

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

frameworks as Futhark¹ [7], Accelerate² [8], AnyDSL³ [9].

In this work we discuss a way to implement GraphBLAS API which combines high-performance computations on GPGPU and the power of high-level programming languages in both application development and possible code optimizations. Our solution is based on metaprogramming techniques: we propose to generate code for GPGPU from a high-level programming language. Namely, we plan to generate OpenCL C from a subset of F# programming language. To translate F# to OpenCL C we use a Brahma.FSharp⁴ which is based on F# quotations metaprogramming techniques⁵. Usage of F# simplifies both implementation of GraphBLAS API, making features of functional programming available, and its utilization in application development with high-level programming language on .NET platform. Moreover, as far as F# is a functional-first programming language, it should make it possible to use advanced optimization techniques and power of type system. Choice of OpenCL C as a target language is motivated by its portability: it is possible to run OpenCL C code on multi-thread CPU, on different GPGPUs (not only Nvidia), and even on FPGA [10], [11]. The utilization of FPGAs may open a way to hardware acceleration of sparse linear algebra and, as a result, of many solutions in different areas such as graph analysis, computational biology, machine learning.

This work in progress, so only tiny not optimized prototype is implemented, but our preliminary evaluation shows that !!!

II. DESIGN PRINCIPLES

In this work we are focused on making development process easier and safer by using !!! Automate optimization. Accurate type-level encoding of domain: monoids, semirings.

Monoids and semirings are closed under operations. Thus, in contrast with GraphBLAS API, all operations in semirings and monoids have the following type: $t \rightarrow t \rightarrow t$ (instead of $t_1 \rightarrow t_2 \rightarrow t_3$ as proposed in the GraphBLAS specification). It makes our definition less flexible, but allows one to generalize some operations, such as closure of relation. We realize, that in some cases such restrictive constraints are not required. Namely, definition of matrix multiplication does not requires a semiring, it just requires two operations \oplus and \otimes with following types: $\otimes : t_1 \rightarrow t_2 \rightarrow t_3$, $\oplus : t_3 \rightarrow t_3 \rightarrow t_3$. But formally, a set with such operations is not a semiring. We think that such case should be investigated separately from semirings, because additional guarantees provided by semirings may be used for code simplification and optimization. For example,

¹Futhark is a purely functional statically typed programming language for GPGPU. Project web page: <https://futhark-lang.org/>. Access date: 12.01.2021.

²Accelerate: GPGPU programming with Haskell. Project web page: <https://www.acceleratehs.org/>. Access date: 12.01.2021.

³AnyDSL is a partial evaluation framework for parallel programming. Project web page: <https://anydsl.github.io/>. Access date: 12.01.2021.

⁴Brahma.FSharp project on GitHub: <https://github.com/YaccConstructor/Brahma.FSharp>. Access date: 12.01.2021.

⁵F# code quotations is a run time metaprogramming technique which allows one to transform written F# code during program execution. Official documentation: <https://docs.microsoft.com/en-us/dotnet/fsharp/language-reference/code-quotations>. Access date: 12.01.2021.

```
1 type RInfinity = R of float | Infinity
2
3 [<Struct>]
4 type MinPlusSemiring =
5     MinPlusSemiring of RInfinity
6 with
7     static member Zero = MinPlusSemiring Infinity
8     static member (+)
9         (MinPlusSemiring x, MinPlusSemiring y) =
10             match x, y with
11             | R x, R y -> System.Math.Min(x,y) |> R
12             | _ -> Infinity
13             |> MinPlusSemiring
14     static member (*)
15         (MinPlusSemiring x, MinPlusSemiring y) =
16             match x, y with
17             | R x, R y -> x + y |> R
18             | _ -> Infinity
19             |> MinPlusSemiring
20     static member op_Implicit (MinPlusSemiring src) =
21         src
```

Listing 1: Example om Min-Plus semiring definition

it may help to solve a problem with explicit zeros⁶ because we should explicitly specify conversion from one semiring to another if required.

Matrices and vectors are equipped with monoid or semiring. Explicit type conversions. Can be automatically removed in some cases during translation time.

We propose to generate OpenCL c code in running time as a way to solve problems with generics: with strong typing all type information become known and can be used to generate kernels for specific types. Moreover, running time code generation is a way to apply advances optimization techniques, such as partial evaluation (or code specialization), which can improve performance of generated code when part of input parameter of kernel becomes known prior its generation [12].

The example of Min-Plus semiring definition is provided in listing 1. Type is defined using discriminated unions (line 1): new set can contains both floats, marked with R and a special value Infinity. Thus floats are extended with infinity as required for accurate definition of Min-Plus semiring. Semiring definition (lines 3–21) includes definition of zero (identity), operations \oplus (lines 8–13) and \otimes (lines 14–19), !!!!

III. IMPLEMENTATION DETAILS

To evaluate ideas described above we start a development of library named GraphBLAS⁷.

We use a Brahma.FSharp library for running time translation of F# code to OpenCL C, and for translated kernels execution. Brahma.FSharp is based on code quotations, thus utilizes strong typing to provide more static code checks, and polymorphic first class functions for general highly abstract code creation. Additionally, Brahma.FSharp provides special workflow builder to simplify heterogeneous programming and automate resource management.

⁶Discussion on zeros removing !!! Access date: 12.01.2021.

⁷Sources of GraphBLAS# on GitHub: <https://github.com/YaccConstructor/GraphBLAS-sharp>. Access date: 12.01.2021.

TABLE I
RESULTS OF ELEMENT-WISE ADDITION EVALUATION

Name	Matrix		SuiteSparse	Math.NET	GraphBLAS#	
	Rows	NNZ			CPU	GPGPU
m1	10	9	8	3	2	1
m2	10	9	8	3	2	1
m3	10	9	8	3	2	1
m4	10	9	8	3	2	1
m5	10	9	8	3	2	1

Abstraction layers which hides details of matrix representation and operations implementation. Currently we are working on COO and CSR formats and respective operations.

IV. EVALUATION

While our implementation of GraphBLAS API is on very early stage, we cannot evaluate it on well-known linear algebra based algorithms. But in order to !!! Elementwise addition.

We perform our experiments on the PC with Ubuntu 18.04 installed and with the following hardware configuration: !!! CPU, !!! RAM, !!!GPGPU with !!!!.

our solution on CPU and GPGPU. For comparison we choose the following libraries.

- SuiteSparse as a ...
- Math.NET Numerics⁸
- GraphBLAST

Dataset description. Matrices form SuiteSparse collection⁹

For each matrix !!!!. For .NET-based implementations *BenchmarkDotNet*¹⁰ is used. Results of performance evaluation are presented in table I. Time is measured in !!!

We can see, that !!!! results analysis and conclusion.

V. CONCLUSION

We present a work in progress that demonstrates a way to utilize both a power of high-level languages and performance of GPGPUs to implement GraphBLAS API. Our preliminary evaluation shows that !!!

In the future, first of all, we should extend our library up to full GraphBLAS API implementation. Moreover, it may be useful for community to implement an analog of LAGraph¹¹ algorithms collection for .NET on the top of our GraphBLAS API implementation.

The next step is evaluation of the solution on real-world cases and comparison with other implementations of GraphBLAS API on different devices and different algorithms. Additionally, it may be interesting to compare our solution

⁸Library which provides numerical computations primitives for .NET: <https://numerics.mathdotnet.com/>. Access date: 12.01.2021.

⁹!!!!

¹⁰*BenchmarkDotNet* allows one to automate benchmarking process for .NET platform. Project web page: <https://benchmarkdotnet.org/>. Access date: 12.01.2021.

¹¹LAGraph is a collection of algorithms implemented using GraphBLAS. Project sources on GitHub: <https://github.com/GraphBLAS/LAGraph>. Access date: 12.01.2021.

with graph analysis libraries and with linear algebra libraries for .NET platform.

Another direction of future work is Brahma.FSharp improvements. First of all, it is necessary to support discriminated unions to make it possible to express custom semirings such as Min-Plus, as presented in listing 1.

Also, it is necessary to add high-level abstractions for both asynchronous programming and for multi-GPU programming. Such mechanisms can be naturally expressed in F# with native primitives for asynchronous programming, and by using high-level abstractions for multiple GPUs management.

Finally, we plan to implement high-level optimizations, like fusion and specialization in Brahma.FSharp.

REFERENCES

- [1] J. Kepner, P. Aaltonen, D. Bader, A. Buluc, 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)*, Sep. 2016, pp. 1–9.
- [2] O. Selvitopi, S. Ekanayake, G. Guidi, G. A. Pavlopoulos, A. Azad, and A. Buluc, "Distributed many-to-many protein sequence alignment using sparse matrices," in *Proceedings of the International Conference for High Performance Computing, Networking, Storage and Analysis*, ser. SC '20. IEEE Press, 2020.
- [3] J. Kepner, M. Kumar, J. Moreira, P. Pattnaik, M. Serrano, and H. Tufo, "Enabling massive deep neural networks with the graphblas," in *2017 IEEE High Performance Extreme Computing Conference (HPEC)*, 2017, pp. 1–10.
- [4] T. A. Davis, "Algorithm 1000: Suitesparse:graphblas: Graph algorithms in the language of sparse linear algebra," *ACM Trans. Math. Softw.*, vol. 45, no. 4, Dec. 2019. [Online]. Available: <https://doi.org/10.1145/3322125>
- [5] A. Buluc and J. R. Gilbert, "The combinatorial blas: Design, implementation, and applications," *Int. J. High Perform. Comput. Appl.*, vol. 25, no. 4, pp. 496–509, Nov. 2011. [Online]. Available: <https://doi.org/10.1177/1094342011403516>
- [6] C. Yang, A. Buluc, and J. D. Owens, "GraphBLAST: A high-performance linear algebra-based graph framework on the GPU," *arXiv preprint*, 2019.
- [7] T. Henriksen, N. G. W. Serup, M. Elsmann, F. Henglein, and C. E. Oancea, "Futhark: Purely functional gpu-programming with nested parallelism and in-place array updates," in *Proceedings of the 38th ACM SIGPLAN Conference on Programming Language Design and Implementation*, ser. PLDI 2017. New York, NY, USA: ACM, 2017, pp. 556–571. [Online]. Available: <http://doi.acm.org/10.1145/3062341.3062354>
- [8] T. L. McDonell, M. M. Chakravarty, G. Keller, and B. Lippmeier, "Optimising purely functional gpu programs," *SIGPLAN Not.*, vol. 48, no. 9, pp. 49–60, Sep. 2013. [Online]. Available: <https://doi.org/10.1145/2544174.2500595>
- [9] R. LeiBa, K. Boesche, S. Hack, A. Pérard-Gayot, R. Membarth, P. Slusallek, A. Müller, and B. Schmidt, "Anydsl: A partial evaluation framework for programming high-performance libraries," *Proc. ACM Program. Lang.*, vol. 2, no. OOPSLA, Oct. 2018. [Online]. Available: <https://doi.org/10.1145/3276489>
- [10] T. Kenter, "Invited tutorial: Opencl design flows for intel and xilinx fpgas: Using common design patterns and dealing with vendor-specific differences," in *FSP Workshop 2019; Sixth International Workshop on FPGAs for Software Programmers*. VDE, 2019, pp. 1–8.
- [11] K. Shagririthaya, K. Kepa, and P. Athanas, "Enabling development of opencl applications on fpga platforms," in *2013 IEEE 24th International Conference on Application-Specific Systems, Architectures and Processors*, 2013, pp. 26–30.
- [12] A. Tyurin, D. Berezun, and S. Grigorev, "Optimizing gpu programs by partial evaluation," in *Proceedings of the 25th ACM SIGPLAN Symposium on Principles and Practice of Parallel Programming*, ser. PPoPP '20. New York, NY, USA: Association for Computing Machinery, 2020, pp. 431–432. [Online]. Available: <https://doi.org/10.1145/3332466.3374507>