# Distillation of Sparse Linear Algebra

Aleksey Tyurin

Saint Petersburg University, Russia

JetBrains Research, Russia

alekseytyurinspb@gmail.com

Ekaterina Vinnik

Saint Petersburg University, Russia

JetBrains Research, Russia

catherine.vinnik@gmail.com

Mikhail Nikoliukin

National Research University
Higher School of Economics, Russia

mnnikolyukin@edu.hse.ru

Daniil Berezun

Saint Petersburg University, Russia

JetBrains Research, Russia

d.berezun@spbu.ru

Semyon Grigorev

Saint Petersburg University, Russia

JetBrains Research, Russia

s.v.grigoriev@spbu.ru

Geoff Hamilton

School of Computing,
Dublin City University, Ireland

geoffrey.hamilton@dcu.ie

Linear algebra (LA) is a common language for many areas, providing high-performance solutions by utilizing the highly parallelizable nature of LA operations. However, intermediate data structures arising during LA expressions evaluation are one of the primary performance bottlenecks, especially when sparse data structures are used. We show that program distillation can be efficiently used to optimize LA-based programs by minimizing occurrence and evaluation of intermediate data structures.

*Keywords*— fusion, high-level synthesis, sparse data processing, linear algebra

## 1 Introduction

Nowadays high-performance processing of a huge amount of data is indeed a challenge not only for scientific computing, but for applied systems as well. Special types of hardware, such as General Purpose Graphic Processing Units (GPGPUs), Tensor Processing Units (TPUs), FPGA-based solutions, along with respective specialized software have been developed to provide appropriate solutions, and the development of new solutions continues. In its turn, sparse linear algebra and GraphBLAS standard in particular, are a way to utilize all these accelerators to provide high-performance solutions in many areas including machine learning [4] and graph analysis [9].

Unfortunately, evaluation of expressions over matrices generates intermediate data structures similar to the well-known example of a pipelined processing of collections: `map g (map f data)`. Suppose `data` is a list, then the first `map` produces a new list which then will be traversed by the second `map`. The same pattern could be observed in neural networks, where initial data flows through network layers, or in linear algebra expressions, where each subexpression produces an intermediate matrix. The last case occurs not only in scientific computations but also in graph analysis [9]. It is crucial that not only the data structures are traversed multiple times, while it is possible to traverse over them only once, but also the intermediate data populates memory (RAM). Extra memory accesses are a big problem for real-world data analysis: the size of data is huge and memory accesses are expensive operations with noticeable latency. While a number of complex real-world cases including stream fusion and dense kernels fusion [10] could be successfully optimized using deforestation [6] and other techniques [5], avoiding intermediate data structures in sparse data processing is still an open problem [9].

## 2 Proposed Solution

The goal of our research is to find out if sparse linear algebra-based programs can be efficiently optimized by program distillation [3] through eliminating intermediate data structures and computations. To answer this question, we have developed a library of matrix operations in POT language: a simple functional language used in distiller.

We use a quad-tree matrix representation [8] since it both avoids indexing and can be implemented via algebraic data type, which itself is very natural for functional programming. Besides, it provides similar compression

| Function | Matrix size | | | | Interpreter | | Reduceron | FHW |
|---|---|---|---|---|---|---|---|---|
| | m1 | m2 | m3 | m4 | Red-s | Reads | Ticks | Ticks |
| seqAdd | $64 \times 64$ | $64 \times 64$ | $64 \times 64$ | $64 \times 64$ | 2.7 | 1.9 | 1.8 | 1.4 |
| addMask | $64 \times 64$ | $64 \times 64$ | $64 \times 64$ | – | 2.1 | 1.8 | 1.4 | 1.4 |
| kronMask | $64 \times 64$ | $2 \times 2$ | $128 \times 128$ | – | 2.2 | 1.9 | 1.4 | 2.7 |
| addMap | $64 \times 64$ | $64 \times 64$ | – | – | 2.5 | 1.7 | 1.7 | 1.5 |
| kronMap | $64 \times 64$ | $2 \times 2$ | – | – | 2.9 | 2.2 | 1.8 | 2.0 |

Table 1: Evaluation results: original program to distilled one ratio of measured metrics is presented

rate to widely adopted CSR and COO, and a natural way to represent both sparse and dense matrices, as well as makes it possible to express basic operations over the representation via recursive functions traversing the tree-like structure. Finally, the quad-tree representation allows to natively exploit divide-and-conquer parallelism in matrices handling functions.

Since general purpose devices like CPUs and GPUs appear to be heavily underutilized when executing sparse routines due to low arithmetic intensity and memory boundness, custom hardware seems to be a promising solution towards mitigating these issues. Distillation provides all the needed optimizations for free for a functional language, so what remains is to translate a functional program into custom hardware. We have opted for two solutions here. The first one is Reduceron [7] — a processor, designed to perform highly performant reductions. The second one is FHW [2] — a project, which is constituted of several compilers that help to translate an arbitrary Haskell progam into System Verilog to eventually provide a bitstream for a custom hardware. It leverages a parallel and pipelined dataflow representation, which is abstracted over, e.g., nodes for memory operations, which makes it possible to come up with a specialized memory for our data structures. While the first case is more typical, the second might provide higher performance for specific tasks.

For now, we propose to use program distillation as the first step of program optimization which, we hope, should reduce memory usage and other unnecessary computations, and then compile a distilled program to the two different hardware platforms by using the respective compiler with platform-specific optimizations.

For evaluation, we have been implementing a library of a subset of linear algebraic routines, only matrix-matrix operations for now: matrix-to-matrix element-wise addition (`mtxAdd`), matrix-to-scalar *apply-to-all* operation (`map`), masking (`mask`), which takes a subset of matrix elements, and Kronecker product (`Kron`).

# 3   Preliminary Evaluation

The following examples, which are a combination of the implemented functions, are used for the evaluation. The examples are fairly practical, for example, one could see a sequence of element-wise additions in a Luby's maximal independent set algorithm, and masking is a one of key operation in GraphBLAS standard.

- `addMask m1 m2 m3 = mask (mtxAdd m1 m2) m3`   •  `kronMask m1 m2 m3 = mask (kron m1 m2) m3`

- `addMap m1 m2 = map f (mtxAdd m1 m2)`   •  `kronMap m1 m2 = map f (kron m1 m2)`

- `seqAdd m1 m2 m3 m4 = mtxAdd (mtxAdd (mtxAdd m1 m2) m3) m4`

We compare original versions of these functions and distilled ones in three ways. First, we use the interpreter of the POT language to measure the number of reductions and memory reads inside `case` expressions. We use the simulator shipped with Reduceron to measure the number of clock ticks necessary to evaluate a program, and Vivado's simulator for FHW-compiled programs to measure the number of clock ticks. It is worth noting that Reduceron has somewhat fixed clock frequency, while frequency for FHW-generated hardware varies for different input programs. Since we do not have external memory at the moment, and all the data lives inside the generated scheme, the logic is not synthesizable for reasonably sized matrices in the case of FHW. We get similar clock frequencies for distilled and non-distilled programs for inputs with smaller matrices and hence assume that clock frequencies are also similar below. Thus, we provide only the number of ticks instead of time.

A set of sparse matrices of appropriate sizes provided at [1] is used. The matrices are converted into boolean ones since POT language lacks the needed primitives at the moment. Average results for several hundreds of different inputs are presented in table 1.

We can see that on average distillation provides up to 3 and 2 times improvement in terms of reductions and memory reads respectively for the interpreter. The number of reductions is also considerably reduced for hardware benchmarks. The lack of matches between ticks for FHW and Reduceron is justified by architecture distinction. All this hopefully makes the proposed solution viable, and we look forward to coming up with full-fledged experiments that would target real hardware and real life competitors like C++ implementations.

## 4  Future Work

In the future, first, we should close a technical debt and make the distiller more stable to handle all the important cases: current implementation can not handle such important functions as matrix-matrix multiplication. Having basic workflow implemented, we should explore how to utilize distillation in the best way for each particular platform. For example, which level of distillation is the best for our particular problem and set of functions? Can we exploit more parallelism using distillation? Can we efficiently exploit the tail-modulo-cons property of the distilled program? What are the limitations of distillation: whether all important cases can be handled?

In addition to it, we plan to improve both FHW and Reduceron and compilers for them in order to make them mature enough to handle real-world examples. The most relevant improvement here, for example, is the support for out-of-chip memory.

## References

[1]  Timothy A. Davis & Yifan Hu (2011): *The University of Florida Sparse Matrix Collection*. ACM Trans. Math. Softw. 38(1), doi:10.1145/2049662.2049663.

[2]  S. Edwards, Martha A. Kim, Richard Townsend, Kuangya Zhai & L. Lairmore (2019): *FHW Project : High-Level Hardware Synthesis from Haskell Programs*.

[3]  Geoffrey Hamilton (2021): *The Next 700 Program Transformers*. arXiv:2108.11347.

[4]  Jeremy Kepner, Manoj Kumar, Jose Moreira, Pratap Pattnaik, Mauricio Serrano & Henry Tufo (2017): *Enabling massive deep neural networks with the GraphBLAS*. 2017 IEEE High Performance Extreme Computing Conference (HPEC), doi:10.1109/hpec.2017.8091098.

[5]  Oleg Kiselyov, Aggelos Biboudis, Nick Palladinos & Yannis Smaragdakis (2016): *Stream Fusion, to Completeness*. CoRR abs/1612.06668. arXiv:1612.06668.

[6]  Simon Marlow & Philip Wadler (1992): *Deforestation for Higher-Order Functions*. In John Launchbury & Patrick M. Sansom, editors: *Functional Programming, Glasgow 1992, Proceedings of the 1992 Glasgow Workshop on Functional Programming, Ayr, Scotland, UK, 6-8 July 1992*, Workshops in Computing, Springer, pp. 154–165, doi:10.1007/978-1-4471-3215-8_14.

[7]  Matthew Naylor & Colin Runciman (2012): *The Reduceron reconfigured and re-evaluated*. Journal of Functional Programming 22(4-5), p. 574–613, doi:10.1017/S0956796812000214.

[8]  I. Simecek (2009): *Sparse Matrix Computations Using the Quadtree Storage Format*. In: *2009 11th International Symposium on Symbolic and Numeric Algorithms for Scientific Computing*, pp. 168–173, doi:10.1109/SYNASC.2009.55.

[9]  Carl Yang, Aydin Buluç & John D. Owens (2019): *GraphBLAST: A High-Performance Linear Algebra-based Graph Framework on the GPU*. CoRR abs/1908.01407. arXiv:1908.01407.

[10]  Zhen Zheng, Pengzhan Zhao, Guoping Long, Feiwen Zhu, Kai Zhu, Wenyi Zhao, Lansong Diao, Jun Yang & Wei Lin (2020): *FusionStitching: Boosting Memory Intensive Computations for Deep Learning Workloads*. CoRR abs/2009.10924. arXiv:2009.10924.