

Parallel Programming in Actor-Based Applications via OpenCL

Paul Harvey and Kristian Hentschel
School of Computing Science
University of Glasgow
Glasgow, U.K.
p.harvey.1@research.gla.ac.uk,
1003734H@student.gla.ac.uk

Joseph Sventek
Dept of Computer and Information Science
University of Oregon
Eugene, OR, USA
jsventek@uoregon.edu

ABSTRACT

GPU and multicore hardware architectures are commonly used in many different application areas to accelerate problem solutions relative to single CPU architectures. The typical approach to accessing these hardware architectures requires embedding logic into the programming language used to construct the application; the two primary forms of embedding are: calls to API routines to access the concurrent functionality, or pragmas providing concurrency hints to a language compiler such that particular blocks of code are targeted to the concurrent functionality. The former approach is verbose and semantically bankrupt, while the success of the latter approach is restricted to simple, static uses of the functionality.

This paper presents an extension to an existing actor-based programming model and runtime to support executing applications on parallel hardware architectures. Besides the glove-like fit of a kernel to the actor abstraction, quantitative code analysis shows that actor-based kernels are always significantly simpler than API-based coding, and generally simpler than pragma-based coding. The structuring of applications in this manner, enables the runtime to automate the initialisation and interaction with these parallel hardware platforms. Performance measurements show that the overheads of actor-based kernels are commensurate to API-based kernels, and range from equivalent to vastly improved for pragma-based annotations, both for sample and real-world applications.

Categories and Subject Descriptors

D.2.8 [Software Engineering]: Metrics—*performance measures*

General Terms

LANGUAGES, MANAGEMENT, PERFORMANCE

Keywords

parallel programming, actors, performance, OpenCL, middleware

1. INTRODUCTION

Due to power consumption, heat dissipation, and clock propagation limits, modern hardware architectures are now designed with many, concurrent processing elements, as opposed to single processing elements with increasing clock rates; examples of such architectures include GPUs and multicore CPUs. These hardware platforms are designed to provide the user with multiple physical threads of execution, thus enabling many computations to occur simultaneously.

Software threads have traditionally been used to enable parallel execution on CPU architectures. However, due to the different nature of GPU hardware architectures, a number of different programming techniques are used. OpenCL is a standardised programming framework available for the main GPU vendors (NVIDIA and AMD), as well as other parallel hardware architectures.

While the OpenCL API enables access to these architectures, there are three main limitations. Firstly, the user is required to write large amounts of boilerplate code to create the OpenCL environment for a particular calculation. Secondly, the programming style requires explicit data movement between the host CPU and the OpenCL device; this requires manually flattening multi-dimensional arrays and structures of non-primitive types. Thirdly, the language and style used to program the device is often different from the programming language being used on the host. A similar argument can be made against the CUDA framework; since CUDA is only available on NVIDIA hardware, this work is focused on the more broadly applicable OpenCL.

OpenACC is a pragma-based approach to concurrent programming, where a developer explicitly annotates sections of code to be parallelised, as well as the data which should be moved between host and device. While this approach abstracts much of the boilerplate code of OpenCL, applying the simple annotations to single threaded code does not guarantee good performance and is not effective for all application classes.

This work describes the application of the actor-model of computation to kernel-based concurrency. In the actor-model

of computation, applications are composed of shared-nothing loci of computation which interact via explicit message passing. This is a strong correlation to the model used in kernel-based applications. We hypothesise that moving from low-level C code to a concurrent, shared-nothing, high-level actor programming model simplifies the use of OpenCL by providing appropriate structuring, thus enabling greater access to high-performance and heterogeneous computing. This structuring enables the runtime to automate the initialisation and communication with parallel hardware architectures. We demonstrate that applications written using actors exhibit less complexity via quantitative metrics compared to handwritten OpenCL in C (C-OpenCL), and that these applications run efficiently on different hardware platforms, with low overhead when compared to C-OpenCL and equivalent or better performance to OpenACC annotated C (C-OpenACC). This is shown for a number of different types of application, including a real-world document ranking example.

The remainder of the paper is arranged as follows: Section 2 describes the OpenCL framework, Section 3 gives an overview of the state of the art, Sections 4 and 5 describe the language and runtime which were used in this work, respectively. Section 6 discusses the extensions which were made to support and abstract the OpenCL framework, Section 7 describes experimental results, and Section 8 summarises the work and discusses future directions.

2. OPENCL

OpenCL is a programming framework for heterogeneous and parallel computing. It is standardised and is managed by the Khronos working group¹. In OpenCL, users are required to think in terms of host and device code, where a host is a coordinator application on the CPU, and a device is an *accelerator*. An accelerator may be a CPU, GPU, FPGA, or co-processor such as the Xeon Phi [12].

2.1 OpenCL Configuration

In OpenCL, the host is tasked with setting up, dispatching, and collecting results from a device. OpenCL is accessed through an API, which enables relatively low-level access to data types and functions in order to program and interact with one or more accelerators.

Creating an OpenCL environment consists of first querying the hardware at runtime to determine the available vendor *platforms* and the *devices* available in each platform. Platforms are essentially drivers provided by the hardware vendor, and the devices represent the actual accelerators. Then, a **context** must be created. A context is an umbrella structure that holds the device(s) to be used, as well as other runtime software constructs. A **command_queue** is then associated with each device and placed within the context. A **command_queue** is used to issue commands to a device. Commands include device queries, memory management operations, and kernel (Section 2.2) invocations. After this, a user creates a **program** with the kernel source file, and compiles it at runtime. The specific function to be executed within the

```

1  __kernel void square(__global float* input,
2                      __global float* output,
3                      const unsigned int count){
4      int i = get_global_id(0);
5      if(i < count)
6          output[i] = input[i] * input[i];
7  }
```

Listing 1: OpenCL Kernel to Compute the Square of An Input Array

compiled source is then used to create the **kernel** object. At this point the OpenCL environment has been constructed.

From here the user allocates memory on the device and then copies host data into this memory. The device memory is then associated with the correct position in the kernel arguments. Then, the number of dimensions upon which the kernel should work is calculated, and the kernel is launched on the device, with this information, via the **command_queue**. Usually, the host then blocks attempting to read data back from the device once it has finished its computation. Once all computation is complete and the device is no longer required, there are appropriate destructor functions. The device itself is treated simply as a functional unit. Data and code are passed to the device, the device executes this code, and the results are read back by the host.

2.2 Kernels

A device runs a special piece of code known as a *kernel*. An OpenCL kernel is written in a C-like syntax and represents the logic of a single thread. The number and groupings of threads are supplied during the configuration stage on the host. These values are known as the **local** and **global** work-sizes, and are used to optimise the allocation of threads to the underlying hardware for a given dataset. Within a kernel, the currently executing thread may be identified via the API. This can be used to customise application logic. The kernel is expressed as a function with parameters. Information for the actual computation is passed to this function as arguments by the host.

The OpenCL model uses a memory hierarchy in which memory is split into **global**, **local**, **private**, and **constant** regions. This is a direct mapping to the hardware configuration of memory found in GPUs, however the same model is applied to all hardware devices. Global memory is shared amongst all threads, local memory is shared between a specified group of threads, and private memory is specific to a thread. Global and local memory are subject to unsynchronised modifications, although there are mechanisms to synchronise access. Constant memory is shared by all threads, but is read only. Listing 1 shows a simple kernel.

3. PARALLEL PROGRAMMING FRAMEWORKS

There are a number of different techniques which are used to program GPUs. These can be split into three different equivalence classes: APIs, fully-automated parallelisation, and semi-automated parallelisation.

3.1 API Approach

¹<https://www.khronos.org/opencl/> - Accessed 5 January 2015

The most low level approach to programming an accelerator is via an API within an existing language. Examples of this include Python [13], Java², and the original implementations in C/C++ of OpenCL [1] and CUDA [17].

While an API is a simple approach, requiring no modifications to the host language, the general drawback of using an API is the need to write large amounts of boilerplate code simply to set up the programming environment, as described in Section 2.1. This boilerplate code often follows the steps required to setup and initialise the relevant framework, and can have very little relation to the programming idioms of the host language - e.g., the Java API requires the use of pointer objects. Similarly, the code required to express the calculations on the accelerator is written in a C-like language, which is embedded in the host language as a string. For non-C-like languages, such as Java and Python, this can be challenging for programmers without experience in C/C++. More generally, this leads to two different programming styles for the host and the accelerator, although it has the potential for the best performance.

3.2 Automated Approach

A different approach is to hide from the developer all the low-level details such as initialisation, data allocation and movement, the specification of work sizes, and when to dispatch kernels. A number of domain-specific languages have been developed to provide such higher level abstractions. Examples of these include **Accelerator** [20], **LIME** [8], **SkelCL** [18], and **Chestnut** [19].

With the exception of Accelerator, which pre-dates OpenCL and CUDA and is no longer under development, the primitives of these languages abstract data movement and computation through operations such as map, reduce, stencil, and other vector operations. In general, the effectiveness of the approach relies on the maturity of the source language, and the quality of the code transformation to the underlying transformation. For example, Chestnut does not support data structures, SkelCL requires stack operations to assign variables to a kernel, and LIME only allows static topologies of tasks and very strongly immutable state.

3.3 Semi-Automated Approach

Much work has been done on semi-automated translation of serial programs to parallel OpenCL or CUDA code. The techniques used include recognition of common parallelisable patterns in the source code, such as nested loops that can be unrolled and executed in parallel. However, most rely on the programmer to provide annotations or some form of refactoring applied to the original code.

AMD's open-source **Aparapi**³ system allows partial offloading of Java code at runtime. Aparapi relies on the programmer to refactor a function or loop into an inner class with a run method containing the computation. The Java bytecode for this method is translated to OpenCL/C code at runtime. The programming model is similar to our work as it requires some refactoring but still allows the kernel code to be written in a subset of the original language, making use

of the primitives (such as classes) provided by that language to express the required meta-information. It also abstracts, to a certain extent, the exact memory layout and management of data movement and provides automatic adjustment to the available devices, rather than having the programmer optimize the code for a specific device. Aparapi offers a trade-off between a simplified programming model, and non-trivial performance cost [7].

OpenACC [22] is a set of explicit annotations and API for C or Fortran code, which enables annotated code to execute on GPUs. Separate **OpenMP** [3] annotations are required for CPUs. With OpenACC, the user must explicitly annotate data to be moved between host and device, as well as sections of code which should be parallelised. Using these hints, the compiler attempts to generate optimised code for the supported patterns. The system enables the use of existing code with fewer modifications than is required to use the OpenCL API. The PGI⁴ compiler provides the best performance, but is proprietary. There are also open-source compilers, with varying quality and reliability. **hiCuda** [10] is a similar approach for the CUDA framework, applying annotations to sequential C.

By their nature, annotations are extraneous to the logic of the underlying application. This has the advantage that existing code can be used as a starting point, with annotations extending or enhancing the functionality. However, this means that there is no intuition as to what is happening *under the hood*. Annotations must be applied to *each* construct to be parallelised, resulting in applications which are increasingly difficult to follow as their size increases. Also, there is no guarantee that the compiler will be able to generate an effective parallel strategy for the annotated section of code. For example, if there is a non-linear data dependency in a for loop, sequential code may be generated instead of parallel. Thus, the approach is unsuitable for all application classes, requiring code to be refactored or rewritten, negating a main advantage of the approach.

The C++ actor framework (**CAF** [4]) is a C++ library providing actor-based concurrency. Using layers of C++ templates, developers perform template metaprogramming to create actors and explicitly send messages between them. The author's goals are similar to this work: using an actor-based abstraction to simplify programming of concurrent systems. Specifically, CAF uses the actor abstraction from their library to abstract OpenCL kernels. Actors are represented by lightweight threads, and the framework supports its own scheduling system. The framework has shown good performance compared to other popular approaches of programming such systems in terms of performance and memory consumption.

By creating a framework in C++, CAF enjoys good performance, and easy integration with existing projects. However, despite the good performance results, there are a number of limitations of this approach compared to this work. Firstly, the use of template metaprogramming is required to take advantage of the popularity of C++, however its use is a non-trivial exercise [16], acting as a barrier to non-

²<http://www.jocl.org/> - Accessed 5 January 2015

³code.google.com/p/aparapi - Accessed 5 January 2015

⁴<http://www.pgroup.com/> - Accessed 5 January 2015

```

1 type Isnd is interface(out integer output)
2 type Ircv is interface(in integer input)
3
4 stage home{
5
6     actor snd presents Isnd {
7         value = 1;
8         constructor() {}
9         behaviour {
10            send value on output;
11            value := value + 1;
12        }
13    }
14
15    actor rcv presents Ircv {
16        constructor() {}
17        behaviour {
18            receive data from input;
19            printString("\nreceived:␣");
20            printInt(data);
21        }
22    }
23
24    boot{
25        s = new snd();
26        r = new rcv();
27        connect s.output to r.input;
28    }
29 }

```

Listing 2: Simple Ensemble Send and Receive Example

expert programmers. Like this work, CAF supports an actor-based abstraction of OpenCL kernels for programming parallel hardware architectures, although unlike this work, CAF specifies kernels in terms of C strings. This means that they do not support kernels containing multi-dimensional arrays or nested references within structures, requiring developers to marshal and demarshal such data types themselves. Also, given the abstraction of the memory allocation and data movement used, it is not possible to leave data on a parallel device: a common optimisation used to greatly reduce execution time. Furthermore, as kernels are compiled by a separate compiler at runtime, error messages are less meaningful and delayed, unlike this work.

4. THE ENSEMBLE LANGUAGE

Ensemble is an actor-based programming language. It is an evolution of the Insense programming language [6] which was originally created to ease the development of software for wireless sensor networks. Unlike traditional actor systems, Ensemble is imperative, not functional. An Ensemble actor has its own private state and a single thread of control. The single thread of control is expressed as a `behaviour` clause. The code within this clause is repeated until explicitly told to stop. All actors execute within a `stage`. Stages represent a memory space, thus many stages may exist per physical machine. Ensemble supports reference counted garbage collection.

Actors have *shared-nothing* semantics, meaning that they share no state. Message passing along uni-directional, typed channels is used to share information between actors. Each channel may have an optional buffer to enable asynchrony during communication. When no buffer is specified or a

buffer is full, the system reverts to synchronous, blocking communication. Channels are either statically associated with a specific actor via an interface, or dynamically created and composed at runtime into one or more 1-1, 1-n, or n-1 topologies.

In actor-based languages, data sent along channels must be duplicated to ensure that shared-nothing semantics are preserved. One novel feature of Ensemble is the ability to declare data *movable* (`mov`). By marking heap allocated memory as movable, only a reference will be sent along a channel, as opposed to a duplicate. Movability is included within the type system of the language, using inter-procedural analysis at compiletime to ensure that once a reference is sent it is not accessed again until it is assigned to, otherwise a compiletime error is generated. Movable types are a simplified form of ownership types [5]. This feature was originally developed to reduce heap memory consumption and fragmentation in small embedded devices, but has been used in this work to enable runtime optimisations of data movement, Section 6.2.3. Listing 2 shows a simple Ensemble program which defines, instantiates and connects two actors, one of which sends linearly increasing values to the other.

Ensemble applications are compiled to Java source code, and then Java class files, Figure 1. These class files are then modified for use with the runtime. This process removes the constant pool and metadata in the class files to save space; the modification process resolves all inter-class references. Although this work has been integrated with Ensemble, applications for the system need not be written in Ensemble. Developers can either create actor-based applications in plain Java, bypassing Ensemble, or generate Java source code from existing systems, as is done in Salsa [21].

5. THE ENSEMBLE RUNTIME

Once compiled to augmented class files, Ensemble application are executed by a virtual machine (VM). The Ensemble VM is a modified version of the JVM, originally designed for programming wireless sensor networks [2]. It supports an actor-based model of computation directly; rather than threads with shared state, the system supports encapsulated actors with explicit message passing. Each VM represents a *stage*.

The VM either executes as a process on Linux-based systems, or as the sole application on InceOS-based systems. InceOS [11] is a language specific operating system (OS) for embedded devices. Within the VM, a thread is created for each actor. When the VM is executing on Linux-based systems, the `Pthread` library [15] is used. When the runtime is executing on InceOS, the InceOS thread library is used. Each actor's thread interprets the bytecode representation of an actor's behaviour clause. The thread executes the interpreter in an infinite loop until explicitly told to `stop` in the language, at which point the actor and its resources are garbage collected, and the thread exits.

Actor scheduling is dictated by inter-actor communication. Should there be no communication, the system reverts to a pre-emptive round robin scheduler.

Currently, hardware interaction is achieved via *system ac-*

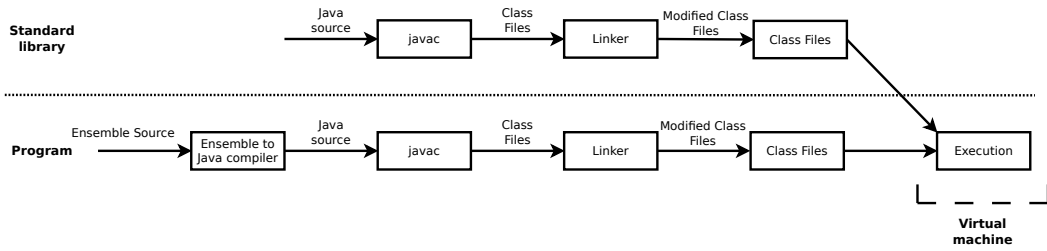


Figure 1: Ensemble Compilation Process

tors. These actors are written in C and provided by the runtime to explicitly access a range of hardware, such as timers, sensors, and explicit networking. User actors interact with these actors via channels, in the same manner as other actors.

6. OPENCCL THROUGH ACTORS

The shared-nothing semantics required by the actor model, and the use of explicit message passing for such actors to communicate, provides the perfect structure to facilitate kernel-based concurrency. Consequently, the existing actor-based framework was extended to enable the initialisation, programming, communication with, and use of parallel hardware platforms via OpenCL. This provides a high-level, single-source solution which leverages OpenCL’s existing low-level, yet portable, split-source infrastructure.

6.1 Language Integration

The similarity between the isolated nature of computation in actor-based programming and the OpenCL model is strong. Both require explicit data movement between loci of computation. In Ensemble, the channel mechanism provides a natural way to facilitate this without the need to go *outside the language* or obfuscate existing code; consequently, OpenCL kernels are represented as actors. The requirement for separate programming of host and device also matches the distinction between actors. The full example of matrix multiplication in the language shown in Listing 3 is referred to throughout this section. The green colour represents additions to the language, and light yellow shows compiler-enforced structure. These highlighted sections represent the only changes to the language. All existing application are compatible with these extensions.

6.1.1 Language Model

In the language, an actor is marked as being a kernel by adding the `openccl` keyword to its definition (line 22). This tells the compiler that the actor’s interface should only contain a single channel, and that a slightly different structure is expected within the behaviour clause. In this way, kernels are integrated into the language, not a separate source file in a different language.

The single channel conveys an `openccl struct` type as defined by the developer (lines 6 - 11). This is a normal `struct` except that the `openccl` keyword tells the compiler that the fields of the `struct` should contain two integer arrays, an `in` channel, and an `out` channel. The arrays convey the local and global worksizes required to dispatch the OpenCL kernel, as discussed in Section 2. The two channels are the

input and output channels for the data to and from the kernel, respectively. The channels are created and configured in the host actor (lines 48 - 52), sent to the kernel actor via the input channel (line 57), and then the data is sent (line 58). The host then waits for data from the actor (line 59). This is an example of the dynamic nature of channels, and enables any type of data to be conveyed to and from the kernel actor safely, without requiring any extra compiler analysis, beyond normal language processing. The use of the channels in this way enables multiple communication patterns.

Within the behaviour loop of the kernel actor it is required that the first two statements are receive statements (lines 26 - 27). The first statement receives an `openccl struct` instance. This enables the runtime to prepare the kernel for dispatch to the device with the appropriate dimensions by using the `worksize` and `groupsize` values. The second statement is used to receive the data that the kernel will process. It is also required that the last statement is a send statement (line 38). This is used to send the processed data onwards. The send statement is also used as a marker, with all statements between the second receive and send statements representing the OpenCL kernel (lines 28 - 37). The standard set of OpenCL calls natively available in a kernel are also available within the OpenCL behaviour clause between these markers (lines 28 - 30), including the math functions. Also, the `global`, `local`, and `private` memory modifiers are available, enabling the developer to take full advantage of the different memory regions which are available in the OpenCL model to improve performance. It is important to note that these are only available within the kernel actor, maintaining the actor memory semantics. Like normal OpenCL, print statements are not allowed in kernels within Ensemble.

One of the key advantages of embedding OpenCL within the language is the ability to preserve multi-dimensional array and structure dereferencing. Currently, OpenCL requires that arrays and structures containing pointers are flattened when being passed to a kernel. In Ensemble, this process is automated, leaving the users with normal dereferencing of such data types within the kernel. This also has the advantage of providing the user with warnings and errors for kernels at compiletime, rather than having to wait until runtime kernel compilation. Furthermore, as all the OpenCL actors are connected by channels, should the user wish to change the device upon which the OpenCL actor should run, the language only requires that the device type be modified in the actor definition. No other change is required. Also, should the developer wish to use a different kernel or a dif-

```

1  type data_t is struct (
2    real [][] a;
3    real [][] b;
4    real [][] result
5  )
6  type settings_t is openc1 struct (
7    integer [] worksizes;
8    integer [] groupsizes;
9    in data_t input;
10   out real [][] output
11 )
12 type dispatchI is interface(
13   out settings_t requests;
14   out data_t dout;
15   in real [][] din
16 )
17 type mulI is interface(
18   in settings_t requests
19 )
20 stage home{
21
22   openc1 <device_index=0, device_type=CPU>
23   actor Multiply presents mulI {
24     constructor() {}
25     behaviour {
26       receive req from requests;
27       receive d from req.input;
28       x = get_global_id(0);
29       y = get_global_id(1);
30       dim = get_global_size(0);
31       c = 0.0;
32
33       for i = 0 .. (dim-1) do {
34         c := c + (d.a[y][i])*(d.b[i][x]);
35       }
36
37       d.result[x][y] := c;
38       send d.result on req.output;
39     }
40   }
41
42   actor Dispatch presents dispatchI{
43     constructor() {}
44     behaviour {
45       s = 1024;
46       ws = new integer[2] of s;
47       gs = new integer[2] of 0;
48       i = new in data_t;
49       o = new out real [][];
50
51       connect dout to i;
52       connect o to din;
53
54       config = new settings_t(ws,gs,i,o);
55       d = generate_data(s);
56
57       send config on requests;
58       send d on dout;
59       receive result from din;
60     }
61   }
62
63   boot{
64     d = new Dispatch();
65     m = new Multiply();
66     connect d.requests to m.requests;
67   }
68 }

```

Listing 3: Ensemble Matrix Multiply Example

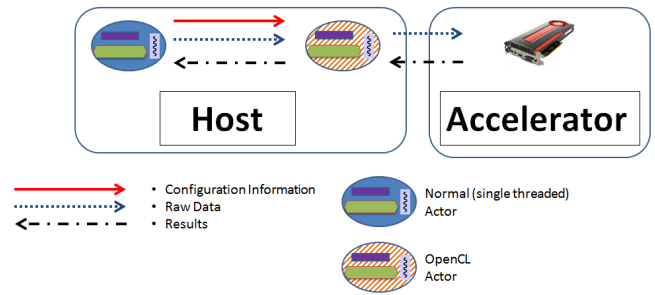


Figure 2: OpenCL in Ensemble Execution Model

ferent device at runtime, all that is required is to reconnect the configuration channel to an appropriate kernel actor's configuration channel.

By tagging an actor as an OpenCL kernel, a user must still write parallel code within the behaviour clause; however, by integrating this into the language model, the process is much simplified compared to a C version, see Section 7.

6.1.2 Compiler Modifications

The compiler was modified to apply slightly different rules to an OpenCL actor, enforcing the structure described previously in Section 6.1.1. The channel operations mirror the explicit data movement required by OpenCL, and it is the channel operations that dictate the creation of device buffers and the movement of memory between the host and the device. `struct` values are flattened so that each field is sent separately, with the compiler generating the appropriate code within the kernel to manage this. This includes pointer-rich structures, such as linked-lists and trees. Multi-dimensional arrays are flattened to single dimensional arrays. Again, the compiler generates appropriate kernel code to manage this. Primitive values are sent as 1D-arrays of one element so as to ensure updates within the kernel are applied to the value. Passing a pointer to the host variable is not an option. A potential optimisation here is to wrap all passed primitive variables in a single array.

6.1.3 Execution Model

Figure 2 shows the execution model when using OpenCL with Ensemble. Rather than the actor running on the device directly, an OpenCL actor is compiled into a bytecode representation of the actor in the normal way. A C representation of the code identified as the kernel is generated, and stored as a string within the actor's bytecode. Should the kernel actor contain functions defined in other code sections, the compiler will generate C equivalents within this string. This is completely hidden from the developer. The bytecode is interpreted as normal, and acts as the host in the traditional OpenCL sense. This actor handles the incoming and outgoing channel communications, as well as preparing, launching, and collecting data from the kernel. This also enables multiple kernels to execute on a single device. This includes, multiple kernels being scheduled for execution at the same time.

6.2 Runtime Integration

6.2.1 Runtime

In order to add OpenCL support into the runtime, a number of modifications were made. Firstly, the use of OpenCL is optional, and conditionally compiled into the runtime. This was necessary to ensure that RAM and ROM on resource-constrained platforms without OpenCL were not consumed. In this way, the same source tree is still available for multiple platforms, reducing maintenance and update effort.

Secondly, during the initialisation of the runtime, a single matrix is created to hold the different platforms and devices available in this system. This is done to ensure that there is only a single `command_queue` per device. This was necessary as race-conditions were observed with multiple `command_queues` per device when reading data. The information passed in the declaration (Listing 3, line 22) of an OpenCL actor is used to index into this matrix at runtime to determine the appropriate `context` and `command_queue`. If no information is given in the declaration, default values are used.

Thirdly, OpenCL wrapper functions were created and made available to the interpreter to simplify and abstract the interaction with the OpenCL API.

6.2.2 Interpreter

Within the interpreter, all OpenCL operations are implemented in C for performance. Each operation described in Section 2.1 is implemented as a custom native operation in the VM. The VM has a special `invokenative` bytecode that is used to provide access to certain primitive operations in the runtime [2]. This includes actor and channel lifecycle operations, and now includes access to the OpenCL wrappers discussed previously. Also, each OpenCL actor is given an `OpenCLEnvironment` variable. This is a runtime structure only visible within the interpreter that is used to store metadata about the platform, device, and device type, as well as the relevant `command_queue` and `context` for a given OpenCL actor. This structure is populated when the actor is created using the information contained in the previously described runtime matrix.

6.2.3 Lazy Evaluation

In OpenCL, a common idiom is to leave data on a device for as long as possible, thus reducing the time spent copying data between the device and host, and ultimately the application's execution time; data movement is often the greatest performance bottleneck. In actor-based languages there is no shared state, hence when messages are sent between actors, a duplicate is created and sent. This ensures no shared state between the actors, and that each actor has a unique copy of the data. While this is correct, it costs time, requires greater memory consumption, and precludes keeping data on the device.

To prevent such duplication, Ensemble supports marking non-primitive types as movable (`mov`), as discussed in Section 4. This approach has been applied to the OpenCL kernels in Ensemble. Marking the `in` channel to the kernel used for accepting data as moveable (`mov`) has two effects. Firstly, once any non-primitive data is copied to the device it is marked as no longer being on the host, and copies of relevant OpenCL data structures are associated with the

runtime representation of the data type. Secondly, the compiler will not generate the code to read this data back from the device. Thus, when the data is sent onwards it will hold a reference to the data on the device.

At this point there are two possibilities for the data that is sent onwards. The first option is that the data successfully arrives at another OpenCL actor without being accessed by the host. In this case, the pointer to the device data is set as the appropriate kernel argument, and the kernel is dispatched. Here, the data was kept on the device at all times. The second option is that the data is either accessed directly by host code, or the data is sent to an OpenCL actor associated with a different context. In both cases, the runtime reads the data back from the device and returns the device memory. As Ensemble uses automatic garbage collection, should the host reference count ever reach zero, both host and device memory will be returned. In this way, the choice to use `mov` gives the user control over memory usage. This optimisation is not possible in the C++ actor framework.

One benefit of lazy evaluation is that the runtime can automatically determine, using the `OpenCLEnvironment`, if the incoming data needs to be moved to the current context, and then do so if required. Currently, OpenCL manages data movement between devices in a single context, but not in different contexts.

7. EVALUATION

This section explores the linguistic complexity and performance cost incurred by the extension to the Ensemble framework by comparison to API and directive-based approaches.

7.1 Applications

A range of applications were chosen to evaluate the linguistic complexity and performance of actor-based OpenCL in Ensemble as compared to equivalent C and OpenACC implementations. The range of applications covered include matrix operations, multiple kernels, parallel reduction and a real world application. All of the source code for the described applications can be found online⁵.

- **Matrix Multiplication** multiplies two 1024^2 matrices in a single kernel.
- **Mandelbrot** computes a 1000 iteration Mandelbrot set in a single kernel.
- **LUD (Lower Upper Decomposition)** factorises a square matrix of 2048 elements, and is a common operation in matrix calculations: this example uses three kernels in series.
- **Matrix Reduction** finds the minimal value in an array of 33,554,432 elements using parallel reduction in a single kernel.
- **Document Ranking** takes a set of documents and using a template determines if these documents are wanted, or unwanted: this is a single kernel.

⁵www.paul-harvey.org/data-and-links.html

7.2 Experimental Setup

The host machine had 16GB of RAM. CPU refers to an Intel Core i5-3550 CPU @ 3.30GHz, and GPU refers to an AMD Radeon R9 290x GPU. All results are the average of one hundred runs of the respective application, and the error bars indicate the standard deviation of the total execution times. All code has been compiled with the `-O2` optimisation. Unless stated otherwise, the same results were observed for all applications and implementations - i.e. all implementations were functionally equivalent. The AMD version of OpenCL 1.2, and version 14.10 of the PGI compiler was used. All configuration information is found within the online code examples. No attempt was made to optimise the explicit kernels or annotated code for any approach to approximate a non-expert or non-computing scientist.

7.3 Code Complexity

Table 1 shows the absolute (arithmetic) difference between the concurrent and non-concurrent code versions for each approach. The number in brackets represents the percentage change. For example, the C-OpenCL version of Matrix Multiplication required 154 (142%) more lines of code than the single threaded C version, whereas the Ensemble-OpenCL version required 8 (10%) fewer lines than the single threaded Ensemble version. As well as the lines of code written, the table also shows McCabe's cyclomatic complexity [14] for applications. This metric quantitatively assesses the number of different paths through a program. Also shown is the assignments, branches, and conditions (ABC) metric that assesses the size/complexity of code [9]. In each case, the number shown is for the entire application; negative values indicate a decrease in the specified metric.

By comparison to C-OpenCL, both Ensemble-OpenCL and C-OpenACC require fewer lines of code, and are simpler by both metrics, with the only exceptions being the matrix multiplication and Mandelbrot cyclomatic complexity examples. The negative values seen are due to the kernel code/actor effectively replacing the outer `for` loop in the single-threaded version, thus reducing code and complexity. In general, annotating code with OpenACC pragmas has little effect on the code size or metrics. The main impact is from having to explicitly specify the sizes of the data to be moved, requiring variables to be accessed.

The implicit kernel `for` loop, plus the actor and channel abstractions accounts for the better results generally shown by Ensemble-OpenCL compared to C-OpenACC. The main discrepancy being the reduction example, which required very different kernel logic to the single-threaded equivalent in both Ensemble and C, however, this mindset is advocated by both approaches, unlike OpenACC.

7.4 Performance

Figures 3a-3e show the comparative performance of Ensemble-OpenCL, C-OpenCL and C-OpenACC for the applications described above. C-OpenACC will refer to both GPU and CPU results, as the PGI compiler was used for both OpenACC and OpenMP pragmas. Each column in each figure displays the time taken to complete the given application, normalised to the Ensemble GPU results. This is done to

highlight relative application performance, rather than absolute execution performance. The columns are split to show the relative amount of time taken to move data to a device, move data from a device, to execute a kernel, and the overhead which represents the total relative application execution time minus these values. This does not apply to C-OpenACC as it was not possible to correctly identify the distinct operations due to the `pragma`-based abstraction.

The general results that can be drawn from these figures is that Ensemble-OpenCL shows comparable performance to C-OpenCL on both CPU and GPU results across all applications. Also, the relative performance of Ensemble-OpenCL to C-OpenACC ranges from equivalent to significantly better on GPUs, and from better to vastly better on CPUs.

Figures 3a and 3b show the results for the matrix multiplication and Mandelbrot applications, respectively. The OpenCL actions (data movement and kernel calculations) are nearly equivalent between the C and Ensemble versions. The higher overheads in the Ensemble version are due to interpreter overhead for the non-kernel code. C-OpenACC shows similar performance on the GPU for Figure 3a, however much worse performance in Figure 3b, even when using the fine-grained `gangs` and `worker` annotations to explicitly specify the `groupsize` and `worksize` to be used on the GPU. Whereas an explicit kernel can take advantage of each thread's position within the 2D architecture of the GPU, the C-OpenACC abstraction cannot, hence the poor performance.

Figure 3c is of particular note as it highlights one of the advantages of using actors and channels. In this application there are three kernels, each performing a different operation. In the Ensemble version, a controller actor *plumbs* the channels of the kernel actors together, creating a pipeline of distinct kernels, Figure 4. The controller then sends the data, and waits for the final result. By comparison, the C-OpenCL version sequentially invokes each kernel from the host. As C-OpenACC is annotated single-threaded code, it also invokes the relevant code sequentially.

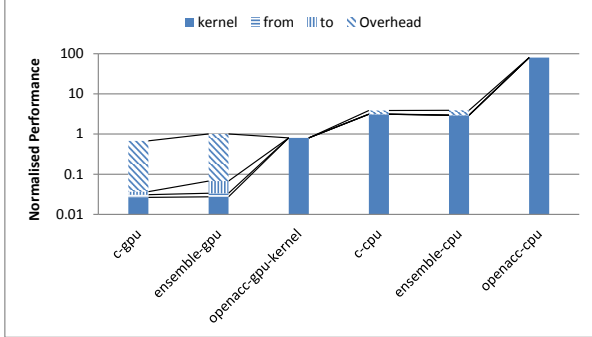
This application also highlights the advantage of movability within the language. Without movability, LUD took approximately 3 minutes (not shown) to complete on the GPU due to all the data movement involved; with movability, it takes approximately five seconds on the GPU. Here we can see that by using movable types, the time taken by Ensemble is comparable to the C version. The higher overheads are caused by the non-OpenCL code being interpreted by the unoptimised VM, not data movement.

In order to obtain comparable performance from OpenACC, annotating the outer loop of the relevant code was not sufficient, requiring use of the non-trivial `gangs` and `worker` annotations. This makes the distinction of host and device explicit, with the added disadvantage of being inline in the code, rather than having distinct sections of code. Again, worse performance is seen from the CPU.

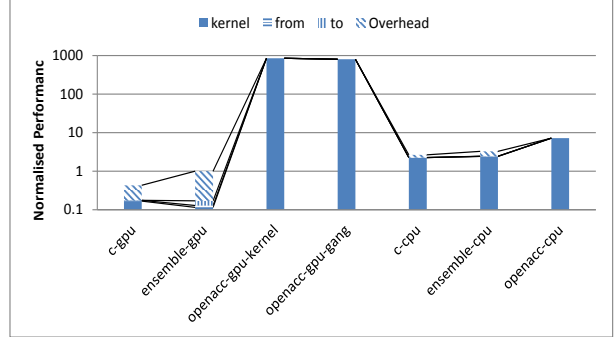
Figure 3d shows that Ensemble-OpenCL closely tracks the performance of C-OpenCL. C-OpenACC performs poorly for this application on both the GPU and CPU due to this

Table 1: Difference Between Single Threaded and Concurrent Code per Approach

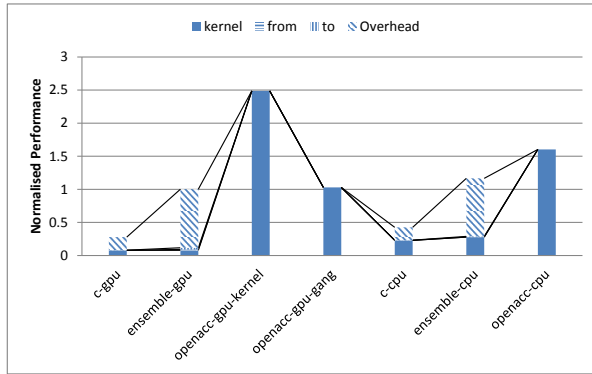
Application	Lines of Code			Cyclomatic Complexity			ABC		
	C	Ensemble	OpenACC	C	Ensemble	OpenACC	C	Ensemble	OpenACC
Matrix Multiplication	154 (142)	-8 (-10)	5 (4)	-1 (-6)	-2 (-25)	0 (0)	134 (638)	2 (12)	1 (5)
Mandelbrot	96 (48)	-4 (-3)	12 (6)	-1 (4)	1 (20)	0 (0)	22 (34)	-6 (-21)	2 (3)
Reduction	266 (157)	72 (73)	3 (2)	19 (105)	4 (80)	1 (5)	103 (163)	10 (66)	0 (0)
LUD	144 (59)	7 (2)	7 (3)	5 (13)	-8 (-28)	1 (2)	200 (183)	-11 (-21)	0 (0)
Document Ranking	45 (9)	-16 (-4)	3 (1)	53 (112)	-1 (4)	0 (0)	405 (111)	-29 (-40)	0 (0)



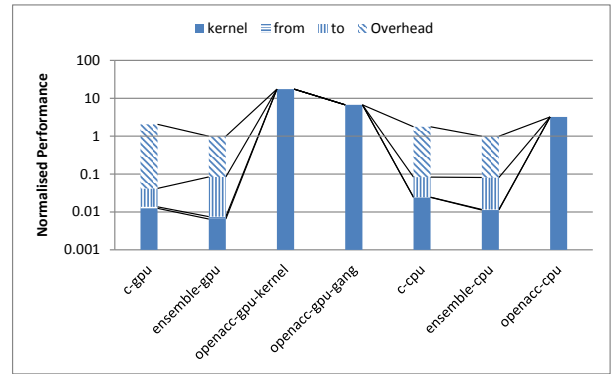
(a) Matrix Multiplication



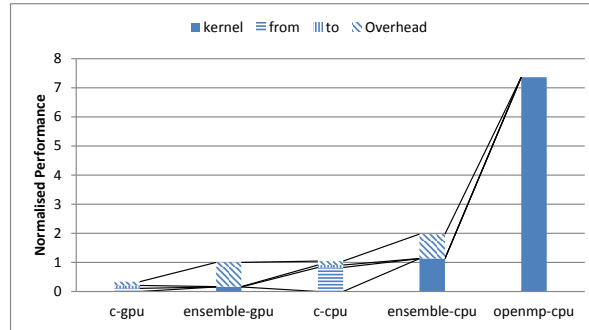
(b) Mandelbrot



(c) LUD



(d) Reduction



(e) Document Ranking

Figure 3: Performance between C-OpenCL, Ensemble-OpenCL and C-OpenACC Normalised to Ensemble GPU

type of application requiring different logic to take advantage of parallel hardware. Again, annotating the sequential version is not enough for this type of application.

Figure 3e shows results for a real world example and indicates that the kernel execution time in Ensemble-OpenCL is greater than in C-OpenCL. This is due to some fundamental

differences between the host languages, rather than the programming models. Firstly, in Ensemble there are no NULL values. This means that all data types must be initialised at the point of creation. This has the advantage of making the language safer, but can lead to increased execution time when the variable is written to before reading - i.e. initialisation was not required. In this application, two arrays

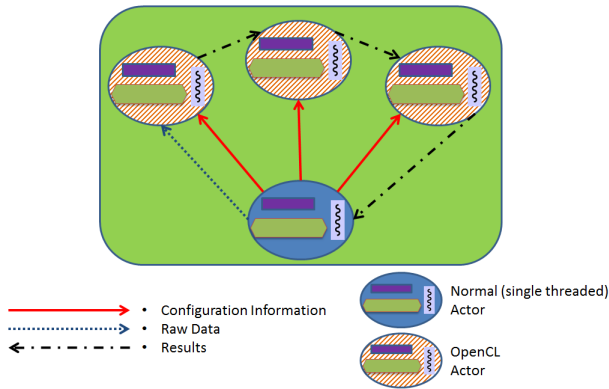


Figure 4: Topology of the LUD Application

are initialised in the kernel, within a loop with many iterations. A similar action is taken in the C version, however the two initialisation loops are combined, effectively *halving* the amount of work done by the C version. Loop unrolling in the code generated from the Ensemble compiler could help with this.

Secondly, the ability of C to use an integer value as both a boolean indicator and numerical value leads to faster code when compared to Ensemble, which has no such overloading. Ensemble uses separate types for numeric and boolean values, which in this application requires a number of control structures to be used, and causes greater execution time.

Thirdly, the C-OpenCL version uses short vector types and operations. Ensemble does not yet support OpenCL specific types such as short vectors, and the associated vector operations upon such types. This limitation is due to time constraints, rather than implementation or theoretical barriers. Such types and operations will be added to the Ensemble type system, and will reduce the execution time. Again, this is a limitation of Ensemble, not actors.

The second observation is the smaller communication time in Ensemble-OpenCL compared to C-OpenCL. This was an unexpected consequence of movability (Section 6). In this application, the kernel execution was run multiple times during each individual run to collect sufficiently large time values. In the C-OpenCL version, the data is copied to and from the device each time. No changes or modifications are made to the data between movements. This is also true in the Ensemble version, however due to the lazy evaluation logic the data on the device is never moved back to the host as it is not required to do so.

The PGI compiler was not able to compile this code, hence no results were obtained for the GPU or CPU from C-OpenACC. The CPU results were generated from the OpenMP pragmas and the gcc compiler. Even with the gcc compiler we still see the slower CPU results by comparison with the other methods.

From the results shown, the general trend is that the performance of hand coded C-OpenCL is comparable to Ensemble-OpenCL. As both C and Ensemble are using the OpenCL

runtime, the main difference in time between these approaches comes from the overhead of the Ensemble VM. There are optimisations that can be applied to the implementation, as discussed previously, however the fact that Ensemble is an interpreted language accounts for the majority of the overhead when compared to the C version, as opposed to issues with the language’s actor-based structure.

Given these results and the previous discussion, we have demonstrated that:

- Ensemble-OpenCL always enables simpler, functionally-equivalent code with many fewer lines of code compared to C-OpenCL, and generally simpler, equivalent code with fewer lines when compared to C-OpenACC;
- Ensemble-OpenCL applications present commensurate performance to C-OpenCL on GPU & CPU;
- Relative performance of Ensemble-OpenCL to C-OpenACC ranges from equivalent to significantly better on GPUs, and from better to vastly better on CPUs.

8. CONCLUSIONS AND FUTURE WORK

8.1 Conclusion

This paper presents an extended runtime and language model for actor-based systems to facilitate and simplify programming kernel-based concurrency. The separation of host and device code, as well as the explicit data movement found in kernel programming, naturally maps to the concept of isolated actors using message-passing communication. Using the information expressed in the language, the runtime is able to automate the initialisation and interaction with the parallel hardware, via OpenCL.

By comparing C-OpenCL with Ensemble-OpenCL, this paper has shown that Ensemble-OpenCL requires less code, is qualitatively simpler, and provides comparable performance with hand written C-OpenCL. Also, by taking advantage of novel features in the language, it is still possible to use common OpenCL optimisation techniques which are in contrast to the normal shared-nothing semantics of actor-based programming. Furthermore, by comparison with C-OpenACC, Ensemble-OpenCL shows performance ranging from better to vastly better, at the cost of slightly larger code size in a small number of cases.

Ultimately, by offering a simple and performant programming model for a range of applications classes, actors provide an ideal abstraction for kernel-based concurrency.

8.2 Future Work

The main bottleneck in performance is the Ensemble VM. To address this, a JIT compiler will be built to remove the overhead of bytecode interpretation. Furthermore, the Ensemble type system will be expanded to include OpenCL types, such as short vectors.

Work is currently in progress to add distribution natively to Ensemble. This includes abstraction of network communication between actors via the channel mechanism, user-driven discovery of runtime resources, remote deployment of actors,

and runtime actor migration. As kernels are represented as actors, they will be able to natively take advantage of these benefits. Once complete, one possible use-case will be to program and load-balance a computer cluster of heterogeneous platforms entirely with actors from within the language.

9. ACKNOWLEDGEMENTS

The authors wish to thank the Carnegie Trust for the Universities of Scotland for funding this work. Also, Wim Vanderbauwhede for his helpful comments and suggestions in the creation of this document.

10. REFERENCES

- [1] The opencl specification. *Khronos OpenCL Working Group*, 1:11–15, 2009.
- [2] C. Cameron, P. Harvey, and J. Sventek. A virtual machine for the insense language. In *MOBILE Wireless MiddleWARE, Operating Systems and Applications (Mobilware), 2013 Intl. Conference on*, pages 1–10, Nov 2013.
- [3] B. Chapman, G. Jost, and R. Van Der Pas. *Using OpenMP: portable shared memory parallel programming*, volume 10. MIT press, 2008.
- [4] D. Charoussat, T. C. Schmidt, R. Hiesgen, and M. Wählisch. Native Actors – A Scalable Software Platform for Distributed, Heterogeneous Environments. In *Proc. of the 4rd ACM SIGPLAN Conference on Systems, Programming, and Applications (SPLASH '13), Workshop AGERE!*, New York, NY, USA, Oct. 2013. ACM.
- [5] D. G. Clarke, J. M. Potter, and J. Noble. Ownership types for flexible alias protection. In *Proceedings of the 13th ACM SIGPLAN Conference on Object-oriented Programming, Systems, Languages, and Applications, OOPSLA '98*, pages 48–64, New York, NY, USA, 1998. ACM.
- [6] A. Dearle, D. Balasubramaniam, J. Lewis, and R. Morrison. A component-based model and language for wireless sensor network applications. In *Proc of the IEEE Intl. Computer Software and Applications Conference*, pages 1303–1308, 2008.
- [7] J. Docampo, S. Ramos, G. Taboada, R. Exposito, J. Tourino, and R. Doallo. Evaluation of java for general purpose GPU computing. In *Advanced Information Networking and Applications Workshops*, pages 1398–1404, March 2013.
- [8] C. Dubach, P. Cheng, R. Rabbah, D. F. Bacon, and S. J. Fink. Compiling a high-level language for GPUs: (via language support for architectures and compilers). In *Proc of the 33rd ACM SIGPLAN Conference on Programming Language Design and Implementation, PLDI '12*, pages 1–12, New York, NY, USA, 2012. ACM.
- [9] J. Fitzpatrick. More c++ gems. chapter Applying the ABC Metric to C, C++, and Java, pages 245–264. Cambridge University Press, New York, NY, USA, 2000.
- [10] T. D. Han and T. S. Abdelrahman. hiCUDA: High-level GPGPU programming. *IEEE Trans. Parallel Distrib. Syst.*, 22(1):78–90, 2011.
- [11] P. Harvey, A. Dearle, J. Lewis, and J. S. Sventek. Channel and active component abstractions for WSN programming - A language model with operating system support. In *SENSORNETS 2012 - Proc of the 1st Intl. Conference on Sensor Networks, Rome, Italy, 24-26 February, 2012*, pages 35–44, 2012.
- [12] J. Jeffers and J. Reinders. *Intel Xeon Phi Coprocessor High Performance Programming*. Newnes, 2013.
- [13] A. Klöckner, N. Pinto, Y. Lee, B. Catanzaro, P. Ivanov, and A. Fasih. PyCUDA and PyOpenCL: A scripting-based approach to GPU run-time code generation. *Parallel Comput.*, 38(3):157–174, Mar. 2012.
- [14] T. McCabe. A Complexity Measure. *Software Engineering, IEEE Transactions on*, SE-2(4):308–320, Dec 1976.
- [15] B. Nichols, D. Buttler, and J. P. Farrell. *Pthreads Programming*. O'Reilly & Associates, Inc., Sebastopol, CA, USA, 1996.
- [16] Z. Porkoláb, J. Mihalicza, and A. Sipos. Debugging c++ template metaprograms. In *Proceedings of the 5th International Conference on Generative Programming and Component Engineering, GPCE '06*, pages 255–264, New York, NY, USA, 2006. ACM.
- [17] J. Sanders and E. Kandrot. *CUDA by example: an introduction to general-purpose GPU programming*. Addison-Wesley Professional, 2010.
- [18] M. Steuwer, P. Kegel, and S. Gorlatch. SkelCL - a portable skeleton library for high-level GPU programming. In *Proc of the IEEE 27th Intl. Symposium on Parallel and Distributed Processing*, pages 1176–1182, 2011.
- [19] A. Stromme, R. Carlson, and T. Newhall. Chestnut: a GPU programming language for non-experts. In *Proc of the 2012 Intl. Workshop on Programming Models and Applications for Multicores and Manycores*, pages 156–167. ACM, 2012.
- [20] D. Tarditi, S. Puri, and J. Oglesby. Accelerator: Using data parallelism to program GPUs for general-purpose uses. In *Proc of the 12th Intl. Conference on Architectural Support for Programming Languages and Operating Systems*, pages 325–335, 2006.
- [21] C. Varela and G. Agha. Programming dynamically reconfigurable open systems with salsa. *SIGPLAN Not.*, 36:20–34, December 2001.
- [22] S. Wienke, P. Springer, C. Terboven, and D. an Mey. OpenACC: First experiences with real-world applications. In *Proc of the 18th Intl. Conference on Parallel Processing*, pages 859–870. Springer-Verlag, 2012.