# SPOC: GPGPU PROGRAMMING THROUGH STREAM PROCESSING WITH OCAML

Mathias Bourgoin, Emmanuel Chailloux, Jean-Luc Lamotte

**HAL Id: hal-00697257**

**https://hal.science/hal-00697257**

Submitted on 15 May 2012

# SPOC : GPGPU PROGRAMMING THROUGH STREAM PROCESSING WITH OCAML

MATHIAS BOURGOIN, EMMANUEL CHAILLOUX and JEAN-LUC LAMOTTE

*Laboratoire d'Informatique de Paris 6 (LIP6 - UMR 7606)*
*Université Pierre et Marie Curie (UPMC - Paris 6)*
*Sorbonne Universités*
*4 place Jussieu, 75005 Paris, France*
*email: Mathias.Bourgoin@lip6.fr,Emmanuel.Chailloux@lip6.fr,Jean-Luc.Lamotte@lip6.fr*

### ABSTRACT

General purpose computing on graphics processing units (GPGPU) consists of using GPUs to handle computations commonly handled by CPUs. GPGPU programming implies developing specific programs to run on GPUs managed by a host program running on the CPU. To achieve high performance implies to explicitly organize memory transfers between devices. Besides, different incompatible frameworks exist making productivity and portability difficult to achieve. In this paper, we describe SPOC, an OCaml library, defining specific data sets in order to automatically manage transfers between GPU and CPU. SPOC also offers a runtime library looking for multiple frameworks and making them usable transparently. We also describe the link between SPOC and the OCaml garbage collector to optimize transfers dynamically. SPOC benchmarks show that SPOC can offer great performance while simplifying GPGPU programming

*Keywords*: OCaml, Stream Processing, GPGPU, High Level Libraries, Data-Parallel

## 1. Introduction

Since 1980, graphics hardware has become more and more complex and efficient, providing many computation units and large dedicated memory. From fixed-function pipelines, Graphics Processing Units (GPU) have become fully programmable through shaders (for custom rendering techniques) and then through specific frameworks for general purpose programming (GPGPU). NVIDIA has been developing the Cuda framework while the Khronos Group offers the OpenCL standard, as a general computing framework for GPGPU programming. OpenCL targets GPUs as well as CPUs and specific accelerators (*e.g.*, Cell-BE processor, FPGAs). Many vendors are now providing an OpenCL implementation with their hardware: currently, every personal computer sold is GPGPU compatible. OpenCL and Cuda are the main frameworks. OpenCL, as a standard, offers some portability while Cuda, which

is proprietary and older, offers a more mature framework with many optimized libraries and tools. However, Cuda is usable on NVIDIA hardware only. Currently, GPGPU software can mainly be written in NVIDIA Cuda assembly, C/C++ for Cuda or OpenCL.

Both frameworks use low-level libraries with explicit hardware and memory management implying complex and error-prone programming. GPGPU developement could easily benefit from a high level language abstracting memory transfers and computations while ensuring type and memory safety. As GPGPU programming aims at high performance, part of the approch is to offer a high-level runtime library abstracting transfers while binding with external GPGPU programs. We followed this approach, choosing OCaml as a composing language, basing memory management (including transfers) on its garbage collector and dynamically compiling GPGPU programs (in Cuda or OpencL). This keeps high performance while improving abstraction, allowing easier and safer programming.

We propose SPOC, an OCaml library which enables to use the power of GPGPU programming with the OCaml language. SPOC uses OCaml as a glue language for the OpenCL and Cuda frameworks with a runtime library allowing automatic management of devices memory and data transfers between CPU and devices. In section 2, we present SPOC and its implementation. Then, we show its efficiency through benchmarks in section 3. In section 4, we present other approaches to offer high-level GPGPU programming. Finally, we conclude on our approach and present our future work.

## 2. Stream Processing with OCaml

### 2.1. *Stream Programming*

GPGPU programming is very hardware dependent. Even with OpenCL, getting high performance implies writing specific programs for different devices. OpenCL and Cuda frameworks use the Stream Processing paradigm based on SIMD (Single Instruction Multiple Data) architectures. A GPGPU device (which we will now call a device) is seen as multiple computational units, which, given a set of data and a computation kernel (a sequence of basic operations), apply this kernel to each element of the set in parallel.
A device offers multiple layers of parallelism:

- computation units which we will call threads,
- threads are grouped into blocks and
- blocks are grouped into a software grid.

The grid is mapped to the device's multi-processors, each multi-processor running one or several blocks. Furthermore, new layers appear when

- running multiple grids in parallel on one device,
- building multi-device computers,

- associating several computers with device(s) into computer grids.

Besides, devices are considered to have a dedicated memory[a], implying the copy of data from the CPU memory to the device's. Device memory is, moreover, divided into several categories:

- global memory (accessible by all threads in the grid)
- shared memory (shared inside a block)
- local memory (local to a thread)
- specific memory banks depending on hardware

Global memory banks can be read/written by both GPU and CPU while shared, local and specific memory banks are only accessible from the GPU. As many levels of parallelism and memory make Stream Programming strongly hardware related, which makes portability and high performance difficult to achieve conjointly. While kernel optimization is mandatory to achieve great performance, memory management and data transfers scheduling also have a high impact on the overall program performance. Having a fast kernel waiting for data can ruin all one's efforts. Low-level APIs, manual memory management and incompatibility between frameworks decrease development productivity, as well as safety and portability of GPGPU programs.

### 2.2. *Why OCaml?*

GPGPU programming is (notoriously) hard (with current language and tools); getting high performance is even harder. Low-level, verbose APIs bring total control over the hardware but also bring confusion and errors. OCaml[1] offers many advantages to make GPGPU programming easier. It is very efficient (mainly through the native compiler provided by Inria) for sequential computation, which is mandatory for programs aiming at having high performance. It is a multi-paradigm language (*e.g.*, functional, imperative, object oriented). It is also highly extensible, among other via the Camlp4 preprocessor (part of the official OCaml distribution) whose main application is to define domain-specific extensions for the OCaml language. Thus, OCaml is a good starting point to experiment various programming methods to deduce the best practices for GPGPU programming. It is a good platform to design new abstractions to describe GPGPU computation easily, while maintaining efficiency. Furthermore, OCaml is fully interoperable with C and offers automatic memory management associated with a configurable garbage collector (GC) which can be used to manage OCaml values as well as external values coming from C bindings. OCaml GC is highly efficient, consisting of a two generation garbage collector combining a Stop&Copy algorithm on the first generation with an incremental Mark&Sweep&Compact algorithm on the second generation. OCaml is also

---

[a]e.g : OpenCL is compatible with CPUs which share memory with themselves

a strongly and statically typed language which as well as the memory manager improves program reliability while increasing development productivity by detecting a lot of common programming errors at compile time. And last, OCaml is compatible with many architectures and operating systems, allowing us to provide a portable solution.

### 2.3. *Abstracting GPUs*

One of the main difficulties in current GPGPU programming comes from hardware heterogeneity: usually a standard CPU architecture linked with GPUs using their own memory banks and programs. Our library abstracts devices by managing most of this heterogeneity automatically. SPOC offers a way to use any current framework transparently. Besides, we built a runtime library offering automatic transfers.

**Abstracting Frameworks** As described previously, GPGPU programming implies choosing a framework/API. We focused our efforts on Cuda and OpenCL. They offer different ways to express computing kernels but also to launch them or to manage memory on devices. Our library unifies these two frameworks allowing to run a program with any type of device, separately or conjointly. We based our library on dynamic linking with the driver APIs of the different frameworks. This allows, during initialization, to detect any device compatible with Cuda or OpenCL on the hardware running the program to make them usable. We unified both APIs into one library to allow portability. Programs written with SPOC can run with any devices. This implies that our library only offers functions shared by both frameworks, while it uses the specificities of each framework in its core to provide the best performance.

**Vectors** Because of transfers, stream processing provides high performance for high computation over data size ratio. However, to benefit from the many computation units on current GPGPU devices, stream processing needs big data sets to compute on. OCaml offers many ways of describing sets of values including arrays and lists. SPOC provides the type vector to represent large data sets. Vectors must be able to contain a large amount of data while being easily transferable to devices through C bindings. As the Bigarray module of OCaml allows us to define large data sets easily interoperable between OCaml and C, we based SPOC vectors on bigarrays. Bigarrays are uniform arrays of predefined types (int, float, complex). To allow OCaml users to define their own vector type we added a Custom vector type. The type of custom vectors must be defined as an OCaml record type including information about the memory size of a vector element as well as specific functions to access vector elements.

**Abstracting Transfers** Efficient copy of data, between CPU and devices, by minimizing useless (and expensive) transfers, is crucial. Kernels are procedures returning

no value but having side effects over data on their global memory. Thus, it is difficult to know which value to transfer back to CPU memory. A solution to manage automatic transfers is to use "cyclic" transfers [2]. When launching a kernel, every vector needed to run the kernel is transferred to the device, and brought back after completion. This ensures that every modified vector is now accessible from the CPU but can trigger unnecessary transfers. Our alternative is to move data only when needed. SPOC vectors store their location (CPU/Devices). When launching a kernel, our runtime library checks if data are present on the GPU or instead on CPU memory. If needed they are transferred (by the CPU) to the GPU.After completion, we do not move any data. When a vector is read or written by the CPU, SPOC checks its location bringing it back if needed. This offers a good optimization over cyclic transfers (especially in programs running multiple kernels sequentially with the same vectors) by limiting useless transfers while assuring that data are always located where they are needed. However, if it does not bring data back to CPU memory and keep them on device memory, SPOC may quickly fill our device memory. This is solved by linking transferred data to the OCaml GC. When vectors become useless to the program or any kernel launched by the program, they are discarded from memory on CPU as well as on device's global memory depending on their current location.

**Managing Kernels** SPOC offers a way to declare external kernels (meaning external sources written in Cuda assembly or OpenCL), which can then be launched from OCaml. To declare an external kernel and make it accessible from the rest of the program, SPOC provides an OCaml syntax extension (using the Camlp4 preprocessor) using keyword "kernel". The programmer provides the OCaml name of the kernel, the types of its parameters and the location of its implementation.

```
1  __kernel void vec_add(__global const float* a,
2                        __global const float* b,
3                        __global float* c, int N) {
4    int nIndex = get_global_id(0);
5    if (nIndex < N)
6      c[nIndex] = a[nIndex] + b[nIndex];
7  }
```

Fig. 1.   OpenCL vector add

For instance, given a kernel written in OpenCL or Cuda (fig. 1) computing the sum of two 32-bit float vectors, he declares a kernel *vector_add* taking three 32-bit float vectors and an integer as parameters (see fig. 2, lines 1 to 3). Both last strings define respectively the kernel file[b] and the name of the kernel function in this file.

---

[b]SPOC will search and compile an OpenCL or Cuda file depending on the hardware

6    *M. Bourgoin, E. Chailloux, J-L Lamotte*

Our extension will statically generate the code instantiating a spoc_kernel OCaml object (defined in the library) providing methods to compile and run it. Kernel arguments are type-checked as well as the complete OCaml program at compile-time.

**Example** The program presented in fig. 2 uses SPOC to sum two vectors.

```
1   kernel vector_add :
2     Spoc.Vector.vfloat32 -> Spoc.Vector.vfloat32 -> Spoc.Vector.vfloat32 ->
3       int -> unit = "kernel_file" "vec_add"
4
5   let example () =
6     let devs = Spoc.Devices.init () in
7     let a = Spoc.Vector.create Spoc.Vector.float32 1024
8     and b = Spoc.Vector.create Spoc.Vector.float32 1024
9     and c = Spoc.Vector.create Spoc.Vector.float32 1024
10    in
11      fill_vectors [a; b; c];
12      let blk = {Spoc.Kernel.blockX = 256;
13                 Spoc.Kernel.blockY = 1;
14                 Spoc.Kernel.blockZ = 1;}
15      and grd = {Spoc.Kernel.gridX = 4;
16                 Spoc.Kernel.gridY = 1;
17                 Spoc.Kernel.gridZ = 1;}
18      in
19        Spoc.Kernel.run devs.(0)(blk,grd) vector_add (a, b, c, 1024);
20        for i = 0 to 1023 do
21          Printf.printf "%g\n" (c.[<i>])
22        done ;;
```

Fig. 2.   OCaml vector add

SPOC is initialized with *Devices.init* (line 6) which returns the array of compatible devices. Then, it defines three 32bit vectors, filled with random values (lines 7 to 11) and describes the dimensions of computations on the device: the size of the grid and blocks of threads that will run the kernel (lines 12 and 17). Bigger dimensions than those physically present on the hardware will force some hardware multi-processors to run several blocks. Here, it describes grid and blocks in order to dedicate one thread for each element of the vectors. Then, it compiles the kernel for a device (kernels are dynamically compiled, with the compiler accessible through the framework implementation corresponding to the used device, ensuring compatibility with every CPU-Device system) (lines 12 to 17). It finally compiles and runs the kernel on the device providing blocks and grid dimensions (line 18), and the parameters needed as a tuple, before printing the obtained result (line 20 to 22). Here, SPOC automatically checks which kind of device will run the kernel and depending on the framework corresponding to this device compiles the correct kernel before executing it (see fig.3). Compiled kernels are cached to prevent multiple compilations in case of multiple use of the same kernel. During this operation, vectors *a*, *b* and *c* are transferred from the CPU (where they were defined) to the
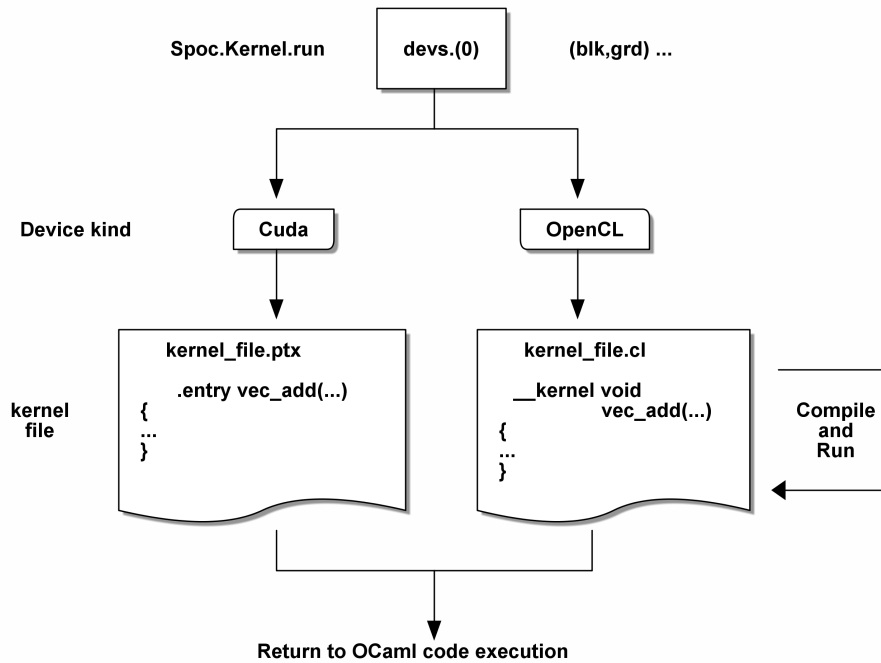
Fig. 3.   Kernel dynamic compilation

device (where the kernel is run). Only vector $c$ is brought back to CPU memory to be accessible in the printing loop. All those transfers are automatic. Vectors $b$ and $c$ are discarded from device memory by the GC.

**External Libraries** We used SPOC to provide a binding to the Cublas V1 library giving access to a set of optimized kernels running on Cuda compatible devices corresponding to the Basic Linear Algebra Subprograms (BLAS) functions. This Cublas module uses SPOC vectors and devices.

**Multi-GPU** During initialization, SPOC searches for any device compatible with Cuda or OpenCL. Then, it becomes easy to use multi-device systems, by dividing data sets and dynamically compile kernels for any device found. SPOC can use any type of device separately or conjointly. Besides, SPOC offers a set of functions, to define sub-vectors, which helps to describe muti-gpu programs (see fig. 4). SPOC's sub-vectors share CPU memory space with their parent. However, they do not share memory space once transferred to devices. Sub-vectors can also be defined to include non contiguous data from their parent vector. This enables to easily divide a vector into many sub-vectors (without extending the CPU memory space

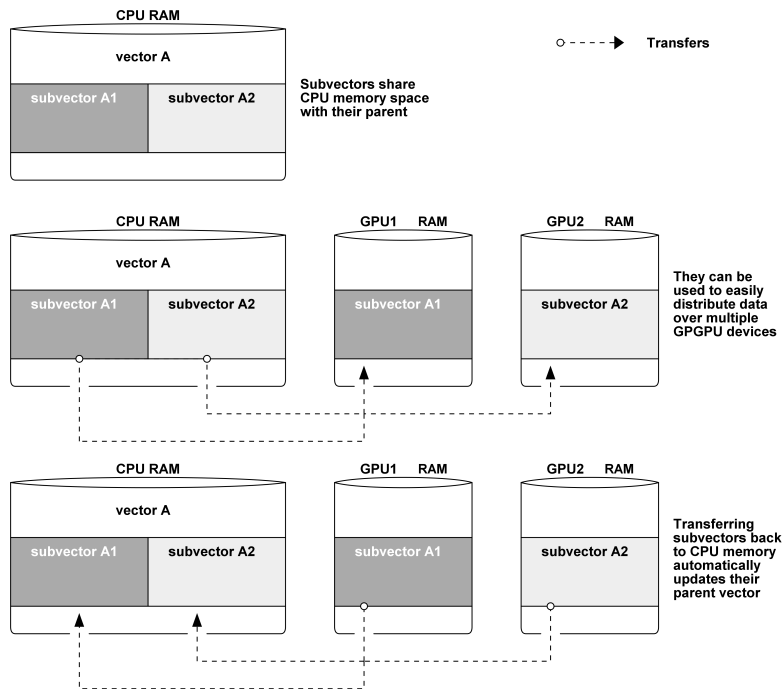8   *M. Bourgoin, E. Chailloux, J-L Lamotte*



Fig. 4.   Subvectors

needs) before sending different sub-vectors to different devices. When transferred back to the CPU memory, sub-vector data are merged with those of the parent vector. Currently, sub-vectors must explicitly be transferred back to CPU memory to enable multiple copies of a same region of a vector to be transferred to different devices without overwriting them multiple times (and spending a lot of time in transfers) when accessing their parents.

## 3. Benchmarks

To check SPOC performance we tested it through two benchmarks (using simple-precision floating-point format). Our test system consists of a quad-core Intel(R) Core(TM) i7 CPU 960 at 3.20GHz (102 GFLOPS simple precision) with two GPUs and 12GB CPU RAM. The GPUs are an NVIDIA Tesla C2070 (1030 GFLOPS, with NVIDIA driver 270.41.19) with 6GB RAM and an AMD Radeon HD 6950 (2.25 TFLOPS, with AMD-APP-SDK-v2.5 (684.213) with CAL 1.4.1353 driver) with 2GB RAM.

**Mandelbrot** This benchmark draws the image resulting from computation of the Mandelbrot set. Measured times include transfers and drawing. This program uses naive kernels to verify that SPOC can easily increase performance of data-parallel programs already written in OCaml. This program is run using one or both GPUs, dividing computations naively (and statically) into two equal sub-computations. We compare our program using OCaml with SPOC to one using OCaml with a C external program for the computations (using OpenMP to use the multicore capabilities of our CPU).

| | OCaml and C | | OCaml and SPOC | | | |
|---|---|---|---|---|---|---|
| | Intel i7 | | AMD 6950 | Tesla C2070 | C2070+6950 | C2070+6950 |
| | 1 Core | 4 Cores | OpenCL | Cuda | Cuda+OpenCL | OpenCL |
| | 892s | 307s | 12.84s | 10.99s | 6.56s | 6.66s |
| Speedup | - | **1** | 23,91 | 27,93 | 46,80 | 46,10 |

Those speedups show that SPOC allows efficient Stream Processing with OCaml (with both GPUs, SPOC provides a 190x speedup to the "same" program written in OCaml only, which takes 1252s (1 thread) to compute). Furthermore, SPOC efficiently distributes computation between multiple GPUs.

**Matrix Multiply** In this benchmark, we multiply two matrices using the Cublas SGEMM function of the Cublas library. Cublas being provided only for NVIDIA GPU we can only use the Tesla C2070 card. SPOC Cublas binding is based on Cublas V1 which cannot launch two Cublas functions on the same device concurrently. The first two multiplications use square matrices of size 21000. This size implies the use of at least 5.2GB $((32/8 * 21000^2 * 3)/10^9 = 5, 2)$ on the device. Matrix_Multiply_1 computes directly over the complete matrices in one time. Matrix_Multiply_2 divides each matrix into 4 squared blocks and implements a block matrix multiplication by using SGEMM. This computation needs several steps implying many transfers between CPU and device but with smaller amounts of data. Our third test computes over bigger matrices (size 25000) needing 7.5GB on the device, forcing our system to trigger at least one garbage collection on the device in order to measure the cost implied by garbage collection. This collection will free memory used during previous steps on the device. Here we use Matrix_Mutliply_2 to provide this out-of-core product with automatic data management.

| | Matrix Multiply 1 | Matrix Multiply 2 | Matrix Multiply 2 |
|---|---|---|---|
| Matrix size | 21000 | 21000 | 25000 |
| Maximum memory needed | 5.2GB | 5.2GB | 7.5GB |
| GFLOPS[c] | 473 | 377 | 421 |

Results showed that using blocks already impacts performance as transferring blocks implies transferring non contiguous data (and as we are not overlaping block transfers with computations). With the third test, each block contains larger contigous

[c]including transfers and computation

data than with the second test, which increases transfer performance. This result showed that the cost of a collection remains moderate, in view of the fact that it automatically provides an out-of-core product.

## 4. Related Works

Many high-level languages now offer bindings to the OpenCL or Cuda frameworks. Some offer higher abstractions to ease GPGPU programming. They mostly propose specific data sets (like SPOC) with specific operators to generate GPGPU code. Some propose embedded DSL to express kernels. It was difficult to compare them directly with SPOC as most of them only work with specific GPUs or are currently in development. In this section, we present some of these approaches and their specificities.

**Hybrid Multicore Parallel Programming (HMPP)** While focusing on C and Fortran, HMPP workbench (CAPS-entreprise) offers directives (very similar to OpenMP) to declare and call codelets which can run on accelerators[3]. This method allows developers to keep their source compatible with CPU-only systems while the HMPP compiler will generate specific code for GPGPU systems. A standard (OpenACC[4]) based on directives (highly compatible with HMPP directives) is currently being developed by the OpenACC organization.

**SkePU** SkePU[5] is a C++ template library with a very similar approach to SPOC. It provides a unified interface to Cuda and OpenCL and uses lazy memory copying to avoid unnecessary transfers. The major difference with SPOC comes from the use of OCaml and its Garbage Collector to implement memory transfers.

**SkelCL** SkelCL[6] is a C++ library providing parallel skeletons for GPGPU programming. Those skeletons allow GPGPU computations over a vector data-type automatically transferred as in our approach. SkeCL skeletons are implemented using OpenCL.

**Accelerator** Accelerator[7] proposes GPGPU computing through the .Net framework (including F#). It allows to declare specific arrays of data combined with overloaded operators which enable to describe computations over those arrays. Computations are optimized and translated to GPGPU kernels (in DirectX shaders).

**JavaCL/ScalaCL** While JavaCL is a binding to the OpenCL framework, ScalaCL[8] (based on JavaCL) goes further, providing specific sets of data, as does Accelerator, with automatic translations of specific operations over classic scala set of data.

**Aparapi** Aparapi[9] is an open source Java library expressing data parallel workloads with a runtime library converting Java bytecode to OpenCL. It provides a

Kernel class whose method *run* must be overridden to express the kernel to run. It currently only works with one AMD device at a time.

**Obsidian**  Obsidian[10] is an embedded language for Haskell to describe operations on arrays. It offers a set of combinators which are translated to Cuda and run on the device.

## 5. Conclusion and Future Work

SPOC[11] is a library which enables launching GPGPU kernels written in OpenCL and Cuda. It relies on specific vector data types associated with a runtime library abstracting memory transfers between CPU and GPGPU devices while dynamically adapting its behavior depending on the compatible hardware (Cuda/OpenCL/Both). It offers great performance and can now be used as a basis to experiment.

**Real world use case**  In order to test our library, measure performance and propose new/better abstractions, we will use SPOC to translate a numerical software (PROP[12]) from FORTRAN+Cuda to (OCaml+SPOC)+Cuda. This work will focus on kernel and memory transfers management as we should be able to reuse the current Cuda kernel implementation as external kernels with SPOC.

**Extensions**  As Cublas V1, we intend to provide bindings to GPGPU Libraries through SPOC starting from Cublas V2. Furthermore, having mainly focused our efforts on our runtime and memory manager, we should now be able to propose extensions to express the computation to run on the device. This could be proposed through a DSL embedded into an OCaml extension, as well as through specific operators over Vectors directly translated to kernels (as many approaches presented section 4).

**Models**  OCaml extensions like BSML[13] and CamlP3L[14] propose models to describe data and task parallel computations. The stream processing paradigm relies heavily on data-parallelism and we could use those models with SPOC. Task-parallel models could also be used especially to describe multi-GPU software.

Most other approaches (presented section 4) use specific data sets, associated with functions/combinators which are then translated to GPGPU kernels. We focused on the runtime library to ensure performance and to offer efficient automatic memory management. OCaml is easily extensible and we can now provide combinators/skeletons which will also generate kernels benefiting from OCaml type-safety.

12   *M. Bourgoin, E. Chailloux, J-L Lamotte*

# References

[1] Xavier Leroy. The Objective Caml system release 3.12 : Documentation and user's manual. Technical report, Inria, 2011. `http://caml.inria.fr`.

[2] T.B. Jablin, P. Prabhu, J.A. Jablin, N.P. Johnson, S.R. Beard, and D.I. August. Automatic CPU-GPU Communication Management and Optimization. In *Proceedings of the 32nd ACM SIGPLAN conference on Programming language design and implementation*, pages 142–151. ACM, 2011.

[3] R. Dolbeau, S. Bihan, and F. Bodin. HMPP: A Hybrid Multi-core Parallel Programming Environment. In *First Workshop on General Purpose Processing on Graphics Processing Units*, 2007.

[4] Nvidia Cray Inc., CAPS Enterprise and The Portland Group. Openacc 1.0 specification, 2011.

[5] Johan Enmyren and Christoph W. Kessler. Skepu: a multi-backend skeleton programming library for multi-gpu systems. In *Proceedings of the fourth international workshop on High-level parallel programming and applications*, HLPP '10, pages 5–14. ACM, 2010.

[6] Michel Steuwer, Philipp Kegel, and Sergei Gorlatch. Skelcl - a portable skeleton library for high-level gpu programming. In *Proceedings of the 25th IEEE International Symposium on Parallel and Distributed Processing Workshops and PhD Forum*, IPDPSW '11, pages 1176–1182. IEEE Computer Society, 2011.

[7] D. Tarditi, S. Puri, and J. Oglesby. Accelerator: Using Data Parallelism to Program GPUs for General-Purpose Uses. *ACM SIGARCH Computer Architecture News*, 34 (5):325–335, 2006.

[8] R. Beck, H.W. Larsen, T. Jensen, and B. Thomsen. Extending Scala with General Purpose GPU Programming. Technical report, Adlborg University, Departement of Computer Science, 2011.

[9] AMD. Aparapi. `http://code.google.com/p/aparapi/`.

[10] Joel Svensson. Obsidian: GPU Kernel Programming in Haskell. Technical Report 77L, Computer Science and Enginering, Chalmers University of Technology and Gothenburg University, 2011.

[11] Mathias Bourgoin, Emmanuel Chailloux, and Jean-Luc Lamotte. Spoc: Stream processing with ocaml, 2012. `http://www.algo-prog.info/spoc`.

[12] N. Stan Scott, M. Penny Scott, Phil G. Burke, Timothy Stitt, V. Faro-Maza, Christophe Denis, and A. Maniopoulou. 2DRMP: A suite of two-dimensional R-matrix propagation codes. *Computer Physics Communications*, 180(12):2424–2449, 2009.

[13] Frédéric Loulergue, Frédéric Gava, and D. Billiet. Bulk Synchronous Parallel ML: Modular Implementation and Performance Prediction. In Vaidy S. Sunderam, Gaétan Dick van Albada, Peter M. A. Sloot, and Jack Dongarra, editors, *International Conference on Computational Science (ICCS)*, LNCS 3515, pages 1046–1054. Springer, 2005.

[14] Roberto Di Cosmo, Zheng Li, Susanna Pelagatti, and Pierre Weis. Skeletal Parallel Programming with OCamlP3l 2.0. *Parallel Processing Letters*, 18(1):149–164, 2008.