

FlexGraph: A Reconfigurable Graph Analytics Accelerator for Heterogeneous CPU-FPGA Architectures

Blaise-Pascal Tine
Georgia Institute of Technology
blaisetine@gatech.edu

Abstract—In recent years, The end of Dennard Scaling [9] and the need for agile development methodology in data centers is pushing the computer architecture community towards designing more specialized and energy efficient systems [24]. This has lead to a new application domain for FPGAs and the emergence of heterogeneous CPU-FPGA computing platforms [12] [10], enabling the design of new energy efficient FPGA accelerators for domain specific applications. Graph Analytics, one the largest applications in production data centers today, faces scaling challenges with the ever increasing size of its workloads and its inherent low locality and memory-bound characteristic. The GraphMat [25] programming framework was proposed that efficiently translates graph algorithms into high performance computation in hardware. We present FlexGraph, a reconfigurable domain-specific accelerator for Graph Analytics targeting commodity heterogeneous CPU-FPGA platforms. FlexGraph is designed using Cocoh [26], a C++ domain specific library for hardware design and simulation, providing a high level programming abstraction with the efficiency of native RTL. FlexGraph’s novel architecture allows an easy implementation of GraphMat kernels using standard C++. To the best of our knowledge it is a first implementation of a graph accelerator on a commodity heterogeneous CPU-FPGA platform. FlexGraph’s preliminary performance is comparable to HLS, we later discuss several optimization strategies that could further improve its speed.

I. INTRODUCTION

The end of Dennard Scaling [9] has moved the focus of computer architecture designers towards power and energy efficient architectures. However, energy efficient solutions such as ASIC designs present a serious limitation in flexibility and production cycle. These constraints have pushed production data centers towards using FPGA-based accelerators [24] for their reconfigurability and energy savings when compared to general-purpose graphics processing units (GPUs). This trend has led to emergence of heterogeneous CPU-FPGA computing platforms [12] [10], enabling the design of new energy efficient FPGA accelerators [7] for domain specific applications. Graph Analytics is one the largest applications in production data centers today [11], its performance suffers from workload imbalance, frequent updates, limited data locality and low compute communication ratio making it memory-bound and energy inefficient. This problem is further exacerbated with the ever increasing size of its dataset, presenting

a scalability challenge for the industry. Several solutions have been proposed to address this problem both at the software level, with better algorithms and programming abstraction [21] [18] [25] [16] [15], and at the hardware level with custom accelerators [13] [1] [19]. GraphMat [25] and many other graph frameworks [25] [16] [15] define a vertex programming model to represent graph computation, but GraphMat is unique in that only the edges transformation of the kernel is accelerated, allowing the programmer to implement its algorithm in C++, instead of a proprietary language, increasing both accessibility and productivity. Flexgraph’s design goals aim at achieving the same balance between efficiency and productivity via hardware acceleration. We developed FlexGraph using the Cocoh’s [26] framework, a C++ domain specific library for hardware design and simulation, providing a high-level programming abstraction with the efficiency of native RTL. Several hardware description languages and libraries [14] [3] [23] [2] [8] [4] [5] have been proposed as alternative to native Verilog or VHDL to improve developer productivity by providing a higher level programming abstraction [3] [23] and execution model [2] [5]. Cocoh’s unique strength is that it provides a uniform development and simulation environment based on the same source. Cocoh’s modules are written directly in standard C++, debugged, tested and simulated using the same source, and later exported to Verilog for FPGA deployment. Another strength of Cocoh’s API is its integration with CHDL [14], allowing gate-level simulation and analysis of the design. Using the Cocoh’s API, developers can write their custom GraphMat kernel for the FlexGraph accelerator and test their design before production. FlexGraph architecture uniquely targets commodity heterogeneous CPU-FPGA computing platforms [12] for data centers. To the best of our knowledge it is a first implementation of a graph accelerator targeting this ecosystem. FlexGraph’s current performance is comparable to HLS, we have identified some performance optimizations that could further improve its performance which we will discuss later in the paper. The remainder of this paper is organized as follows: Section 2 describes GraphMat computing model and its SPMV kernel, section 3 describes FlexGraph’s architecture, section 4 describes the experimental setup, section 5 describes our results analysis,

section 6 describes the related work, section 7 discusses future work and finally, section 8 summarizes our main contribution and results.

II. FLEXGRAPH ARCHITECTURE

In this section, we describe the overall architecture of the FlexGraph Accelerator.

A. FlexGraph Computing Model

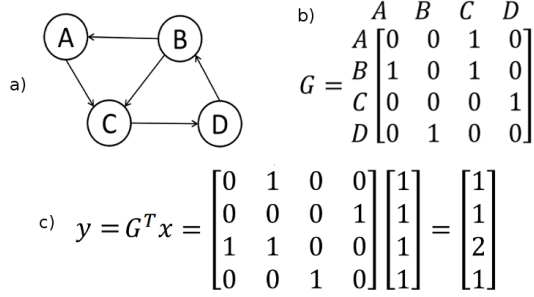


Fig. 1. Graph Vertex Programming

FlexGraph's architecture is based on GraphMat [25]'s vertex computing paradigm where a graph algorithm is described as a linear algebra set of operations operating on sparse matrix with input vertices. Figure 1 shows an example where graph (a) is converted into matrix (b), allowing us to execute operations on the graph using simple matrix vertex multiplication (c). This above example will generate the number of incoming edges for each vertex. All GraphMat algorithms, Breadth First Search (BFS) [27], PageRank [22], Single Source Shortest-Path (SSSP) [29], etc., are expressed in this fashion such that computation can be efficiently accelerated in hardware. In FlexGraph, we accelerate the matrix vertex computation on the FPGA, instead of executing it on the multi-core host processor.

B. The Doubly Compressed Sparse Matrix

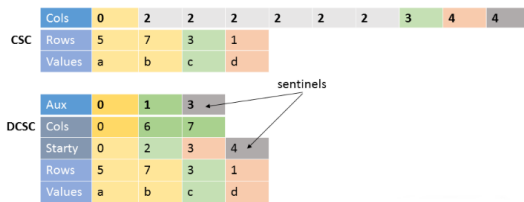


Fig. 2. DCSC hyperspace matrix format

Most Data Analytics graphs are represented in a sparse matrix format to save memory space due to the large amount of empty edges that exist. FlexGraph's matrices are stored using a Doubly Compressed Sparse Column (DCSC) [6] format. The format allows minimal traversal into the sparse matrix structure, saving necessary memory bandwidth when fetching empty columns. Figure 2 shows a sample matrix $A = \{(0,5,a), (0,7,b), (6,3,c), (7,1,d)\}$ with four non-zero edges (src, dst, weight) encoded using DCSC

versus the conventional CSC [30] format. The traversal over CSC requires accessing all consecutive column ranges ($cols[i+2] - cols[i]$) even though most of the distances are empty. DCSC saves on bandwidth by using an additional indirection buffer inside its structure to only encode non-zero columns. It is important to note that the saving this format introduced additional space and compute overhead that only pays off when the matrix is very large and sparse, a prevalent characteristic of Graph Analytics datasets.

C. The Sparse Matrix Multiplication Kernel

FlexGraph architecture implements a custom Sparse Matrix Vertex Multiplication (SPMV) kernel. Its architecture diverges from conventional implementations in two ways: First, it uses the dirty masks from the input vertices to select the edges to process and return a new dirty masks capturing the intersection of the active input vertices and the non-zero edges actually visited. Second, its matrix data structure (DCSC) contains an additional indirection buffer for accessing the matrix columns. These two properties present unique performance challenges when designing the accelerator. Listing 1 shows the pseudo-code of the SPMV kernel. The program iterates through each column and fetches the corresponding input vertex active mask to check if the column should be processed, then access all the non-zero rows of the matrix for that column to evaluate each edge. In the code's comments on the right are highlighted the different types of performance hogs present: semi-random memory accesses (*), fully-random memory accesses (!) and control branches (?). The memory accesses for the columns and rows data (lines 4 and 9) are semi-random because only the access cannot be predicted and after that the address is simply incremented. FlexGraph's architecture attempts to address some of those performance hogs using several optimizations detailed in section 4.

```

1 def SPMV_kernel(x_values, x_activemask):
2     y_values = {0}, y_activemask = {0}
3     for col in (c_start, c_end):
4         (a_x, r_start, r_end) = coldata[col] # *
5         x_value = x_values[a_x] # !
6         x_active = x_activemask[a_x] # !
7         if (x_active): # ?
8             for row in (r_start, r_end):
9                 (a_y, a_value) = rowdata[row] # *
10                y_values[a_y] += a_value * x_value # !
11                x_activemask[a_y] = true # !
12    return (y_values, y_activemask)

```

Listing 1. Pseudo-code for SPMV kernel

D. FlexGraph Microarchitecture

Externally, FlexGraph input and output signals implement an Accelerator Functional Unit (AFU) interface defined by Intel's Accelerator Abstraction Layer (AAL) [12]. AFUs implementing the interface are able to bind with the FPGA's board support package (BSP) and seamlessly communicate with the AAL software running on the host processor. Listing 2 shows AAL device interface implementation using Cocoh's API. Our FlexGraph Accelerator's

class in Cocoh simply derives from this interface and override the *initialize()* function to provide its implementation and Cocoh takes care of generating the corresponding verilog module. It implements three input signals (*start*, *qpi_in*, *ctx*) and two output signals (*qpi_out*, *done*). The AFU starts execution when the *start* signal is asserted and communicate completion by asserting the *done* signal. The *ctx* signal provides application's specific data like constants and buffers address in the case of FlexGraph. The *qpi* in/out signals implement Intel Quick Path Interconnect (QPI) interface [28], providing single channel read/write ports for accessing external shared memory. The FPGA socket hosts a 64-byte cache coherent interface (CCI) [20] that connects to the host processor's last level cache (LLC) via QPI.

```

1 class aal_device {
2 public:
3     virtual out_t initialize() (
4         const ch_logic& start,
5         const qpi::in_t& qpi_in,
6         const afu_ctx_t& ctx,
7         qpi::out_t& qpi_out,
8         ch_logic& done
9     ) const = 0;
10 };

```

Listing 2. AAL Device Interface in Cocoh C++

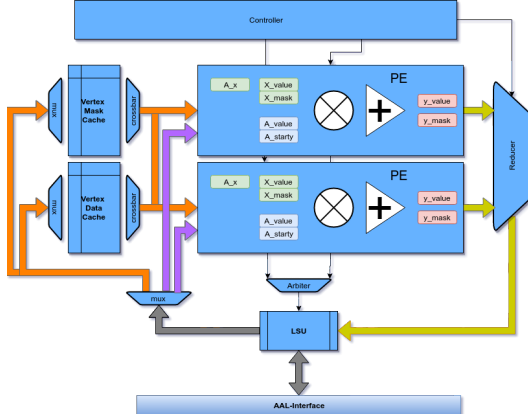


Fig. 3. FlexGraph Microarchitecture

Figure 3 illustrates FlexGraph accelerator microarchitecture. It is comprised of four main module types, a controller, processing elements (PEs), a load store unit (LSU) and vertex caches. The main controller is responsible for starting the accelerator, scheduling tasks for execution, reporting hardware counters and terminating the execution. The processing elements execute partition tasks assigned to them by the controller and communicate with the LSU to access the matrix and vertex data for their partition. They are also responsible for sending their final output result back to memory. The LSU is the module responsible for managing external communication between the accelerator and memory via the QPI interface. It directly binds to the QPI ports defined in listing 2. The vertex caches store intermediate vertex data and active masks for sharing between the processing elements. FlexGraph processing pipeline

slightly resembles the SPMV execution steps illustrated in listing 1. We made several important modifications to it to improve performance.

III. FLEXGRAPH PARALLELIZATION

The first optimization we performed in FlexGraph early during the design phase was the support of multiple processing units. The reasoning behind it was to extract the maximum bandwidth out of the LSU such that the QPi is always busy processing a request when some processing element are stalled waiting for their data to return. the other advantage of parallelization for FlexGraph is to increase the overall accelerator throughput by processing multiple a larger chunk of the workload per unit of time. The scheduling of the partitions for execution on the processing elements is controlled by a dispatch unit inside the main controller module. The dispatch unit fetches partitions data from the LSU and pass down each partition in a first come first serve fashion to the processing elements.

A. DCSC Matrix partitioning

| | c0 | c1 | | c2 | | c3 |
|----|-----|----|-----|-----|--|----|
| P0 | x | | | x | | |
| | x x | | | | | |
| | | | | | | |
| P1 | xxx | x | | xxx | | |
| | x | xx | xxx | | | xx |
| | | | | | | |
| | | | | | | |

Fig. 4. DCSC Matrix Partitioning

To enable the parallelization of FlexGraph tasks, the DCSC sparse matrix is first partitioned into aligned partitions of 32 consecutive rows containing non-zero edges. we choose a partition size of 32 mainly to match the 32-bit size of the bitmasks encoding the vertices that are active. These masks are used by the software on the host processor to determine the regions of the acceleration's output buffer that have been updated. Figure 4 illustrates sample matrix partitioning in which the non-zero horizontal regions have been broken into two partitions P0 and P1. The non-zero column ranges (*c0*, *c1*, *c2*, *c3*) represent the selected chunks covered by each partition. Partition P0 will only contain ranges *c0* and *c2* while P2 has non-zeros in all four ranges. The partitioning scheme is not ideal because of the workload in-balance that might exist when the number of non-zero varies disproportionally between partitions.

B. Synchronising Memory Accesses

An implementation challenge we faced when supporting multiple processing elements in FlexGraph was the synchronization of write accesses for the output active masks. Because the LSU data transfer granularity is 64 byte (matching CCI [20]), FlexGraph need to manage shared

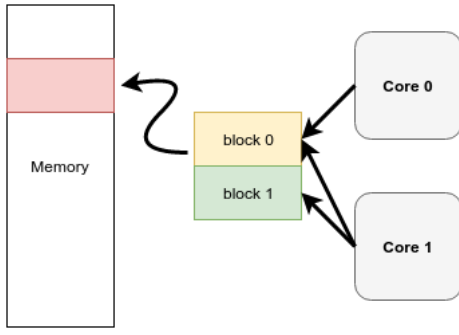


Fig. 5. Active Masks Write Synchronization

write accesses to the same block to avoid having a processing element override the content of another one. Luckily this doesn't pose any problem for writing out partition output values because the 32 rows within a partition occupy $32 * 4 \text{ bytes} = 2 * 64 \text{ bytes}$ blocks. Because the partitioning reinforces 32 rows alignment, each processing element can safely write into their assigned blocks. However, writing out the active masks poses a challenge because a single 4-byte mask encodes the active states of all 32 rows in a partition. This causes the processing elements to potentially share a single 64-byte block when writing the masks out to memory. To alleviate the contention, FlexGraph keep N active 64-byte blocks locally assigned to either one of the N processing elements. It assigns an ownership bit mask to each block to track their reference count and the processing element assigned to them. It also tracks the current address value assigned to each block. When a processing element is ready to write its 4-byte active mask, it passes the mask address and value to the LSU. The LSU goes to last active block the processing element previous wrote to and check if the block address matches. If so, it simply adds the new content to the existing block. If there is no match, it clears its ownership bit and flushes the block to memory if the ownership mask goes to zero. Then it looks up the other active block if anyone already has the address to use it, otherwise if it acquires the last unused block. Figure 5 shows a simplified illustration of the scheme with two processing elements. We don't need to keep more than N blocks in local storage for this scheme to work because the block address references are always incremental and never regress, meaning that there is always going to be free block available to use. It is important to also point out that for this scheme to work, the operation has to be atomic. FlexGraph has a single communication channel between all processing elements and the LSU which does some round robin arbitration and blocks the write mask request until it is committed. We also investigated an alternative solution to avoid this synchronization, which is simply increasing the partition size to $16 * 32$ rows, allowing the aggregated active masks to occupy a full 64-byte block. They were two major issues with this approach. Firstly, the total size of all resident partitions will be too large if we support multiple

processing elements. Secondly, it will worsen the workload imbalance that is already present with 32 rows.

IV. FLEXGRAPH OPTIMIZATIONS

This section describes some optimizations we added to FlexGraph to improve its performance.

A. Optimizing Memory Accesses via Stream Buffers



Fig. 6. Stream Buffers

As stated before, FlexGraph's LSU transfer granularity is 64 bytes, which is much larger than the 4 bytes of elements accessed by the processing elements. To save on bandwidth, the LSU uses stream buffers that fetch a 64-byte blocks of data but extracts 4-byte elements at the time. It is implemented using a fifo structure in the back-end where 64-byte responses from QPI go into. In the front-end, there is a temporary 64-byte block that is extracted from the fifo on demand to deliver 4-byte element at the time. We employ a shift register to extract the 4-byte element to pass it to the processing element. Figure 6 shows an illustration of the stream buffer concept.

B. Caching Vertex Values and Masks

Another optimization we implemented in FlexGraph is the caching of vertex values and masks. Looking at the SPMV pseudo-code in listing 1, we can observe in lines 5 and 6 that there is a random memory access to buffers x_values and $x_activemask$. FlexGraph employs a stream buffer like concept to consume 64-byte vertex data once they arrive inside the processing element to each iteration loop. However, the same block can be referenced by another processing element and sharing them inside the cache structure can save unnecessary memory traffic. We implemented two small fully associative caches with first-in-First-out replacement to hold active vertex data during execution. The caches are implemented inside the LSU which is responsible for managing them. When a vertex request arrives from a processing element, the LSU first looks up the cache if the block is already present and return it. If the block is not there, it sends the request to QPI. QPI output interface support a 14-bit metadata field that is used to identify the block once it is returned. We use that field to store the index of the block such that we can compute the cache tag when it arrives. Because the cache is small, we implemented a one cycle tag lookup for the LSU to know if the block is present and accessible.

C. Executing Non-blocking Memory writes

FlexGraph memory writes are all non-blocking, this applies to both the output values and active masks. This allows the processing elements to push their write request

and resume execution while the LSU is processing them. Upon write responses from the QPI channel, the LSU internally keeps track of the count of outstanding requested to ensure that all writes have completed. The QPI interface exposes two write response ports by which the memory replies could be sent. This allows servicing two responses simultaneously. The LSU implements a mechanism to process both channels simultaneously when they are active. At completion time when all processing elements are done executing the current run, the controller waits for all outstanding write requests to complete before asserting the *done* signal.

V. EXPERIMENTAL SETUP

In this work, we simulated our evaluation environment entirely in software using the Cocoh [26] C++ framework. Figure 7 shows an overview of what the simulation environment contains. At the top we have the accelerator implementation with its AAL interface, both modelled using Cocoh’s API. At the bottom, we have the Cocoh’s simulation runtime that bind the accelerator and the back-end host CPU simulator. Inside our host CPU simulator we modelled a QPI memory simulator which managed the shared memory space and served all the memory requests coming from the accelerator. In the current design we used a 1 GB shared memory space for holding the matrix and vertex data, including the output result from the accelerator. On the actual Xeon platform [12], the shared memory space limit is 2 GB.

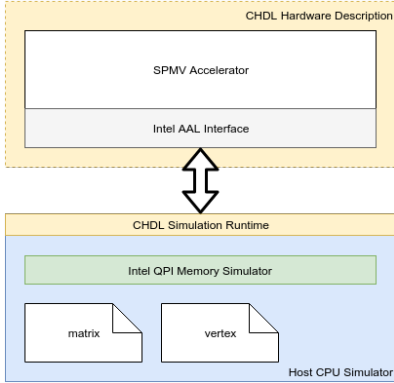


Fig. 7. FlexGraph Simulation Environment

| Simulation Parameters | | Accelerator Parameters | |
|-----------------------|------------|------------------------|--------------|
| Shared Memory | 1 GB | Clock Frequency | 200 Mhz |
| QPI Max Requests | 90 | Vertex Data Cache | 8x64 Bytes |
| CCI Hit Rate | 80% | Vertex Mask Cache | 8x64 Bytes |
| LLC Hit Rate | 50% | Cols Data Buffers | 2x2x64 Bytes |
| CCI Latency | 4 cycles | Rows Data Buffers | 2x5x64 Bytes |
| LLC Latency | 32 cycles | Output Buffer per PE | 2x64 Bytes |
| Memory Latency | 100 cycles | Number of PEs | 2 |

Fig. 8. FlexGraph Simulation Parameters

A. Intel QPI Simulation

We simulated the Intel QPI memory transaction using an analytic model capturing the latency of memory transfers when the request hits local CCI [20] cache or the host processor LLC. All outgoing memory requests from the accelerator to CCI take two cycles and the responses from CCI hits take another two cycles, giving a minimum round trip latency of four cycles for read requests. Figure 8 shows some configuration parameters we use for both the simulation and accelerator modelling.

B. Graph Datasets

We used the Graph500 [17] synthetic datasets for our evaluation. We configured graph generator similar to GraphMat [25] setting the default RMat paramters (A=0.57, B=C=0.15 and D=1-A-B-C). We used five scaling factors for our graphs (6, 8, 10, 12, 14) with corresponding number of non-zero ranging from 1024 to 262,144 accordingly. These datasets are relatively small mainly due to simulation time constraints, but the setup could have supported up to a scaling factor of 20 (about 16,777,216 non-zeros). Due to time constraints we did not test any real world dataset, that will be discussed in our future work section. Figure 9 lists down the number of vertices and non-zeros for our dataset.

| dataset | m6 | m8 | m10 | m12 | m14 |
|--------------|------|------|-------|-------|--------|
| scale factor | 6 | 8 | 10 | 12 | 14 |
| num vertices | 64 | 256 | 1024 | 4096 | 16384 |
| Non-zeros | 1024 | 4096 | 16384 | 65536 | 262144 |

Fig. 9. Graph500 dataset

VI. RESULTS ANALYSIS

In this section, we discuss our evaluation results.

A. Throughput Performance

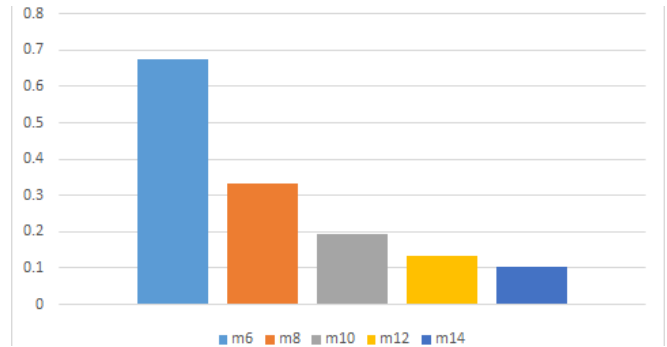


Fig. 10. FlexGraph Throughput Performance

Figure 10 shows the throughput performance of the accelerator in GFlops. We estimated the performance using the following equation:

$$Perf = \frac{2 * nonzeros}{latency}$$

The performance varies from 0.67 Gflops to 0.1 Gflops as the size of the dataset increases. This result is comparable to our implementation of a similar kernel using the Altera OpenCL HLS, also showing similar performance degradation as the workload size increases. In prior work, GraphOps [19] peak performance for their FPGA implementation was 0.18 GFlops with a graph size of 512K vertices, about x3 our throughput. The target platform has peak memory bandwidth is 6.0 GB/s, this shows that there is room for improvement. than We identified several optimization opportunities that we will discuss later in this section.

B. The Impact of Parallelism

Figure 11 shows the throughput performance of the accelerator in GFlops for a single-core system versus the dual-core implementation. We can see a considerable drop in performance (by almost half) across all datasets when a single processing element is executing. This is happening because there is much less memory level parallelism that can be exploited when a single core is running. In a dual-core system, the LSU receives a lot more memory requests from the processing element and has much less idle time when QPI is taking more cycles to return the blocks. When a processing element is stalled on a request, another may be able to make progress. We believe that increase the number of processing element beyond two will have a possible impact in performance for FlexGraph, possibly doubling the current throughput.

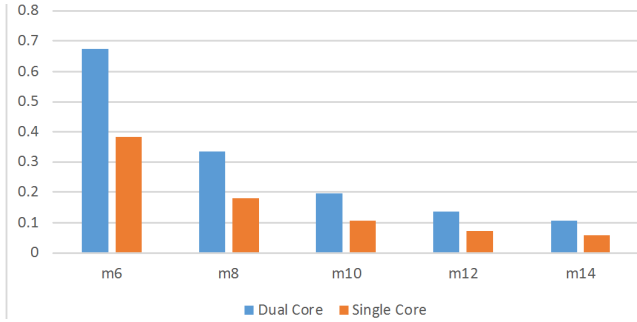


Fig. 11. FlexGraph Single Core Performance

C. The Impact of Memory Optimizations

Figure 12 shows the accelerator performance with all memory optimizations disabled in GFlops. Once again the graph shows a similar performance degradation as the dataset size decreases. The memory optimizations give the accelerator an average performance speed-up of 45%, which is significant considering the fact the the system is memory bound.

D. Summary Discussion

In summary, FlexGraph current average performance is 0.1 Gflops (using the lower bound of the largest dataset

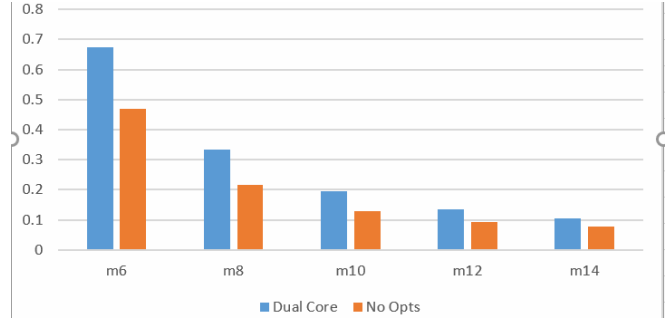


Fig. 12. FlexGraph Performance with memory optimizations

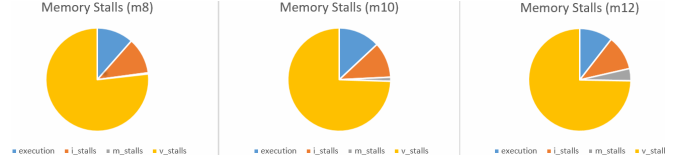


Fig. 13. FlexGraph Memory Stalls

in our experiment). This is in part with our HLS implementation of the same kernel. We have identified several optimization opportunities to boost the accelerator performance.

We added some performance counters into FlexGraph accelerator to monitor the memory stalls in each processing elements and identify bottlenecks. Figure 13 shows the memory stalls distribution over three datasets *m8*, *m10* and *m12*. The blue region represents the execution cycles, the orange region represents the stall cycles occurring when we access the matrix column indices, this is *coldata* buffer in the SPMV pseudo-code (line 4 in listing 1). The gray region represents the stall cycles when accessing the vertices active masks (line 6 in listing 1). We can significantly reduce the stall cycles by prefetching the blocks before they are needed to hide the memory latency. Using a simple next-line prefetcher will suffice since we know the access pattern on *coldata*.

Another performance optimization is in reducing the impact of the branch operation when checking if a vertex is active of not (line 7 in listing 1). We will investigate adding a preliminary test on the entire 32-bit mask to see if it is zero and skip the whole block altogether.

Lastly, another performance opportunity is increasing the number of preprocessing elements. Figure 11 showed the

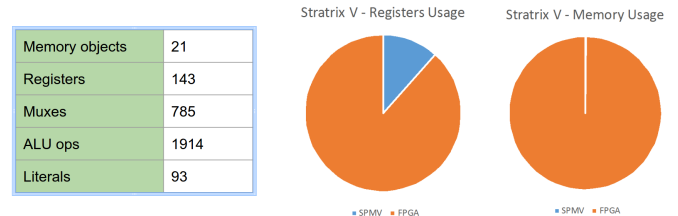


Fig. 14. Cocoh's hardware cost evaluation

performance gain of using two cores versus a single one. Having multiple processing elements will improve the system memory level parallelism by overlapping the execution of some processing element while others are stalled on memory accesses. We ran some estimate of the current design hardware cost using Cocoh's framework (see Figure 14) and it showed that our memory usage is very small (about 1% of the target Stratix V FPGA capacity). Also, FlexGraph accelerator only currently consumes 15% of the registers capacity. With this reserve, we anticipate extending the number of processing elements to 16. A challenge that we will have to resolve when adding additional processing elements is the synchronization of write accesses for the active masks. Another challenge will be the arbitration latency in the LSU when communicating with all nodes, since the single LSU will now become a bottleneck in the system.

VII. RELATED WORK

In this section, we discuss prior works related to this project.

A. Graphicionado

Graphicionado [13] is a much recent project that is similar to FlexGraph. In this work, the authors implemented a graph accelerator targeting Intel Xeon platform. Graphicionado also uses a vertex programming model of graphs computation similar to GraphMat [25]. It is a more specialized accelerator compared to FlexGraph in that it implements the entire graph algorithm in hardware, sacrificing flexibility for efficiency. Although the design allows some reconfigurability using FPGA, it still add many inconveniences compared to software. FlexGraph on the other hand only accelerate the graph edges computation in hardware, most of the algorithm is still written expressed in C++ and runs on the host processor like in GraphMat. FlexGraph reconfigurability is minimal and restricted to the SPMV module (changing the matrix datatype or the reduce operation). Graphicionado main architecture strength is their efficient optimizations of the graph access patterns to reduce stalls inside the pipeline.

B. GraphOps

GraphOps [19] is another graph accelerator also targeting FPGAs like FlexGraph. In this work, the authors propose an accelerator architecture that efficiently process a graph in memory encoded using a proposed locality-optimized scheme. Their main contribution is the graph storage representation which provides an efficient memory access pattern. A drawback of their proposal is that fact that the host processor has to pre-process the graph before the accelerator can consume it, which can add some considerable latency, considering that Graph Analytics algorithms tend to update the graph frequently.

VIII. FUTURE WORK

As stated in the results discussion, FlexGraph performance can be greatly improved once we implement the proposed optimizations. Doing this first step will be our priority for future work. The Cocoh C++ library we used in this work did not yet have support for verilog generation, we plan to add that support such that we could install the accelerator on the actual FPGA board for evaluation. We also anticipate to do some power and energy efficiency analysis on the final design, hopefully identifying all the components in the system consuming the most energy. We also plan to add some real world datasets to our evaluation, popular benchmarks include Flickr, Facebook, Wikipedia, Twitter, Netflix and LiveJournal.

IX. CONCLUSION

In this paper we presented FlexGraph, a reconfigurable Graph Analytics Accelerator targeting commodity heterogeneous CPU-FPGA platforms. FlexGraph was implemented using Cocoh C++ library, providing a productive and efficiency development support for modifying the accelerator when needed to accommodate all graph algorithms. FlexGraph architecture implements a customized SPMV kernel adapted for graphs vertex computation. It implements some memory optimizations and tasks parallelization to improve its performance. FlexGraph current performance is not ideal and we identified several strategies to fix that, including doing some memory prefetching and adding additional processing elements.

REFERENCES

- [1] J. Ahn, S. Hong, S. Yoo, O. Mutlu, and K. Choi. Tesseract: A scalable processing-in-memory accelerator for parallel graph processing. In *Proceedings of the 42Nd Annual International Symposium on Computer Architecture, ISCA '15*, pages 105–117, New York, NY, USA, 2015. ACM.
- [2] Arvind. Bluespec: A language for hardware design, simulation, synthesis and verification invited talk. In *Proceedings of the First ACM and IEEE International Conference on Formal Methods and Models for Co-Design, MEMOCODE '03*, pages 249–, Washington, DC, USA, 2003. IEEE Computer Society.
- [3] J. Bachrach, H. Vo, B. Richards, Y. Lee, A. Waterman, R. Avienis, J. Wawrzyniek, and K. Asanovi. Chisel: Constructing hardware in a scala embedded language. In *Design Automation Conference (DAC), 2012 49th ACM/EDAC/IEEE*, pages 1212–1221, June 2012.
- [4] P. Bellows and B. Hutchings. Jhdl-an hdl for reconfigurable systems. In *FPGAs for Custom Computing Machines, 1998. Proceedings. IEEE Symposium on*, pages 175–184, Apr 1998.
- [5] G. Berry and G. Gonthier. The esterel synchronous programming language: Design, semantics, implementation. *Sci. Comput. Program.*, 19(2):87–152, Nov. 1992.
- [6] A. Buluc and J. R. Gilbert. On the representation and multiplication of hypersparse matrices. In *2008 IEEE International Symposium on Parallel and Distributed Processing*, pages 1–11, April 2008.
- [7] Y.-k. Choi, J. Cong, Z. Fang, Y. Hao, G. Reinman, and P. Wei. A quantitative analysis on microarchitectures of modern cpu-fpga platforms. In *Proceedings of the 53rd Annual Design Automation Conference, DAC '16*, pages 109:1–109:6, New York, NY, USA, 2016. ACM.
- [8] J. Decaluwe. Myhdl: A python-based hardware description language. *Linux J.*, 2004(127):5–, Nov. 2004.

- [9] H. Esmailzadeh, E. Blem, R. S. Amant, K. Sankaralingam, and D. Burger. Dark silicon and the end of multicore scaling. In *2011 38th Annual International Symposium on Computer Architecture (ISCA)*, pages 365–376, June 2011.
- [10] Y. L. Fei Chen. Fpga acceleration in a power8 cloud. <https://openpowerfoundation.org/blogs/fpga-acceleration-in-a-power8-cloud>.
- [11] D. Gage. The new shape of big data. <http://www.wsj.com/articles/SB10001424127887323452204578288264046780392>, 2013.
- [12] P. Gupta. Intel xeon+fpga platform for the data center. <http://reconfigurablecomputing4themas.net/files/2.2%20PK.pdf>.
- [13] T. J. Ham, W. Lisa, S. Narayanan, S. Nadathur, and M. Margaret. Graphicionado: A high-performance and energy efficient accelerator for graph analytics. In *the 49th Annual IEEE/ACM International Symposium on Microarchitecture (MICRO)*, 2016.
- [14] C. D. Kersey. Chdl: Opensource c++ hardware design library. <https://github.com/cdkersey/chdl>, 2012.
- [15] Y. Low, D. Bickson, J. Gonzalez, C. Guestrin, A. Kyrola, and J. M. Hellerstein. Distributed graphlab: A framework for machine learning and data mining in the cloud. *Proc. VLDB Endow.*, 5(8):716–727, Apr. 2012.
- [16] G. Malewicz, M. H. Austern, A. J. Bik, J. C. Dehnert, I. Horn, N. Leiser, and G. Czajkowski. Pregel: A system for large-scale graph processing. In *Proceedings of the 2010 ACM SIGMOD International Conference on Management of Data*, SIGMOD '10, pages 135–146, New York, NY, USA, 2010. ACM.
- [17] R. C. Murphy, K. B. Wheeler, B. W. Barrett, and J. A. Ang. Introducing the graph 500. <http://www.graph500.org/>, 2010.
- [18] D. Nguyen, A. Lenharth, and K. Pingali. A lightweight infrastructure for graph analytics. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*, SOSP '13, pages 456–471, New York, NY, USA, 2013. ACM.
- [19] T. Oguntebi and K. Olukotun. Graphops: A dataflow library for graph analytics acceleration. In *Proceedings of the 2016 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, FPGA '16, pages 111–117, New York, NY, USA, 2016. ACM.
- [20] N. Oliver, R. R. Sharma, S. Chang, B. Chitlur, E. Garcia, J. Grecco, A. Grier, N. Ijhi, Y. Liu, P. Marolia, H. Mitchel, S. Subhaschandra, A. Sheiman, T. Whisonant, and P. Gupta. A reconfigurable computing system based on a cache-coherent fabric. In *2011 International Conference on Reconfigurable Computing and FPGAs*, pages 80–85, Nov 2011.
- [21] L. Page, S. Brin, R. Motwani, and T. Winograd. Apache spark's api for graph and graph-parallel computation. <http://spark.apache.org/graphx/>.
- [22] L. Page, S. Brin, R. Motwani, and T. Winograd. The pagerank citation ranking: Bringing order to the web. Technical Report 1999-66, Stanford InfoLab, November 1999. Previous number = SIDL-WP-1999-0120.
- [23] P. R. Panda. Systemc - a modeling platform supporting multiple design abstractions. In *System Synthesis, 2001. Proceedings. The 14th International Symposium on*, pages 75–80, 2001.
- [24] A. Putnam, A. M. Caulfield, E. S. Chung, D. Chiou, K. Constantinides, J. Demme, H. Esmailzadeh, J. Fowers, G. P. Gopal, J. Gray, M. Haselman, S. Hauck, S. Heil, A. Hormati, J.-Y. Kim, S. Lanka, J. Larus, E. Peterson, S. Pope, A. Smith, J. Thong, P. Y. Xiao, and D. Burger. A reconfigurable fabric for accelerating large-scale datacenter services. *SIGARCH Comput. Archit. News*, 42(3):13–24, June 2014.
- [25] N. Sundaram, N. Satish, M. M. A. Patwary, S. R. Dulloor, M. J. Anderson, S. G. Vadlamudi, D. Das, and P. Dubey. Graphmat: High performance graph analytics made productive. *Proc. VLDB Endow.*, 8(11):1214–1225, July 2015.
- [26] B. Tine. Cocoh, a c++ domain specific library for hardware design and simulation. <http://casl.gatech.edu/research/cocoh>, 2017.
- [27] wikipedia. Breadth-first search. https://en.wikipedia.org/wiki/Breadth-first_search.
- [28] wikipedia. Intel quickpath interconnect. https://en.wikipedia.org/wiki/Intel_QuickPath_Interconnect.
- [29] wikipedia. Shortest path problem. https://en.wikipedia.org/wiki/Shortest_path_problem.
- [30] wikipedia. Sparse matrix. https://en.wikipedia.org/wiki/Sparse_matrix.