Obsidian: GPGPU Programming in Haskell

Joel Svensson

Chalmers University of Technology

May 15, 2009

Joel Svensson Obsidian: GPGPU Programming in Haskell

・ 同 ト ・ ヨ ト ・ ヨ ト …

General-purpose computations of GPUs

Why are GPUs interesting for non-graphics programming:

- Cost efficient highly parallel computers
- ~500 processing cores (NVIDIA 295GTX)
- Taste of the future, today!
 - Programming a 500 core machine is very different from a 2,4,8,16 core machine.





GPUs are made to draw Triangles. As many as possible as quickly as possible.



프 🕨 🗆 프

Screenshot



Figure: Fallout 3, Bethesda Softworks.

Joel Svensson Obsidian: GPGPU Programming in Haskell

Vertex programs:

- Executed per vertex.
- 3d to 2d transformations.

Fragment programs:

- Executed per fragment (potential pixel).
- Computes a color value.





Vertex programs:

- Executed per vertex.
- 3d to 2d transformations.

Geometry programs: Fragment programs:

- Executed per fragment (potential pixel).
- Computes a color value.





Unified Architecture



Figure: The NVIDIA 8800GTX GPU architecture, with 8 pairs of multiprocessors. Diagram courtesy of NVIDIA.

ヘロア 人間 アメヨア 人口 ア

ъ

Unified Architecture

In each multiprocessor:

- 16Kb Shared memory.
- 8192 32-bit registers.

Memory:

- Uncached device memory.
- Read-only constant memory.
- Read-only texture memory.



イロン 不同 とくほ とくほ とう

NVIDIA Compute Unified Device Architecture:

- NVIDIA's Hardware architecture + programming model.
- Provides compiler and libraries.
 - Extension to C.
- Enables/Simplifies implementing general purpose algorithms on the GPU

ヘロン 人間 とくほ とくほ とう

э.

```
global static void sum(int * values,int n)
 extern ____shared___ int shared[];
 const int tid = threadIdx.x;
 shared[tid] = values[tid];
 for (int j = 1; j < n; j *= 2) {
 syncthreads();
   if ((tid + 1) % (2*j) == 0)
     shared[tid] += shared[tid - j];
  }
 values[tid] = shared[tid];
```

(《문》《문》 문

The code listing on the previous slide defines a Kernel:

- A kernel is executed by a block of threads:
 - Up to 512 threads per block.
 - Run on one multiprocessor.
- The threads are further divided into Warps
 - The scheduled unit.
 - Group of 32 threads.
- Many blocks can be executed concurrently:
 - Referred to as a *Grid* of Blocks.

・聞き ・ヨキ ・ヨト

synchtreads()

Barrier synchronisation primitive:

- Barrier across all the threads of a block.
- Used to coordinate communication
 - Within a block.

```
if (even(tid)) {
    /* do something */
    ___syncthreads();
}
```

More ____syncthreads()



Joel Svensson Obsidian: GPGPU Programming in Haskell

◆臣→

Obsidian Outline:

- Embedded in Haskell.
- Tries to stay in the spirit of Lava.
 - Combinator library.
- Higher level of abstraction compared to CUDA.
 - While still assuming knowledge of architecture characteristics in the programmer.

・ 同 ト ・ ヨ ト ・ ヨ ト …

æ

- Generate efficient code for GPUs from short and clean high level descriptions.
- Make design decisions easy.
 - Where to place data in the memory hierarchy.
 - What to compute where, and when.

個人 くほん くほん

- Generate efficient code for GPUs from short and clean high level descriptions.
- Make design decisions easy.
 - Where to place data in the memory hierarchy.
 - What to compute where, and when.
- We are not there yet.

・ 同 ト ・ ヨ ト ・ ヨ ト ・

(ロ) (同) (目) (日) (日) (の)

Running an Obsidian program on the GPU

Obsidian> execute GPU (sumUp 4) [0..15] [120]

▲□ ▶ ▲ 三 ▶ ▲ 三 ▶ ● 三 ● ● ● ●

```
global void generated(unsigned int* input, unsigned int* result) {
unsigned int tid = (unsigned int)threadIdx.x;
extern shared unsigned int s data[];
unsigned int attribute ((unused)) *sm1 = &s data[0];
unsigned int __attribute__((unused)) *sm2 = &s_data[8];
sm2[tid] = ((unsigned int)((input[(tid << 1)] + input[((tid << 1) + 1)]));</pre>
syncthreads();
if (tid < 4) {
sml[tid] = ((unsigned int)((((int)(sm2[(tid << 1)])) + ((int)(sm2[((tid << 1) + 1)]))));</pre>
syncthreads();
if (tid < 2) {
sm2[tid] = ((unsigned int)((((int)(sm1[(tid << 1)])) + ((int)(sm1[((tid << 1) + 1)]))));</pre>
syncthreads();
if (tid < 1) {
sml[tid] = ((unsigned int)((((int)(sm2[(tid << 1)])) + ((int)(sm2[((tid << 1) + 1)]))));</pre>
___syncthreads();
if (tid < 1) {
result[tid] = ((unsigned int)(sm1[tid]));
```

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

Key parts of Obsidian:

Arrays

```
data Arr a = Arr (IxExp -> a) Int
```

Obsidian programs

```
data a :-> b
 = Pure (a -> b)
 | Sync (a -> Arr FData) (Arr FData :-> b)
```

• Collection of combinators and functions.

two, ->-, sync, pair, unpair, zipp, unzipp, etc

▲冊 ▶ ▲ 臣 ▶ ▲ 臣 ▶ ● ○ ○ ○ ○ ○

The Obsidian sync has many roles :

- Values are stored in shared memory.
 - Enables sharing of computed results between threads.
- Introduces parallelism.
- Assigns work to threads.
 - The length of the input array specifies the number of threads.

```
sync :: (Flatten a) => Arr a :-> Arr a
```

instances of Flatten have functions toFData and fromFData defined on them.

・ 同 ト ・ ヨ ト ・ ヨ ト …

sync :: Flatten a => Arr a :-> Arr a
sync = Sync (fmap toFData) (Pure (fmap fromFData))

(ロ) (同) (三) (三) (三) (○)

◆□▶ ◆□▶ ◆ □▶ ◆ □▶ ● □ ● ○ ○ ○

```
addOne :: Arr IntE :-> Arr IntE
addOne = Pure (fmap (+1)) ->- sync
```

```
Obsidian> execute GPU addOne [0..7] [1,2,3,4,5,6,7,8]
```

```
Obsidian> execute GPU addOne' [0..7]
[(1,2),(3,4),(5,6),(7,8)]
```

[0, 1, 3, 6, 10, 15, 21, 28]

(ロ) (同) (三) (三) (三) (0)

Drawing a Sklansky



(ロ) (四) (モ) (モ) (モ) (モ)

```
global void generated (unsigned int* input, unsigned int* result) {
unsigned int tid = (unsigned int)threadIdx.x:
extern shared unsigned int s data[];
unsigned int attribute ((unused)) *sm1 = &s data[0];
unsigned int attribute ((unused)) *sm2 = &s data[8];
sm2[tid] = ((unsigned int))((((tid & 0xfffffff)) < 1))?
((int)(input[tid])) :
(((int)(input[(tid & 0x6)])) + ((int)(input[tid])))));
___syncthreads();
sml[tid] = ((unsigned int)((((tid & Oxfffffffb) < 2) ?</pre>
((int)(sm2[tid])) :
(((int)(sm2[((tid \& 0x4) | 0x1)])) + ((int)(sm2[tid])))));
syncthreads();
sm2[tid] = ((unsigned int))(((tid < 4))?
((int)(sm1[tid])) :
(((int)(sm1[3])) + ((int)(sm1[tid])))));
syncthreads();
result[tid] = ((unsigned int)(sm2[tid]));
```

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □

- Obsidian is work in progress.
 - Changing rapidly.
- Promising: for some applications we are generating quite efficient code.
 - More needs to be done.
 - Generalise.

・ 同 ト ・ ヨ ト ・ ヨ ト …

Applications:

- Sorting.
- Parallel prefix.
- Reduction.

Joel Svensson Obsidian: GPGPU Programming in Haskell

< 一型

More Questions?

