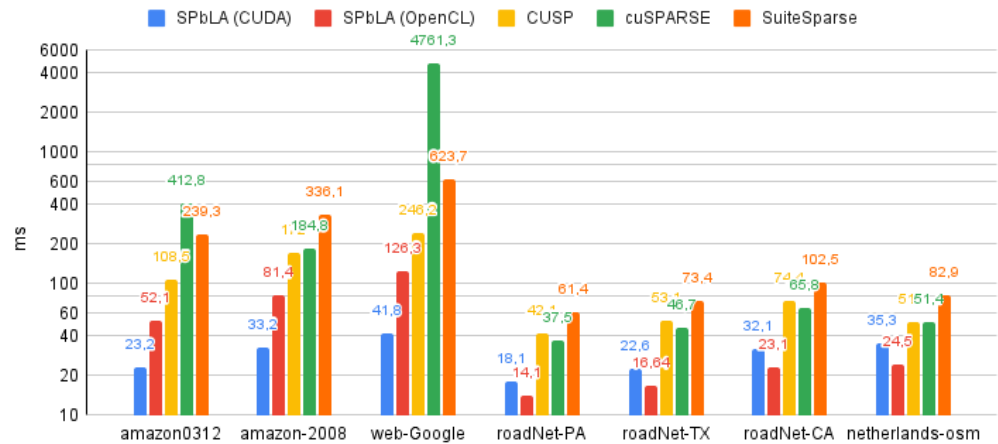


Figure 1: Matrix-matrix multiplication time consumption. Time in milliseconds.

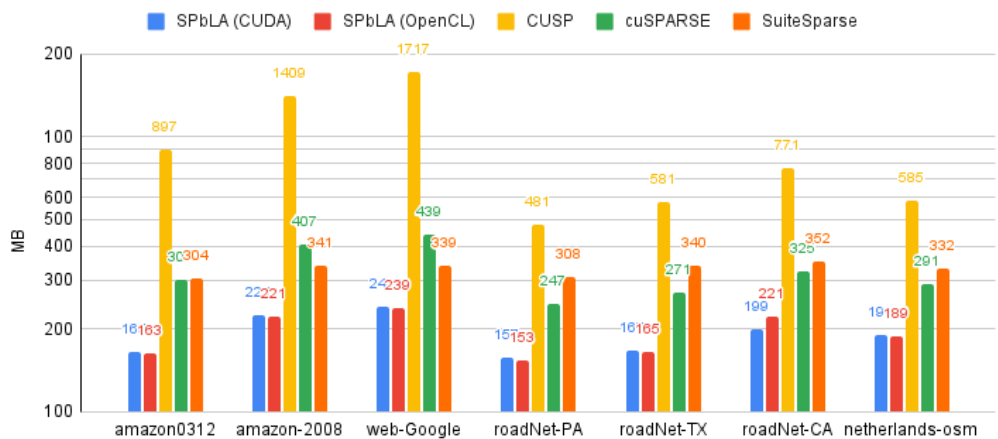


Figure 2: Matrix-matrix multiplication memory consumption. Memory in megabytes.

## Future research

First direction of the future research is experiments with virtual memory, so-called Shared Virtual Memory (SVM) in OpenCL and Unified Memory (Managed memory) in Nvidia Cuda. Our results show, that processed data can easily exceed available VRAM. Thus, it is worth to study if implemented algorithms can process more data without significant modifications in matrices structure or workload balance.

Moreover, it is necessary to further extend library to work in multi-GPU environment. This step introduces a number of issues, such as memory management among computational units as well as proper workload dispatch and granularity of parallel tasks. Potential solution is to

87 use a hybrid sparse matrix format, such as quadtree, utilize virtual memory and expose more  
88 control over expressions evaluations to the user in order to support matrix and expression level  
89 granularity among computational units.

90 Finally, we plan to generalize computational kernels and primitives in order to support arbitrary  
91 types and operations, defined by user. This step will allow defining custom elements and  
92 functions, which will be executed on GPU similarly as it is done for predefined Boolean values.

## 93 Acknowledgements

94 Work on the SPbLA project was supported by a grant from JetBrains Research.

## 95 References

- 96 Anderson, J., Novák, A., Sükösd, Z., Golden, M., Arunapuram, P., Edvardsson, I., & Hein, J.  
97 (2013). Quantifying variances in comparative RNA secondary structure prediction. *BMC*  
98 *Bioinformatics*, 14, 149. <https://doi.org/10.1186/1471-2105-14-149>
- 99 Azimov, R., & Grigorev, S. (2018). *Context-free path querying by matrix multiplication*. 1–10.  
100 <https://doi.org/10.1145/3210259.3210264>
- 101 Dalton, S., Bell, N., Olson, L., & Garland, M. (2014). *Cusp: Generic parallel algorithms for*  
102 *sparse matrix and graph computations*. <http://cusplibrary.github.io/>
- 103 Davis, T. (n.d.). *SuiteSparse matrix collection (the university of florida sparse matrix collec-*  
104 *tion)*. <https://sparse.tamu.edu/>
- 105 Davis, T. A. (2019). Algorithm 1000: SuiteSparse:GraphBLAS: Graph algorithms in the  
106 language of sparse linear algebra. *ACM Trans. Math. Softw.*, 45(4). [https://doi.org/10.](https://doi.org/10.1145/3322125)  
107 [1145/3322125](https://doi.org/10.1145/3322125)
- 108 Greathouse, J. L., Knox, K., Poła, J., Varaganti, K., & Daga, M. (2016). CISPARE: A  
109 vendor-optimized open-source sparse BLAS library. *Proceedings of the 4th International*  
110 *Workshop on OpenCL*. <https://doi.org/10.1145/2909437.2909442>
- 111 Kepner, J., Aaltonen, P., Bader, D., Buluc, A., Franchetti, F., Gilbert, J., Hutchison, D.,  
112 Kumar, M., Lumsdaine, A., Meyerhenke, H., McMillan, S., Yang, C., Owens, J. D.,  
113 Zalewski, M., Mattson, T., & Moreira, J. (2016). Mathematical foundations of the  
114 GraphBLAS. *2016 IEEE High Performance Extreme Computing Conference (HPEC)*, 1–9.  
115 <https://doi.org/10.1109/HPEC.2016.7761646>
- 116 Liu, W., & Vinter, B. (2015). A framework for general sparse matrix-matrix multiplication  
117 on GPUs and heterogeneous processors. *J. Parallel Distrib. Comput.*, 85(C), 47–61.  
118 <https://doi.org/10.1016/j.jpdc.2015.06.010>
- 119 Orachev, E., Epelbaum, I., Azimov, R., & Grigorev, S. (2020). *Context-free path querying by*  
120 *kronecker product* (pp. 49–59). [https://doi.org/10.1007/978-3-030-54832-2\\_6](https://doi.org/10.1007/978-3-030-54832-2_6)
- 121 Orachev, E., Karpenko, M., Alimov, P., & Grigorev, S. (2021). *SPbLA: Sparse boolean linear*  
122 *algebra for CPU, cuda and OpenCL computations*. <https://pypi.org/project/pyspbld/>
- 123 *Pygraphblas: A python wrapper around the GraphBLAS API*. (2021). Github. [https://github.](https://github.com/Gragegon/pygraphblas)  
124 [com/Gragegon/pygraphblas](https://github.com/Gragegon/pygraphblas)
- 125 *Sparse matrix library in cuda*. (n.d.). <https://docs.nvidia.com/cuda/cusparse/>
- 126 Yang, C., Buluç, A., & Owens, J. D. (2019). GraphBLAST: A high-performance linear algebra-  
127 based graph framework on the GPU. *arXiv Preprint*.

- 128 Zhang, Q., Lyu, M. R., Yuan, H., & Su, Z. (2013). Fast algorithms for dyck-CFL-reachability  
129 with applications to alias analysis. *SIGPLAN Not.*, 48(6), 435–446. [https://doi.org/10.](https://doi.org/10.1145/2499370.2462159)  
130 [1145/2499370.2462159](https://doi.org/10.1145/2499370.2462159)
- 131 Zhang, X., Feng, Z., Wang, X., Rao, G., & Wu, W. (2015). Context-free path queries on  
132 RDF graphs. *CoRR*, abs/1506.00743. [https://doi.org/10.1007/978-3-319-46523-4\\_38](https://doi.org/10.1007/978-3-319-46523-4_38)

DRAFT