

Work Scheduling on Heterogeneous Resources

Gabriele Keller

Edward Pierzhalski

Supervisor

February 12, 2015

1 Introduction

There are two notable trends in computer hardware. First, CPUs are gaining cores per chip while having reduced gains in clock speed. Second, GPUs are developing into generic computing units, with APIs that expose more of the power of their architecture. As more computing problems are found to be parallelizable, there is an incentive to try and move work that would have been done on the CPU to the GPU, particularly operations on vectors and streaming data. At the same time, research suggests that the gap between CPU and GPU processing rates is not as high as one expects. [1, 2] This gives an incentive to design languages and frameworks to ease the division of work over both the CPU and GPU.

1.1 Why GPU Programming is Painful

Two of the most popular frameworks for GPU programming are CUDA and OpenCL, which are both low-level imperative languages. In theory, the problem of partitioning work over the CPU and GPU can be solved at the level of these languages. In practice, even for GPU-focused tasks such as a monoidal scan, sorts, or folds, efficient implementations tend to require intimate knowledge of the memory and process architecture of the GPU. Harris et al. [3], Satish et al. [4] This is an added cognitive strain on users of these frameworks, and makes composing GPU tasks difficult: a map followed by a scan is substantially different to just a map, or just a scan.

These languages are also impractical for use as a single source language over multiple types of processors - for instance, while CUDA and OpenCL can target CPU architectures, the resulting binaries are substantially less performant than code originally written with a CPU architecture in mind.¹

The full semantics of the CUDA and OpenCL languages allow for arbitrary write-effects from any kernel to any array, and include operations such as global barriers. This makes reasoning about kernels difficult for the same reason that hidden state impedes compositional reasoning in many imperative languages. The effect is compounded when kernels are expected to run on separate devices; the result is that these low-level languages make for a poor foundation when trying to attack this problem.

1.2 Why Scheduling is Painful

There are several obstacles to effectively partitioning, let alone scheduling, a workload over different work processors. The major issues both revolve around memory, namely copy speed and access patterns.

¹See Jie Shen et al. [5] for a discussion focusing on OpenCL.

Memory Transfer Speed

Inside of a single GPU, data transfer rates can reach upwards of 300 GB/s. [6] However, CPUs and GPUs are typically connected over a PCIe bus. Modern PCIe busses have a bandwidth cap of around 8GB/s. [7] This places strong constraints on how frequently data can be shuttled between components. In the current Accelerate architecture, producer/consumer fusion reduces the number of intermediate structures required to perform work, which helps mitigate this issue. However, this fusion can impede attempts to partition work.

Memory Access Patterns

Parallelizable computations may have varying degrees of contiguity in their memory accesses. When computing a simple map operation, adjacent output locations depend on adjacent input locations, and so we can partition the work by index almost arbitrarily. Similarly, a tree-traversal monoidal fold can, at any given depth, be split over processors. However, any output index of a backpermute can be affected by the value at any index of the input. Although we can parallelise the work, we cannot partition the memory access. (Expand on examples?)

1.3 Making Things Better

A large body of work has been produced on generation of parallel code for the GPU and for the CPU, [8] however less has been done on the problem of doing so automatically for a general-purpose program. Previous work has investigated dynamic scheduling of a task composed of kernels, however the kernels had to have both CPU and GPU versions provided by the user of the library. (cite) Ideally, we should be able to separate the declaration of our problem from the generation of parallel code to run on whatever hardware happens to be available. Correctly partitioning and allocating work in the light of the above issues will be a key component of this thesis.

2 Background

Phrase the problem as follows: you have some source code in some language, describing a computation you want performed. It would be nice if, given this source specification, we could partition and schedule the work on to any processing unit available.

2.1 Combinators

The frameworks described previously are distinguished by being *imperative*: they describe, step by step, the actions performed by the GPU. Alternatively, we can describe our program using *combinators*, which abstract the combination of computations. For example, map is a combinator that takes a function `f : a -> b`, along with some list `as : List a`, and produces a new list, `map f as : List b`, which contains the result of applying `f` to every

element in `as`. We don't particularly care how `map` performs this: we just care that the result is what we expect.

This separation between *what* we want to happen, and *how* we want it to be done, makes combinators quite useful as a source language for computations to be performed, assuming the combinators have 'sufficiently sane' semantics. (yeah this needs a better definition) Once the user describes their program, we can interpret the description into whatever implementation we need to in order to satisfy the semantics of the given program.

2.2 Describing Semantics

The design of a language framework for a specialised purpose (for instance, automatically parallelised vector programming) results in a Domain Specific Language, or DSL. Many DSLs are implemented *inside* another host language, commonly called 'embedding'. There are essentially two flavours of Embedded DSLs (EDSLs):

Shallow Embeddings

Shallow embeddings are almost synonymous with constructing a library. Terms in the DSL correspond to library functions that, when evaluated, directly perform whatever computations that they represent.

Deep Embeddings

Alternatively, we can use the host languages' data definitions to construct a direct representation of the combinators. Say we define a new version of `map`, call it `map'`. Instead of a function, define it as a data constructor: an object containing a function and a list, i.e `map' : (a -> b) -> List a -> DSL b`, where `DSL` is a type representing our DSL embedding. The representation data type is typically called the Abstract Syntax Tree, or AST.

This representation has issues: we can't nest it, for instance, since `map'` takes a `List` yet produces a `DSL`. However, if we change `map'` to take a `DSL`, we will have completely removed lists from our list operation description language! The solution is to introduce a way to 'lift' lists into our DSL, which we can do by including a data constructor `lift : List a -> DSL a`. We can then give `map'` the type `(a -> b) -> DSL a -> DSL b`.

In order to do anything with these data structures, we need to interpret or evaluate them. For instance, we can define a recursive evaluation function `eval`:

```
eval : DSL a -> List a
eval dsl = dsl match
  lift list          -> list
  map' f anotherDsl -> map f (eval anotherDsl)
```

Which essentially 'replaces' `map'` with `map` and lifted lists with their underlying ones.

Benefits of Deep Embeddings

This may seem like an unnecessary and over-complicated overhead, however it gives us the flexibility to do other things with our description of list computations. One of the most important is that once a user has constructed a term in the DSL, we can manipulate it in the host language. For array DSLs, many of these manipulations are for optimisation.

As an example, consider the evaluation of `map f (map g as)`. Semantically, the final result is the application of the composition of `f` and `g` on the elements of `as`. However, since `map` is a simple function, the two calls to `map` will produce two intermediate lists. On a CPU, this is a relatively small (but not ignorable!) overhead; yet if we were to naively send off the work to the GPU as two separate mapping steps (presumably over arrays instead of lists), we would suddenly run into memory bottlenecks transferring the data over the PCIe bus.

Compare this with the equivalent deep embedding, `map' f (map' g (lift as))`. Since this is just a nested data structure, we can ‘pull out’ the functions and compose them, producing a new data structure `map' (f . g) (lift as)`. When we finally evaluate this new DSL value, it will be converted into a call to `map` that only makes a single list, avoiding memory bottlenecks.

2.3 Accelerate

Accelerate is an EDSL targeting general-purpose GPU programming. It is embedded in Haskell, a high-level general purpose functional programming language. Accelerate models general GPU programming using a deeply-embedded AST, which is then transformed through a variety of intermediate representations and functions over them. These transformations perform a wide variety of changes aimed at converting the AST from the high-level, user-facing description into a low-level description appropriate for compilation to machine code for the relevant platform.

All levels of the Accelerate framework make heavy use of advanced features of the Haskell compiler. Most prominently, GADTs and type families are used to encode many properties of the AST, such as the shape of arrays, and levels of parallelism (these advanced type annotations are left out of this report for brevity). Other features include the use of De Bruijn indices to capture lambda abstraction, and also to facilitate term reuse. However, despite how interesting all these properties are, they mostly form a backdrop to our main focus, the compilation phases.

The final, machine-code-generating compilation step is intentionally kept separate from the AST transformations, to make changing compilation targets as straightforward as possible. (does this add overhead? could we do cool things if we screwed with the tree earlier, or provided different KernIR primitives specialised for the CPU?) Recently, a new backend kit has been published for Accelerate, cleanly separating these compiler stages and exposing a simplified kernel-level representation compared to previous versions of Accelerate.

One of these changes introduced a new pair of combinators to the language, `split` and `concat`. They have the intuitive semantics that their names imply, of splitting and concatenating arrays. The combinators are carried through to the KernIR stage? only in SimpleAcc?

used by the compilation phases to ‘fragment’ arrays and operations on them. For instance, `map f arr` may be translated² into

```
let (x, y) = split arr in
    concat (map f x) (map f y)
```

When later stages in the compiler encounter the `concat` node, they may schedule the two arguments to `concat` to run on different processors before combining them.

Take up two pages describing the internals of Accelerate. Note: get a hand on R/T’s code so we can see exactly how they go from fissioning the SimpleAcc AST to actually running code on separate processing units.

Aside: A Frustrating Operation

There is one family of array operations in the Accelerate DSL that is particularly non-performant for heterogeneous scheduling. The `permute` and `backpermute` operations use a function from one set of array indices to another in order to permute an array. For instance, `backpermute f old` will produce an array `res` such that `res ! i = old ! (f i)`.

Since the range of `f` is unrestricted, the value at any index of the output array may depend on the value at any index of the input array. This complicates the partitioning of the work: although we can split the task of constructing the resulting array, we can’t partition the input array between processors.

In many GPU use-cases, the strongest constraint on throughput is memory transfer speed, which makes copying the entire input array to each processor infeasible. The result is that these permutations are *de facto* synchronisation points: every operation before them must all have their results returned and combined before the permutation. This has a potentially interesting interaction with data fission: normally, a map composed with a backpermute would be fused together in some fashion. We may ask the question: when is it worth fissioning data for the map before recombining for the permutation, and when is the fused operation preferred?

3 Previous Work

3.1 Heterogeneous Scheduling

Previous work on scheduling a parallelisable workload onto heterogeneous processors has explored solutions in low-level contexts, and has encountered a variety of obstacles.

²This is pseudocode: the details of the Accelerate DSL are omitted for brevity. Example drawn from Newton et al. [9].

Binary Analysis

In work by Lee, Samadi, Park, and Mahlke [8], the authors construct a framework that takes a users' data parallel OpenCL kernel and partitions the work across devices. This is achieved through hooking into the OpenCL API layer to expose a single, large device.

Kernels are split at runtime using a decision tree heuristic, using approximations on the transfer costs of moving data between processors, and also the variance of performance of processors. These parameters are estimated by the partitioning system, however it is not mentioned what the method of estimation is.

To determine whether a workload can be safely separated to several devices, data flow analysis is used on array indexing. If the data flow analysis does not sufficiently guarantee safe array access, the kernel is considered un-parallelisable. This is an example of where a higher-level abstraction may have assisted in splitting the workload. These are both potentially interesting directions for research on improving Accelerate, however they are beyond the current scope.

Profiling

3.2 Array DSLs

3.3 Heterogeneous Scheduling in Accelerate

A recent (at the time of writing, unpublished) paper by Newton, Holk, and McDonell [9] describes a first foray into heterogeneous scheduling in Accelerate. This is also the work that led to the previously described improvements to Accelerate - the separation of AST transformation phases, the improved final kernel representation, and the separation and abstraction of the hardware backends. In addition to these contributions, in their paper Newton et al. introduce data fissioning as described previously, and explore two methods of scheduling tasks over the fissioned data.

Bulk-Synchronous Parallel Tasks were fissioned into some constant number of segments, which were in turn round-robin allocated onto processors. Scheduling was performed in batches, and tasks were fissioned in dependency order.

This first method had no notion of device affinity, and so ran into memory bottlenecks since data would be frequently copied from processors to the host before the next batch was allocated. In addition, the batch allocation added significant synchronisation overhead between rounds. These deficiencies resulted in extreme performance degradation, and the authors do not pursue the method further (which is a hint that we should not either).

Single Program, Multiple Data As an alternative strategy, the authors use the dependency graph to ensure that at least some of the dependencies for a fragment are all allocated

to the same device. Data fissioning occurred in the same manner as with the bulk-synchronous method.

The enforcement of the scheduling semantics of fissioning were split between both the fissioning stage of AST processing, and in the final device scheduler. Namely, the fissioning algorithm is what made decisions about work ratios and so expected, for example, the first fraction of the segments to go on one device and the rest to another. This was separately reflected in the scheduler as well. This latter method forms the foundation of this project.

4 What We're Going To Do

There is much potential work within the scope of what has been discussed so far. For the purposes of this project, the high-level goal is to have Accelerate able to fission a program written in its DSL (already done by Trevor and Ryan - but not automatically?) into 'jobs' that can be scheduled onto a heterogeneous set of workers (CPUs, GPUs) with a demonstrable performance boost, while also addressing some of the issues expressed in previous work. (Need to clarify what I'm adding that isn't just 'making the scheduler better') In addition, the following are potential points of focus to provide guidance in achieving the desired result.

4.1 Dynamic Fissioning

The data fissioning described previously was always static: the SimpleAcc term would be fissioned into some finite and static number of array segments, and each work component would be sent off to the same processing unit (determined by index).

We can improve on this by determining the degree of fissioning and the allocation of fissioned segments dynamically, depending on factors such as the relative processing speeds of the processors.

4.2 Avoiding Sharing

Currently doesn't try and avoid copying data, just hoping that the thing works

4.3 Promoting Sharing

Currently can't have several devices working on a 'different part of the program' at the same time. Need to clarify.

4.4 Small Array Optimisations

One minor issue not addressed by existing approaches to heterogeneous resource scheduling is the cost of operations on small arrays. Loading, running, and using the result of a GPU kernel has significant overhead. Below a certain array size, the processing time of operations on an

array is dominated by the high latencies of both loading the GPU kernel and data transfer between the CPU and GPU.

In these cases, it would be more effective to perform the operation on the CPU. The existing Accelerate framework does not differentiate workloads by array size, and so does not take advantage of this. Once we implement dynamic fissioning, we can look into including this optimisation.

4.5 Extensions

Below is a list of potentially valuable additions to the goals of this project. Since they are not critical to the goal, they will be pursued only if time permits.

Improved Concatenation Currently, `concat` nodes are converted into array generators using conditional indexing on the segments. For instance if we have `b = concat b1 b2`, and if `b1` and `b2` are split at some index `s`, this will be translated into

```
b = generate (\i -> if (i < s) then (b1 ! i) else (b2 ! i))
```

Although there are situations where the intervening `generate` may be removed by inlining or through indexing with a constant value, performance may be improved otherwise by adding memory-copying primitives to the lower-level array representations.

LLVM Support The current LLVM backend for Accelerate does not support some operations, such as stencils. Expanding support for this feature will allow more comprehensive comparisons of fissioning transformations.

Device Affinity Dynamic fissioning would be improved if we had a notion of ‘device affinity’, so that data dependencies in the task graph are not unnecessarily copied between processors.

Processor Affinity For some tasks such as sorting and non-monoidal folds, CPUs may be as fast as GPUs. For this reason, after fissioning we may prefer that a given segment of work be performed on a particular kind of processor. Handling this form of ‘task affinity’ may provide additional improvements to performance.

(Find more examples. Is this relevant? Does Accelerate contain enough operations that are performant on CPUs?) (Note: find out why Ryan/Trevor didn’t add the benchmark for the CPU+GPU performance on the Black Scholes stuff.)

4.6 Initial Technical Plan

The existing framework relies on the scheduling backend reacting to `split` and `concat` nodes in the final representation, and splitting work accordingly. Currently, the process of introducing these fragmenting nodes occurs in a single phase, after high-level map fusion. Investigating the

interaction of data fissioning with other operations and optimisations will involve separating out the fusion steps, so as to insert fissioning between them. Adding information about available processors will require providing some data to most of the compilation steps: we can achieve this with something like the reader monad, if needed. Adding affinities may involve adding metadata to the fissioning nodes themselves.

5 Work Schedule

As this is partially a development thesis, we must form a development plan. The below schedule details the hopes and expectations for when parts of the development phase will be completed. Included are milestones to be used as checks on whether these expectations have been met. The final report is expected to be written in tandem with the development stages: each week and each milestone has an implicit requirement, ‘have something written down about this’.

(Need to clarify what parts are research and what parts are development. Also need to clarify what granularity is required other than ‘mess with Accelerate to learn about it, then mess with Accelerate to add something to it’.)

5.1 Timetable

Week 1 - 4: Familiarise myself with Accelerate.

Week 1: User-facing DSL.

Milestone: a simple Accelerate program, say a Newtonian gravity simulation.

Week 2 - 3: Compilation phases.

Milestone: implement a simple phase?

Week 3 - 4: Accelerate backends (mainly LLVM)

Week 5 - 9: Write scheduler.

Week 5 - 6: Have fissioning working on the AST.

Milestone: produce index-partitioned ASTs at kernel generation step.

Milestone: produce ASTs with optimisations for small array sizes.

Week 6 - 7: Emit to both CPU and GPU-facing backends.

Milestone: have above kernels produced and split between CPU and GPU targets, at some static fraction.

Milestone: implement small-array CPU optimisation as part of scheduling.

Week 7 - 9: Implement dynamic scheduling.

Milestone: a runtime utility to determine what backends and processors are available, along with their processing speeds and bus transfer limits (both between and within processors).

Milestone: (Extension) have task fissioning and task allocation take into account memory dependencies and relative processing speed of processors.

Week 10 - 12: Benchmarks.

Milestone: Compare performance on common problems (Black Scholes, n -body) with and without scheduler.

Milestone: Compare vs. previous scheduler implementation from Newton et al. [9], hopefully seeing improvement.

Milestone: Show improvement due to separate stages of scheduler optimisation: small-array CPU scheduling, task partitioning, preventing permute fusion.

Milestone: (Extension) Show improvement due to dynamically informed allocation vs. static strategy.

Week 13 - Due: Buffer time for the inevitable. Clean-up for final report.

5.2 Measures of Success

For those milestones directly related to construction of the scheduler, there are two ‘parameters’ for success: depth and breadth.³ ‘Depth’ in this context refers to successfully processing an AST through all the relevant compiler passes. ‘Breadth’ refers to having depth completion across Accelerates’ many array operations (maps, folds, permutes, and recently iteration constructs). Clearly depth completion is a necessary prerequisite for most of the desired results concerning scheduling, however breadth is important for assessing the interaction between task fissioning and compilation optimisations like fusion.

³This conceptual division is not due to us: we rephrase it here, but it is originally from Newton et al. [9].

References

- [1] Victor W. Lee, Changkyu Kim, Jatin Chhugani, Michael Deisher, Daehyun Kim, Anthony D. Nguyen, Nadathur Satish, Mikhail Smelyanskiy, Srinivas Chennupaty, Per Hammarlund, and others. Debunking the 100x GPU vs. CPU myth: an evaluation of throughput computing on CPU and GPU. In *ACM SIGARCH Computer Architecture News*, volume 38, pages 451–460. ACM, . URL <http://dl.acm.org/citation.cfm?id=1816021>.
- [2] Chris Gregg and Kim Hazelwood. Where is the data? why you cannot debate CPU vs. GPU performance without the answer. In *Performance Analysis of Systems and Software (ISPASS), 2011 IEEE International Symposium on*, pages 134–144. IEEE. URL http://ieeexplore.ieee.org/xpls/abs_all.jsp?arnumber=5762730.
- [3] Mark Harris, Shubhabrata Sengupta, and John D. Owens. Parallel prefix sum (scan) with CUDA. 3(39): 851–876. URL <http://www.eecs.umich.edu/courses/eecs570/hw/parprefix.pdf>.
- [4] Nadathur Satish, Mark Harris, and Michael Garland. Designing efficient sorting algorithms for manycore GPUs. URL <http://www.nvidia.com/docs/io/67073/nvr-2008-001.pdf>.
- [5] Jie Shen, Jianbin Fang, H. Sips, and A. L. Varbanescu. Performance traps in OpenCL for CPUs. pages 38–45. IEEE. ISBN 978-1-4673-5321-2, 978-1-4673-5321-2, 978-0-7695-4939-2. doi: 10.1109/PDP.2013.16. URL <http://ieeexplore.ieee.org/lpdocs/epic03/wrapper.htm?arnumber=6498531>.
- [6] NVIDIA. GeForce GTX 780-TI specifications, 2015. URL <http://www.geforce.com/hardware/desktop-gpus/geforce-gtx-780-ti/specifications>.
- [7] PCI-SIG. PCI Express 3.0 frequently asked questions, 2015. URL https://www.pcisig.com/news_room/faqs/pcie3.0_faq.
- [8] Janghaeng Lee, Mehrzad Samadi, Yongjun Park, and Scott Mahlke. Transparent CPU-GPU collaboration for data-parallel kernels on heterogeneous systems. In *Proceedings of the 22nd international conference on Parallel architectures and compilation techniques*, pages 245–256. IEEE Press, . URL <http://dl.acm.org/citation.cfm?id=2523756>.
- [9] Ryan R. Newton, Eric Holk, and Trevor L. McDonell. Converting data to task-parallelism by rewrites. URL <http://www.cse.unsw.edu.au/~tmcdonell/papers/acc-multidev-icfp2014-sub.pdf>.