

HSTREAM: A directive-based language extension for heterogeneous stream computing

Suejb Memeti
Department of Computer Science
Linnaeus University
Växjö, Sweden
suejb.memeti@lnu.se

Sabri Pllana
Department of Computer Science
Linnaeus University
Växjö, Sweden
sabri.pllana@lnu.se

Abstract—Big data streaming applications require utilization of heterogeneous parallel computing systems, which may comprise multiple multi-core CPUs and many-core accelerating devices such as NVIDIA GPUs and Intel Xeon Phi. Programming such systems require advanced knowledge of several hardware architectures and device-specific programming models, including OpenMP and CUDA. In this paper, we present HSTREAM, a compiler directive-based language extension to support programming stream computing applications for heterogeneous parallel computing systems. HSTREAM source-to-source compiler aims to increase the programming productivity by enabling programmers to annotate the parallel regions for heterogeneous execution and generate target specific code. The HSTREAM runtime automatically distributes the workload across CPUs and accelerating devices. We demonstrate the usefulness of HSTREAM language extension with various applications from the STREAM benchmark. Experimental evaluation results show that HSTREAM can keep the same programming simplicity as OpenMP, and the generated code can deliver performance beyond what CPUs-only and GPUs-only executions can deliver.

Index Terms—stream computing, heterogeneous parallel computing systems, source-to-source compilation

I. INTRODUCTION

Nowadays, a huge amount of data is generated throughout various mechanisms, such as scientific measurement and experiments (including genetics, physics, and astronomy), social media (including Facebook, and Twitter), and health-care [6, 10]. The current challenges of big data include storing and processing very large files.

However, in big data applications, not necessarily the entire data has to be processed at once. Furthermore, in most of the big-data applications, the data may be streamed, which means flowing in real-time (for instance data coming from different sensors in the Internet of Things), and therefore, it may not be available entirely. In such cases, the data needs to be processed in chunks and continuously [13].

Heterogeneous parallel computing systems comprise multiple non-identical processing units (PU), including CPUs on the host and accelerating devices (such as GPU, Intel Xeon Phi, and FPGA). Most of the top supercomputers in the world [20] comprise multiple nodes with heterogeneous processing units. For instance, the nodes of the current number

one supercomputer in the TOP500 list consist of two IBM POWER9 CPUs and six NVIDIA Volta V100 GPUs.

While the combination of such heterogeneous processing units may deliver high performance, scalability, and energy efficiency, programming and optimizing such systems is much more complex [3, 12, 4]. Different manufacturers of accelerating devices prefer to use different programming frameworks for offloading (which means transferring the data and control from the host to the device). For instance, OpenMP is used to offload computations to Intel Xeon Phi accelerators, whereas CUDA and OpenCL are used for offloading computations to GPUs.

The programming complexity of such systems leads to system underutilization. For example, applications designed for multi-core processing are able to utilize the available resources on the host CPUs. However, while the host CPUs are performing the actual work, the accelerating devices remain idle. On the other hand, most of the applications designed for accelerating devices use the CPU resources just for performing the data transfers and initiation of kernels, which is often performed by a single thread. Modern multi-core CPUs comprise a larger number of cores/threads, hence most of the CPU resources remain idle.

Researchers have proposed different techniques to address challenges of heterogeneous parallel programming and big data. For instance, source-to-source compilation techniques are proposed to ease programming of data-parallel applications [22, 8, 1, 18]. Similarly, approaches that are based on C++ template library are used for stream and data parallel computing systems [5, 7]. Pop and Cohen [16] propose a language extension to OpenMP for stream computing on multi-core architectures. To the best of our knowledge, there does not exist any source-to-source compiler that supports stream computing for heterogeneous parallel computing systems.

In this paper, we present HSTREAM, a compiler directive-based language extension that supports heterogeneous stream computing. HSTREAM aims to keep the same simplicity as programming with OpenMP and to enable programmers to easily utilize the available heterogeneous parallel computing resources on the host (CPU threads) and device (GPUs, or Intel Xeon Phi). The overview of the HSTREAM solution is depicted in Fig. 1. The HSTREAM source-to-source compiler

arXiv:1809.09387v1 [cs.PL] 25 Sep 2018

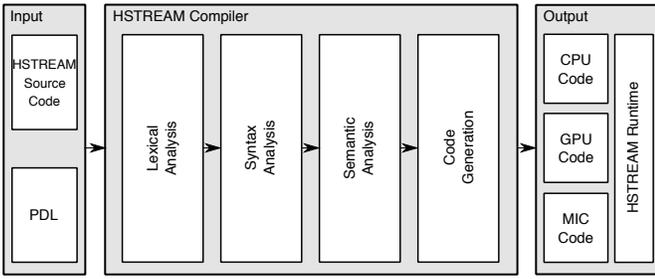


Fig. 1: Overview of the proposed approach.

performs several analysis steps (including lexical, syntactical, and semantical) and generates target specific code from a given source code annotated with HSTREAM compiler directives and a PDL file that describes the hardware architecture. HSTREAM supports code generation for multi-core CPUs using OpenMP, GPUs using CUDA, and Intel Xeon Phi (also known as MIC) using Intel Language Extensions for Offloading (LEO). The HSTREAM runtime is responsible for scheduling the workload across the heterogeneous PUs.

We use the HSTREAM source-to-source compiler to generate the heterogeneous version of the STREAM benchmark [9]. We evaluate the generated heterogeneous STREAM benchmark with respect to programming productivity and performance. The experimental results show that HSTREAM keeps the same simplicity as OpenMP, and the code generated for execution on heterogeneous systems delivers higher performance compared to CPUs-only and GPUs-only execution.

Major contributions of this paper include:

- HSTREAM compiler - a source-to-source compiler for generating target specific code from high-level directive-based annotated source code.
- HSTREAM runtime - a runtime system for scheduling the workload across various non-identical processing units.
- Evaluation of the usefulness of HSTREAM using applications from the STREAM and STREAM2 benchmarks.

The rest of the paper is structured as follows. The design, execution model, and implementation aspects of HSTREAM are described in Section II. Section III describes the experimental environment (including the system configuration, STREAM benchmark, and evaluation metrics) and experimental results. We compare and contrast our work with the current state-of-the-art in Section IV. Section V concludes this paper and provides information about future work.

II. HSTREAM: LANGUAGE EXTENSION TO SUPPORT HETEROGENEOUS STREAM COMPUTING

In this section, we describe the design, the execution model, and the implementation aspects of HSTREAM.

A. Design

OpenMP 4.5 supports offloading of computations to accelerators. However, a single loop is usually offloaded to a single device, and while one device is performing some computations, the other PUs (including host CPUs and accelerators)

remain idle. Distributing the data and computations of the same loop across multiple accelerating devices and host CPUs requires additional programming investment.

HSTREAM enables the automatic distribution of data and computations across different PUs. It enables programmers to easily exploit the available resources in heterogeneous parallel computing systems. The source-to-source code generation helps to reduce the programming time investment and errors that may come when explicitly handling communication and synchronization between various processing units.

While HSTREAM is designed for heterogeneous computing of stream applications, it can also be used for data-parallel applications. In the context of HSTREAM, data-parallel applications process data that is entirely available in memory, whereas stream computing process some data that is streamed from an I/O device. Streams are read in chunks, stored in local buffers first, processed by multiple heterogeneous PUs, and transformed back to a stream or stored somewhere in memory.

Listing 1 shows the syntax of the HSTREAM compiler directive, which starts with the `#pragma` keyword followed by the `hstream` keyword that stands for heterogeneous stream. Thereafter, multiple clauses can occur in any order. Details about each of the directive clauses are provided below.

Listing 1: An example of the HSTREAM compiler directive.

```

1 #pragma hstream in (...) out (...) inout (...) device (...)
2   scheduling (...)
3 {
4   //body
  }
```

In clause: The *in clause* is used to indicate the data that should be transferred to the accelerating devices for processing. The syntax for the *in clause* is inspired from the Intel LEO and looks as follows: `in(variable_ref [, variable_ref ...])`. The *variable_ref* can be a simple data type, array, or stream. For example, `in(a)` where *a* may be either a simple data type or an array, and `in(a : double)` where *a* is a stream of data of type double. The *in clause* may accept multiple variables of different types. For instance, in `in(a, b, c : int)`, the first variable (*a*) is an array, *b* is a scalar value, and *c* is a stream of integers. The *in clause* can be used multiple times within the same directive. For instance, `in(a, b) in(c : int)` is equivalent to the example above.

Out clause: The *out clause* is used to indicate the variables that need to be transferred from the accelerators to the host memory. The syntax for the *out clause* is similar to the *in clause* and looks as follows: `out(variable_ref [, variable_ref ...])`.

InOut clause: The *inout clause* is used to indicate the variables that need to be transferred to and back from the accelerating devices. The syntax looks as follows: `inout(variable_ref [, variable_ref ...])`. The *inout clause* combines the functionality of the *in* and *out clause*. For example, `inout(a)` has the same functionality as using the *in clause* and the *out clause* separately (`in(a) out(a)`).

Device clause: In comparison to the existing OpenMP *device clause*, which is used to specify only one accel-

erating device id as offloading target, the *device clause* of the HSTREAM language extension allows providing a list of PU ids that will collaboratively process the input data elements provided using the *in*, *out*, and *inout* clauses. The syntax of the *device clause* looks as follows: *device(device_id [, device_id ...])*. The following examples describe different scenarios of the use of *device clause*: (1) *device(*)* is the default value of the *device clause*, which means that all PUs should be used including the CPUs and accelerating devices; (2) *device(1)* means that only PU with id 1 should be used; (3) *device(0,1,2)* means that PUs with id 0, 1 and 2 will work together to process some data. Details for each PU (such as, id, number of cores, cache size, core frequency, and global memory) are extracted from the platform description file (.pdl), which needs to be provided as input to our compiler. The PDL is a platform description language for the explicit description of heterogeneous parallel computing systems [19].

Scheduling clause: The *scheduling clause* is used to determine the sizes of data chunks for each PU. There are different scenarios on how to use the *scheduling clause*: (1) Device specific distribution, where the programmer will explicitly set the chunk size for each processing unit. The syntax looks as follows, *scheduling(device_id : chunk_size [, device_id : chunk_size ...])*; (2) Uniform distribution, where the programmer explicitly sets a constant uniform distribution for all PUs. The syntax looks as follows, *scheduling(chunk_size)*; and (3) Automatic distribution, where the HSTREAM runtime system automatically determines the chunk sizes for each of the PU. The syntax for automatic distribution looks as follows, *scheduling(AUTO)*.

B. The execution model

Figure 2 depicts an overview of the HSTREAM execution model. The main components of our solution are, the data producer, data processor, and data store. The data producer reads the data in batches from an input stream (that can be a file, data coming from the network, ...). The data processor stores the batches locally and then applies a specified function to each data item. Once the data is processed (consumed), the output is sent to the third component (named data store) to either write the data to a file or print it.

Please note that this process can be overlapped, which means that while data producer is reading one batch, the data consumer can process another batch, and at the same time the data store can write another batch to a file. Figure 3 shows an example of overlapping the reading, processing, and writing of data. To achieve this result, we need three separate threads, T1 is responsible to read the data, T2 is responsible to initiate the data processing, and T3 writes the data. First, the data producer starts reading the first batch and then notifies T2 that there is data available for processing. While T2 reads data, T1 can continue reading the next batch. When T2 finishes processing the first batch, it will notify T3 that there is data ready to be written. Then, T3 will start writing the output data to a file. While T3 is writing the first batch, T2 may start processing

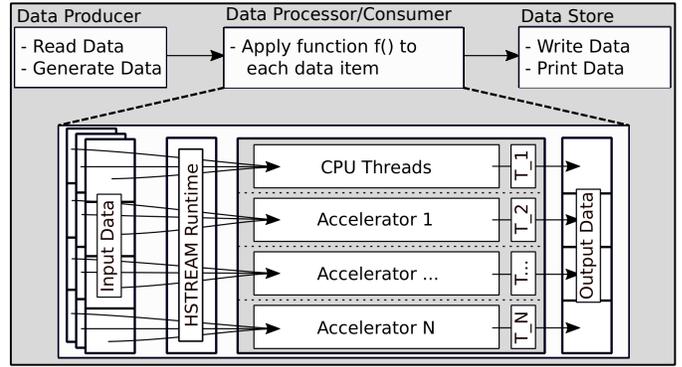


Fig. 2: The execution model of HSTREAM solution for stream computing in heterogeneous parallel computing systems.

the second batch (only if T1 has finished reading batch 2), and T1 can start reading batch 3.



Fig. 3: An example of overlapping data read, data process, and data write.

The data processor in Fig. 2 uses all the available PUs (or specific ones, depending on what the developer has provided in the scheduling clause) to process the input data. There is a separate CPU thread controlling each PU, including (remaining) CPU threads, GPUs, and Intel Xeon Phis. Each CPU thread (except the one that controls the CPU threads) is responsible to transfer the corresponding data chunk from the host local storage to the memory of the accelerating device. The data will be processed in the accelerating devices, and then transferred back to the host. The thread that controls the remaining CPU threads does not need to do any explicit data transfer, because the controlling thread and the processing threads share the same DRAM. Please note that the same overlapping strategy used for reading, processing, and writing data can be used within the data processing component, such that the data transfer (host-to-device, and device-to-host) is overlapped with the data processing.

C. Implementation

In this section, we will describe the tools and techniques used to implement our source-to-source compiler. Thereafter, throughout an example, we will describe the transformation process from a high-level C++ code annotated with HSTREAM compiler directives, to C++ code with OpenMP directives for execution on host CPUs, Intel LEO for execution on Intel Xeon Phi coprocessors, and CUDA for execution on GPU accelerators.

1) *Source-to-Source Compiler:* Figure 4 depicts an overview of the implementation steps, including the HSTREAM language definition, front-end, and back-end.

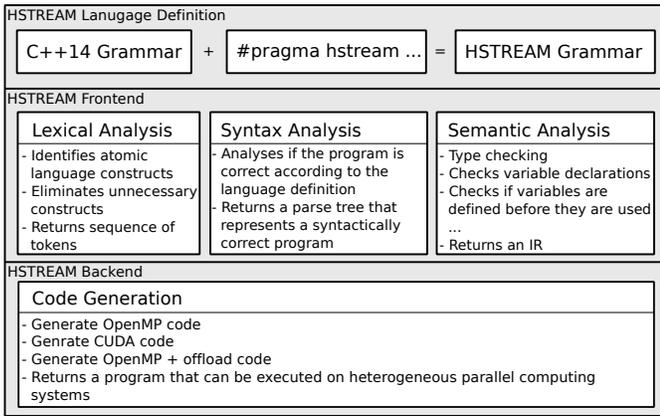


Fig. 4: Overview of the implementation of our source-to-source compiler.

We use ANTLR4 (Another Tool for Language Recognition) [14] to write the grammar for our compiler, that we call *HSTREAM grammar*. As a basis for our work, we used the C++14 grammar from ANTLR Project GitHub repository of a collection of ANTLR4 grammars [17]. We extend the C++14 grammar with the *HSTREAM* compiler directives to allow developers to annotate parts of the code that need to be executed in heterogeneous parallel computing systems. Listing 2 shows an excerpt of the *HSTREAM* grammar written in ANTLR4.

Listing 2: An excerpt of the *HSTREAM* Grammar.

```

1 directives : directive* ;
2 directive : PRAGMA HSTREAM clauses body ;
3 clauses   : clause* ;
4 clause   : inclause
5           | outclause
6           | inoutclause
7           | deviceclause
8           | schedulingclause
9           ;

```

The ANTLR tool is used to generate the parser, which is used for lexical, syntax, and semantic analysis. The lexical analyzer takes the input file, identifies atomic language constructs (with the help of regular expressions and pattern rules), removes any unnecessary constructs, such as comments and white spaces, and returns a sequence of tokens. In the case of illegal language constructs (tokens) errors will be generated.

The sequence of tokens is passed to the syntax analyzer (also known as the parser) to check (with the help of context-free grammars) if the input source code is correct according to the language definition. The parser generates a parse tree, which represents a syntactically correct program.

The generated parse tree is used by the semantic analyzer to judge whether the syntax derives any meaning. Some of the main tasks associated with the semantic analysis include type checking, scope resolution, checking for variable declarations, and checking whether variables are defined before they are used. If there are semantic errors, appropriate feedback will

be provided, such as type mismatch, undeclared variables, multiple variable declarations within a scope, and variable access is out of scope. For instance, an *HSTREAM* directive with two or more *scheduling* or *device* clauses is syntactically correct, but semantically wrong, because only one such clause can be provided. If no semantic errors are identified, then an intermediate representation (IR) of the source code will be created, which may be used for optimization and translation.

We do not perform any optimization of the code during source-to-source compilation, but we use the IR to generate the target-specific code. We use the String Template (ST) Library [15] to generate the target source code. A simple example of the ST library that is used to generate the CUDA memory transfer statements is shown in Listing 3.

Listing 3: An example of the string template library that is used to generate the CUDA memory transfer statements.

```

cuda_memcpy_host_to_device(from, to, type) ::=
  "cudaCheckError(cudaMemcpy($from$,
    d_$$to$, sizeof($type$)*myN,
    cudaMemcpyHostToDevice));"

```

Based on the type of the system architecture and based on the type of PUs provided in the *device clause* the corresponding functions will be generated. At the time of writing this paper, *HSTREAM* source-to-source compiler supports code generation for OpenMP, CUDA, and Intel LEO.

2) *Code transformation example*: In this section, throughout an example we describe the transformation process from high-level C++ code annotated with *HSTREAM* directives, to target specific source code. The source code for the *TRIAD* function written in C++ with *HSTREAM* annotations is shown in Listing 4. Line 1 shows the *HSTREAM* pragma directive, which includes the *hstream* keyword, the *in* and *out clause* for memory management, and the *device* and *scheduling clause* for scheduling of the workload. The *TRIAD* function takes three array variables and one scalar variable as input and outputs the result to one of the arrays. In this example, we want to use all available PUs (*device(*)*), and we want a uniform distribution of the workload for each of the PUs (*scheduling(4096)*).

Listing 4: *HSTREAM* TRIAD function of the STREAM benchmark.

```

1 #pragma hstream in(b,c,a,scalar) out(a) device(*)
   scheduling(4096)
2 {
3   a = b+scalar*c;
4 }

```

We assume that we have a heterogeneous platform that comprises CPUs, GPUs, and Intel Xeon Phi co-processor. Please note that we describe our hardware architecture using XML-based platform description language [19]. The *HSTREAM* compiler will generate three types of functions, each designed for execution on one of the PUs. For example, an OpenMP based function will be generated for execution on multi-core CPUs, a CUDA kernel will be generated for execution on

GPU devices, and another pragma-based function that uses Intel LEO for execution on Intel Xeon Phi.

Listing 5 shows the TRIAD code generated for execution on host CPUs.

Listing 5: TRIAD function of the STREAM benchmark designed for execution on multi-core CPUs.

```

1 #pragma omp parallel for
2 for (int i=start; i<finish; i++)
3 {
4     a[i] = b[i]+scalar*c[i];
5 }

```

The corresponding function for execution on the GPU accelerators is shown in Listing 6. Please note that Listing 6 shows only the CUDA kernel, whereas the code for memory management, such as allocation, the transfer from host to device and vice-versa, and the memory deallocation is handled by our runtime scheduler (shown in Algorithm 1).

Listing 6: TRIAD function of the STREAM benchmark designed for execution on GPUs.

```

1 __global__ void GPU_Triad( double *b, double *c, double *a,
2     double scalar, int len) {
3     int idx = threadIdx.x + blockIdx.x * blockDim.x;
4     if (idx < len)
5     {
6         a[idx] = b[idx]+ scalar*c[idx];
7     }
8 }

```

Listing 7 shows the generated source code for execution on the Intel Xeon Phi (also known as MIC), which corresponds to the TRIAD function from the STREAM benchmark.

Listing 7: Triad function of the STREAM benchmark designed for execution on Intel Xeon Phi.

```

1 #pragma offload target(mic: cpu_thread_id) in(a[my_start:
2     my_finish]) in(c[my_start:my_finish]) out(c[my_start:
3     my_finish])
4 {
5     #pragma omp parallel for
6     for (int i = my_start; i < my_finish; i++)
7     {
8         a[i] = b[i]+scalar*c[i];
9     }
10 }

```

The pseudo-code for the runtime scheduler, which is responsible for distributing the workload across the heterogeneous processing units, is shown in Algorithm 1. The generated runtime class has information for the hardware architecture, which is derived from the provided platform description file (.pdl). This class, together with the information provided in the *device clause*, are used to determine how many threads we need to engage. Since one thread controls a separate PU, we create as many threads as there are PUs. In the initialization step (see Line 1) we create an instance to the runtime class that has information about the system and create two shared variables that keep track of the current state (start and finish positions) of the processed data.

Each CPU thread controls a PU, and each thread is responsible to determine the *start* and *finish* index of its own chunk of data (see Line 5). To avoid multiple PUs processing the same amount data, we need to perform this process in a critical

section, which means that while one thread is determining its *start* and *finish* positions, no other thread can do the same. Furthermore, thread private variables of the *start* and *finish* variables are created, to enable other threads to pick other chunks of data while another one is processing its chunk of data.

Once the *start* and *finish* positions are determined, the data is ready for processing. Using a single *if-then-else* statement we check the type of the PU (see Line 6, 13, and 15) and perform the corresponding steps for each type. For instance, for GPU accelerated devices, we need explicit data management, such as device memory allocation, transferring data from host to device and vice-versa, and memory deallocation (see Line 7-12). Similarly, for Intel Xeon Phi accelerators, explicit data management is performed through the *in* and *out* clauses of the Intel LEO directive (see Listing 7 Line 1). For CPUs, there is no need for explicit data transfer because the controlling thread and the processing threads share the same DRAM (see Listing 5).

ALGORITHM 1: The algorithm used for runtime scheduling and workload distribution.

Data: List of PUs, workload

Result: Distribute the workload across these PUs and process it simultaneously

```

1 initialization;
2 create as many threads (T) as PUs;
3 foreach T do
4     while not reached the end of data do
5         determine start and finish positions;
6         if pu.type is GPU then
7             /* explicit data management is performed
8              using CUDA API calls */
9             create device variables;
10            allocate device memory;
11            copy chunks of data from host to device;
12            execute CUDA kernel ; // Example: Listing 6
13            copy data back to host memory;
14            free memory ;
15        else if pu.type is MIC then
16            /* explicit data management is done using
17             LEO in and out clauses */
18            offload data and control to MIC device ;
19            // Example: Listing 7
20        else if pu.type is CPU then
21            /* no explicit data management is needed
22             */
23            execute the CPU code ; // Example Listing 5

```

III. EVALUATION

In this section, we first describe the experimentation environment, including the hardware configuration, application benchmarks, and the data-set used for evaluation of our approach. Thereafter, we discuss the results of the study.

A. Experimentation environment

For evaluation of our approach, we have used applications from the industry standard STREAM and STREAM2 benchmarks [9]. We vary the stream size, chunk size, and

the number of available resources. Details about the system configuration, application benchmark, considered data-sets, and the measurement metrics will follow.

1) *System configuration*: We used our heterogeneous system named *DISA*, which comprises two Intel Xeon Gold CPUs, and four Quadro P4000 GPUs. Please note that in our experiments the CPU hyper-threading is disabled. Table I lists the details of our systems. For experimental evaluation, we vary the number of resources used in the *DISA* system, such as CPU only, 1GPU, 2GPUs, 3GPUs, 4GPUs, CPU+1GPU, CPU+2GPUs, CPU+3GPUs, and CPU+4GPUs.

TABLE I: The system configuration details for *DISA*.

Specs	Intel Xeon	NVIDIA GPU
Type	Gold 6148	Quadro P4000
Frequency (GHz)	2.4 - 3.7	1.48
# of Cores	20	1792
# of Threads	40	/
Cache (MB)	27.5	/
Memory (GB)	768	8
TDP (W)	150	105

2) *Benchmark application and data-set*: We use applications from the *STREAM* and *STREAM2* benchmarks to evaluate our approach. We use the *HSTREAM* source-to-source compiler to generate the heterogeneous version of the *STREAM* benchmark [9]. We use the *COPY*, *SCALE*, *SUM*, and *TRIAD* functions from the *STREAM* benchmark, and the *FILL* and *DAXPY* functions from the *STREAM2* benchmark. Details and information about the considered functions of the *STREAM* benchmark are available on-line ¹.

The stream array size is varied between 256, 512, 1024, 2048, 4096, and 8192MB to simulate various scenarios of the workload. Furthermore, when we split the workload among the available PUs, we vary the chunk sizes between 1, 2, 4, 8, 16, 32, and 64MB. The chunk size indicates the amount of data that should be sent for processing to a specific PU.

To address the variability of the results, we repeat each experiment 10 times and report the average value.

3) *Metrics*: In our experiments, we consider two aspects: (1) the programming productivity and (2) the performance. We measure the programming productivity with respect to the total lines of code and the lines of code that are specific to a programming framework required to parallelize the code. We use our tool, named *CodeStat* [11] to measure the programming productivity. With respect to the performance, we measure the throughput, which reflects the amount of data that can be processed within a time unit.

B. Results

In this section, we first describe the evaluation results with respect to programming productivity. Thereafter, we describe the results with respect to performance.

TABLE II: The programming effort expressed in lines of code (LOC) required to program the *STREAM* benchmark using different programming frameworks.

	Total	OpenMP	CUDA	Other
sequential (C)	131	0	0	0
multi-core (OpenMP)	214	8	0	0
accelerator (CUDA)	190	0	55	0
multi-core + acc (HSTREAM)	210	8	0	57
HSTREAM generated	1195	69	131	0

1) *Programming productivity*: Table II shows the programming productivity expressed in lines of code (LOC). We consider the total LOC, and LOC specific to a programming language, such as OpenMP, and CUDA. For *HSTREAM* input we also show the LOC required to create the platform description file. We compare different versions of implementations of the *STREAM* benchmark, including sequential, multi-core version using OpenMP, accelerated version using CUDA, and the heterogeneous version using *HSTREAM*. The last row shows the LOC of the *HSTREAM* generated version.

We may observe that the multi-core version of *STREAM* benchmark requires 8 OpenMP specific lines of code to be added, whereas the accelerated version requires about 55 CUDA specific lines of code. To parallelize the *STREAM* benchmark with *HSTREAM* language extension there are needed only 8 *HSTREAM* specific LOC, exactly the same as with OpenMP. For *HSTREAM* the developer needs to provide the PDL file as well, the LOC for which depends on the system. For the *DISA* system, the PDL file has 57 lines of code. The generated code, which is what a programmer would need to write manually for execution on heterogeneous systems that comprise host CPUs, and accelerating devices such as GPUs and Intel Xeon Phi, requires in total 1195 LOC, of which 69 are OpenMP and 131 are CUDA specific lines of code.

We may conclude that *HSTREAM* maintains the same level of programming complexity as OpenMP, at the cost of providing an XML-based description of the platform.

2) *Performance*: Figure 5 depicts the throughput (MB/s) of the selected functions from the *STREAM* and *STREAM2* benchmark, including *COPY*, *SCALE*, *ADD*, *TRIAD*, *FILL*, and *DAXPY*. We vary the stream size between 256, 512, 1024, 2048, 4096, and 8192MB. From our experiments, we have observed that the stream size does not impact the throughput, therefore we show the results only for the largest stream size (that is 8192MB). We also vary the chunk size (CS) between 1, 2, 4, 8, 16, 32, and 64MB. We vary the number of processing units engaged for computation, such as CPU only, 1GPU, 2GPUs, 3GPUs, 4GPUs, CPU+1GPU, CPU+2GPUs, CPU+3GPUs, and CPU+4GPUs. However, due to space limitation, we only show results for CS 2, 8, and 32MB, and system configurations CPU, 4GPUs, and CPU+4GPUs.

We may observe that in all cases, the execution that uses all available resources (CPU+4GPUs) results with the highest throughput. With respect to the chunk size, we may observe

¹<http://www.cs.virginia.edu/stream/>

that the throughput depends on the type of processing units engaged in the execution, and the type of computations. For example, the more complex functions such as TRIAD and DAXPY benefit more when executed in accelerated devices, whereas simpler functions such as FILL and ADD perform better on CPUs. The highest throughput observed in our experiments is when executing the DAXPY function on CPU+4GPUs with a chunk size of 32MB.

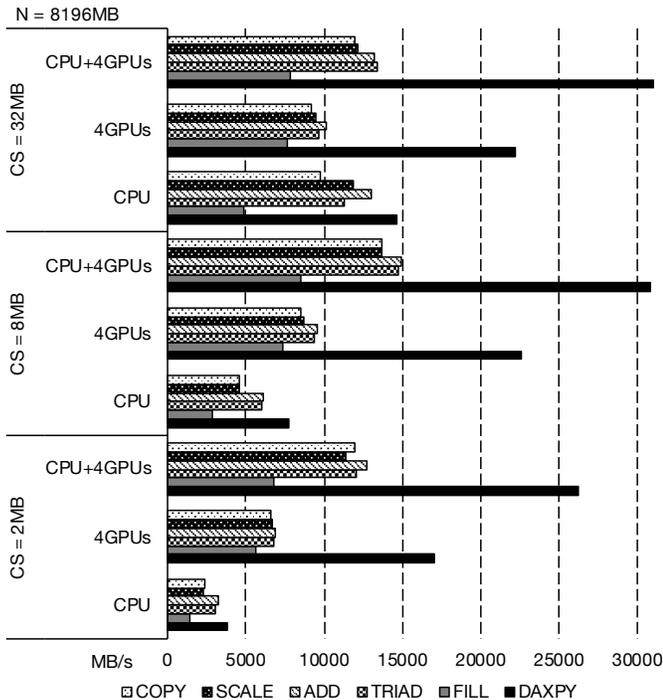


Fig. 5: The performance of heterogeneous STREAM benchmark when using host CPUs only, 4GPUs, and CPU+4GPUs on our DISA system. We vary the chunk size between 2, 8, and 32MB. The input size used in this experiment is 8192MB.

IV. RELATED WORK

Recent versions of OpenMP and OpenACC support offloading computations to accelerating devices [21]. Offloading to a single accelerating device can be easily achieved by using simple compiler-directives, whereas offloading to multiple devices requires additional programming effort. HSTREAM enables developers to use high-level compiler directives (similar to OpenMP and OpenACC) to develop applications that can be executed on multiple heterogeneous devices, such as multi-core CPUs, NVIDIA GPUs, and Intel Xeon Phi accelerators.

Zhang and Mueller [24] present GStream, a framework for data-streaming on systems accelerated with GPUs. GStream provides an application interface that software developers can use to express the parallelism of their streaming application without explicitly writing MPI messaging or CUDA memory copy instructions. Similar to our approach, the aim is to reduce the development time. Rather than exposing system developers to a new application programming interface, we have decided

to introduce a new simple OpenMP like directive that developers can use to parallelize their streaming applications.

Yan et al. [22] propose a language extension to OpenMP for data-parallel processing. Their approach distributes data and computations of parallel loops across multiple PUs that may be homogeneous or heterogeneous. In comparison, we describe an OpenMP language extension for stream computing.

Pop and Cohen [16] propose a stream computing extension to the OpenMP programming model. Their extension decomposes programs into tasks and in an explicit manner provides instructions on how data should flow among these tasks. In comparison to their work, which targets multi-core architectures, we provide stream computing support for heterogeneous parallel computing systems that may comprise multiple multi-core processors and many-core accelerating devices.

Del Rio Astorga et al. [5] and Ernstsson et al. [7] propose pattern based high-level application programming interface for stream and data computing on heterogeneous parallel computing systems. While template and skeleton based libraries may be helpful for generic applications, they are not recommended in case programmers want to use data-structures or algorithms optimized for a particular type of problem.

Zhang et al. [23] propose an auto-tuning approach for stream applications running on systems accelerated with Intel Xeon Phi. The authors exploit the pipeline parallelism through temporal sharing, which means that computations are overlapped with communication (data transfer from host to device, and vice-versa). We employ similar techniques in our solution, but in comparison we support heterogeneous systems accelerated with GPUs as well.

With the aim to alleviate the programming of heterogeneous systems, several source-to-source compilers [2, 8, 1, 18] are proposed that can generate target specific code from a high-level representation. In comparison, our approach targets streaming applications. Furthermore, rather than introducing a new programming language, our solution extends a well-established programming model, such as OpenMP, to enable acceleration of streaming applications.

V. CONCLUSION AND FUTURE WORK

We presented our HSTREAM language extension to support stream computing on heterogeneous parallel computing systems. HSTREAM source-to-source compiler can automatically generate device-specific code, such as OpenMP for CPUs, CUDA for GPUs, and Intel Language Extension for Offloading for Intel Xeon Phi, from a high-level source code annotated with OpenMP-like compiler directives. The HSTREAM runtime is responsible to distribute the workload across the engaged processing units accordingly. We have evaluated the usefulness of our HSTREAM solution for stream computing in heterogeneous parallel computing systems with the STREAM benchmark. We have observed that while HSTREAM keeps the same programming simplicity as OpenMP, the generated code outperforms the CPU and GPU only versions of the code.

Future work may focus on extending the HSTREAM runtime to support dynamic and adaptive workload scheduling.

Furthermore, we aim to extend the source-to-source compiler to support additional accelerating devices (such as FPGAs) and programming frameworks (such as OpenCL).

REFERENCES

- [1] Jason Ansel, Cy Chan, Yee Lok Wong, Marek Olszewski, Qin Zhao, Alan Edelman, and Saman Amarasinghe. *PetaBricks: a language and compiler for algorithmic choice*, volume 44. ACM, 2009.
- [2] Thomas Henry Beach and Nicholas J Avis. An intelligent semi-automatic application porting system for application accelerators. In *Proceedings of the combined workshops on UnConventional high performance computing workshop plus memory access workshop*, pages 7–10. ACM, 2009.
- [3] Siegfried Benkner, Sabri Pllana, Jesper Larsson Träff, Philippas Tsigas, Andrew Richards, Raymond Namyst, Beverly Bachmayer, Christoph Kessler, David Moloney, and Peter Sanders. The PEPPER Approach to Programmability and Performance Portability for Heterogeneous many-core Architectures. In *ParCo*, 2011.
- [4] Paweł Czarnul. Benchmarking performance of a hybrid intel xeon/xeon phi system for parallel computation of similarity measures between large vectors. *International Journal of Parallel Prog.*, 45(5):1091–1107, Oct 2017. ISSN 1573-7640. doi: 10.1007/s10766-016-0455-0.
- [5] David Del Rio Astorga, Manuel F Dolz, Javier Fernández, and J Daniel García. A generic parallel pattern interface for stream and data processing. *Concurrency and Computation: Practice and Experience*, 29(24), 2017.
- [6] Ciprian Dobre and Fatos Xhafa. Parallel programming paradigms and frameworks in big data era. *Int. J. Parallel Program.*, 42(5):710–738, 10 2014. ISSN 0885-7458. doi: 10.1007/s10766-013-0272-7.
- [7] August Ernstsson, Lu Li, and Christoph Kessler. Skepu2: Flexible and type-safe skeleton programming for heterogeneous parallel systems. *International Journal of Parallel Programming*, 46(1):62–80, Feb 2018. ISSN 1573-7640. doi: 10.1007/s10766-017-0490-5.
- [8] Alcides Fonseca and Bruno Cabral. AeminiumGPU: An Intelligent Framework for GPU Programming. In *Facing the Multicore-Challenge III*, pages 96–107. Springer, 2013.
- [9] John D. McCalpin. Memory bandwidth and machine balance in current high performance computers. *IEEE Computer Society Technical Committee on Computer Architecture (TCCA) Newsletter*, pages 19–25, 12 1995.
- [10] Sueb Memeti and Sabri Pllana. Analyzing large-scale DNA Sequences on Multi-core Architectures. In *Computational Science and Engineering (CSE), 18th International Conference on*, pages 208–215. IEEE, 2015.
- [11] Sueb Memeti, Lu Li, Sabri Pllana, Joanna Kolodziej, and Christoph Kessler. Benchmarking OpenCL, OpenACC, OpenMP, and CUDA: Programming Productivity, Performance, and Energy Consumption. In *Proceedings of the 2017 Workshop on Adaptive Resource Management and Scheduling for Cloud Computing, ARMS-CC '17*, pages 1–6, New York, NY, USA, 2017. ACM. ISBN 978-1-4503-5116-4. doi: 10.1145/3110355.3110356.
- [12] Sueb Memeti, Sabri Pllana, Alécio Binotto, Joanna Kolodziej, and Ivona Brandic. Using meta-heuristics and machine learning for software optimization of parallel computing systems: a systematic literature review. *Computing*, Apr 2018. ISSN 1436-5057. doi: 10.1007/s00607-018-0614-9.
- [13] Eric Mizell and Roger Biery. *Introduction to GPUs for Data Analytics*. O’Reilly, 2017.
- [14] Terence Parr. *The definitive ANTLR 4 reference*. Pragmatic Bookshelf, 2013.
- [15] Terence John Parr. Enforcing strict model-view separation in template engines. In *Proceedings of the 13th international conference on World Wide Web*, pages 224–233. ACM, 2004.
- [16] Antoniu Pop and Albert Cohen. A stream-computing extension to openmp. In *Proceedings of the 6th International Conference on High Performance and Embedded Architectures and Compilers*, pages 5–14. ACM, 2011.
- [17] Antlr Project. Grammars written for ANTLR v4. <https://github.com/antlr/grammars-v4/wiki>, 2018. [Online; accessed 06-June-2018].
- [18] Christopher J Rossbach, Yuan Yu, Jon Currey, Jean-Philippe Martin, and Dennis Fetterly. Dandelion: a compiler and runtime for heterogeneous systems. In *Proceedings of the Twenty-Fourth ACM Symposium on Operating Systems Principles*, pages 49–68. ACM, 2013.
- [19] Martin Sandrieser, Siegfried Benkner, and Sabri Pllana. Explicit platform descriptions for heterogeneous many-core architectures. In *Parallel and Distributed Processing Workshops and Phd Forum (IPDPSW), 2011 IEEE International Symposium on*, pages 1292–1299. IEEE, 2011.
- [20] TOP500. TOP500 Supercomputer Sites. <http://www.top500.org/>, 2016. Accessed: July, 2018.
- [21] R. Xu, S. Chandrasekaran, and B. Chapman. Exploring programming multi-gpus using openmp and openacc-based hybrid model. In *2013 IEEE International Symposium on Parallel Distributed Processing, Workshops and Phd Forum*, pages 1169–1176, May 2013. doi: 10.1109/IPDPSW.2013.263.
- [22] Yonghong Yan, Jiawen Liu, Kirk W Cameron, and Mariam Umar. Homp: Automated distribution of parallel loops and data in highly parallel accelerator-based systems. In *Parallel and Distributed Processing Symposium (IPDPS), International*, pages 788–798. IEEE, 2017.
- [23] Peng Zhang, Jianbin Fang, Tao Tang, Canqun Yang, and Zheng Wang. Auto-tuning streamed applications on intel xeon phi. 2018.
- [24] Yongpeng Zhang and Frank Mueller. Gstream: A general-purpose data streaming framework on gpu clusters. In *Parallel Processing (ICPP), 2011 International Conference on*, pages 245–254. IEEE, 2011.